University of Windsor

Scholarship at UWindsor

Electronic Theses and Dissertations

Theses, Dissertations, and Major Papers

1980

OPTIMIZATION PROBLEMS IN STATISTICS.

KAI SANG. WONG University of Windsor

Follow this and additional works at: https://scholar.uwindsor.ca/etd

Recommended Citation

WONG, KAI SANG., "OPTIMIZATION PROBLEMS IN STATISTICS." (1980). *Electronic Theses and Dissertations*. 2951.

https://scholar.uwindsor.ca/etd/2951

This online database contains the full-text of PhD dissertations and Masters' theses of University of Windsor students from 1954 forward. These documents are made available for personal study and research purposes only, in accordance with the Canadian Copyright Act and the Creative Commons license—CC BY-NC-ND (Attribution, Non-Commercial, No Derivative Works). Under this license, works must always be attributed to the copyright holder (original author), cannot be used for any commercial purposes, and may not be altered. Any other use would require the permission of the copyright holder. Students may inquire about withdrawing their dissertation and/or thesis from this database. For additional inquiries, please contact the repository administrator via email (scholarship@uwindsor.ca) or by telephone at 519-253-3000ext. 3208.



Canadian Theses on Microfiche Service

Bibliothèque nationale du Canada Direction du développement des collections

Service des thèses canadiennes sur microfiche

NOTICE

AVIS

The quality of this microfiche is heavily dependent upon the quality of the original thesis submitted for microfilming. Every effort has been made to ensure the highest quality of reproduction possible.

If pages are missing, contact the university which granted the degree.

Some pages may have indistinct print especially if the original pages were typed with a poor typewriter ribbon or if the university sent us a poor photocopy.

Previously copyrighted materials (journal articles, published tests, etc.) are not filmed.

Reproduction in full or in part of this film is governed by the Canadian Copyright Act, R.S.C. 1970, c. C-30. Please read the authorization forms which accompany this thesis.

THIS DISSERTATION
HAS BEEN MICROFILMED
EXACTLY AS RECEIVED

La qualité de cette microfiche dépend grandement de la qualité de la thèse soumise au microfilmage. Nous avons tout fait pour assurer une qualité supérieure de reproduction.

S'il manque des pages, veuillez communiquer avec l'université qui a conféré le grade.

La qualité d'impression de certaines pages peut laisser à désirer, surtout si les pages originales ont été dactylographiées à l'aide d'un ruban usé ou si l'université nous a fait parvenir une photocopie de mauvaise qualité.

Les documents qui font déjà l'objet d'un droit d'auteur (articles de revue, examens publiés, etc.) ne sont pas microfilmés.

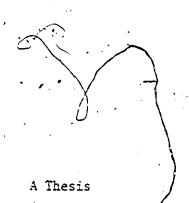
La reproduction, même partielle, de ce microfilm est soumise à la Loi canadienne sur le droit d'auteur, SRC 1970, c. C-30. Veuillez prendre connaissance des formules d'autorisation qui accompagnent cette thèse.

LA THÈSE À ÉTÉ MICROFILMÉE TELLE QUE NOUS L'AVONS REÇUE

Ottawa, Canada K1A 0N4 OPTIMIZATION PROBLEMS IN STATISTICS

 $\left[\mathsf{C}\right]$

by KAI SANG WONG



Submitted to the Faculty of Graduate Studies through the Department of Mathematics in Partial Fulfillment of the Requirements for the Degree of Doctor of Philosophy at The University of Windsor.

Windsor, Ontario

© KAI SANG WONG'. 1980

740380

ABSTRACT

The theory of symbolic matrix derivatives is connected to the theory of differentials. It is shown that symbolic matrix derivatives are nothing but linear transformations of the representations of certain differentials. Representations of various differential rules are obtained and compared with those obtained by various authors. As illustrations, particular attention is given to the product rule. The theory of monotone operators is used to find the optimal solutions of various optimization problems in statistics. Some algebraic results which might be of interest by themselves are obtained to prove the main results. Optimal control models of regression experiments are presented to illustrate optimization problems with solutions on the boundary of the region of concern.

ACKNOWLEDGEMENTS

I would like to take this opportunity to thank Professor Cervin, Professor Gillen, Professor Ponnapalli, and Professor Traynor for their interest in my work.

I would like to express my sincere thanks to Professor Lemire and Professor Atkinson, who, as past and present chairmen of the Department of Mathematics, have been giving me help and encouragement during my four and a half year stay in the University.

Finally, I would like to express my deep sense of gratitude to my supervisor, Professor Chi Song Wong, who constantly gave me guidance, encouragement, help and rigorous training in mathematical statistics.

I acknowledge the financial support received from the Natural Sciences and Engineering Research Council.

TABLES OF CONTENTS

		Page
ABSTRACT		iii.
ACKNOWLED	GEMENTS	i⊽.
CHAPTER	• •	
0. INT	RODUCTION	1
. T DED		
	RESENTATIONS OF DIFFERENTIALS	. 5
1.1	Preliminaries	5
1.2	Differentials	. 8
1.3	Differential Rules	- 9
1.4	Quadratic Differentials	12
1.5	Representations	13
1.6	Examples ?	30
II. MIN	IMA AND MAXIMA IN LINEAR MODELS AND	
MUL	FIVARIATE ANALYSIS	34
2.1	Introduction (م 34
2.2	Preliminaries	34
2.3	Contingency Tables	38
2.4	Maximum Likelihood Estimates of Multivariate Normal Model	43
2.5	Multivariate Regression Models	44
2.6	Multivariate Linear Hymotheses	7.6

TABLE OF CONTENTS (CONT'D),

•	•		Page
	2.7	Quadratic Estimates	50
	2.8	Minimum Distances and Principal Components	52
	2.9	Factor Analysis	55
	2.10	Growth Curve Models	• , 57
	2.11	Simultaneous Equation Models	59
III.	OPTI	MAL CONTROL OF A REGRESSION EXPERIMENT	62
	3.1	Introduction .	62
	3.2	Optimal Control of a Regression Experiment	62.
	3.3	D-, A-, and D _s - Optimal Models	69
,	3.4	Reproducing Kernel Hilbert Spaces	78
REFER	ENCES		93

VITA AUCTORIS

CHAPTER ZERO

Introduction

Symbolic matrix derivatives emerge as means to handle statistical and mathematical problems when the variables in a hypothesized model are many. As modern science and technology advance, the need to take more factors than one into consideration while setting up any mathematical or statistical model becomes increasingly pronounced. Multivariate analysis, originated with the paper of Hotelling (1931, see, e.g., Anderson (1958)), has become the central and one of the most important branches of statistics and data analysis. Though multivariate models are usually simple in. form, the computation of test statistics and estimates is often difficult. Estimation, for instance, is usually confined to either maximum likelihood or Bayesian methods depending, partially, on one's philosophical belief. The maximum likelihood estimates are values which maximize a given likelihood function. The maximum likelihood method is more familiar to the practitioners in the field presumably because Bayesian methods are more recent and less known, and because it is employed in the classical book by Anderson (1958), or more recently, in the book by Kshirsagar (1972). The maximum likelihood method has the additional advantage that it is connected directly to hypothesis testing. A third and the oldest method, the least square method, which is prominent in the theory of linear models, is equivalent to the maximum likelihood method under the normal theory. This might be another reason for the dominance of the maximum likelihood method in multivariate analysis.

The need for matrix differentiation was pointed out in the paper by Dwyer and MacPhail (1948). In this paper, the authors define two kinds of derivatives symbolically, and apply it to the problem of least squares, canonical correlations and orthogonal regression. They also give two tables to illustrate various derivatives of real-valued functions with matrix variables or vice versa. Following the above approach, Dwyer (1967) explores further applications of matrix derivatives and works out various formulae which are useful in multivariate analysis. The connection between matrix derivatives and Jacobians was examined. Results similar to those of Deemer and Olkin (1951) (originated by P.L. Hsu) are also obtained. Later on, Tracy and Dwyer (1969) consider the derivatives of vectors with respect to vectors and use it to represent derivatives of matrix-valued functions with respect to matrices. This also leads to the second order derivatives of real-valued functions with respect to vectors or matrices. They give several applications in multivariate analysis with an attempt to justify that the critical values of the matrix-valued functions they find give the absolute minima (maxima). Following these presentations, Tracy and Singh (1972) and Singh (1972) generalize certain results to the case of partitioned matrices. In the period between the publication of these papers, the need for matrix differentiation is now well recognized and is referred to by Anderson (1958), Rao (1973) and Graybill(1969).

Econometricians also need matrix differentiation in the development of their theory. For example, Neudecker wrote three papers in the period of 1967 to 1969. In the investigation of

matrix differentiation, he uses the differential notions algebraically. In the 1969 paper, he also puts certain elements of a matrix in vector form to represent the derivatives of matrix-valued functions of matrix variables.

The paper by Vetter in 1970 is worth mentioning. In the paper, he derives a chain rule and differential rules for matrix product and Kronecker product and gives several examples to illustrate the applications of his differential rules. His work is aimed at the applications of matrix differentiation to system and control theory (see, for examples, Athans and Schweppe (1965), Athans and Tse (1967) and Athans (1967)). In 1973, McDonald and Swaminathan presented a system of matrix calculus and labelled them as McD.-S. calculus. In their paper, they give their own definitions of matrix derivatives and derive a chain rule and various product rules. Later, MacRae (1974), McDonald (1976), Swaminathan (1976) and Bentler and Lee (1975, 1978) all try to formulate and develop matrix derivatives further in this direction.

However, while the techniques of matrix derivatives are applied to various optimization problems in statistics or other disciplines such as econometrics, there is a lack of justification for the optimality of the solutions obtained by using matrix derivatives. The formulae given by various authors are long, complicated and difficult to remember. Another disadvantage with the existing methods is that there is no unity in the matrix calculus developed by various authors. Each researcher has his own basic definitions and formulae. Such a situation could lead to confusion and jeopar-

dize the development, understanding and aplications of the theory. One purpose of this dissertation is to connect the theory of matrix calculus to the familiar theory of multidimensional calculus (see, e.g., Apostal (1957) and Fleming (1977)) and linear algebra which can be treated and referred to as finite dimensional functional analysis. We shall show that the symbolic matrix derivatives mentioned above are nothing but linear transformations of the representations of certain differentials. The theory of monotone operators developed in the sixties (see, for example, Opial (1967)) will be used to find the optimal solutions of various optimization problems in statistics (Wong (to appear), Wong and Wong (1979, to appear)). We obtain some new algebraic results which are interesting in their own right. Optimal control models of regression experiments related to Chang (1979), Dorogovcev (1971) and Kiefer (1974) are presented here to illustrate optimization problems with solutions on the boundary of the region of concern. A problem raised in Chang and Wong (1979) in this connection is solved.

Chapter one will be devoted to matrix differentials and its representations. Chapter two will be devoted to the applications of differentials to the maximum likelihood theory and certain other optimization problems in statistics. Chapter three will be devoted to the solutions of certain problems of optimal control of a regression experiment.

For the sake of completeness, we include certain related results of Dr. Chi Song Wong, some of which are published and some to be published.

CHAPTER ONE .

Representations of Differentials

1.1 Preliminaries

Let L, L₁, L₂, ..., L_n be vector spaces over the real field R. Recall that the product Π L_i of L_i's is the linear space Ω of all functions (x_i) on $\{1, 2, \ldots, n\}$ such that $x_i \in L_i$, i = 1, 2, ..., n. Each L_i will be considered as a linear subspace of Ω by identifying each $x_i \in L_i$ with $f_i(x_i)$ in Ω , where f_i is the isomorphism of L_i into Ω such that $(f_i(x_i))(j) = 0$ if $j \neq i$, $f_i(x_i)(i) = x_i$. Thus Ω is the direct sum L₁ \oplus L₂ \oplus ... \oplus L_n of L₁, L₂, ..., L_n. Let Λ be a function of Ω into L. Λ is said to be multilinear (bilinear when n=2) if Λ is coordinatewise linear on each L_i, i.e., Λ is linear in x_i when all other x_k 's are fixed. $\Lambda((x_i))$ will be written as $x_1 \wedge x_2 \wedge \ldots \wedge x_n$.

Let I be a nonempty finite set. R^I will denote the family of all functions of I into R and will be equipped with the usual pointwise scalar multiplication and addition. R^I is a finite dimensional vector space over R. When $I = \{1, 2, \ldots, m\}$ x $\{1, 2, \ldots, n\}$, R^I is the linear space M_{mxn} of all mxn matrices (over R). When n = 1, M_{mxn} will be denoted by R^m . In general, f in $R^J \times K$ is called a $J \times K$ matrix. $\{u_i\}_{i \in I(L)}$, $\{u_{\ell}\}_{j \in I(L_{\ell})}$

will be bases of L, L₂ respectively. Let $x \in L$. Then $x = \sum_{i=1}^{n} for$ some unique x_i 's in R. The function $[x] \equiv \{x_i\}$ on I(L) is called the <u>linear representation</u> of x with respect to $\{u_i\}$.

Theorem 1.1.1. [] above is an isomorphism of L onto $R^{I(L)}$.

 Ω will be equipped with the basis $\{u_{\ell,j}\}_{j \in I(L_{\ell})}$, $\ell=1,2,3\ldots,n$. Suppose that Λ is multilinear. Then there exists a function $[\Lambda]\equiv (a_{i,j_1,j_2,\ldots,j_n})$ of the Cartesian product $I=I(L)\times\prod_{\ell=1}^n I(L_{\ell})$ into R such that for any $u_{\ell,j_{\ell}}$'s.

$$\mathbf{u}_{1,j_1} \wedge \mathbf{u}_{2,j_2} \wedge \dots \wedge \mathbf{u}_{n,j_n} = \sum_{i \in I(L)} \mathbf{a}_{i,j_1,j_2,\dots,j_n} \mathbf{u}_{i,j_n}$$

[Λ] is called the <u>linear representation</u> of Λ (with respect to the given bases).

Theorem 1.1.2. [] above is an isomorphism of the family of all multilinear functions of Ω into L onto $R^{\rm I}$.

Suppose that n=1. Then $\Lambda \in \mathcal{Z}(L_1, L)$, i.e. Λ is a linear transformation of L_1 into L. For any $f \in R^{I \times J}$, $g \in R^{J \times K}$, the $(\underline{\text{matrix}})\underline{\text{product}}$ fg of f,g is defined as an element in $R^{I \times K}$ such that each

$$(fg)((i,k)) = \sum_{j \in J} f((i,j))g((j,k)).$$

Theorem 1.1.3.

- (a) $[\Lambda(x)] = [\Lambda][x], x \in L_1, \Lambda \in \mathcal{L}(L_1/L).$
- (b) For any $\Lambda_1 \in \mathcal{L}(L, L_2)$ and $\Lambda_2 \in \mathcal{L}(L_1, L)$, $[\Lambda_1 \circ \Lambda_2] = [\Lambda_1] [\Lambda_2],$

where o is the composition for functions.

(c) [] is an isomorphism of the linear space $\mathcal{L}(L_1,L)$ onto $R^{\mathrm{I}(L)} \times \mathrm{I}(L_1)$. Hence $R^{\mathrm{I} \times \mathrm{I}}$ is an algebra which is isomorphic to $\mathcal{L}(L,L)$.

Theorem 1.1.4. Let

$$(f,g) = \sum_{i \in I} f(i)g(i), f, g \in R^{I}.$$

Then (,) is an inner product and $R^{\rm I}$ with (,) is a Hilbert space. Every Hilbert space will be equipped with an orthonormal basis. The usual orthonormal basis for $R^{\rm I}$ is $\{e_i\}$:

$$e_{i}(j) = \delta_{ij}, i,j \in I,$$

where \hat{o}_{ij} 's are the Kronecker signs. When $I=J\times J$, (,) above is called the <u>trace inner product</u> for R^I and the norm induced by (,) is called the <u>trace norm</u> for R^I .

1.2 Differentials

Let L,M be nontrivial finite dimensional Hilbert spaces (over R). Let $\{u_i\}_{i \in I(L)}$ be an orthonormal basis of L. Let C be an open set of L and f be a function of C into M. Let $x \in C$. f is said to be <u>differentiable</u> at x and has <u>differential</u> df(x) if there exists a linear transformation $df(x) \in \mathcal{L}(L,M)$ such that

$$\lim_{h\to 0} \frac{f(x+h) - f(x) - df(x)(h)}{th!} = 0 \quad (\epsilon M) .$$

Here $\{x,y\}$ is the norm induced by the inner product in L. Let $\{v_j\}_{j\in I(M)}$ be an orthonormal basis of M. Then $f(x)=\sum\limits_{j\in I(M)}f_j(x)v_j$ for some unique $f_j(x)$, s. f(x) will be denoted by $(f_j(x))$ and f will be denoted by $(f_j(x))$. $\frac{\partial f_j(x)}{\partial x_j}$ will denote

$$\lim_{t\to 0} \frac{f_j(x + tu_i) - f_j(x)}{dt} \quad (\varepsilon R)$$

and is called the <u>partial derivative</u> of f_i at the point x in the direction of u_i . Suppose that df(x) exists. Then

$$(df(x)(u_i), v_j) = \frac{\partial f_j(x)}{\partial x_i}$$

and so the representation [df(x)] of df(x) is

$$[df(x)] = \left(\frac{\partial f_j(x)}{\partial x_i}\right)_{(j,i)} \in I(M) \times I(L)$$

Let A be a subset of C. We say that $f \in A^{(1)}$ if all partial derivatives $\frac{\partial f_j(x)}{\partial x_i}$ of f are continuous for all $x \in A$; $f \in A^{(2)}$ if

 $\int_{0}^{\infty} df \, \epsilon \, A^{(1)}$. d(df) will be denoted by d²f.

1.3 Differential Rules

In the following, we shall—present a series of rules which are direct generalizations of the corresponding ones in real variable calculus. With the usual calculus, we could prove easily all of the following rules algebraically.

Theorem 1.3.1., (Linear rule). Let L,M be nontrivial finite dimensional Hilbert spaces. Let A be a subset of L, f,g be functions into M such that f,g ϵ A⁽¹⁾. Then, for every x ϵ A, dx ϵ L, α , β ϵ R,

 $d(\alpha f + \beta g)(x)(dx) = \alpha df(x)(dx) + \beta dg(x)(dx).$

Theorem 1.3.2. (Chain rule). Let L,M,H be nontrivial finite dimensional Hilbert spaces. Let A be a subset of L,B be a subset of M,g,f be functions into H and M respectively such that $f \in A^{(1)}$, $g \in B^{(1)}$. Then, for every $dx \in L$, $x \in A$ with $y = f(x) \in B$,

d(g(f(x)))(dx) = dg(y)(df(x)(dx)).

Theorem 1.3.3. (Rule for linear functions). Let L,M be nontrivial finite dimensional Hilbert spaces and f be a linear transformation of L into M. Then, for every x, dx ϵ L,

$$df(x)(dx) = f(dx).$$

Let L, L₁, L₂, . . . , L_n be montrivial finite dimensional Hilbert spaces, $\Omega = L_1 \oplus L_2 \oplus \ldots \oplus L_n$, A Ω , f be a function into L such that f ϵ A⁽¹⁾. Let $\mathbf{x} = (\mathbf{x_i})$ be an element of the domain of f,

$$g_i(u) = f((x_1, \ldots, x_{i-1}, u, x_{i+1}, \ldots, x_n))$$

The differential of g_i at $u = x_i$ will be denoted by $\partial_{x_i} f(x)$.

