

**Matrix Differential Calculus
with Applications in Statistics
and Econometrics**

WILEY SERIES IN PROBABILITY AND STATISTICS

Established by *Walter E. Shewhart and Samuel S. Wilks*

Editors: *David J. Balding, Noel A. C. Cressie, Garrett M. Fitzmaurice, Geof H. Givens, Harvey Goldstein, Geert Molenberghs, David W. Scott, Adrian F. M. Smith, and Ruey S. Tsay*

Editors Emeriti: *J. Stuart Hunter, Iain M. Johnstone, Joseph B. Kadane, and Jozef L. Teugels*

The *Wiley Series in Probability and Statistics* is well established and authoritative. It covers many topics of current research interest in both pure and applied statistics and probability theory. Written by leading statisticians and institutions, the titles span both state-of-the-art developments in the field and classical methods.

Reflecting the wide range of current research in statistics, the series encompasses applied, methodological, and theoretical statistics, ranging from applications and new techniques made possible by advances in computerized practice to rigorous treatment of theoretical approaches. This series provides essential and invaluable reading for all statisticians, whether in academia, industry, government, or research.

A complete list of the titles in this series can be found at
<http://www.wiley.com/go/wsps>

Matrix Differential Calculus with Applications in Statistics and Econometrics

Third Edition

Jan R. Magnus

*Department of Econometrics and Operations Research,
Vrije Universiteit Amsterdam, The Netherlands*

and

Heinz Neudecker[†]

*Amsterdam School of Economics,
University of Amsterdam, The Netherlands*

WILEY

This third edition first published 2019
© 2019 John Wiley & Sons

Edition History

John Wiley & Sons (1e, 1988) and John Wiley & Sons (2e, 1999)

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, except as permitted by law. Advice on how to obtain permission to reuse material from this title is available at <http://www.wiley.com/go/permissions>.

The right of Jan R. Magnus and Heinz Neudecker to be identified as the authors of this work has been asserted in accordance with law.

Registered Offices

John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, USA
John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester, West Sussex, PO19 8SQ, UK

Editorial Office

9600 Garsington Road, Oxford, OX4 2DQ, UK

For details of our global editorial offices, customer services, and more information about Wiley products visit us at www.wiley.com.

Wiley also publishes its books in a variety of electronic formats and by print-on-demand. Some content that appears in standard print versions of this book may not be available in other formats.

Limit of Liability/Disclaimer of Warranty

While the publisher and authors have used their best efforts in preparing this work, they make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives, written sales materials or promotional statements for this work. The fact that an organization, website, or product is referred to in this work as a citation and/or potential source of further information does not mean that the publisher and authors endorse the information or services the organization, website, or product may provide or recommendations it may make. This work is sold with the understanding that the publisher is not engaged in rendering professional services. The advice and strategies contained herein may not be suitable for your situation. You should consult with a specialist where appropriate. Further, readers should be aware that websites listed in this work may have changed or disappeared between when this work was written and when it is read. Neither the publisher nor authors shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

Library of Congress Cataloging-in-Publication Data applied for
ISBN: 9781119541202

Cover design by Wiley
Cover image: © phochi/Shutterstock
Typeset by the author in L^AT_EX

10 9 8 7 6 5 4 3 2 1

Contents

<i>Preface</i>	xiii
--------------------------	------

Part One — Matrices

1 <i>Basic properties of vectors and matrices</i>	3
1 Introduction	3
2 Sets	3
3 Matrices: addition and multiplication	4
4 The transpose of a matrix	6
5 Square matrices	6
6 Linear forms and quadratic forms	7
7 The rank of a matrix	9
8 The inverse	10
9 The determinant	10
10 The trace	11
11 Partitioned matrices	12
12 Complex matrices	14
13 Eigenvalues and eigenvectors	14
14 Schur's decomposition theorem	17
15 The Jordan decomposition	18
16 The singular-value decomposition	20
17 Further results concerning eigenvalues	20
18 Positive (semi)definite matrices	23
19 Three further results for positive definite matrices	25
20 A useful result	26
21 Symmetric matrix functions	27
<i>Miscellaneous exercises</i>	28
<i>Bibliographical notes</i>	30
2 <i>Kronecker products, vec operator, and Moore-Penrose inverse</i>	31
1 Introduction	31
2 The Kronecker product	31
3 Eigenvalues of a Kronecker product	33
4 The vec operator	34
5 The Moore-Penrose (MP) inverse	36