Theorem 1.3.4. (Leibniz's rule). Let L, L₁, ..., L_n be non-trivial finite dimensional Hilbert spaces. Let $A \subseteq \Omega \equiv \prod_{i=1}^n L_i$ (= L₁ \oplus L₂ \oplus ... \oplus L_n) and f be a function into L such that f ϵ A⁽¹⁾. Then, for every $x = (x_i) \epsilon A$, $dx = (dx_i) \epsilon \Omega$,

$$df(x) (dx) = \sum_{i=1}^{n} \partial_{x_i} f(x) (dx_i).$$

The following result follows from Theorems 1.3.3 and 1.3.4.

Theorem 1.3.5. (Rule for multilinear functions). Let L, L_1 , L_2 , ..., L_n be nontrivial finite dimensional Hilbert spaces. Let Λ be a multilinear function of $\Omega \equiv \prod_{i=1}^n L_i$ into L. Then for every $(\mathbf{x_i})$, $(d\mathbf{x_i}) \in \Omega$,

$$d(\mathbf{x}_1 \wedge \mathbf{x}_2 \wedge \dots \wedge \mathbf{x}_n)((d\mathbf{x}_i)) = \sum_{i=1}^n \mathbf{x}_1 \wedge \mathbf{x}_2 \wedge \dots \wedge \mathbf{x}_{i-1} \wedge d\mathbf{x}_i \wedge \mathbf{x}_{i+1} \wedge \dots \wedge \mathbf{x}_n.$$

The following rule follows from Theorems 1.3.2 and 1.3.5. Theorem 1.3.6. (Multivariate product rule). Let M, L_1 , ..., L_n be nontrivial finite dimensional Hilbert spaces. Let Λ be a multilinear function of $\Omega = \bigcap_{i=1}^n L_i$ into L. Let $A \subseteq M$ and f_i be functions into L_i such that $f_i \in A^{(1)}$, $i=1,2,\ldots,n$. Then, for every $x \in A$, $dx \in M$,

$$d(f_1(x)\wedge \dots \wedge f_n(x))(dx) = \sum_{i=1}^n f_1(x)\wedge f_2(x)\wedge \dots \wedge f_{i-1}(x)\wedge (df_i(x))$$

$$(dx))\wedge f_{i+1}(x)\wedge \dots \wedge f_n(x).$$

The following result is a special case of Theorem 1.3.6.

Theorem 1.3.7. (Product rule). In Theorem 1.3.6, suppose that n=2. Then

$$df_1(x) \wedge f_2(x)(dx) = df_1(x)(dx) \wedge f_2(x) + f_1(x) \wedge df_2(x)(dx).$$

: :

We shall now generalize the usual Hadamard product $\stackrel{*}{\approx}$ and the Kronecker product $\stackrel{*}{\approx}$ for matrices: $\stackrel{*}{\approx}$ is nothing but the pointwise product for R^I ; $\stackrel{*}{\otimes}$ is a function of $R^I \times J \times R^K \times L$ into $R^{(I \times K) \times (J \times L)}$ such that for any $f \in R^I \times J$, $g \in R^K \times L$, each $R^I \times R^K \times L$ into $R^I \times R^K \times L$ such that for any $f \in R^I \times L$, $g \in R^K \times L$, each $R^I \times R^K \times L$.

* and \otimes are obviously bilinear. The products \oplus and $\widehat{\mathfrak{M}}$ defined in Singh (1972), Tracy and Singh (1972), the product \odot defined in Khatri and Rao (1968) and the product \Longrightarrow defined in Swaminathan (1976) are all bilinear. The usual matrix product, inner product and many others are also bilinear. So the above product rule can be applied to all of them.

1.4 Quadratic Differential Forms

Let L be a nontrivial finite dimensional Hilbert space and A be an open subset of L. Let f be a real-valued function such that $f \in A^{(2)}$. Let $x \in A$. We can calculate $d^2f(x)$ from df(x). However, as we shall see in Chapter Two, it is more convenient to calculate $d^2f(x)$ through the quadratic form Q(T) of $T = d^2f(x)$:

$$Q(T)(dx) = (dx, T(dx)), dx \varepsilon L.$$

Theorem 1.4.1. Let L be a nontrivial finite dimensional Hilbert space, A be an open subset of L. Let f be a real-valued function such that $f \in A^{(2)}$. Let $\{u_i\}_{i \in I(L)}$ be an orthonormal basis for L. Let $x \in A$, $dx \in L$. Then

(i)
$$\left[d^2f(x)\right] = \left(\frac{\partial^2f(x)}{\partial x_i \partial x_j}\right)$$
 and $\left[d^2f(x)\right]$ is symmetric,

where each
$$\frac{\partial^2 f(x)}{\partial x_i \partial x_j}$$
 denotes $\frac{\partial}{\partial x_i} \left(\frac{\partial f(x)}{\partial x_j} \right)$.

(ii)
$$Q(d^2f(x))(dx) = \partial_{dx}(df(x)(dx))(dx).$$

 $[d^2f(x)] \ \ above \ is \ called \ the \ \underline{Hessian \ matrix} \ \ of \ f \ at \ x \ with \\ respect \ to \ \{u_i\}.$

1.5 Representations

Theorem 1.5.1 (Linear rule). In Theorem 1.3.1,

$$[d(\alpha f + \beta g)(x)] = \alpha[df(x)] + \beta[dg(x)],$$

i.e.,

$$[d(\alpha f + \beta g)(x)(dx)] = \alpha[df(x)][dx] + \beta[dg(x)][dx].$$

Theorem 1.5.2. (Chain rule). In Theorem 1.3.2,

$$[dg(f(x))] = [dg(y)][df(x)],$$

i.e.

$$[dg(f(x))(dx)] = [dg(y)][df(x)][dx].$$

Theorem 1.5.3. (Rule for linear functions). In Theorem 1.3.3,

$$[df(x)] = [f],$$

i.e.

$$[df(x)(dx)] = [f][dx].$$

Theorem 1.5.4. (Leibniz's rule). In Theorem 1.3.4,

$$[df(x)] = [[\partial_{x_i} f(x)]]$$

(a partitioned $I(L) \times (\bigcap_{i=1}^{n} I(L_{i}))$ matrix), i.e.,

$$[df(x)(dx)] = \sum_{i=1}^{n} [\partial_{x_i} f(x)] [dx_i].$$

Theorem 1.5.5. (Rule for multilinear functions). In Theorem 1.3.5,

$$[d(\mathbf{x}_1 \wedge \mathbf{x}_2 \wedge \dots \wedge \mathbf{x}_n)] = [[\mathbf{x}_1 \wedge \mathbf{x}_2 \wedge \dots \wedge \mathbf{x}_{i-1} \wedge \dots \wedge \mathbf{x}_{i+1} \wedge \dots \wedge \mathbf{x}_n]]$$

(a partitioned I(L) x ($\prod_{i=1}^{n}$ I(L_i)) matrix), i.e.,

$$[d(\mathbf{x}_1 \wedge \mathbf{x}_2 \wedge \dots \wedge \mathbf{x}_n)(d\mathbf{x})] = \sum_{i=1}^n [\mathbf{x}_1 \wedge \dots \wedge \mathbf{x}_{i-1} \wedge \dots \wedge \mathbf{x}_{i+1} \wedge \dots \wedge \mathbf{x}_n][d\mathbf{x}_i].$$

Theorem 1.5.6. (Multivariate product rule). In Theorem 1.3.6,

$$[d(f_1(x) \wedge f_2(x) \wedge ... \wedge f_n(x))] = \sum_{i=1}^{n} [f_1(x) \wedge ... \wedge f_{i-1}(x) \wedge ... \wedge f_{i+1}(x)$$
$$\wedge ... \wedge f_n(x)][df_i(x)],$$

i.e.,

$$\left[d(f_1(x)\wedge f_2(x)\wedge \ldots \wedge f_n(x))(dx)\right] = \sum_{i=1}^n \left[f_1(x)\wedge \ldots \wedge f_{i-1}(x)\wedge \cdot \wedge f_{i+1}(x)\right]$$
$$\wedge \ldots \wedge f_n(x)\left[df_i(x)\right][dx].$$

Theorem 1.5.7. (Product rule). In Theorem 1.3.7,

$$[d(f_1(x) \land f_2(x))] = [.\land f_2(x)][df_1(x)] + [f_1(x) \land .][df_2(x)],$$

i.e.,

$$[d(f_1 \wedge f_2)(x)(dx)] = [. \wedge f_2(x)][df_1(x)][dx] + [f_1(x) \wedge .][df_2(x)][dx].$$

Let us use the term "theory of differentials" to denote the usual theory of differentials and the related topics in linear algebra, and use the term "theory of matrix derivatives" to denote those results on matrix derivatives obtained in various papers in statistics, such as Dwyer and MacPhail (1948), Dwyer (1967), Dwyer and Tracy (1969), Neudecker (1967, 1968, 1969), Vetter (1970), MacDonald and Swaminathan (1973), MacRae (1974). Some results in the theory of differentials have been presented in section 1.1.1 - 1.5.7. Most of these results are reformulations of familiar results from undergraduate analysis (Theorem 1.3.5 is exercise 2-14 of Spivak (1965))

and are presented in such a way that all of them appear to be trivial. Our main contribution here is to connect the theory of differentials to the theory of matrix derivatives. Such a connection justifies, corrects, simplifies and generalizes the theory of matrix derivatives. As an illustration, we shall pay particular attention to Theorem 1.5.7.

Lemma 1.5.8. Let Λ be a bilinear function of $M_{m_1 \times m_1} \times M_{m_2 \times m_2}$ into $\mathbb{R}^{1 \times J}$. Let $[\Lambda] = (a(r,s),(k,2)) = (a(r,s),((i,j),(k,2)))$ be the representation of Λ with respect to the usual basis. Then

- (i) $a(r,s) = \delta_{ri}\delta_{jk}\delta_{\ell s}$, if Λ is the usual matrix product with $n_1=m_2$.
- (ii) $a_{(i,j),(k,\ell)}^{(1,1)} = \delta_{ik}\delta_{j\ell}$, if Λ is the trace inner product.

(For simplicity, we shall write $a_{(i,j),(k,\ell)}$ for $a_{(i,j),(k,\ell)}$)

- (iii) $a(r,s) = \delta_{rk} \delta_{\ell s} \delta_{ir} \delta_{sj}$, if Λ is the Hadamard product.
- (iv) $a_{(i,j),(k,\ell)}^{(i',k'),(j',\ell')} = \hat{o}_{ii}, \hat{o}_{jj}, \hat{o}_{kk}, \hat{o}_{\ell\ell}, \quad \text{if } \Lambda \text{ is the Kronecker}$ product.

Proof.

(i) Let $\{e^{(i,j)}\}$, $i=1, 2, ..., m_1, j=1, 2, ..., n_1, \{f^{(k,\ell)}\}$, $k=1, 2, ..., n_1, \{\ell=1, 2, ..., n_2, k=1, 2, ..., n_1, k=1, 2, ..., n_1, k=1, 2, ..., n_2, k=1, 2,$

$$e^{(i,j)}f^{(k,\ell)} = \sum_{r,s} a^{(r,s)}_{(i,j),(k,\ell)}g^{(r,s)}.$$

So

$$(e^{(i,j)}f^{(k,\ell)})_{\alpha\beta} = \sum_{r,s} a^{(r,s)}_{(i,j),(k,\ell)} g^{(r,s)}_{\alpha\beta}$$
$$= \sum_{r,s} a^{(r,s)}_{(i,j),(k,\ell)} \delta_{r\alpha} \delta_{s\beta}$$
$$= a^{(\alpha,\beta)}_{(i,j),(k,\ell)}.$$

On the other hand,

$$(e^{(i,j)}f^{(k,\ell)})_{\alpha\beta} = \sum_{u} e^{(i,j)}f^{(k,\ell)}_{u\beta}$$

$$= \delta_{i\alpha}\delta_{jk}\delta_{\alpha\beta} .$$



S٥

$$\delta_{ri}\delta_{jk}\delta_{\ell s} = a^{(r,s)}_{(i,j)(k,\ell)}$$

(ii) Let $\{e^{\left(i,j\right)}\}$ i=1, 2,..., m, j=1, 2,..., n be the usual basis of M_{mxn} . Then

$$a_{(i,j)(k,\ell)} = e^{(i,j)}e^{(k,\ell)}$$

$$= \sum_{\alpha,\beta} e^{(i,j)}e^{(k,\ell)}$$

$$= \delta_{i\alpha}\delta_{j\beta}\delta_{k\alpha}\delta_{\ell\beta}$$

$$= \delta_{ik}\delta_{j\ell}.$$

(iii) Let $\{e^{(i,j)}\}$ i=1, 2,... m, j=1, 2,..., n be the usual basis of M_{maxn} . Then

$$(e^{(i,j)} * e^{(k,\ell)})_{\alpha\beta} = e^{(i,j)}_{\alpha\beta} e^{(k,\ell)}_{\alpha\beta}$$
$$= \tilde{o}_{i\alpha} \delta_{j\beta} \delta_{k\alpha} \delta_{\ell\beta},$$

and

$$(e^{(i,j)} \div e^{(k,\ell)})_{\alpha\beta} = \sum_{r,s} a^{(r,s)}_{(i,j)(k,\ell)} e^{(r,s)}_{\alpha\beta}$$
$$= a^{(r,s)}_{(i,j)(k,\ell)} \delta_{r\alpha} \delta_{s\beta}.$$

a(r,s) $a(i,j)(k,\ell) = \delta_{ir}\delta_{js}\delta_{kr}\delta_{\ell s}$.

(iv) Let $\{e^{(i,j)}\}\ i=1,\,2,\ldots,\,m_1,\,j=1,\,2,\ldots,\,n_1,\,\{f^{(k,\ell)}\}\ i=1,\,2,\ldots,\,m_2,\,j=1,\,2,\ldots,\,n_2\,$ and $\{g^{(j'',\ell'')}\}\ i''=1,\,2,\ldots\,m_1,\,k''=1,\,2,\ldots\,m_2,\,j''=1,\,2,\ldots\,n_1,\,k''=1,\,2,\ldots\,n_2\,$ be the usual basis of $M_{m_1}\times n_1,\,M_{m_2}$ and $R^{(\{1,\,2,\ldots,m_1\}}\times\{1,2,\ldots,m_2\})\times(\{1,\,2,\ldots,n_1\}\times\{1,2,\ldots,n_2\})$ respectively. Then

- >

$$(e^{(i,j)} \otimes f^{(f,\ell)})_{((i',k'),(j',\ell'))}$$

$$= \sum_{(i'',k''),(j'',\ell'')} a^{(i'',k''),(j'',\ell'')} g^{(i'',k''),(j'',\ell'')} g^{(i'',k''),(j'',\ell'')}$$

$$= a^{((i',k')(j',\ell'))}_{(i,j)(k,\ell)}$$

and

$$(e^{(i,j)} \otimes f^{(k,\ell)})_{((i',k')(j',\ell'))} = \delta_{ii}, \delta_{jj}, \delta_{kk'}, \delta_{\ell\ell'}.$$

The required result follows.

i.e.d.

Before presenting the representations of the differentials of various matrix products, we shall rewrite Theorem 1.5.7 in such a form that Lemma 1.5.8 can readily be used. Let

A
$$\varepsilon$$
 R^{I₁xI₂x...xI_k},

B ε R^{I_i}, A \mathfrak{I} B = $(\sum_{j_1 \varepsilon I_i} a_{(j_1, j_2, \dots, j_k)} b_{j_i})$
 ε R<sup>I₁x...xI_{i-1}xI_{i+1}x...xI_k,

B \mathfrak{I} A = A \mathfrak{I} B.</sup>

For example, let $A = (a_{ijk})$, i,j,k = 1,2, with $a_{111} = 1$, $a_{112} = 2$, $a_{121} = 3$, $a_{122} = 4$, $a_{211} = 5$, $a_{212} = 10$, $a_{221} = 7$, $a_{222} = 9$, $B = (b_2)$. $\ell = 1$, 2 with $b_1 = 1$, $b_2 = -1$. Then $B ① A = (c_{jk})$ with $c_{11} = -4$,

 $c_{12} = -8$, $c_{21} = -4$ and $c_{22} = 15$. Also B ② A = (d_{1k}) with $d_{11} = -2$, $d_{12} = -2$, $d_{21} = -2$, $d_{22} = 1$. B ① A and B ② A could be put into matrix form as follows:

$$B \textcircled{1} A = \begin{pmatrix} -4 & -8 \\ -4 & -5 \end{pmatrix}$$

and

$$B(2)A = -2 -2$$

$$-2 -1$$

Lemma 1.5.9. Let Λ be a bilinear function of $R^{I_1^{\times}I_2} \times R^{I_3^{\times}I_4}$ into $R^{I_{\times}J}$. Let $A \in R^{I_1^{\times}I_2}$, $B \in R^{I_3^{\times}I_4}$. Then

$$[A \land .] = A ② [\Lambda]$$
 , $[. \land B] = [\Lambda] ③ B$.

Proof.

Let $\mathcal{B}_1 = \{d^{k\ell}\}, \mathcal{B}_2 = \{e^{rs}\}, \mathcal{B} = \{g^{uv}\}$ be the usual bases of $R^{I_1 \times I_2}$, $R^{I_3 \times I_4}$ and $R^{I \times J}$ respectively. Write

$$[\Lambda] = (a_{(u,v),(k,\hat{z}),(r,s)}) = (a_{(k,\hat{z}),(r,s)}^{(u,v)}),$$

$$\dot{A} = \sum_{k,\ell} a_{kl} d^{k\ell}, \quad B = \sum_{r,s} b_{rs} e^{rs},$$

where $a_{k\hat{\mathcal{L}}}$, b_{rs} ϵR . Then

$$d^{k\ell} \wedge B = \sum_{r,s} b_{rs} (d^{k\ell} \wedge e^{rs})$$

$$= \sum_{r,s} b_{rs} \sum_{u,v} a^{(u,v)}_{(k,\ell),(r,s)} g^{uv}$$

$$= \sum_{u,v} (\sum_{r,s} b_{rs} a^{(u,v)}_{(k,\ell),(r,s)}) g^{uv}.$$

So

$$[.\Lambda B] = (\sum_{r,s} b_{rs} a_{(k,\ell),(r,s)}^{(u,v)}((u,v),(k,\ell))$$
$$= [\Lambda] \ \mathfrak{B}.$$

Similarly,

$$A \wedge e^{rs} = \sum_{k,\ell} a_{k\ell} [d^{k\ell} \wedge e^{rs}]$$

$$= \sum_{k,\ell} a_{k\ell} \sum_{u,v} a_{(k,\ell),(r,s)}^{(u,v)} g^{uv}$$

$$= \sum_{u,v} (\sum_{k,\ell} a_{k\ell} a_{(k,\ell),(r,s)}^{(u,v)}) g^{uv}.$$

So

$$[A\Lambda.] = (\sum_{k,\ell} a_{k\ell} a_{(k,\ell)}^{(u,v)}, (r,s))((u,v),(k,\ell))$$
$$= A ② [\Lambda].$$
q.e.d.

We close this section with a series of corollaries showing certain specializations of Theorem 1.5.10. The usual bases for $\mathbb{R}^{\mathbb{T}}$'s will be assumed.

Corollary 1.5.11. (Inner product rule). In Theorem 1.5.10, suppose that Λ is the trace inner product. (Whence $L_1 = L_2$ and L = R.)

Then

$$[d(f_1(x) \cdot f_2(x))] = f_1(x)[df_2(x)] + f_2(x)[df_1(x)].$$

Proof.

Let Y ε L₁. By Lemma 1.5.8,

$$[\Lambda] \begin{tabular}{l} \begin{$$

Similarly,

$$Y @ [\Lambda] = (\sum_{i,j} Y_{ij} \delta_{ik} \delta_{j\ell})$$
$$= (Y_{k\ell})$$
$$= Y.$$

Therefore, by Theorem 1.5.10,

$$[d(f_1(x) \cdot f_2(x))] = f_1(x)[df_2(x)] + f_2(x)[df_1(x)].$$
 q.e.d.

An equivalent version of corollary 1.5.11 can be found in Bentler and Lee (1978).

Corollary 1.5.12. (Matrix product rule). In Theorem 1.5.10, suppose that Λ is the matrix product on $M_{m_1 \times m_1} \times M_{n_1 \times n_2}$ Then

$$[d(f_1(x)f_2(x))] = (I_{m_1} \otimes (f_1(x))')[df_1(x)] + (f_1(x) \otimes I_{m_2})$$
$$. [df_2(x)].$$

Proof.

Let Y ϵ M_{m₁xn₁}, Z ϵ M_{n₁xn₂}. Then by Lemma 1.5.8,

$$Z \textcircled{2} [\Lambda] = (\sum_{i,j} Z_{ij} \delta_{ri} \hat{\delta}_{jk} \hat{\delta}_{\ell s})$$
$$= (Z_{rk} \delta_{\ell s})$$
$$= Z \textcircled{3} I_{n_2}.$$

So by Theorem 1.5.10,

$$[d(f_1(x)f_2(x))] = (I_m \otimes (f_2(x))')[df_1(x)] + (f_1(x) \otimes I_{n_2})$$

$$[df_2(x)].$$

An equivalent version of corollary 1.5.12 can be found in McDonald and Swaminathan (1973).

For any $A \in R^{i}$, $B \in R^{i-1}$, the <u>Hadamard product</u> $A \stackrel{*}{\Rightarrow} B$

. (or write as B_j^*A) is defined as the element in $R^{i=1}$ i such that

$$A = i_j B = (a_j b_{i_1, i_2, \dots, i_{j-1}, i_j, i_{j+1}, \dots, n}).$$

For example, let

$$I_1 = \{1,2\} \times \{1,2\}, I_2 = \{1,2\} \times \{1,2,3\},$$

$$C = (c_{(i,j),(k,2)}) \in R^{I_1 \times I_2}$$

$$D = (d_{i,j}) \in R^{I_1}.$$

with

$$c_{(i,j)(k,\ell)} = ijk\ell,$$

$$d_{ij} = i+j.$$

Then

$$D *_1 C = ((i+j)(ijk\varrho)) \cdot \epsilon R^{I_1 \times I_2}$$

and can be written as the following 4 x 6 matrix by arranging the indices of the rows and columns of D $\stackrel{*}{=}_1$ C in lexicographical order:

$$\begin{pmatrix}
2 \times 1 & 2 \times 2 & 2 \times 3 & 2 \times 2 & 2 \times 4 & 2 \times 6 \\
6 \times 1 & 6 \times 2 & 6 \times 3 & 6 \times 2 & 6 \times 4 & 6 \times 6 \\
6 \times 1 & 6 \times 2 & 6 \times 3 & 6 \times 2 & 6 \times 4 & 6 \times 6 \\
16 \times 1 & 16 \times 2 & 16 \times 3 & 16 \times 2 & 16 \times 4 & 16 \times 6
\end{pmatrix}$$

i.e.,

Corollary 1.5.13. (Hadamard product rule). In Theorem 1.5.10, suppose that Λ is the Hadamard product on M $_{mxn}$ x M $_{mxn}$. Then

$$[d(f_1(x)*f_2(x))] = f_1(x)*_1[df_2(x)] + f_2(x) *_1[df_2(x)],$$

i.e.,

$$[d(f_1(x)*f_2(x))] = [f_1(x)][df_2(x)] + [f_2(x)][df_1(x)],$$

where [A] is defined to be the mn x mn diagonal matrix whose diagonal elements $d_{ii}=a_i$, the entry of a 1 x mn vector $A=(a_i)$ at the position i. (See Bentler and Lee (1978)).

Proof.

Let Y ϵ M_{mxn}. By Lemma 1.5.7,

$$[\Lambda] \ \ \Im \ \ Y = (\sum_{k,\ell} Y_{k\ell} \delta_{ik} \delta_{j\ell} \delta_{ir} \delta_{sj})$$

$$= (Y_{ij}\delta_{ir}\delta_{sj})_{((r,s),(i,j))}.$$

Similarly,

$$Y \ge [\Lambda] = (Y_{k\ell} \delta_{ir} \delta_{sj}).$$

So by Theorem 1.5.10 ,

$$[d(f_1(x)*f_2(x))] = f_1(x)*_1[df_2(x)] + f_2(x)*_1[df_1(x)].$$
 q.e.d.

For any A ϵ R $(I_1xI_2)x(I_3xI_4)$, B ϵ R I_5xI_6 , define

$$A \otimes_{1} B = (a_{(i_{1},i_{2}),(i_{3},i_{4})}^{b_{(i_{5},i_{6})}}) \varepsilon R^{((I_{1}xI_{5})x(I_{2}xI_{6}))x(I_{3}xI_{4})}$$

$$A \bar{x}_{2} B = (a_{(i_{1},i_{2}),(i_{3},i_{4})}^{b_{(i_{5},i_{6})}}) \epsilon R^{(I_{1}xI_{2})x((I_{3}xI_{5})x(I_{4}xI_{6}))}$$

$$\mathbb{E} \otimes_{1} \mathbb{A} = (a_{(i_{1}, i_{2}), (i_{3}, i_{4})} b_{(i_{5}, i_{6})}) \in \mathbb{R}^{((I_{5} \times I_{1}) \times (I_{6} \times I_{2})) \times (I_{3} \times I_{4})},$$

$$B_{\tilde{x}_{2}} = (a_{(i_{1},i_{2}),(i_{3},i_{4})}b_{(i_{5},i_{6})}) \epsilon_{R}^{(I_{1}xI_{2})x((I_{5}xI_{3})x(I_{6}xI_{4}))}$$

Corollary 1.5.14. (Kronecker product rule). In Theorem 1.5.10,

suppose that Λ is the Kronecker product on R $^{\rm I}{}_{\rm 1}{}^{\rm xI}{}_{\rm 2}$ $^{\rm I}{}_{\rm 3}{}^{\rm xI}{}_{\rm 4}$. Then

$$[df_{1}(x) \otimes f_{2}(x)] = ((I_{m_{1}} \otimes I_{n_{1}}) \otimes_{1} f_{2}(x))[df_{1}(x)] + (f_{1}(x) \otimes_{1} (I_{m_{2}} \otimes I_{n_{2}}))[df_{2}(x)].$$

(1.5.1)

Proof.

Let $\{e^{ij}\}$ i=1, 2,... m_1 , j=1, 2,... n_1 , $\{h^{k\ell}\}$ k=1, 2,... m_2 , $\ell=1$, 2,... $\ell=1$, 2,

$$a_{(i,j),(k,\hat{\ell})}^{(r,s)} = (e^{ij} \otimes h^{k\hat{\ell}})_{rs}.$$

Let Y ϵ R^I1^{xI}2, Z ϵ R^I3^{xI}4. Then

$$\begin{split} [\wedge] \ \ & \exists \ Z = (\sum\limits_{k,\hat{z}} Z_k \ \delta_{ii}, \delta_{jj}, \delta_{kk}, \delta_{i\ell},) \\ & = (Z_k, \delta_{ii}, \delta_{jj},) \\ & = (I_{m_1} \otimes I_{n_1}) \otimes_1 Z. \end{split}$$

Similarly,

$$Y @ [\Lambda] = Y \otimes _{1}(I_{m_{2}} \otimes I_{n_{2}})[df_{2}(x)].$$

By Theorem 1.5.10, we obtain

$$[df_{1}(x) \otimes f_{2}(x)] = ((I_{m_{1}} \otimes I_{n_{1}}) \otimes_{I} f_{2}(x)) [df_{1}(x)] +$$

$$(f_{1}(x) \otimes_{I} (I_{m_{2}} \otimes I_{n_{2}})) [df_{2}(x)].$$
q.e.d.

We shall now compare the Kronecker product rule obtained by Bentler and Lee (1978) with corollary 1.5.14. It can be proved that the derivative of a matrix-valued function f with respect to

a matrix variable x given by Bentler and Lee (1978) is, in terms of our notations, the transpose of [df(x)]. The Kronecker product rule in Bentler and Lee (1978) states that, in terms of our notations in corollary 1.5.14,

$$[d(f_{1}(x) \otimes f_{2}(x))] = [[df_{1}(x)](I_{m_{1}n_{1}} \otimes \overline{f_{2}(x)}) + [df_{2}(x)]$$

$$(\overline{f_{1}(x)} \otimes I_{m_{2}n_{2}})][I_{m_{1}} \otimes E^{n_{2}m_{1}} \otimes I_{n_{2}}],$$

$$(1.5.2)$$

where \bar{A} is the 1 x mm row vector $[A_1, A_2, \ldots, A_m]$ with each A_i the i^{th} row of an mxm matrix A, E^{mr} is the mr x mr matrix such that, for $1 \le g \le mr$ and $1 \le h \le mr$, $e_{gh} = 1$, if g = r(j-1) + k, h = m(k-1) + j (0 < $j \le m$; 0 < $k \le r$), and $e_{gh} = 0$ otherwise. E^{mr} is called a <u>commutation matrix</u> in Magnus and Neudecker (1979) and also appears in Tracy and Dwyer (1969) and MacRae (1974) where the notation $I_{(m,r)}$ is used. As an example, let m = 2, r = 3. Then

$$\mathbf{E}^{2\times3} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

In order to compare (1.5.1) with (1.5.2), we consider $(\mathrm{I}_{\mathrm{m_1}\mathrm{n_1}} \otimes \, \overline{\mathrm{f_2(x)}}) (\mathrm{I}_{\mathrm{m}} \, \otimes \, \mathrm{E}^{\mathrm{n_2}\mathrm{m_1}} \otimes \, \mathrm{I}_{\mathrm{n}} \,) \text{ and } (\mathrm{I}_{\mathrm{m_1}} \otimes \, \mathrm{I}_{\mathrm{n_1}}) \otimes_{\mathrm{l}} \mathrm{f_2(x)}.$

Since

$$((\mathbf{I}_{\mathbf{m}_1} \otimes \mathbf{I}_{\mathbf{n}_1}) \otimes \mathbf{f}_2(\mathbf{x}))' = (\delta_{ij} \delta_{k\ell} \mathbf{f}_{2\alpha\beta})((j,\ell),((i,k)(\alpha,\beta)))$$

and

$$(I_{m_1} \otimes \overline{f_2(x)})(I_{m_1} \otimes E^{n_2m_1} \otimes I_{n_2})$$

$$= (\delta_{jj}, \delta_{\ell\ell}, f_{2\alpha'\beta'})((j,\ell)((j',\ell'), (\alpha',\beta'))$$

$$= \begin{cases} (\delta_{ij}, \delta_{\ell'k} \delta_{\alpha\alpha'} \delta_{\beta\beta'})(((j', \ell')(\alpha', \beta')), ((i,k), (\alpha,\beta))) \\ \sum_{j', \ell', \alpha', \beta'} \delta_{jj}, \delta_{\ell\ell'} f_{2\alpha'\beta'}, \delta_{ij'}, \delta_{\ell'k} \delta_{\alpha\alpha'} \delta_{\beta\beta'} ((j, \ell), ((i,k), (\alpha,\beta))) \end{cases}$$

$$(\delta_{ij}\delta_{kl}f_{2\alpha\beta}).$$

$$((\mathbf{I}_{\mathbf{m}_{1}} \otimes \mathbf{I}_{\mathbf{n}_{1}}) \otimes_{\mathbf{I}} \mathbf{f}_{2}(\mathbf{x}))' = (\mathbf{I}_{\mathbf{m}_{1}} \otimes \overline{\mathbf{f}_{2}(\mathbf{x}}))(\mathbf{I}_{\mathbf{m}_{1}} \otimes \mathbf{E}^{\mathbf{n}_{2}\mathbf{m}_{1}} \otimes \mathbf{I}_{\mathbf{n}_{2}}).$$

(1.5.3)

Similarly,

$$(f_1(x) \otimes_1 (I_{m_2} \otimes I_{n_2}))' = (\overline{f_1(x)} \otimes I_{m_2 n_2})(I_{m_1} \otimes E^{n_2 m_1} \otimes I_{n_2}).$$

$$(1.5.4)$$

(1.5.3) and (1.5.4) give a shorter and rigorous proof of the result (1.5.2) of Bentler and Lee (1978). It is more important to realize that although (1.5.1) and (1.5.2) represent the very same product rule, our (1.5.1) is simpler than (1.5.2).

In this section, we have used various product rules, especially the Kronecker product rule to illustrate the simplicity of our general, rigorous and practical approach. We can carry out similar conclusions for many rules other than the Kronecker product rule.

1.6 Examples

Example 1.6.1. Let $X, dX \in M_{pxq}$, f(X) = X'. Then

$$df(X)(dX) = (dX)$$

and

$$[df(X)] = (\delta_{is}\delta_{jr})((r,s),(i,j))$$

Example 1.6.2. Let $X, dX \in M_{pxp}$, f(X) = [X]. Then, df(X)(dX) = tr(adjX)'dX.

Hence if X^{-1} exists, then

$$df(X) = |X| trX^{-1}dX$$

and

$$[df(X)] = |X| (X^{-1})'.$$

The above result can easily be proved by Leiberz's rule and the rule for linear functions. The main result in Golberg (1972) follows from Example 1.6.2 and the chain rule.

Example 1.6.3. Let $dX \in M_{qxp}$, $A \in M_{pxq}$, $B \in M_{pxr}$. Then, since $X \to AXB$ is linear,

$$d(AXB)(dX) = A(dX)B.$$

Now by corollary 1.5.12 ,

$$[d(AXB)] = (I_q \otimes B') [d(AX)]$$

$$= (I_q \otimes B') ((A \otimes I_p)[dX])$$

$$= (I_q \otimes B') (A \otimes I_p)$$

$$= A \otimes B'.$$

Example 1.6.4. Let X, dX ϵM_{pxp} such that X^{-1} exists. Then

$$d(\ln |X|)(dX) = trX^{-1}(dX)$$

and

$$[dln | X |] = (X^{-1})^{\dagger}.$$

(ii)
$$dX^{-1}(dX) = -X^{-1}(dX)X^{-1}$$

and

$$[dx^{-1}] = -x^{-1} \otimes (x^{-1})$$
.

Proof.

We shall prove (ii). By Cramer's rule, it is easy to see that the inverse function is differentiable on its domain. Since $-xx^{-1} = I_p$,

$$0 = d(XX^{-1})(dX)$$

$$= dX(dX)X^{-1} + X(dX^{-1}(dX))$$

$$= (dX)X^{-1} + X(dX^{-1}(dX)).$$

Hence

$$dx^{-1}(dx) = -x^{-1}(dx)x^{-1}$$
.

By Corollary 1.5.12,

$$0 = (I_p \otimes X')[dX^{-1}] + X^{-1} \otimes I_p$$

Hence

$$[dX^{-1}] = -(I_{p} \otimes X')^{-1}(X^{-1} \otimes I_{p})$$

$$= -(I_{p} \otimes (X')^{-1})(X^{-1} \otimes I_{p})$$

$$= -X^{-1} \otimes (X^{-1})'.$$
q.e.d.

The function f in the following example is mentioned in Bentler and Lee (1978) to challenge the existing matrix differentiation methods which had been developed up to that date.

Example 1.6.5. (Bentler and Lee (1978, p.225)). Let

$$X = (X_1, X_2, X_3, X_4, X_5), \text{ where } X_i \in M_{m_i \times m_i}, i = 1, 2, ... 5,$$

 $Y = (X_1 \otimes X_2)X_3(X_4 + X_5)$ such that Y^{-1} exist. Let $f(X) = trY^{-1}$. By Leibniz's rule,

$$df(X)(dX) = \sum_{i=1}^{5} \partial_{X_{i}} f(X)(dX_{i}).$$

By the rules for linear functions and inverses,

$$\partial_{X_{1}} f(X) (dX_{1}) = -trY^{-1} \partial_{X_{1}} Y(dX_{1}) Y^{-1}$$

$$e^{-trY^{-1}(dX_1 \otimes X_2)X_3(X_4 + X_5)Y^{-1}}$$

Likewise,

$$\partial_{X_2} f(X) (dX_2) = -trY^{-1}(X_1 \otimes dX_2)X_3(X_4 + X_5)Y^{-1},$$

$$\partial_{\mathbb{X}_{3}} f(\mathbb{X}) (d\mathbb{X}_{3}) = -tr \mathbb{Y}^{-1} (\mathbb{X}_{1} \otimes \mathbb{X}_{2}) (d\mathbb{X}_{3}) (\mathbb{X}_{4} * \mathbb{X}_{5}) \mathbb{Y}^{-1},$$

$$\partial_{X_4} f(X) (dX_4) = -trY^{-1}(X_1 \otimes X_2)X_3(dX_4 * X_5)Y^{-1},$$

$$\partial_{\mathbb{X}_{5}} f(\mathbb{X}) (d\mathbb{X}_{5}) \; = \; - \text{trY}^{-1} (\mathbb{X}_{1} \otimes \mathbb{X}_{2}) \mathbb{X}_{3} (\mathbb{X}_{4} \dot{\times} d\mathbb{X}_{5}) \mathbb{Y}^{-1}.$$

CHAPTER TWO

Maxima And Minima in Linear Models and Multivariate Analysis

2.1 Introduction

Optimization problems in statistics can be divided up into two types: one which has the optimal solutions appeared in the interior of the region of concern and the other one which has the optimal solutions appeared on the boundary of the region of concern. For the first type, matrix differentiation is an important tool. In this chapter, we shall illustrate how matrix differentiation could be applied in various situations.

2.2 Preliminaries

In order to prove the optimality results in the following sections, we need several results which are of interest in their own right.

 P_n will denote the set of all nxn positive definite matrices; S_n will denote the Hilbert space of all nxn symmetric matrices over R; SpA will denote the range of the spectrum, i.e., the set of all eigenvalues of A. It is well-known that A ϵ P_n implies that SpA \subset (o,∞) . Moreover, Sp(AB) \subset $[o,\infty)$ if A ϵ P_n, B ϵ M_{nxn} and B is nonnegative definite.