6	Existence and uniqueness of the MP inverse	37
7	Some properties of the MP inverse	38
8	Further properties	39
9	The solution of linear equation systems	41
	<i>Miscellaneous exercises</i>	43
	<i>Bibliographical notes</i>	45
3	<i>Miscellaneous matrix results</i>	47
1	Introduction	47
2	The adjoint matrix	47
3	Proof of Theorem 3.1	49
4	Bordered determinants	51
5	The matrix equation $AX = 0$	51
6	The Hadamard product	52
7	The commutation matrix K_{mn}	54
8	The duplication matrix D_n	56
9	Relationship between D_{n+1} and D_n , I	58
10	Relationship between D_{n+1} and D_n , II	59
11	Conditions for a quadratic form to be positive (negative) subject to linear constraints	60
12	Necessary and sufficient conditions for $r(A : B) = r(A) + r(B)$	63
13	The bordered Gramian matrix	65
14	The equations $X_1A + X_2B' = G_1, X_1B = G_2$	67
	<i>Miscellaneous exercises</i>	69
	<i>Bibliographical notes</i>	70

Part Two — Differentials: the theory

4	<i>Mathematical preliminaries</i>	73
1	Introduction	73
2	Interior points and accumulation points	73
3	Open and closed sets	75
4	The Bolzano-Weierstrass theorem	77
5	Functions	78
6	The limit of a function	79
7	Continuous functions and compactness	80
8	Convex sets	81
9	Convex and concave functions	83
	<i>Bibliographical notes</i>	86
5	<i>Differentials and differentiability</i>	87
1	Introduction	87
2	Continuity	88
3	Differentiability and linear approximation	90
4	The differential of a vector function	91
5	Uniqueness of the differential	93
6	Continuity of differentiable functions	94

7	Partial derivatives	95
8	The first identification theorem	96
9	Existence of the differential, I	97
10	Existence of the differential, II	99
11	Continuous differentiability	100
12	The chain rule	100
13	Cauchy invariance	102
14	The mean-value theorem for real-valued functions	103
15	Differentiable matrix functions	104
16	Some remarks on notation	106
17	Complex differentiation	108
	<i>Miscellaneous exercises</i>	110
	<i>Bibliographical notes</i>	110
6	<i>The second differential</i>	111
1	Introduction	111
2	Second-order partial derivatives	111
3	The Hessian matrix	112
4	Twice differentiability and second-order approximation, I	113
5	Definition of twice differentiability	114
6	The second differential	115
7	Symmetry of the Hessian matrix	117
8	The second identification theorem	119
9	Twice differentiability and second-order approximation, II	119
10	Chain rule for Hessian matrices	121
11	The analog for second differentials	123
12	Taylor's theorem for real-valued functions	124
13	Higher-order differentials	125
14	Real analytic functions	125
15	Twice differentiable matrix functions	126
	<i>Bibliographical notes</i>	127
7	<i>Static optimization</i>	129
1	Introduction	129
2	Unconstrained optimization	130
3	The existence of absolute extrema	131
4	Necessary conditions for a local minimum	132
5	Sufficient conditions for a local minimum: first-derivative test	134
6	Sufficient conditions for a local minimum: second-derivative test	136
7	Characterization of differentiable convex functions	138
8	Characterization of twice differentiable convex functions	141
9	Sufficient conditions for an absolute minimum	142
10	Monotonic transformations	143
11	Optimization subject to constraints	144
12	Necessary conditions for a local minimum under constraints	145
13	Sufficient conditions for a local minimum under constraints	149
14	Sufficient conditions for an absolute minimum under constraints	154

15	A note on constraints in matrix form	155
16	Economic interpretation of Lagrange multipliers	155
	<i>Appendix: the implicit function theorem</i>	157
	<i>Bibliographical notes</i>	159