Theorem 2.2.1. Let A, B ϵ P_n such that A \neq B. Then

(i)
$$Sp((A-B)(A^{-1}-B^{-1})) = (-\infty, 0],$$

(ii)
$$Sp((A-B)(A^{-1}-B^{-1})) \neq \{o\}.$$

Hence,

$$tr(A-B)(A^{-1}-B^{-1}) < 0.$$

Recently, Ky Fan generalized the above result to the following:

Let A be a non-singular Hermitian matrix of order n, with p positive eigenvalues and q negative eigenvalues, p+q=n. Let B be a positive definite Hermitian matrix of order n, and let

$$C = (A^{-1}-B^{-1})(A-B)$$
.

Then all the eigenvalues of C are real, p of them non-positive, and q of them are positive. Furthermore, if exactly r of the eigenvalues of $B^{-1/2}AB^{-1/2}$ (or the similar matrix AB^{-1}) are equal to 1, then C has exactly r eigenvalues equal to zero.

Theorem 2.2.2. Let A, B, W ϵ P_n. Then

$$tr(AWA-BWB)(A^{-1}-B^{-1}) < 0.$$

Hence

$$tr(A^2-B^2)(A^{-1}-B^{-1}) < 0.$$

.Proof.

Since A, B are positive definite, there exist a non-singular matrix P and a diagonal matrix D such that

$$A = PP^{\dagger}$$
.

So

Let $E = I - D^{-1} - D^2 + D$. Then E is diagonal. Let a_i be the entry of D at (i,i). Then the entry b_i of E at (i,i) is

$$b_{i} = 1 - \frac{1}{a_{i}} - a_{i}^{2} + a_{i}$$

$$= -\frac{1}{a_{i}} (a_{i}^{3} - a_{i}^{2} - a_{i} + 1)$$

$$= -\frac{1}{a_{i}} (1 - a_{i})^{2} (1 + a_{i}) \le 0.$$

= tr $P'WP(I-D^{-1}-D^{2}+D)$.

At least one of b_i < 0, otherwise all a_i = 1, i.e., A = B, a contradiction to $A \neq B$. Since (c_{ij}) = P'WP is positive definite, all $c_{ii} > 0$, So

$$tr(AWA-BWA)(A^{-1}-B^{-1}) = \sum_{i=1}^{n} c_{ii}b_{i} < 0.$$
 q.e.d.

The following result is important in justifying the differentiability of certain useful functions in statistics and is overlooked by various authors.

Theorem 2.2.3. P_n is open in S_n .

Proof.

Let $C_0 = (c_{oij}) \in P_n$, $B = \{Z \in \mathbb{R}^n : \|Z\| = 1\}$, $x \in B$. Consider $f_x : f_x(A) = x'Ax$, $A \in S_n$. Then f_x is linear on S_n and is therefore continuous.

$$|f_{x}(A)| = |\sum_{i,j} x_{i}a_{ij}x_{j}|$$

$$\leq ||A|| \sum_{\infty} \sum_{j=1}^{n} |x_{i}||x_{j}|$$

$$\leq ||A|| n^{2}.$$

Let $M = n^2$. Then $|f_x(A)| \le M |A|$ for all $x \in B$ and $A \in S_n$. Since $x \to x'C_0x$ is continuous on R^n and B is compact, $m \equiv \min \{x'C_0x : x' \in S_n\}$

 $x \in B$ exists. Since $C_0 \in P_n$, m > 0. Let $e = \frac{m}{2M}$, $C \in S_n$ with $C - C_0 \parallel_{\infty} < e$, $x \in B$. Then $f_x(C) = f_x(C - C_0) + f_x(C_0) \ge -M \parallel C - C_0 \parallel_{\infty} + m > -Me + m > 0$. So $C \in P_n$ and hence P_n is open in S_n .

q.e.d.

2.3 Contingency Tables

Consider the problem of maximum likelihood estimation and testing using likelihood ratio procedures regarding contingency tables. For simplicity, we demonstrate the case of rxk tables. The case of higher dimensions can be done similarly. The maximum likelihood function based on a sample $\{x_{ij}\}$ that forms a rxk table is

$$L(\theta) = c \prod_{i=1}^{r} \prod_{j=1}^{k} (\theta_{ij})^{x_{ij}},$$

where

$$c = \frac{n!}{\frac{r}{r} \frac{k}{k}}$$

$$. \prod_{i=1}^{n} \prod_{j=1}^{n} (x_{ij}!)$$

$$\theta \ \epsilon \ \widehat{\underline{\mathbb{H}}} = \{(\theta_{ij})_{(i,j)\neq \ (r,k)} : \theta_{ij} > o, \sum_{(i,j)\neq (r,k)} \theta_{ij} < 1\},$$

Here x_{ij} is the number of sample values belonging to the cell I_{ij} and $n\theta_{ij}$ is the theoretical number of sample values belonging to

 I_{ij} . We wish to find the maximum likelihood estimate $\widehat{\theta}((\mathbf{x}_{ij}))$ of θ , i.e., $\widehat{\theta}((\mathbf{x}_{ij}))$ ε Θ with $L(\widehat{\theta}((\mathbf{x}_{ij}))) = \max L(\Theta)$. It can be proved that $\widehat{\theta}((\mathbf{x}_{ij}))$ exists if and only if all \mathbf{x}_{ij} 's are positive. We shall assume that all \mathbf{x}_{ij} 's are positive. Note that Σ \mathbf{x}_{ij} is equal to the sample size n. Let $f = -\ln(\frac{L}{c})$. Then $f(\widehat{\theta}((\mathbf{x}_{ij}))) = \min f(\Theta)$. Θ may be considered as an open subset of \mathbb{R}^{rk-1} . Let $\theta \in \Theta$. Then

$$f(\theta) = -\sum_{i=1}^{r} \sum_{j=1}^{k} x_{ij} \ln \theta_{ij}.$$

Let $d\theta = (d\theta_{ij}) \in \mathbb{R}^{rk-1}$. Then by Leibniz's rule,

$$df(\theta)(d\theta) = \left(-\frac{x_{ij}}{\theta_{ij}} + \frac{x_{rk}}{\theta_{rk}}\right) \cdot (d\theta_{ij}),$$

where is the usual inner product. Let $\theta_1 = (\theta_{1ij})$, $\theta_2 = (\theta_{2ij})$ such that $\theta_1 \neq \theta_2$. Then

$$\begin{aligned} (\mathrm{df}(\theta_1) - \mathrm{df}(\theta_2)) (\theta_1 - \theta_2) &= \sum_{(i,j) \neq (r,k)} ((-\frac{x_{ij}}{\theta_{1ij}} + \frac{x_{rk}}{\theta_{1rk}}) - \\ (-\frac{x_{ij}}{\theta_{2ij}} + \frac{x_{rk}}{\theta_{2rk}})) (\theta_{1ij} - \theta_{2ij}) \end{aligned}$$

$$= \sum_{(i,j) \neq (r,k)} \frac{x_{ij} (\theta_{1ij} - \theta_{2ij})^2}{\theta_{1ij} \theta_{2ij}} +$$

$$x_{rk}(\frac{1}{\theta_{lrk}} - \frac{1}{\theta_{2rk}}) \sum_{(i,j)\neq(r,k)} (\theta_{lij} - \theta_{2ij}).$$

Since
$$\sum_{i=1}^{r} \sum_{j=1}^{k} \theta_{1ij} = \sum_{i=1}^{r} \sum_{j=1}^{k} \theta_{2ij} = 1$$
,

$$(df(\theta_1) - df(\theta_2))(\theta_1 - \theta_2) = \sum_{(i,j) \neq (r,k)} \frac{x_{ij}(\theta_{1ij} - \theta_{2ij})^2}{\theta_{1ij}\theta_{2ij}} +$$

$$=\sum_{i=1}^{r}\sum_{j=1}^{k}\frac{x_{ij}(\theta_{1ij}-\theta_{2ij})^{2}}{\theta_{1ij}\theta_{2ij}}>0.$$

Thus, df is strictly monotone on \bigoplus (see, e.g., Opial (1967), p.84). So f is strictly convex on \bigoplus . Let $df(\theta) = 0$. Then

$$\frac{\mathbf{x}_{ij}}{\theta_{ij}} = \frac{\mathbf{x}_{rk}}{\theta_{rk}} = \frac{\mathbf{x}_{rk}}{\mathbf{x}_{ij}} = \mathbf{x}_{ij}$$

$$\frac{\mathbf{x}_{ij}}{\mathbf{x}_{rk}} = \frac{\mathbf{x}_{rk}}{\mathbf{x}_{ij}} = \mathbf{x}_{ij}$$

$$\mathbf{x}_{i=1}$$

$$\mathbf{x}_{ij} = \mathbf{x}_{rk}$$

i.e.,
$$\hat{\theta}((x_{ij})) = (\frac{x_{ij}}{n})$$
.

We now consider the problem of testing the hypothesis of independence, i.e., we want to test

$$H_{o}: \theta_{ij} = \theta_{i}, \theta_{.j} \text{ for all } i=1, 2, \ldots, r, \ j=1, 2, \ldots, k.$$
 against

$$H_1: \theta_{ij} \neq \theta_{i}, \theta_{ij}$$
 for some i, j,

where

$$\theta_{i} = \sum_{j=1}^{k} \theta_{ij}, \quad \theta_{ij} = \sum_{i=1}^{r} \theta_{ij}.$$

Extend $\hat{\theta}$ such that $\hat{\theta}((x_{ij})) = (\frac{x_{ij}}{n})$ for all nonegative intergers with $\sum x_{ij} \le n$, $(x_{ij}) \in \mathbb{R}^{rk-1}$. We first assume that all x_{ij} 's are positive. Then, under H_0 , the likelihood function L_1 is:

$$L_{1}(\alpha) = c \prod_{i=1}^{r} \prod_{j=1}^{k} (\theta_{i}, \theta_{\cdot j})^{x_{ij}}$$

where

$$\alpha = ((\theta_{i}^{p}), (\theta_{-j})) \in \bigoplus_{0} = \{(a,b) \in \mathbb{R}^{r-1} \times \mathbb{R}^{k-1} : a_{i} > 0, b_{i} > 0, \\ \sum_{i=1}^{r-1} a_{i} < 1, \sum_{j=1}^{k-1} b_{j} < 1\},$$

$$\sum_{i=1}^{r} \theta_{i} = 1 = \sum_{j=1}^{k} \theta_{i}.$$

المكاف

Note that Θ_o is open and convex in R^{r+k-2} . Let $f = -\ln(\frac{L_1}{c})$, $d\alpha = ((d\theta_{i}), (d\theta_{i})) \in \Theta_o$. Then

$$\mathrm{df}(\alpha)(\mathrm{d}\alpha) = \left(\left(-\frac{x_{\underline{i}}}{\theta_{\underline{i}}} + \frac{x_{\underline{r}}}{\theta_{\underline{r}}}\right), \left(-\frac{x_{\underline{j}}}{\theta_{\underline{j}}} + \frac{x_{\underline{k}}}{\theta_{\underline{k}}}\right)\right). \left(\left(\mathrm{d}\theta_{\underline{i}}\right), \left(\mathrm{d}\theta_{\underline{j}}\right)\right).$$

Let
$$\alpha_1 = ((\theta_{1i},)(\theta_{1\cdot j})), \quad \alpha_2 = ((\theta_{2i},)(\theta_{2\cdot j})) \in \bigoplus_0 \text{ with } \alpha_1 \neq \alpha_2.$$
Then

$$(\mathrm{df}(\alpha_1) - \mathrm{df}(\alpha_2))(\alpha_1 - \alpha_2) = \sum_{i=1}^{r} \frac{x_i \cdot (\theta_{1i} - \theta_{2i} \cdot)^2}{\theta_{1i} \cdot \theta_{2i}} +$$

$$\sum_{j=1}^{k} \frac{x_{\cdot j} (\theta_{i \cdot j} - \theta_{2 \cdot j})^2}{\theta_{1 \cdot j} \theta_{2 \cdot j}}$$

> 0.

Therefore f is strictly convex, L is strictly concave and $(\widehat{\theta}_{\underline{i}}),$ $(\widehat{\theta}_{\underline{i}})$ with

$$\hat{\theta}_{i}.((x_{i}.)) = \frac{x_{i}.}{n}, \hat{\theta}_{i}.((x_{i}.)) = \frac{x_{i}}{n}$$

are the maximum likelihood estimates of (θ_i) , (θ_j) under H_0 . Extend $(\hat{\theta}_i)$, $(\hat{\theta}_j)$ such that

$$(\widehat{\theta}_{i})((x_{i})) = (\frac{x_{i}}{n}),$$

$$(\widehat{\theta}_{i})((x_{i})) = (\frac{x_{i}}{n}),$$

for all x_i , $x_{.j} \in \{0, 1, 2, ..., n\}$ such that $\sum_{i} \le n$, $\sum_{i} \le n$. Let

$$\lambda((\mathbf{x}_{ij})) = \frac{L(\widehat{\theta}((\mathbf{x}_{ij})))}{L_1((\widehat{\theta}_{i\cdot})((\mathbf{x}_{i\cdot})), (\widehat{\theta}_{\cdot j})((\mathbf{x}_{\cdot j})))}$$

Then

$$\lambda((\mathbf{x}_{ij})) = \frac{\prod_{i=1}^{r} \prod_{j=1}^{k} \left(\frac{\mathbf{x}_{ij}}{n}\right)^{\mathbf{x}_{ij}}}{\prod_{i=1}^{r} \prod_{j=1}^{k} \left(\frac{\mathbf{x}_{i}}{n} \cdot \frac{\mathbf{x}_{\cdot j}}{n}\right)^{\mathbf{x}_{ij}}}$$

$$= \prod_{i=1}^{r} \prod_{j=1}^{k} \left(\frac{n\mathbf{x}_{ij}}{\mathbf{x}_{i} \cdot \mathbf{x}_{\cdot j}}\right)^{\mathbf{x}_{ij}}.$$

It is intuitively clear that we shall reject H_0 when $\lambda((x_{ij}))$ is large, say larger than a constant c. The test with the above region of rejection will be denoted by ϕ_c . The first type risk of ϕ_c is determined by c (but not vice versa). ϕ_c is often referred to as a maximum (supremum) likelihood ratio procedure (Lehman (1959), p.15).

2.4 Maximum likelihood estimates of multivariate normal model Let C, D ϵ $P_{\rm n}$,

$$f(C) = 1/2 \text{ Nin } C[-1/2 \text{ trCD}, C \in P_n.$$

In finding the maximum likelihood estimate of (μ, Σ) based on a random sample (x_{α}) of a normal population of mean μ and dispersion matrix Σ , it is necessary to prove that the maximum value of $f(P_n)$ is $f(ND^{-1})$ (Anderson (1958)). Let

 $g(C) = 1/2 \text{ Nln} | C | - 1/2 \text{ trCD}, C \in M_{nxn}, (C) > 0.$

44

Let h be the identity function on S_n , as a function of S_n into $M_{n\times n}$. Then

$$f(C) = g(h(C)), C \in P_n$$

By Theorem 2.2.3 and the chain rule

$$df(C)(dC) = 1/2 tr(NC^{-1}-D)(dC), C \epsilon P_n, dC \epsilon S_n$$

Let C_1 , $C_2 \in P_n$ such that $C_1 \neq C_2$. Then by Theorem 2:2.1,

$$(df(C_1)-df(C_2))(C_1-C_2) = \frac{N}{2} tr(C_1^{-1} - C_2^{-1})(C_1-C_2)$$

< 0.

So f is strictly concave on P_n . Since $df(ND^{-1}) = 0$, $f(ND^{-1}) = \max f(P_n)$.

The role of h is important because P_n is open in S_n but has empty interior in $M_{n\times n}$. One may compare the above proof with that of Anderson (1958, p.45-47), Smith (1978) and Watson (1964). Our proof is constructive, rigorous and simple. Note that by the result of Dykstra (1971), the maximum likelihood estimator $\widehat{\Sigma}$ of Σ exists if and only if $N \ge n$.

2.5 Multivariate regression models

The likelihood function based on a sample $x_j \in N(BZ_j, \Sigma_{n\times n})$ j=1, 2,..., N is

$$L(B,\Sigma) = (2\pi)^{-\frac{Nn}{2}} |\Sigma^{-1}|^{\frac{N}{2}} \exp \left[-1/2 \sum_{j=1}^{N} (x_{j}^{-}BZ_{j}) | \Sigma^{-1}(x_{j}^{-}BZ_{j})\right],$$

where B ϵ M $_{\rm nxq}$, Σ ϵ P $_{\rm n}$ are unknown parameters. We want to find the maximum likelihood estimate of (B, Σ). (See, for example, Anderson (1958), Chapter 8).

Let

$$f(B,C) = \ln \left(\frac{L(B,C^{-1})}{\frac{Nn}{2}} \right)$$
, $B \in M_{nxq}$, $C \in P_n$.

We proceed to find \hat{B} first. Fix C ϵ P_n and let f_C(B) = f(B,C). Then

$$f_{C}(B) = \frac{N}{2} \ln |C| - 1/2 \sum_{j=1}^{N} (x_{j} - BZ_{j})C(x_{j} - BZ_{j}), \quad B \in M_{mxq},$$

and

$$df_{C}(B)(dB) = \sum_{j=1}^{N} (x_{j}-BZ_{j})'C(dB)Z_{j}.$$

Let B_1 , $B_2 \in M_{nxq}$. Since C is positive definite,

$$(df_{C}(B_{1})-df_{C}(B_{2}))(B_{1}-B_{2}) = -\sum_{j=1}^{N} Z_{j}'(B_{1}-B_{2})'C(B_{1}-B_{2})Z_{j}$$

$$\leq 0.$$

So f_C is concave on P_{n} . Let $df_C(B) = 0$, $dB \in M_{nxp}$. Then $df_C(B)(dB) = 0$, i.e.,

$$\sum_{j=1}^{N} (x_j - BZ_j)'C(dB)Z_j = tr(\sum_{j=1}^{N} Z_j(x_j - BZ_j)')C(dB) = 0.$$

Hence

$$\sum_{j=1}^{N} Z_{j}(x_{j}-BZ_{j})'C = 0.$$

Since $C \in P_n$, $\widehat{B} = (\sum_{j=1}^{N} x_j Z_j')(\sum_{j=1}^{N} Z_j Z_j')^{-1}$ is the solution for $df_C(B) = 0$. Since f_C is concave on P_n , $f_C(\widehat{B}) = \max_{j=1}^{N} f_C(M_{nxq})$. Let $g(C) = f(\widehat{B}, C)$,

$$D = \sum_{j=1}^{N} (x_j - \widehat{B}Z_j) (x_j - \widehat{B}Z_j)'.$$

Then, by the result in section 2.4, $g(ND^{-1}) = \max g(P_n)$. Thus $(\widehat{B}, \frac{1}{N}, D)$ is the maximum likelihood estimate of (B, C).