Part Three — Differentials: the practice

8	<i>Some important differentials</i>	163
1	Introduction	163
2	Fundamental rules of differential calculus	163
3	The differential of a determinant	165
4	The differential of an inverse	168
5	Differential of the Moore-Penrose inverse	169
6	The differential of the adjoint matrix	172
7	On differentiating eigenvalues and eigenvectors	174
8	The continuity of eigenprojections	176
9	The differential of eigenvalues and eigenvectors: symmetric case	180
10	Two alternative expressions for $d\lambda$	183
11	Second differential of the eigenvalue function	185
	<i>Miscellaneous exercises</i>	186
	<i>Bibliographical notes</i>	189
9	<i>First-order differentials and Jacobian matrices</i>	191
1	Introduction	191
2	Classification	192
3	Derisatives	192
4	Derivatives	194
5	Identification of Jacobian matrices	196
6	The first identification table	197
7	Partitioning of the derivative	197
8	Scalar functions of a scalar	198
9	Scalar functions of a vector	198
10	Scalar functions of a matrix, I: trace	199
11	Scalar functions of a matrix, II: determinant	201
12	Scalar functions of a matrix, III: eigenvalue	202
13	Two examples of vector functions	203
14	Matrix functions	204
15	Kronecker products	206
16	Some other problems	208
17	Jacobians of transformations	209
	<i>Bibliographical notes</i>	210
10	<i>Second-order differentials and Hessian matrices</i>	211
1	Introduction	211
2	The second identification table	211
3	Linear and quadratic forms	212
4	A useful theorem	213

5	The determinant function	214
6	The eigenvalue function	215
7	Other examples	215
8	Composite functions	217
9	The eigenvector function	218
10	Hessian of matrix functions, I	219
11	Hessian of matrix functions, II	219
	<i>Miscellaneous exercises</i>	220

Part Four — Inequalities

11	<i>Inequalities</i>	225
1	Introduction	225
2	The Cauchy-Schwarz inequality	226
3	Matrix analogs of the Cauchy-Schwarz inequality	227
4	The theorem of the arithmetic and geometric means	228
5	The Rayleigh quotient	230
6	Concavity of λ_1 and convexity of λ_n	232
7	Variational description of eigenvalues	232
8	Fischer's min-max theorem	234
9	Monotonicity of the eigenvalues	236
10	The Poincaré separation theorem	236
11	Two corollaries of Poincaré's theorem	237
12	Further consequences of the Poincaré theorem	238
13	Multiplicative version	239
14	The maximum of a bilinear form	241
15	Hadamard's inequality	242
16	An interlude: Karamata's inequality	242
17	Karamata's inequality and eigenvalues	244
18	An inequality concerning positive semidefinite matrices	245
19	A representation theorem for $(\sum a_i^p)^{1/p}$	246
20	A representation theorem for $(\text{tr } A^p)^{1/p}$	247
21	Hölder's inequality	248
22	Concavity of $\log A $	250
23	Minkowski's inequality	251
24	Quasilinear representation of $ A ^{1/n}$	253
25	Minkowski's determinant theorem	255
26	Weighted means of order p	256
27	Schlömilch's inequality	258
28	Curvature properties of $M_p(x, a)$	259
29	Least squares	260
30	Generalized least squares	261
31	Restricted least squares	262
32	Restricted least squares: matrix version	264
	<i>Miscellaneous exercises</i>	265
	<i>Bibliographical notes</i>	269

Part Five — The linear model

12	<i>Statistical preliminaries</i>	273
1	Introduction	273
2	The cumulative distribution function	273
3	The joint density function	274
4	Expectations	274
5	Variance and covariance	275
6	Independence of two random variables	277
7	Independence of n random variables	279
8	Sampling	279
9	The one-dimensional normal distribution	279
10	The multivariate normal distribution	280
11	Estimation	282
	<i>Miscellaneous exercises</i>	282
	<i>Bibliographical notes</i>	283
13	<i>The linear regression model</i>	285
1	Introduction	285
2	Affine minimum-trace unbiased estimation	286
3	The Gauss-Markov theorem	287
4	The method of least squares	290
5	Aitken's theorem	291
6	Multicollinearity	293
7	Estimable functions	295
8	Linear constraints: the case $\mathcal{M}(R') \subset \mathcal{M}(X')$	296
9	Linear constraints: the general case	300
10	Linear constraints: the case $\mathcal{M}(R') \cap \mathcal{M}(X') = \{0\}$	302
11	A singular variance matrix: the case $\mathcal{M}(X) \subset \mathcal{M}(V)$	304
12	A singular variance matrix: the case $r(X'V^+X) = r(X)$	305
13	A singular variance matrix: the general case, I	307
14	Explicit and implicit linear constraints	307
15	The general linear model, I	310
16	A singular variance matrix: the general case, II	311
17	The general linear model, II	314
18	Generalized least squares	315
19	Restricted least squares	316
	<i>Miscellaneous exercises</i>	318
	<i>Bibliographical notes</i>	319
14	<i>Further topics in the linear model</i>	321
1	Introduction	321
2	Best quadratic unbiased estimation of σ^2	322
3	The best quadratic and positive unbiased estimator of σ^2	322
4	The best quadratic unbiased estimator of σ^2	324
5	Best quadratic invariant estimation of σ^2	326