2.6 Multivariate Linear Hypotheses

It is well known that the model discussed in section 2.5 includes certain models in regression analysis as well as analysis of variance which have various applications, especially in econometrics and psychometrics. The following related problem assumes

patterns in the dispersion matrix Σ_{nxn} of the population. It could be viewed as a general approach to the problem of analysis of variance components (see, for example, Rao (1973)). For details of the problem, one could consult Roger and Young (1978) and the references there. The log-likelihood function L based on the data Y is given by

$$L(\beta, \varphi) = -\frac{N}{2} \left\{ \text{n ln} 2\pi - \text{ln} \left[\Sigma^{-1} \right] + \frac{1}{N} \text{tr}(Y - H\beta) \cdot (Y - H\beta) \Sigma^{-1} \right\},$$

$$f_{\phi}(\beta) = L(\beta, \phi), \quad \beta \in M_{axa}$$

Let $d\beta \in M_{qxn}$. Then

$$df_{0}(\beta) (d\beta) = tr (Y - H\beta) \Sigma^{-1} (H(d\beta))'$$

Let β_1 , $\beta_2 \in M_{qxn}$. Since Σ^{-1} is positive definite,

Hence f_{φ} is concave on M_{qxn} . $df_{\varphi}(\beta) = 0$ is equivalent to

$$H (Y-H\beta) \cdot \Sigma^{-1} = 0$$

٥r

$$H'H\beta = H'Y$$

which is called the <u>normal equation</u> for β .

Since $H^+ = (H'H)^+ H'$,

$$\hat{\beta}_{\cdot} = H^{\dagger}Y + (I_{q} - H^{\dagger}H) Z,$$

where Z is any qxn matrix and H is the Moore - Penrose inverse of H. As before, let $f_{\widehat{\beta}}(\psi) = L(\widehat{\beta}, \phi)$. Note that $\psi + \phi$ is a one-to-one mapping and hence the maximum likelihood estimate $\widehat{\phi}$ of ϕ corresponds to the maximum likelihood estimate $\widehat{\psi}$ of ψ . Now, let $d\psi = (d\psi_g) \in G$. Then, by Leibniz's rule,

$$d \Sigma^{-1}(d\psi) = \sum_{g=1}^{m} \partial_{\psi} \Sigma^{-1}(d\psi_{g})$$
$$= \sum_{g=1}^{m} G_{g} d\psi_{g}$$

ঔ্ট্য

and so

$$\mathrm{d}\hat{\beta}(\psi),\;(\mathrm{d}\psi)\;=\;-\;\frac{N}{2}\;\left[\mathrm{tr}\;\Sigma\;\left(\stackrel{m}{\Sigma}\;\mathsf{G}_{g}\;\mathrm{d}\psi_{g}\right)\;-\;\frac{1}{N}\;\mathrm{tr}\;\mathsf{B}\;\left(\stackrel{m}{\Sigma}\;\mathsf{G}_{g}\mathrm{d}\psi_{g}\right)\right],$$

where

$$B = (Y - H\hat{\beta})' (Y - H\hat{\beta})$$

$$= Y' (I - HH^{+}) Y.$$

Let $\psi_1 = (\psi_{1g})$, $\psi_2 = (\psi_{2g}) \in G$ such that $\psi_1 \neq \psi_2$ and

$$\Sigma_{i} = \sum_{g=1}^{m} \phi_{ig} G_{g}, i = 1, 2.$$
 Then

$$\begin{split} (\mathrm{d} f_{\widehat{\beta}}(\psi_1) \; - \; \mathrm{d} f_{\widehat{\beta}}(\psi_2)) \; \; (\psi_1 \; - \; \psi_2) \; &= \; - \; \frac{N}{2} \; [\, \mathrm{tr}(\Sigma_1 - \Sigma_2) \; (\; \stackrel{m}{\Sigma} \; \; G_g(\psi_1 - \psi_2 - \psi_2 - \psi_2)) \,] \\ &= \; - \; \frac{N}{2} \; (\, \mathrm{tr}(\Sigma_1 - \Sigma_2) \; (\Sigma_1^{-1} - \Sigma_2^{-1}) \,) \end{split}$$

So $\hat{f}_{\hat{\beta}}$ is strictly concave. Let $d\hat{f}_{\hat{\beta}}(\psi)=0$. Then $d\hat{f}_{\hat{\beta}}(\psi)$ $(d\psi)=0$,

for every dψ ε G and

i.e.,

tr
$$(\sum_{h=1}^{m} \psi_h G_h)G_g = \frac{1}{N}$$
 tr BG_g for all $g = 1, 2, ..., m$.

Since $\operatorname{tr} G_h G_g = \delta_{hg}$, $\frac{1}{N} \operatorname{tr} BG_g$ is the solution of $\operatorname{df}_{\widehat{\beta}} (\psi) = 0$ and, hence, is the maximum likelihood estimate $\widehat{\psi} = (\widehat{\psi}_g)$ of ψ .

2.7 Quadratic Estimates

Consider the regression model $Y=X\beta+\epsilon$, where X is a known nxp matrix of rank p, ϵ is a normal random vector with mean 0 and dispersion matrix σ^2I . We wish to find a quadratic estimate Y'AY of σ^2 that has the smallest mean square error $\frac{\pi}{\delta}$ $E((Y'AY-\sigma^2)^2)$, $A \in S_n$. Theil and Schweitzer (1961) showed that it is equivalent to finding $A \in S_n$ such that A minimizes

$$f(A) = 2 \operatorname{tr} A^2 + (1 - \operatorname{tr} A)^2$$
, $A \in S_n$, $AX = 0$.

Calvert and Seber (1978) found the desired A using nearest point projections in Hilbert spaces. We shall obtain the desired A in a simple, constructive and rigorous way. Let

$$g(A) = f(A) - tr N'AX,$$

where N' is a Lagrange multiplier. Let $dA \in S_n$. Then dg(A) (dA) = 4 trA(dA) - 2(1-trA) trdA - trN'(dA)X $= tr (4A - 2(1-trA)I_n - XN')dA.$

Let A_1 , A_2 & S_n such that $A_1 \neq A_2$. Since $A_1 - A_2$ is symmetric, $(dg(A_1) - dg(A_2)) (A_1 - A_2) = tr [4(A_1 - A_2) + 2 tr (A_1 - A_2)] (A_1 - A_2)$ $= 4tr (A_1 - A_2)^2 + 2 (tr(A_1 - A_2))^2$ > 0

Therefore g is strictly convex on S_n . Let dg(A) = 0. Then

$$4A - 2(1-trA)I_n - NX' = 0.$$
 (2.7.1)

Multipling X on both sides of (2.7.1), we obtain, by imposing AX=0,

$$N = -2 (1 - trA) X (X'X)^{-1}. (2.7.2)$$

Substituting (2.7.2) in (2.7.1), we have

$$4A = 2(1 - trA) [I_n - X(X'X)^{-1} X'] = 0.$$

So by taking the trace,

$$trA = \frac{n-p}{n-p+2} .$$

Therefore

$$\hat{A} = \frac{1}{n-p+2} [I_n - X(X'X)^{-1} X']$$

satisfies dg(A) = 0. Since g is convex on S_n , $g(\widehat{A}) = \min g(S_n)$. Now by the theory of Lagrange multipliers (see, for example, Lehman (1959), p. 87), $f(\widehat{A}) = \min f(S_n)$. Note that A is never positive definite.

2.8 Minimum Distance and Principal Components

Let $A \in P_h$. We want to find the absolute maximum (minimum) value of f(x) = x'x, $x \in R^h$ subject to the side condition x'Ax = a, where a is a fixed, positive real number. This problem is a mathematical abstraction of the principal component method in multivariate analysis (Anderson (1958, Chapter II)). In statistics, $A = \Sigma^{-1}$, where Σ is the dispersion matrix of a normal random vector. The problem is to compare the Euclidean norm with the norm $\|\cdot\|_{\Sigma}$ of the reproducing kernel Hilbert space associated with Σ . Euclidean norm is the norm of a reproducing kernel Hilbert space associated with an identity matrix. In general, let C be an nxn nonnegative definite matrix,

$$(x,y)_{C} = x' C^{\dagger} y, \qquad x,y \in \mathbb{R}^{h}.$$

Then $(\ ,\)_C$ is a pseudo inner product, i.e., $(\ ,\)_C$ is bilinear, symmetric and $(x,x)\geq 0$ for all $x\in R^h$. $(\ ,\)_C$ is an inner product if and only if C^{-1} exists. Let $B=C^+$. Then B is nonnegative definite. Now to compare $\|\ \|_C$ with $\|\ \|_{\Sigma}$, we consider

$$f(x) = x' B x subject to x'Ax = a.$$

The problem raised earlier is this problem with $B = I_h$. Using the method of Lagrange multipliers, we consider

$$g(x) = f(x) + \lambda (a - x'Ax)$$

$$= x'Bx + \lambda (a - x'Ax),$$

where λ is a Lagrange multiplier. Let $dx \ \epsilon \ R^{\mbox{\scriptsize h}}.$ Then

$$dg(x)(dx) = 2x'Bdx - 2\lambda x'A(dx)$$
$$= 2x'(B-\lambda A)(dx).$$

Let dg(x) = 0. Then $Bx = \lambda Ax$ or $A^{-1}Bx = \lambda x$. Suppose that we are interested in finding the absolute maximum value of f subject to x'Ax = a. Let $\lambda_{max} = \max Sp (A^{-1}B)$, e_{max} be a nonzero eigenvector of $A^{-1}B$ corresponding to λ_{max} , $c = e'_{max} A e_{max}$, $x_{max} = (\frac{a}{c})^{1/2} e_{max}$. Then $x'_{max} A x_{max} = a$, $B x_{max} = \lambda A x_{max}$ and $B - \lambda_{max} A$ is nonnegative definite. Let $x, y \in R^h$. Then

$$(dg(x)-dg(y)) (x-y) = 2(x-y)' (B-\lambda_{max}A) (x-y)$$

 $\leq 0.$

Thus g is concave on R^h and $g(x_{max}) = \max g(R^h)$. So x_{max} maximizes x'Bx subject to x'Ax = a, x ϵ R^h . Similarly, we may use $\lambda_{min} = \min Sp (A^{-1}B)$ to find the absolute minimum value of f subject to x'Ax = a.

The above argument can be carried out for some other problems. For example, Bush and Olkin (1959), Rao (1973) and Goldberger (1964) considered the problem of minimizing x'Ax subject to B'x = u, where $A \in P_h$ and u belongs to the column space of B'. As before, let

$$f(x) = x'Ax + \lambda'(u-B'x)$$

with λ a Lagrange multiplier. Then, for every dx ϵ R^h,

$$df(x)(dx) = 2x'A(dx) - \lambda'B'(dx)$$
$$= (2x'A - \lambda'B')(dx).$$

Let x, $y \in R^h$ such that $x \neq y$. Since A is positive definite,

$$(df(x) - df(y)) (x-y) = 2 (x-y)^{\dagger}A(x-y)$$

> 0.

So f is strictly convex on R^h . Choose $\lambda = 2(B'A^{-1}B)^-u$, $x = A^{-1}B(B'A^{-1}B)^-u$. Then df(x) = 0 and B'x = u. Hence x above minimizes x'Ax subject to B'x = u.

In addition to principal components, we shall give another applications of the result $g(x_{max}) = max \ g(\ R^h)$ obtained above. Using the earlier notations, let us consider

$$h(x) = \frac{x'Bx}{x'Ax} , \quad x \in \mathbb{R}^h \setminus \{0\}.$$

Then

$$\max \quad h(x) = \max \quad \max \quad \frac{x'Bx}{a}$$

$$x \in \mathbb{R}^{h} \setminus \{0\}$$

$$= \frac{1}{e'_{\max} A} \frac{A^{-1}B}{e_{\max}} e'_{\max} B e_{\max}$$

$$= \frac{e'_{\max} A A^{-1}B}{e'_{\max} A e_{\max}}$$

Hence for $x \in \mathbb{R}^h \setminus \{0\}$,

$$\lambda_{\min} \le \frac{x'Bx}{x'Ax} \le \lambda_{\max}$$

and the inequalities are sharp.

2.9 Factor Analysis

In the following, we shall give an example that the likelihood function of parameters in a given model is not concave and hence the method of monotone operators fail. This pathological example shows that statisticians have to find some other method to make sure that the maximum likelihood estimates they

١,

-

find are indeed the absolute maximum value of the likelihood function, especially when the estimates are obtained numerically by iteration methods.

Consider a data matrix $x=(x_{\alpha i})$ of N observations on p response variables. x is given to be an observation of a random matrix $X=(x_{\alpha i})$ with N independent rows, each having a multivariate normal distribution with the same dispersion matrix Σ of the form

$$\Sigma = \wedge \wedge^{\dagger} + \Psi$$

where Ψ s P_p and is diagonal. This model is a familiar factor analysis model (see, for example, Lawley and Maxwell (1963)). We want to estimate Ψ . The log-likelihood function L_1 , as a function of Ψ only, is

$$L_{1}(\Psi) = c - \frac{N}{2} \ln |\Sigma| - \frac{N}{2} \operatorname{tr} (A\Sigma^{-1}),$$

where c is a constant, A is nonnegative definite and does not depend on Σ . Let $\Psi=(\delta_{ij}\Psi_i)$. Then Ψ can be regarded as (Ψ_i) in \mathbb{R}^p . Let $d\Psi$ s \mathbb{R}^p . Then

$$\begin{split} dL_{1}(\Psi)(d\Psi) &= -\frac{N}{2} \text{ tr } \Sigma^{-1} d\Sigma(d\Psi) - \frac{N}{2} \text{ trAd} \Sigma^{-1}(d\Psi) \\ &= -\frac{N}{2} \text{ tr } \Sigma^{-1}(d\Psi) + \frac{N}{2} \text{ trA} \Sigma^{-1}(d\Psi) \Sigma^{-1} \\ &= -\frac{N}{2} \text{ tr } (\Sigma^{-1} - \Sigma^{-1} A \Sigma^{-1})(d\Psi) \,. \end{split}$$

Let $\Psi_1^{},\;\Psi_2^{}$ be two distinct values for $\Psi.$ Then

$$\begin{split} \Delta & \equiv \left(\text{d} \mathbb{L}_1(\Psi_1) - \text{d} \mathbb{L}_1(\Psi_2) \right) \ (\Psi_1 - \Psi_2) \ = \ - \frac{N}{2} \left[\text{tr} (\Sigma_1^{-1} - \Sigma_2^{-1}) (\Psi_1 - \Psi_2) \right] \\ & - \text{tr} (\Sigma_1^{-1} \text{A} \ \Sigma_1^{-1} - \Sigma_2^{-1} \text{A} \ \Sigma_2^{-1}) (\Psi_1 - \Psi_2) \right] \\ & = \ - \frac{N}{2} \left[\text{tr} (\Sigma_1^{-1} - \Sigma_2^{-1}) (\Sigma_1 - \Sigma_2) \right. \\ & \left. - \text{tr} (\Sigma_1^{-1} \text{A} \ \Sigma_1^{-1} - \Sigma_2^{-1} \text{A} \ \Sigma_2^{-1}) (\Sigma_1 - \Sigma_2) \right]. \end{split}$$

By Theorems 2.2.1, and 2.2.2,

$$tr(\Sigma_1^{-1}-\Sigma_2^{-1}) (\Sigma_1-\Sigma_2) < 0$$

and

- tr
$$[\Sigma_1^{-1}A \Sigma_1^{-1}-\Sigma_2^{-1}A \Sigma_2^{-1}) (\Sigma_1-\Sigma_2)] \ge 0.$$

Since Δ is continuous in A, Δ > 0 for A near 0. Thus with probability greater than 0, L is not concave.

2.10 Growth Curve Models

Consider the data matrix $x=(x_{\alpha i})$ in Section 2.9. Instead of assuming that the dispersion matrix Σ is in some structural form, we assume that

$$E(X) = A \xi K$$

where A is an Nxh matrix of rank h and K is a gxp matrix of rank g, both being fixed matrices with h \leq N and g \leq p. ξ is the unknown hxg parameter matrix we wish to estimate. We also assume that Σ is known. This is the familiar growth curve model considered by Potthoff and Roy (1964), (See also Joreskog (1970)). The log-likelihood function, as a function of ξ , is

$$L(\xi) = c - \frac{N}{2} \ln |\Sigma| - \frac{N}{2} \operatorname{trA} \Sigma^{-1},$$

Where c is a constant and

$$A = \frac{1}{N} (x - A \xi K)' (x - A \xi K).$$

Let $d\xi \in M_{hxg}$. Then

$$\begin{split} \mathrm{d}L(\xi)(\mathrm{d}\xi) &= -\,\frac{N}{2}\,\,\mathrm{tr}\,\,\mathrm{d}A(\mathrm{d}\xi)\Sigma^{-1} \\ &= N\,\,\mathrm{tr}\,\,(x-A\xi K)^{\,\prime}\,\,A\,\,(\mathrm{d}\xi)K\Sigma^{-1} \,. \end{split}$$

Since Σ is positive definite,

$$Q(d^{2}L(\xi))(d\xi) = -N \operatorname{tr} (A(d\xi)K)' A (d\xi)K\Sigma^{-1}$$

$$\leq 0.$$

Thus L is concave on M_{hxg} . Let $dL(\xi) = 0$. Then

$$K\Sigma^{-1} (x-A\xi K)' A = 0$$

or

$$A'A\xi K \Sigma^{-1} K' = A'x \Sigma^{-1} K'.$$

Hence $(A'A)^{-1} Ax \Sigma^{-1} K' (K' \Sigma^{-1} K)^{-1}$ is the solution of $dL(\xi) = 0$ and is, therefore, the maximum likelihood estimate $\hat{\xi}$ of ξ .

2.11 Simultaneous Equation Models

Let $U=(u_{ij})$ be an Nxt matrix of normally distributed random variables u_{ij} with E(U)=0 and $E(\frac{1}{t}UU')=\Sigma$. Consider the model

BY +
$$\Gamma Z = U$$
,

ţ .

where B and Γ are NxN and NxA matrices of parameters. Assume that Σ is known. The problem is to obtain the maximum likelihood estimators of B and Γ . This is the familiar problem of full information maximum likelihood estimation of the structural parameters of a simultaneous linear structural equation model considered by Fisk (1967, Chapter 4), Neudecker (1967), Koopmans (1950) and Tracy and Singh (1972). The log-likelihood function of (B,Γ) is given by

$$L(B,\Gamma) = c + t \ln B + \frac{t}{2} \ln \Sigma^{-1} - \frac{t}{2} \operatorname{tr} \Sigma^{-1} AMA',$$

where

$$M = \frac{1}{t} \begin{bmatrix} Y \\ Z \end{bmatrix} (Y', Z'),$$

$$A = (B, \Gamma) \epsilon M_{Nx(N+\Lambda)}$$

c is a constant and Σ , B are positive definite. Note that Neudecker (1967) and Tracy and Singh (1972) do not assume that B is positive

definite and they do not prove that the optimal solution they found are indeed the maximum likelihood estimate of B.

Let $dA \in M_{N\times(N+\Lambda)}$, $dA \neq 0$, dB = P(dA). Since projection P of A to B is linear,

$$dL(A)(dA) = t trB^{-1}(dB)'-t tr\Sigma^{-1}AM(dA)'$$

Since B, M and Σ are positive definite,

$$Q(d^{2}L(A))(dA) = -t tr B^{-1}dBB^{-1}(dB)' - t tr \Sigma^{-1}(dA)M(dA)'$$
< 0.

Thus L is strictly convex on $M_{Nx(N+\Lambda)}$. Let dL(A)=0. Then (and only then)

$$(B^{-1},0) = \Sigma^{-1}AM.$$
 (2.11.1)

The solution of (2.11.1) is, therefore, the maximum likelihood estimate of (B,Γ) .

For simplicity, let us use the above example to demonstrate certain advantages of our approach. One reason that our approach is so simple is because the functions $A \to L(A)$, $dA \to dL(A)$ (dA), $dA \to Q(dL^2(A))(dA)$ we are dealing with are real-valued.

our notations, up to one-to-one linear transformations,

Þ

$$[d^{2}L(A)] = -t \cdot \begin{pmatrix} [(B')^{-1},0] & \otimes & B_{1}^{-1} \\ & & & \\ [(B'')^{-1},0] & \otimes & B_{p}^{-1} \end{pmatrix} + t'MA'\Sigma^{-1}AM \otimes \Sigma^{-1}$$

$$t \left(\begin{array}{c} \Sigma^{-1}AM \otimes \{MA'\Sigma^{-1}\}_{1} \\ \\ \Sigma^{-1}AM \otimes \{MA'\Sigma^{-1}\}_{q} \end{array}\right) - t M \otimes \Sigma^{-1}.$$

Tracy and Singh shows that, in terms of our notations, up to one-toone transformations,

Here C_i is the ith row of C and A E B is defined in Tracy and Singh (1972). One can see from the above results that their approaches are by no means simple. Also, they use the usual theory of calculus without connecting it to their own theories of matrix calculus. Moreover, we show that L is strictly concave. So (2.11.1) cannot have more than one solution. Thus one can use the conventional numerical methods, such as the method of Fletcher and Powell (1963), to solve (2.11.1).

E CHAPTER III

Optimal Control of a Regression Experiment

3.1 Introduction

We have seen in the preceding chapter that using the first and second order differentials, many optimization problems in statistics can be solved practically and rigorously. However, for optimal problems of the second type, often, some other methods are needed. Linear programming techniques may be used to optimize a linear function with linear inequality constraints. Variational techniques may be used to optimize a functional with linear or nonlinear inequality constraints. In this chapter, we shall use a linear regression model to illustrate certain methods of solving the optimal problems of the second type.

3.2 Optimal Control of a Regression Experiment

Consider the regression model m(f):

$$Y(t) = f(t) \theta + \varepsilon(t) , t \varepsilon T, \qquad (3.2.1)$$

where T is a nonempty topological space, $Y(t) = (y_1(t), ..., y_n(t))'$,

$$\mathbf{f}(\mathbf{t}) = (\mathbf{f}_{ij}(\mathbf{t})) \; \boldsymbol{\epsilon} \; \mathbf{M}_{nxk}, \; \; \boldsymbol{\epsilon}(\mathbf{t}) = (\boldsymbol{\epsilon}_1(\mathbf{t}), \dots, \boldsymbol{\epsilon}_n(\mathbf{t}))', \; \; \boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k)'.$$

 ϵ R^k . $\{\epsilon(t)\}_{t\in T}$ is a vector stochastic process on some probability

space (Ω, \mathcal{Q}, P) such that each $E(\epsilon(t)) = 0$ ϵ R^n , each $r(i,s,j,t) = E(\epsilon_i(s)\epsilon_j(t))$ is known and each r(i,j,t) is continuous on TxT. We shall be interested in the control problem of choosing f in some chosen set X to minimize a certain functional $\Phi(\cdot)$ of the dispersion matrix $\Sigma(f)$ of the least square estimator $\hat{\theta}$ of θ . Chang and Wong (to appear) consider this problem for the case n=1 and Dorogovcev (1971) considers this problem for the case n=1 and k=2. Related problems also appear in Chang (1979), Chang and Wong (to appear) and Mehra (1974). When T is a singleton, the above problem is familiar in optimal design theory. Seé, for example, Federov (1972) and Kiefer (1974).

In practice, n represents the number of observations. It turns out that the general case (n≥1) can be reduced to the case n=1 so that the optimality results in Chang and Wong (to appear) and Chan and Wong (to appear) can be applied. For simplicity, we shall assume that $T = \{a,b\}$, where a,bs R and a
b. We can easily generalize the obtained results to more general settings. Let W be an m-dimensional linear subspace of $\mathcal{L}^2(T)$ and $\{\eta_1,\eta_2,\ldots,\eta_m\}$ be an orthonormal basis in W. We shall assume that each η_1 is continuous on T, and in (3.2.1), $f = (f_1,\ldots,f_k)$, where each $f_j = (f_{ij})$, $f_{ij} \in W$. Thus $f_{ij} = \sum_{\alpha} f_{(\alpha,i),j} \eta_{\alpha}$ for some unique $f_{(\alpha,i),j}$'s. Unless otherwise stated, we shall keep all notations. Note that all $f_i \in H \equiv W^n$.

Theorem 3.2.1. $\{f_j\}$ is independent in H if and only if $F = \{f_{(\alpha,i),j}\}_j \text{ is independent in } M_{mxn}.$

We shall assume that $\{f_j\}$ is independent in H. F above will be treated as an element in $R^{(\{1,2,\ldots,m\}x\{1,2,\ldots n\})x\{1,2,\ldots,k\}}$. Let g, h ϵ H, ϵ

$$(g,h) = \sum_{i=1}^{h} (g_i,h_i).$$

Then H with (,) is a Hilbert space. $\| \ \|$ will denote the norm induced by the above inner product (,). (,) is the usual product inner product for H. X will denote $\{(f_{ij}): all\ f_{ij} \in \mathbb{W}, f_j$'s are independent in H, $\| f_j \| \le 1$, $j=1,2,\ldots,k\}$, 1>0. For simplicity, we shall take l=1.

Let w $\in \Omega$. Then the least square estimate of θ , denoted by $\widehat{\theta}$ (w), is defined to be the θ which minimizes

$$\parallel \Upsilon(^{\cdot})(\omega) - F\theta \parallel^2$$
 , $\theta \in \mathbb{R}^k$. (3.2.2)

M(f) will denote F'F, where F' = $\{c_{j,(\alpha,i)}\}$ with each $c_{j,(\alpha,i)}$ = $f_{(\alpha,i),j}$. Thus M(f) ϵM_{kxk} . Each U= (t_{ij}) ϵ R^{IxJ} can be considered as a linear transformation [U] of R^J into R^I such that $[U](x_i)$ =

 $U(x_j)(=(\Sigma t_{ij}x_j) \in \mathbb{R}^I)$, $(x_j) \in \mathbb{R}^J$. Thus the theory of linear transformations can be applied to \mathbb{R}^{IxJ} . Now r(M(f)) = r(F) = k. So $M(f)^{-1}$ exists.

Theorem 3.2.2.

$$\widehat{\theta}(\mathbf{w}) = \mathbf{M}(\mathbf{I})^{-1}((\mathbf{Y}(\mathbf{Y}(\mathbf{w}), \mathbf{f}_{\mathbf{j}})), \quad \mathbf{w} \in \Omega.$$
 (3.2.3)

Proof. Let $w \in \Omega$. Let j=1,2,...,k, Y = Y(')(w). Then

$$(Y - f \hat{\theta}(\omega), f_j) = 0$$

i.e.,

$$(Y,f_{j}) = (f \hat{\theta}(w),f_{j})$$

$$= \sum_{i,\ell} (f_{i\ell},f_{ij}) \hat{\theta}_{\ell}(w).$$

Since

$$\begin{split} (f_{ij}, f_{il}) &= \int_T f_{il}(t) f_{ij}(t) dt \\ &= \sum_{\alpha} f_{(\alpha,i)}, \mathcal{L}^f(\alpha,i), j, \end{split}$$

$$(F'F)\hat{\theta}(\omega) = ((Y,f_j)).$$

Therefore

$$\widehat{\theta}(\omega) = (F'F)^{-1}((Y,f_j))$$

$$= M(\cdot)^{-1}((Y,f_j)). \qquad q.e.d.$$

 $\Sigma(f)$ will denote the dispersion matrix of $\widehat{\vartheta}.$

Theorem 3.2.3.

$$\Sigma(f) = M(f)^{-1}(F'RF)M(f)^{-1},$$

where

$$R = (r_{(\alpha,i),(\beta,j)}) = (\int_{T} \int_{T} r(i,s,j,t) \eta_{\alpha}(s) \eta_{\beta}(t) ds dt) \in \mathbb{R}^{L \times L}$$

with

$$L = \{1,2,\ldots,n\} \times \{1,2,\ldots,m\}.$$

Proof.

$$\Sigma(f) = M(f)^{-1}\Sigma_{\Xi} M(f)^{-1},$$

where Σ_{Ξ} is the dispersion matrix of Ξ with $\Xi(w) = ((Y(\cdot)(w),f_j))$

Since each $E(\epsilon_i(t)) = 0$, by Fubini's theorem, $E((\epsilon(\cdot), f_j)) = 0$. Thus by Fubini's theorem,

$$\begin{split} & \Sigma_{Z} = (\mathbb{E}((\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{v}))) \\ & = (\sum_{i, \hat{\mathcal{L}}} \mathbb{E}((\varepsilon_{i}(\cdot), f_{iu})(\varepsilon_{\hat{\mathcal{E}}}(\cdot), f_{\hat{\mathcal{E}}v}))) \\ & = (\sum_{i, \hat{\mathcal{L}}} \mathbb{E}(\int_{T} \int_{T} \varepsilon_{i}(s) \varepsilon_{\ell}(t) f_{iu}(s) f_{\ell v}(t) ds dt)) \\ & = (\sum_{i, \hat{\mathcal{L}}} \mathbb{E}(\int_{T} \int_{T} \mathbb{E}(\varepsilon_{i}(s) \varepsilon_{\ell}(t)) f_{\alpha u i} f_{\beta v \hat{\mathcal{E}}} \eta_{\alpha}(s) \eta_{\beta}(t) ds dt) \\ & = (\sum_{i, \hat{\mathcal{L}}, \alpha, \beta} f_{\alpha u i}(\int_{T} \int_{T} \mathbb{E}(i, \hat{\mathcal{L}}, \hat{\mathcal{L}}, \hat{\mathcal{L}}) \eta_{\alpha}(s) \eta_{\beta}(t) ds dt) f_{\beta v \ell}) \\ & = \mathbb{E}(\mathbb{E}((\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{v}))) \\ & = \mathbb{E}((\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{v})) \\ & = \mathbb{E}((\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{v}))) \\ & = \mathbb{E}((\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{v})) \\ & = \mathbb{E}((\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{v}))) \\ & = \mathbb{E}((\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{v}))) \\ & = (\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{u})(\varepsilon(\cdot), f_{v})(\varepsilon(\cdot), f_{v})(\varepsilon(\cdot),$$

Therefore

$$\Sigma(f) = M(f)^{-1}(F'RF)M(F)^{-1}.$$
 q.e.d.

With the above estimator $\widehat{\theta}$ of θ , our objective is to choose an f in X such that for a certain function ϕ , $\phi(\Sigma(f)) = \min \phi(\Sigma(g))$. The following are three important ϕ 's in-optimum degeX sign theory. These and together with some other criterions can be found in Federov (1972). Let ϕ be a convex functional on P_k . Then m(f) is called a ϕ -optimal model if $\phi(\Sigma(f)) = \min \phi(\Sigma(g))$. A geX ϕ -optimal model m(f) is said to be

(i) D-Optimal model if

$$\phi(C) = 1C1^{-1}$$
 , $C \in P_k$; (3.2.4)

(ii) A-Optimal if

$$\phi(C) = tr C^{-1}$$
; (3.2.5)

(iii) D_S-Optimal model if

$$\phi(C) = \{C_s^{-1}\}, \quad C \in P_{\nu}, \quad (3.2.6)$$

where C_s is obtained from C by deleting the last (k-s) rows and columns. $\phi_1, \ \phi_2, \phi_3$ will denote respectively ϕ in (i) (ii) and (iii).

Lemma 3.2.3. Let h be a function of P_k into R such that h ϵ (P_k) (2) ($P_k \subset S_k$). Let

$$h_{\mathcal{E}}(A) = \exp(\mathcal{E}(A) + h(A)), \quad A \in P_k,$$

 $\ell \in S_k^\#$, the conjugate space of S_k^- . Then the following two conditions are equivalent:

- (i) h is (strictly) convex.
- (ii) h_{ℓ} is (strictly) convex for all $\ell \in S^{\#}$.

Theorem 3.2.4. $\phi_1, \ \phi_2$ are strictly convex.

Proof. Let $A \ \epsilon \ P_k$, $dA \ \epsilon \ S_k$. Then

$$d\phi_1(A)(dA) = - tr A^{-1}(dA)A^{-1}$$

= - tr A^{-2}dA.

So for any distinct A, B ϵ /P $_k,$ by Theorem 2.2.2,

$$(d\phi_1(A) - d\phi_1(B))(A-B) = -tr(A^{-2}-B^{-2})(A-B)$$

> 0

Hence ϕ_1 is strictly convex on P_k . Let

$$\phi(A) = -\ln |A| , \quad A \in P_k.$$

Let A ϵ P_k, dA ϵ S_k. Then

$$d\phi(A)(dA) = - trA^{-1}(dA)$$
.

So for any distinct A, B ϵ $\mbox{ P}_k,$ by Theorem 2.2.1,

$$(d\phi(A) - d\phi(B))(A-B) = -tr(A^{-1}-B^{-1})(A-B) > 0.$$

Therefore Φ is strictly convex on P_k . By Lemma 3.2.3 with $\hat{z}=0$, Φ_2 is strictly convex. q.e.d.

.Since $C \rightarrow C_s$ is linear, by Theorem 3.2.4, ϕ_3 is convex.

3.3 D-, A- and D_s -Optimal models

To find the D-, A- and D_s-optimal models m(f), one has to solve for each i=1,2,3, the optimization problem $\phi_i(f)=\max_{g\in X}\phi_i(g)$. Using results in Chang and Wong (to appear) and Chan and Wong (to appear), the above problem can be solved.

Let H_1 be an mnxk matrix obtained by re-arranging the "rows" $f_{(\alpha,i)}$ of the "matrix" F in any preassigned order. Let R_1 be the mnxmn matrix obtained by re-arranging the "rows" and "columns" of R in the corresponding way. Since multiplication of matrices does not depend on these re-arrangements,

$$\Sigma(f) = (H_1'H_1)^{-1}H_1'R_1H_1(H_1'H_1)^{-1}.$$

Since \mathbf{R}_1 is positive definite, there exist an orthogonal matrix \mathbf{P}_1 such that

$$PR_1P' = \begin{pmatrix} d_1 & & 0 \\ & \cdot & \\ & & \cdot \\ 0 & & d_{mn} \end{pmatrix} \equiv D,$$

where $0 < d_1 \le d_2 \le \ldots \le d_{mn}$. It is important to stress here that if one changes the above arrangement of rows and columns, then R_1 will be replaced by a similar matrix R_2 whence R_1 but not D will be changed. Let $G = PH_1$. Then

$$\Sigma(f) = (G'G)^{-1}G'DG(G'G)^{-1}.$$

Moreover, the columns G_{j} of G has norm less than or equal to 1:

$$\|G_{\beta}\|^{2} = \sum_{\alpha,j} f_{(\alpha,j),\beta}^{2}$$

$$= \sum_{j} \|f_{j\beta}\|^{2}$$

$$= \|f_{\beta}\|^{2}$$

≤ 1.

Thus the results in Chang and Wong (to appear) and Chan and Wong (to appear) for m(f) with n=1, m replaced by mn can be applied. We shall assume that D is positive definite. The following results follow trivially from Theorem 3 in Chan and Wong (to appear).

Theorem 3.3.1. Let

$$c = \sum_{j=1}^{k} d_j^{1/2},$$

$$D_{1} = \begin{pmatrix} d_{1} & & 0 \\ & \ddots & \\ & & d_{k} \end{pmatrix}$$

$$G = \frac{1}{e^{-1/2}} \begin{pmatrix} D_1^{-1/4} \\ 0 \end{pmatrix} Q^*$$

Where Q is an orthogonal matrix such that Q $D_1^{-1/2}Q^*$ has equal diagonal elements. Then m(f) is A-optimal and

$$\Sigma(f) = c Q D_1^{1/2} Q^{T}$$
.

Of course, in order to find an optimal model m(f), one must first obtain F from G and then obtain f from F.

Theorem 3.3.2. The following conditions are equivalent:

- a) m(f) is D-optimal.
- b) $G = Q(\frac{U}{0})$,

for some orthogonal matrix $\mathbb Q$ and sxs diagonal matrix $\mathbb U$ whose diagonal elements all belong to $\{-1,1\}$.

Theorem 3.3.3. The following conditions are equivalent:

- a) m(f) is D -optimal.
- $\mathbf{P}) \qquad \mathbf{G} = \mathbf{G} \begin{pmatrix} \mathbf{Q} & \mathbf{Q} \\ \mathbf{Q} & \mathbf{A} \\ \mathbf{Q} & \mathbf{Q} \end{pmatrix} \quad ,$

where U is an sxs diagonal matrix such that each diagonal element of U is -1 or 1, V is a (k-s)x(k-s) nonsingular matrix such that each column has Euclidean norm equal or less than 1.

For illustration, we shall give an example. Earlier notations will be kept. Suppose that for m(f), n=2=m=k, a=0, b=1, all $\epsilon_1(s)$, $\epsilon_2(t)$ are independent normal random variables of mean 0 and

$$r(i,s,j,t) = r(s,t)\delta_{ij}$$

where for s,t in (0,1),

$$r(s,t) = (1-s) (1-t) \min \left\{ \frac{s}{1-s}, \frac{t}{1-t} \right\},$$

$$r(1,t) = r(s,1) = r(1,1) = 0.$$

Since
$$\frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}}$$
 (2t-1), $\frac{1}{\sqrt{2}} - \frac{\sqrt{3}}{\sqrt{2}}$ (2t-1) are orthonormal in $\chi^2[0,1]$,

we may take

$$\eta_1 = \frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}} (2t-1),$$

$$\eta_2 = \frac{1}{\sqrt{2}} - \frac{\sqrt{3}}{\sqrt{2}} (2t-1).$$

We shall use Theorem 3.3.1 to find an A-optimal model m(f). Now

$$r_{(1,1),(1,1)} = \int_0^1 \int_0^1 \left(\frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}} (2t-1) \right) \left(\frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}} (2s-1) \right)$$

$$(1-s)(1-t) \min \{\frac{s}{1-s}, \frac{t}{1-t}\} dsdt$$

$$= \int_0^1 \int_s^1 \left(\frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}} (2t-1) \right) \left(\frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}} (2s-1) \right) (1-s) t dt ds$$

$$+ \int_0^1 \int_0^s \left(\frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}} (2t-1) \right) \left(\frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}} (2s-1) \right) (1-t) s ds dt$$

$$= \frac{1}{10} . . .$$

$$r(1,1,),(2,1) = \int_{0}^{1} \int_{0}^{1} \left(\frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}} (2s-1)\right) \left(\frac{1}{\sqrt{2}} - \frac{\sqrt{3}}{\sqrt{2}} (2t-1)\right).$$

$$(1-s)(1-t) \min \left\{ \frac{s}{1-s}, \frac{t}{1-t} \right\} dsdt$$

$$= \int_{0}^{1} \int_{s}^{1} \left(\frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}} (2s-1)\right) \left(\frac{1}{\sqrt{2}} - \frac{\sqrt{3}}{\sqrt{2}} (2t-1)\right) (1-s) t dt ds$$

$$+ \int_{0}^{1} \int_{0}^{s} \left(\frac{1}{\sqrt{2}} + \frac{\sqrt{3}}{\sqrt{2}} (2t-1)\right) \left(\frac{1}{\sqrt{2}} - \frac{\sqrt{3}}{\sqrt{2}} (2s-1)\right) (1-t) s ds dt$$

$$= \frac{-1}{15} .$$

Similarly,

$$r(2,1),(2,1) = r(2,2),(2,2) = r(1,2),(1,2) = \frac{-1}{10},$$
 $r(2,1),(1,1) = r(2,2),(1,2) = r(1,2),(2,2) = \frac{-1}{15},$
 $r(\alpha,i),(\beta,j) = 0$ elsewhere.

We shall use the order:

$$(1,1) \longrightarrow 1, (1,2) \longrightarrow 3, (2,1) \longrightarrow 2, (2,2) \longrightarrow 4.$$

R is then reduced to

$$R_{1} = \begin{pmatrix} \frac{1}{10} & -\frac{1}{15} & 0 & 0 \\ -\frac{1}{15} & \frac{1}{10} & 0 & 0 \\ 0 & 0 & \frac{1}{10} & -\frac{1}{15} \\ 0 & 0 & -\frac{1}{15} & \frac{1}{10} \end{pmatrix}$$

It is easy to diagonalize

$$\begin{pmatrix} \frac{1}{10} & -\frac{1}{15} \\ -\frac{1}{15} & \frac{1}{10} \end{pmatrix} \qquad \text{to} \qquad \begin{pmatrix} \frac{1}{30} & 0 \\ 0 & -\frac{1}{6} \end{pmatrix}$$

bу

$$\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ 1 & -1 \end{pmatrix}$$

So R_1 can be diagonalized to

$$\begin{pmatrix}
\frac{1}{30} & & & & & & & \\
& & \frac{1}{6} & & & & \\
& & & \frac{1}{30} & & \\
& & & & \frac{1}{6}
\end{pmatrix}$$

by .

$$A = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & -1 & 1 \end{pmatrix}$$

Let P be the matrix obtained from A by interchanging the second and third columns. Then

$$P = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 0 & 1 & 0 \\ 1 & 0 & -1 & 0 \\ 0 & -1 & 0 & 1 \\ 0 & -1 & 0 & 1 \end{pmatrix}.$$

and

$$D \equiv P^*R_1P = \begin{pmatrix} \frac{1}{30} & 0 & 0 & 0 \\ 0 & \frac{1}{30} & 0 & 0 \\ 0 & 0 & \frac{1}{6} & 0 \\ 0 & 0 & 0 & \frac{1}{6} \end{pmatrix}$$

Thus
$$d_1 = d_2 = \frac{1}{30}$$
, $d_3 = d_4 = \frac{1}{6}$, $c = \frac{(d_1^{\frac{1}{2}} + d_2^{\frac{3}{2}})}{2} = \frac{1}{\sqrt{30}}$,

$$D_1 = \frac{1}{30} I_2$$
.

Let Q be a 2x2 orthogonal matrix. Then Q D_1 Q' has equal diagonal

elements. In fact,
$$Q D_1^{\frac{1}{2}} Q' = \frac{1}{30} I$$

Şc

$$G = \frac{1}{C^{\frac{1}{2}}} \begin{pmatrix} D_1^{\frac{1}{2}} \\ 0 \end{pmatrix} Q'$$

$$= \begin{pmatrix} I_2 \\ 0 \end{pmatrix} Q'.$$

$$H_1 = P'G$$

$$= \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & & 1 \\ 0 & & 0 \\ 1 & & -1 \\ 0 & & 0 \end{pmatrix} Q'$$

For illustration, let $Q = I_2$. Then

$$H_{1} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 0 & 0 \\ 1 & -1 \\ 0 & 0 \end{pmatrix}$$

Using the earlier order, we obtain $F = (f_{(\alpha,i),j})$ with

$$f_{(1,1),1} = \frac{1}{\sqrt{2}} = f_{(1,1),2}$$

$$f_{(1,2),1} = \frac{1}{\sqrt{2}} = -f_{(1,2),2}$$

$$f_{(\alpha,i),j} = 0$$
 elsewhere.

Let t ϵ [0,1]. Then

$$f_{11}(t) = f_{(1,1),1} \eta_1(t) + f_{(2,1),1} \eta_2(t)$$

$$=\frac{1}{\sqrt{2}}\eta_1(t),$$

i.e.

$$f_{11}(t) = \frac{1}{2}((1-\sqrt{3}) + 2\sqrt{3} t).$$

Similarly

$$f_{12}(t) = f_{11}(t) = f_{21}(t),$$

 $f_{12}(t) = -\frac{1}{\sqrt{2}} \eta_1(t).$

Thus

$$f(t) = \begin{pmatrix} \frac{1}{2} ((1-\sqrt{3}) + 2\sqrt{3} t) & \frac{1}{2} ((1-\sqrt{3}) + 2\sqrt{3} t) \\ \frac{1}{2} ((1-\sqrt{3}) + 2\sqrt{3} t) & -\frac{1}{2} ((1-\sqrt{3}) + 2\sqrt{3} t) \end{pmatrix},$$

ts[0,1]. With this A-optimal model m(f),

$$\Sigma_{\widehat{\theta}} = \Sigma(f)$$

$$= c QA^{\frac{1}{2}} Q'$$

$$= c D_1^{\frac{1}{2}}$$

$$= \frac{1}{30} I_2.$$

i.e. $\hat{\theta}_1$, $\hat{\theta}_2$ are independent random variables of variance $\frac{1}{30}$. D-optimal and D_s-optimal models m(\hat{f}) can be found numerically by using Theorems 3.3.2 and 3.3.3.

3.4. Reproducing Kernel Hilbert spaces

The notion of least squares in section 3.3 does not depend on the covariance function R. There appears to have no sufficient reason to use the Euclidean norm to measure a random distance. In this section, we shall use a distance associated with R. Instead of assuming that D is positive definite, we shall assume that $R(i,s,j,t)=R_i(s,t)\delta_{ij}$, where all R_i are continuous. So for $i\neq j$, all $y_i(s)$; $y_j(t)$ are independent. Chang (1979) deals with the case n=1. For completeness, let us define the notion of positive definiteness mentioned above. Let K be a continuous real-valued function on TxT. Note that all R_i are examples of K. Let $f \in \mathcal{L}(T)$.

$$\hat{K}(f)(t) = \int_a^b f(s)K(s,t)ds, \quad t \in [\hat{a},b].$$

By the continuity of K and the compactness of [a,b], $\widehat{K}(f)$ is continuous on [a,b]. So \widehat{K} is a linear function of $\mathcal{L}(T)$ into $\mathcal{L}^{\infty}(T)$. Since

$$\mathcal{Z}^{\infty}(T) \subset \mathcal{Z}^{2}(T) \subset \mathcal{Z}(T)$$
,

 \hat{K} , restricted to $\mathcal{L}^2(T)$, is a linear function of $\mathcal{L}^2(T)$ into itself. We shall assume that \hat{K} is restricted to $\mathcal{L}^2(T)$. For simplicity, let $f \in \mathcal{L}^2(T)$, $s \in T$. Since $K(s,) \in \mathcal{L}^2(T)$, by Schwartz's inequality,

$$|\hat{K}(f)(s)| \le ||K(s,)|| ||f||.$$

Thus

where $\| K \|$ is the norm of K in $\mathcal{L}^2(TxT)$. Let $B(\mathcal{L}^2(T))$ be the algebra of all continuous linear functions of $\mathcal{L}^2(T)$ into $\mathcal{L}^2(T)$ with the operator norm $\| \cdot \|$. Then $\hat{K} \in B(\mathcal{L}^2(T))$ and as it was shown above,

$$\|\widehat{K}\| \le \|K\|$$

Since Λ is one-to-one, we may identify \widehat{K} with K and write K for \widehat{K} . Since $K \in B(\mathcal{Z}^2(T))$, the notions and theory of operators in Hilbert spaces can be applied, e.g., K is said to be <u>positive (nonnegative)</u> definite if K is self adjoint, i.e. K=K', the adjoint of K, and for any nonzero f in the Hilbert space $\mathcal{Z}^2(T)$, $(Kf,f) > 0 (\geq 0)$. The adjoint K' of K is defined as $K' \in B(\mathcal{Z}^2(T))$ such that

$$(K'f,g) = (f,Kg)$$
, $f,g \in \mathcal{L}^2(T)$.

Thus K=K' if and only if K is symmetric. So all R_i are self adjoint. We shall now let i=1,2,...,n, K= R_i and prove that K is nonnegative definite. Let f $\varepsilon \mathcal{L}^2(T)$. Then

$$(Kf,f) = \int_{a}^{b} (\int_{a}^{b} f(s)K(s,t)ds)f(t)dt$$

$$= \int_{a}^{b} \int_{a}^{b} f(s)f(t) E (\varepsilon_{i}(s)\varepsilon_{i}(t))dsdt$$

$$= E (\int_{a}^{b} \int_{a}^{b} f(s)f(t)\varepsilon_{i}(s)\varepsilon_{i}(t)dsdt)$$

$$= \frac{1}{7} E ((\int_{a}^{b} f(s)\varepsilon_{i}(s)ds)^{2})$$

$$> 0.$$

So K is nonnegative definite. For convenience, we shall assume that all R_i are positive definite. This assumption is not serious and amounts to: For any $i=1,2,\ldots,n$ and any $A\in \mathcal{A}$ with P(A)=1, $\mathcal{L}^2(T)$ is the closed linear span of $\{\epsilon(.)(w):w\in A\}$. It can be shown by the theory of Wiener processes that r in the example of section 3.3 is nonnegative definite.

The background of the above discussions can be found in Riesz and Nagy (1955), especially the Chapters IV and VI, respectively, on integral equations and completely continuous symmetric transformations of Hilbert space.

Now by Mercer's theorem, for each i, there exist an orthogonal basis $\{\phi_{iu}\}$ for $\chi^2(T)$ and a sequence $\{\lambda_{iu}\}$ of positive real numbers such that

$$R_{i}(s,t) = \sum_{u=1}^{\infty} \lambda_{iu} \phi_{iu}(s) \phi_{iu}(t)$$
 (3.4.1)

in $\mathcal{L}^2(T)$ and $\{\sum_{u=1}^{\mathbf{Y}} \lambda_{\underline{i}u} \phi_{\underline{i}u}(s) \phi_{\underline{i}u}(t)\}$ converges uniformly to $R_{\underline{i}}$ on TxT. Here $\{\lambda_{\underline{i}u}\}$ is the spectrum of $R_{\underline{i}}$, (3.4.1) is the spectral representation of $R_{\underline{i}}$ and each $\phi_{\underline{i}u}$ is an eigenvector of $R_{\underline{i}}$ corresponding to $\lambda_{\underline{i}u}$. Let

$$H(R_i) = \{f \in \mathcal{L}^2(T) : \sum_{u=1}^{\infty} \frac{1}{\lambda_{iu}} (f, \phi_{iu})^2 < \infty \},$$

where (,) is the usual inner product for $\mathcal{L}^2(T)$. For g, h ϵ H(R_i), we define $g_{iu} = (g, \phi_{iu})$ (h_{iu} = (h, ϕ_{iu})) and

$$(g,h)_{R_i} = \sum_{u=1}^{\infty} \frac{s_{iu}^{h}_{iu}}{\lambda_{iu}}$$

Then $H(R_i)$ is a Hilbert space and will be referred to as the <u>reproducing kernel Hilbert space (RKHS) induced by R_i .</u> Let $R = (R_i)$ and, as in section 3.2, let H(R) be the product of $H(R_i)$'s equipped with the product inner product $(,)_R$, i.e.,

$$(g,h)_{R} = \sum_{i=1}^{n} (g_{i},h_{i})_{R_{i}}, \qquad g_{i},h_{i} \in H(R_{i}),$$

$$g = (g_{i}), h = (h_{i}).$$

Then H(R) is a Hilbert space and is called the <u>reproducing kernel</u>. Hilbert space (RKHS) induced by R. $\|\cdot\|_{R}$ will denote the norm induced by $(\cdot,\cdot)_{R}$.

$$\|g\|_{R}^{2} = (g,g)_{R}$$
 , $g \in H$. (3.4.3)

We assume that f_j 's in m(f) are continuous on T and are independent in H(R). Let $w \in \Omega$. Ξ in R^k , denoted by $\widetilde{\theta}(w)$, is the Gauss - Markov estimator of θ based on the observation Y(.)(w) if

$$\|Y(.)(\omega) - f z\|_{R}^{2} = \min_{\Theta \in R^{k}} \|Y(.)(\omega) - f \Theta\|_{R}^{2}$$

Since f_1 , f_2 , ..., f_k are independent in H(R) and $\{f\theta:\theta\in\mathbb{R}^k\}$ is the finite dimensional linear space S spanned by f_1 , f_2 , ..., f_k , $\widetilde{\theta}(w)$ uniquely exists and with Y = Y(.)(w),

$$(Y - f \widetilde{\Theta}(\omega), f_j)_R = 0$$
 , $j = 1, 2, ..., k$.

Thus

$$(Y,f_j)_R = (f \tilde{\Theta}(w), f_j)_R$$

$$= \sum_{\ell} \tilde{\Theta}_{\ell}(w) (f_{\ell},f_j)_R.$$

Hence

$$(\sum_{\ell=1}^{k} \widetilde{\theta}_{\ell}(\mathbf{w}) (\mathbf{f}_{\ell}, \mathbf{f}_{j})_{R}) = ((\mathbf{Y}, \mathbf{f}_{j})_{R}).$$

Let

$$\mathtt{M}(\mathtt{f}) = ((\mathtt{f}_{\underline{\ell}},\mathtt{f}_{\underline{j}})_{\mathtt{R}}) \; .$$

Then $M(f)^{-1}$ exists and

$$\xi M(f) \widetilde{\Theta}(\omega) = ((\Upsilon, f_j)_R)$$
.

Therefore

$$\widetilde{\Theta}(\mathbf{w}) = \mathbf{M}(\mathbf{f})^{-1} \left((\mathbf{Y}, \mathbf{f}_{\mathbf{j}})_{\mathbf{R}} \right). \tag{3.4.4}$$

Let W be a random vector of (Ω, \mathcal{Q}, P) into \mathbb{R}^k . W is said to be a <u>continuous linear estimator</u> of θ if there exist linear continuous functionals L_1 , L_2 , \ldots , L_k on $H(\mathbb{R})$ such that

$$W(w) = (L_j(Y(.)(w)))$$
, $w \in \Omega$.

Theorem 3.4.1. $\widetilde{\theta}(w)$ is the best unbiased linear continuous estimator (BLUE) of θ .

Proof. It suffices to prove that $\widetilde{\theta}(w)$ is unbiased for θ and that $\widetilde{\theta}(w)$ has minimum variance among all unbiased linear continuous estimator of θ . Let $w \in \Omega$, and W be a continuous linear estimator of $\overline{\theta}$. Then by Riesz's representation theorem,

$$W(w) = ((Y(.)(w),g_{j})_{R}),$$

for some g_1, g_2, \dots, g_k in H(R). Write $g_j = (g_{ij})$. Then

$$E(W) = \left(\int_{\Omega} \sum_{i=1}^{n} \left(Y_{i}(.)(\omega), g_{ij}\right)_{R_{i}} dP(\omega)\right)$$

$$= (\sum_{i=1}^{n} \int_{\Omega} (Y_{i}(.)(\omega), g_{ij})_{R_{i}} dP(\omega))$$

$$=(\sum_{i=1}^{n}\int_{\Omega}\sum_{u=1}^{\infty}\frac{g_{iju}}{\lambda_{iu}}\int_{a}^{b}Y_{i}(t)(\omega)\phi_{iu}(t)dtdP(\omega))$$

$$= (\sum_{i=1}^{n} \sum_{u=1}^{\infty} \frac{g_{iju}}{\lambda_{iu}} \int_{a}^{b} \sum_{\ell=1}^{k} f_{i\ell}(t) \theta_{\ell} \phi_{iu}(t) dt)$$

$$= (\sum_{\ell=1}^{k} \Theta_{\ell} \sum_{i=1}^{n} \sum_{u=1}^{\infty} \sum_{k=1}^{k} \frac{\hat{s}_{iju}}{\lambda_{iu}} f_{i\ell u})$$

$$= (\sum_{\ell=1}^{k} \Theta_{\ell} \sum_{i=1}^{n} (g_{ij}, f_{i\ell})_{R_{i}})$$

$$= (\sum_{\ell=1}^{k} \theta_{\ell} (g_{j}, f_{\ell})_{R}).$$

So

$$E(W) = ((f_{\ell}, g_{j})_{R})e.$$

(3.4.5)

Now by (3.4.5),

$$E(\widetilde{\theta}(\omega)) = E(M(f)^{-1}((Y(.),f_j)_R))$$

$$= M(f)^{-1}((f_\ell,f_j)_R) \theta$$

$$= \theta.$$

i.e., $\widetilde{\theta}(\omega)$ is unbiased for $\theta.$ Let $\Sigma_{\widetilde{W}}$ be the dispersion matrix of W. Then

$$\begin{split} E_{W} &= E((\varepsilon(.)(\cdot),g_{j})_{R},(\varepsilon(.)(\underline{\cdot}),g_{k})_{R}) \\ &= (\sum_{i=1}^{n} \sum_{i'=1}^{n} \sum_{u=1}^{\infty} \sum_{u'=1}^{\infty} \frac{1}{\lambda_{iu}\lambda_{iu'}} g_{iju}g_{i'ku'} \int_{a\cdot a}^{b} \int_{a\cdot a}^{b} \\ & E(\varepsilon_{i}(s)\varepsilon_{i},(t))\phi_{iu}(s)\phi_{i'u'}(t)dsdt) \\ &= (\sum_{i=1}^{n} \sum_{u=1}^{\infty} \sum_{u'=1}^{\infty} \frac{1}{\lambda_{iu}\lambda_{iu'}} g_{iju}g_{iku'} \int_{a\cdot a}^{b} \int_{a}^{b} R_{i}(s,t) \end{split}$$

$$\phi_{iu}(s)\phi_{iu},(t)dsdt)$$

$$= (\sum_{i=1}^{\infty} \sum_{u'=1}^{\infty} \sum_{u'=1}^{\infty} \frac{1}{\lambda_{iu}\lambda_{iu'}} g_{iju}g_{iku'} \sum_{u''=1}^{\infty} \frac{1}{\lambda_{iu''}} \delta_{uu''}\delta_{u''u''})$$

$$= (\sum_{i=1}^{n} \sum_{u=1}^{\infty} \frac{1}{\lambda_{iu}} g_{iju} g_{iku}^{\dagger}).$$

So

$$\Sigma_{W} = ((g_{j}, g_{k})_{R}).$$
 (3.4.6)

By (3.4.6),

$$\Sigma(f) = M(f)^{-1} \Sigma_{(Y(.),f_j)_R} (M(f)^{-1})^{-1}$$

$$= M(f)^{-1} M(f) M(f)^{-1}.$$

S٥٠

$$\Sigma(f) = M(f)^{-1}. \tag{3.4.7}$$

By (3.4.5), W is unbiased for θ if and only if

$$((\xi_j, f_2)_R) = I_k.$$
 (3.4.8)

It can be proved by (3.4.6) - (3.4.8) that if W is unbiased for θ , then $\Sigma_W - \widetilde{\Sigma}_{\widetilde{\theta}(\omega)}$ is nonnegative definite. Hence $\widetilde{\theta}(\omega)$

has minimum variance among all unbiased continuous estimators of θ .

q.e.d.



To obtain D-, D_s- and A-optimal models m(f), we need to specify the family X of all eligible f. Let L ϵ (0, ∞) and let X be the family of $f=(f_{ij})$ in (3.2.1) such that $\|f_j\|_R \leq L$, $j=1,2,\ldots,k$.

Theorem 3.4.2. (a)
$$\min_{g \in X} |\Sigma(g)| = \frac{1}{L^{2k}}$$

(b) m(f) is a D-optimal model if and only if f_1, f_2, \ldots, f_k are orthogonal in H(R) with norm L.

Proof. (a) Since

$$|\Sigma(g)| = |M(g)|^{-1}$$
, $g \in X$,

it suffices to prove that

$$\max_{g \in X} |M(g)| = L^{2k}$$
.

Let g ε X and let $\{\lambda_j\}$ be the spectrum of M(g). Then

$$\sum_{j=1}^{k} \lambda_{j} = \operatorname{tr} M(g) = \sum_{j=1}^{k} (f_{j}, f_{j})_{R}$$

$$= \sum_{j=1}^{k} \|f_j\|_{R}^{2}$$

$$\leq \sum_{i=1}^{k} L^2$$

$$= k L^2$$
. (3.4.9)

So, by-the familiar arithematic mean - geometric mean inequality,

$$|M(g)| = \prod_{j=1}^{k} \lambda_{j}$$

$$\leq \left(\sum_{j=1}^{k} \frac{\lambda_{j}}{k}\right)^{k}$$

$$\leq \left(\frac{1}{k} \cdot k L^{2}\right)^{k}$$

$$= L^{2k}$$

Now suppose that f_1, f_2, \ldots, f_k are orthogonal in H(R) with norm L. Then

$$|M(f)| = |((f_i, f_j)_R)|$$

= $|L^2 I_k|$
= L^{2k} ,

proving (a).

(b) From (a), it suffices to prove that $\min |M(g)| = |M(f)| = L^{2k}$ $g \in X$ implies that f_1 , f_2 , ..., f_k are orthogonal in H(R) with norm L. Let $\{\lambda_j\}$ be the spectrum of M(f). Since $\{\lambda_j\}$ minimizes $\prod_{i=1}^p \lambda_i$ with $\prod_{i=1}^p \lambda_i > 0$

and

$$\sum_{j=1}^{k} \lambda_{j} \leq k L^{2}, \qquad (3.4.10)$$

 $\lambda_1 = \lambda_2 = \dots = \lambda_p = L^2$. Since M(f) is positive and has identical eignevalues, $\lambda_j = L^2$, M(f) = L^2 I_k. Thus f₁, f₂, ..., f_k are orthogonal in H(R) and have norm L. q.e.d.

Note here that even for n=1, the above result is more general than Theorem 1 (i) in Chang (1979). Also our proof is different from the one given by Chang who uses Hadamard's inequality to prove (a).

Theorem 3.4.3.

- (a) $\min_{g \in X} \operatorname{tr} \Sigma(g) = \frac{k}{L^2}$.
- (b) m(f) is an A-optimal model if and only if f_1 , f_2 , ..., f_k are orthogonal in H(R) with norm L.

Proof. (a) Let geX. Then

tr
$$\Sigma(g) = \operatorname{tr} M(g)^{-1}$$
.

Let $\{\lambda_j^{\cdot}\}$ be the spectrum of M(g). Then

$$\operatorname{tr} \Sigma(g) = \sum_{i=1}^{k} \frac{1}{\lambda_i}.$$
 (3.4.11)



To minimize tr $\Sigma(g)$ is equivalent to minimize (3.4.11) subject to (3.4.10), which, again, is equivalent to minimize (3.4.11) subject to

all
$$\lambda_{i} > 0$$
 , $\sum_{i=1}^{k} \lambda_{i} = k L^{2}$. (3.5.12)

Let $\lambda = (\lambda_i)$ and

$$\phi(\lambda) = \sum_{i=1}^{k} \frac{1}{\lambda_i} - \mu(\sum_{i=1}^{k} \lambda_i - k L^2),$$

where μ is a Lagrange multiplier. Let $d\lambda = (d\lambda_i) \ \epsilon \ R^k$. Then

$$d\phi(\lambda)(d\lambda) = (-\frac{1}{\lambda_i^2} - \mu) d\lambda.$$

Let $\alpha = (\alpha_i)$, $\beta = (\beta_i) \in \mathbb{R}^k$ with all α_i , $\beta_i > 0$. Then

$$(d\phi(\alpha)-d\phi(\beta))(\alpha-\beta) = ((\frac{1}{\beta_i^2} - \frac{1}{\alpha_i^2}))((\alpha_i - \beta_i))$$

$$=\sum_{i=1}^{k}\frac{(\alpha_{i}-\beta_{i})^{2}(\alpha_{i}+\beta_{i})}{\alpha_{i}^{2}\beta_{i}^{2}}$$

≥ 0.

So φ is convex on the open convex set B of all λ in $\,R^k$ with each $.\lambda_{\hat{\tau}}\,>\,0\,.\,$ Let $d\varphi(\lambda)\,=\,0\,.\,$ Then

$$\lambda_{i} = (-\mu)^{-\frac{1}{2}}$$
 , $i = 1, 2, ..., k$.

5

Choose μ so that (3.5.12) holds. Then

$$\lambda_{i} = L^{2}$$
 , $i = 1, 2, ..., k$, $\mu = -\frac{1}{L^{4}}$,

and so λ_0 with each $\lambda_{0i} = L^2$ gives the minimum value of $\phi(B)$. By the theory of Lagrange's multipliers, λ_0 minimizes (3.4.11) subject to (3.5.12) and

$$\operatorname{tr} \Sigma(f) \geq \sum_{i=1}^{\infty} \frac{1}{\lambda_{0i}} = \frac{k}{L^2}$$

Now choose f such that f_1 , f_2 , ..., f_k are orthogonal in H(R) with norm L. Then f ϵ X and $\Sigma(f) = \frac{1}{L^2} \, I_k$. Thus

$$\operatorname{tr} \Sigma(f) = \frac{k}{L^2} ,$$

proving (a).

(b) By (a), it suffices to prove that tr $\Sigma(f) = \frac{k}{L^2}$ implies that f_1, f_2, \ldots, f_k are orthogonal in H(R) with norm L. From the proof of (a), the spectrum $\{\lambda_{0j}\}$ of M(f) minimizes (3.5.11) subject to (3.5.12). So

$$\lambda_{01} = \lambda_{02} = \dots = \lambda_{0k} = L^2,$$

i.e., $M(f) = L^2 F_k$. Hence f_1 , f_2 , ..., f_k are orthogonal in H(R) with norm L.

Again, even for n=1, Theorem 3.4.3 is more general than Theorem 1 (ii) in Chang (1979). Also, our proof is different from the one given by Chang who uses the Gram - Schmidt process to prove (a). Combining Theorem 3.4.2 and Theorem 3.4.3 we have the following nice result.

Theorem 3.4.4. Let $f \in X$. Then m(f) is D-optimal if and only if it is A-optimal.

From the view point of section 3.3, where D-optimal models are not equivalent to A-optimal models, Theorem 3.4.4 is a very strong and, yet, a desirable result. One then wonders which estimator and optimal model should be used. The Gauss-Markov estimator is the BLUE of θ . However, for finding this estimator, one needs the spectral representation of all r_i , which are often difficult to obtain. Also, we seldom know all r precisely and if we do not know all r_i precisely, then a repeated use of r_i 's will probably increase the uncertainty of the chosen model and the chosen estimator. Since the optimal models m(f) obtained in this section depend on a certain RKHS, practically, the design f may be difficult to control (or, say, construct). On the other hand, it is relatively easy to control the optimal models in section 3.3 and compute the underlying least square estimates. Also, in terms of motivations and conclusions, A-optimal models m(f) in section 3.3 tend to ignore larger eigenvalues of R. Indeed, with the assumption of this section, if $n \ge k$, then only the smallest eigenvalue of R

contributes to the optimal models. On the other hand, the A-optimal models m(f) in this section are obtained by more or less treating all eigenvalues of r_i 's as equal. To conclude, like many other statistical models, it is up to the workers to observe the reality and then decide which optimal model to use.

REFERENCES

- ANDERSON, T.W. (1958). Introduction to Multivariate Statistical Analysis. John Wiley and Sons, Inc., New York.
- APOSTAL, T.M. (1957). Mathematical Analysis. Addison Wesley, Mass.
- ATHANS, M. (1967). The Matrix Minimum Principle. Information and Control, 2, 592-606.
- ATHANS, M. and SCHWEPPE, F.C. (1965). Gradient Matrices and Matrix Calculation. M.I.T. Lincoln Lab. Lexington, Mass. Technical Report.
- ATHANS, M. and TSE, E. (1967). A Direct Derivation of the Optimal Linear Filter Using the Maximum Principle. IEEE Trans. Automatic Control, AC-12, 690-698.
- BENTLER, P.M. and LEE, Sik-Yum (1975). Some Extension of Matrix Calculus. General Systems, 20, 145-50.
- BENTLER, P.M. and LEE, Sik-Yum (1978). Matrix Derivatives with Chain Rules and Rules for Simple Hadamard and Kronecker Products. J. Math. Psychology, 17, 255-262.
- BUSH, K.A. and OLKIN, I. (1959). Extrema of Quadratic Forms with Applications to Statistics. Biometrika, 46, 483-486.

- CHAN, N.N. and WONG, Chi Song (To Appear). Existance of an A-Optimal Model for a Regression Experiment. J. Math. Analysis and Its Appl.
- CHANG, D.S. (1979). Design of Optimal Control for a Regression Problem. Ann. Statist., 7, 1078-1085.
- CHANG, D.S. and WONG, Chi Song (To Appear). Corrections to "Design of Optimal Control for a Regression Problem." Ann. Statist.
- CHANG, D.S. and WONG, Chi Song (To Appear). A General Approach to Optimal Control of a Regression Experiment. J. Multivariate Analysis.
- CHANG, D.S. and WONG, Chi Song (1979). Optimal Control of a Regression Experiment. Seminar on Functional Analysis, Institute of Mathematics, National Tsing Hua University, Taiwan, (Editor: Chi Song Wong).
- CALVERT, B. and SEBER, G.A. F. (1978). Minimization of Functions of a Positive Semidefinite Matrix A Subject to AX = 0. J. Multivariate Analysis, 8, 173-180.
- DEEMER, W.L. and OLKIN, I. (1951). The Jacobians of Certain Matrix Transformations Useful in Multivariate Analysis. Biometrika, 38, 345-67.

- DORÓGOVCEV, A. Ja. (1971). Problems of Optimal Control of a Regression Experiment. Selected Transl. in Math. Statist. and Probability, 10, 35-41.
- DWYER, P.S. and MacPHAIL, M.S. (1948). Symbolic Matrix Derivatives.

 Ann. Math. Statist., 19, 517-34.
- DWYER, P.S. (1967). Some Applications of Matrix Derivatives in Multivariate Analysis. J. Amer. Statist. Assoc., 62, 607-625.
- DYKSTRA, R.L. (1970). Establishing the Positive Definiteness of the Sample Covariance Matrix. Ann. Math. Statist. 41, 5153-54.
- FEDEROV, V.V. (1972). Theory of Optimal Experiments. Academic Press, New York.
- FISK, P.R. (1967). Statistically Dependent Equations. Charles Griffin and Company. London.
- FLEMING, W.H. (1977). Functions of Several Variables, Second Edition. Addison-Wesley, Reading.
- .FLETCHER, R. and POWELL, M.J.D. (1963). A Rapidly Convergent Descent Method for Minimization. Computer Journal, 6, 163-168.
 - GOLBERG, M.A. (1972). The Derivative of a Determinant. Amer. Math. Monthly, 79, 1124-1126.

- GOLDBERGER, A.S. (1964). Econometric Theory. John Wiley and Sons Inc., New York.
- GRAYBILL, F.A. (1969). Introduction to Matrices With Applications in Statistics. Wadsworth, California.
- HOTELLING, H. (1931). The Generalization of Student's Ratio. Ann. Math. Statist., 2; 360-378.
- *JENNRICH, R.I. (1973). Standard Error for Obliquely Rotated Factor Loadings. Psychometrika, 38, 593-603.
- JORESKOG, K.G. (1970). A General Method for Analysis of Covariance Structures. Biometrika, 57, 239-251.
- KIEFER, J. (1974). General Equivalence Theory for Optimal Designs.
 Ann. Statist., 2, 849-79.
 - KHATRI, C.G. and RAO, C.R. (1968). Solutions to Some Functional Equations and Their Applications to Characterization of Probability Distributions. Sankhya, Series A, 30, 167-180.
 - KOOPMANS, T.C. (1950). Statistical Inference in Dynamic Economic Models. Cowles Commission for Research in Economics, Monograph No. 10. John Wiley, New York.

- KSHIRSAGAR, A.M. (1972). Multivariate Analysis. Mercel Dekker, New York.
- LEHMAN, E.L. (1959). Testing Statistical Hypothesis. John Wiley and Son Inc., New York.
- LAWLEY, D.N. and MAXWELL, A.E. (1963). Factor Analysis as a Statistical Method. Butterworths, London.
- MacRAE, E. (1974). Matrix Derivatives With an Application to an Adaptive Linear Decision Problem. Ann. Statist., 2, 337-46.
- MAGNUS, J.R. and NEUDECKER, H. (1979). The Commutation Matrix: Some Properties and Applications. Ann. Statist., 7, 381-394.
- McDONALD, R.P. and SWAMINATHAN, H. (1973). A Simple Matrix Calculus With Applications to Multivariate Analysis. General System, 18, 37-54.
- McDONALD, R.P. (1976). The McDonald-Swaminathan Calculus: Clarifications, Extensions and Illustrations. General System, 21, 87-94.
- MEHRA, R.K. (1974). Optimal Input Signals for Parameter Estimation in Dynamic System -- Survey and New Results. IEEE Trans. Automatic Control, AC-19, 753-68.

- *NEL, D.G. (1978). A Review of Matrix Differentiation in Statistics. Technical Report No. 38, Dept. of Math. Statist., University of Orange Free State, South Africa.
- NEUDECKER, H. (1967). On Matrix Procedures for Optimizing Differentiable Scalar Functions of Matrices. Statist. Neerlandica, 21, 101-107.
- NEUDECKER, H. (1968). The Kronecker Matrix Product and Some of its Applications in Econometrics. Statist. Neerlandica, 22, 69-82.
- NEUDECKER, H. (1969). Some Theorems on Matrix Differentiation With Special Reference to Kronecker Matrix Products, J. Amer. Statist. Assoc., 64, 953-63.
- OLKIN, I. (1953). Note on "The Jocabians of Certain Matrix Transformations Useful in Multivariate Analysis." Biometrika, 40, 43-46.
- OPIAL, Z. (1967). Lecture Notes on Nonexpansive and Monotone
 Mapping in Banach Spaces. Center for Dynamical Systems, Brown
 University, Providence.
- POTTHOFF, R.F. and ROY, S.N. (1964). A Generalized Multivariate Analysis of Variance Model Useful Especially for the Growth Curve Problems. Biometrika, 5, 313-16.

- RAO, C.R. (1973). Linear Statistical Inference and its Applications, Second Edition. John Wiley and Sons Inc., New York.
- RIESZ, F. and SZ-NAGY, B. (1955, 3rd Edition). Functional Analysis. (Translated by Boron, L.F.), Frederick Ungar Publishing Co.
- ROGER, G.S. and YOUNG, D.L. (1978). On Testing a Multivariate Linear Hypothesis When the Covariance Matrix and its Inverse Have the Same Pattern. J. Amer. Statist. Assoc., 73, 203-7.
- STNGH, R.P. (1972). Some Generalizations in Matrix Differentiation With Applications in Multivariate Analysis. Doctoral Dissertation. University of Windsor.
- SMITH, D.W. (1978). A Simplified Approach to the Maximum Likelihood Estimation of the Covariance Matrix. Amer. Statistician, 32, 28-29.
- SPIVAK, M. (1965). Calculus of Manifolds. Benjaman, New York.
- SWAMINATHAN, H. (1976). Matrix Calculus for Functions of Partitioned Matrices. General System, 21, 95-99.
- THEIL, H. and SCHWEITZER, A. (1961). The Best Quadratic Estimator of the Residual Variance in Regression Analysis. Statist.

 Neerlandica, 15, 19-23.

- TRACY, D.S. and DWYER, P.S. (1969). Multivariate Maxima and Maxima With Matrix Derivatives. J. Amer. Statist. Assoc., 64, 1576-94.
- TRACY, D.S. and SINGH, R.P. (1972). A New Matrix Product and its Application in Partitioned Matrix Differentiation. Statist. Neerlandica, 26, 143-57.
- TRACY, D.S. and SINGH, R.P. (1975). Some Applications of Matrix Differentation in the General Analysis of Covariance Structures. Sankhya, Series A, 37, 269-280.
- VETTER, W.J. (1970). Derivative Operators on Matrices. IEEE Trans.
 Automatic Control, AC-15, 241-244.
- WATSON, G.S. (1964). A Note on Maximum Likelihood. Sankhya, Series A, 26, 303-4.
- WONG, Chi Song (To Appear). Matrix Derivatives and Its Applications in Statistics. J. Math. Psychology.
- WONG, Chi Song (To Appear). Mathematical Statistics. Tamkang Chair, Tamkang College of Arts and Sciences, Taiwan.
- WONG, Chi Song and WONG, Kai Sang (1979). A First Derivative Test for the Maximum Likelihood Estimates. Bull. Inst. Math. Acad. Sinica, 7, 313-21.

WONG, Chi Song and WONG, Kai Sang (To Appear). Minima and Maxima in Multivariate Analysis. Canadian J. Statist.

*These references have been suggested by some examiners in the examination committee. They were not available to the author at the time of writing this dissertation.

VITA AUCTORIS

Born

February, 22, 1952.

Chung San, China

M.Sc.

May, 1976 😁

University of Windsor

Lecturer (Part Time)

February - Joly, 1979

National Tsing Hua University, Taiwan

At present, Assistant - Management Sciences, Management Sciences Division, Bell Canada.