6	The best quadratic and positive invariant estimator of σ^2	327
7	The best quadratic invariant estimator of σ^2	329
8	Best quadratic unbiased estimation: multivariate normal case	330
9	Bounds for the bias of the least-squares estimator of σ^2 , I	332
10	Bounds for the bias of the least-squares estimator of σ^2 , II	333
11	The prediction of disturbances	335
12	Best linear unbiased predictors with scalar variance matrix	336
13	Best linear unbiased predictors with fixed variance matrix, I	338
14	Best linear unbiased predictors with fixed variance matrix, II	340
15	Local sensitivity of the posterior mean	341
16	Local sensitivity of the posterior precision	342
	<i>Bibliographical notes</i>	344

Part Six — Applications to maximum likelihood estimation

15	<i>Maximum likelihood estimation</i>	347
1	Introduction	347
2	The method of maximum likelihood (ML)	347
3	ML estimation of the multivariate normal distribution	348
4	Symmetry: implicit versus explicit treatment	350
5	The treatment of positive definiteness	351
6	The information matrix	352
7	ML estimation of the multivariate normal distribution: distinct means	354
8	The multivariate linear regression model	354
9	The errors-in-variables model	357
10	The nonlinear regression model with normal errors	359
11	Special case: functional independence of mean and variance parameters	361
12	Generalization of Theorem 15.6	362
	<i>Miscellaneous exercises</i>	364
	<i>Bibliographical notes</i>	365
16	<i>Simultaneous equations</i>	367
1	Introduction	367
2	The simultaneous equations model	367
3	The identification problem	369
4	Identification with linear constraints on B and Γ only	371
5	Identification with linear constraints on B , Γ , and Σ	371
6	Nonlinear constraints	373
7	FIML: the information matrix (general case)	374
8	FIML: asymptotic variance matrix (special case)	376
9	LIML: first-order conditions	378
10	LIML: information matrix	381
11	LIML: asymptotic variance matrix	383
	<i>Bibliographical notes</i>	388

17	<i>Topics in psychometrics</i>	389
1	Introduction	389
2	Population principal components	390
3	Optimality of principal components	391
4	A related result	392
5	Sample principal components	393
6	Optimality of sample principal components	395
7	One-mode component analysis	395
8	One-mode component analysis and sample principal components	398
9	Two-mode component analysis	399
10	Multimode component analysis	400
11	Factor analysis	404
12	A zigzag routine	407
13	A Newton-Raphson routine	408
14	Kaiser's varimax method	412
15	Canonical correlations and variates in the population	414
16	Correspondence analysis	417
17	Linear discriminant analysis	418
	<i>Bibliographical notes</i>	419

Part Seven — Summary

18	<i>Matrix calculus: the essentials</i>	423
1	Introduction	423
2	Differentials	424
3	Vector calculus	426
4	Optimization	429
5	Least squares	431
6	Matrix calculus	432
7	Interlude on linear and quadratic forms	434
8	The second differential	434
9	Chain rule for second differentials	436
10	Four examples	438
11	The Kronecker product and vec operator	439
12	Identification	441
13	The commutation matrix	442
14	From second differential to Hessian	443
15	Symmetry and the duplication matrix	444
16	Maximum likelihood	445
	<i>Further reading</i>	448
	<i>Bibliography</i>	449
	<i>Index of symbols</i>	467
	<i>Subject index</i>	471

Preface

Preface to the first edition

There has been a long-felt need for a book that gives a self-contained and unified treatment of matrix differential calculus, specifically written for econometricians and statisticians. The present book is meant to satisfy this need. It can serve as a textbook for advanced undergraduates and postgraduates in econometrics and as a reference book for practicing econometricians. Mathematical statisticians and psychometricians may also find something to their liking in the book.

When used as a textbook, it can provide a full-semester course. Reasonable proficiency in basic matrix theory is assumed, especially with the use of partitioned matrices. The basics of matrix algebra, as deemed necessary for a proper understanding of the main subject of the book, are summarized in Part One, the first of the book's six parts. The book also contains the essentials of multivariable calculus but geared to and often phrased in terms of differentials.

The sequence in which the chapters are being read is not of great consequence. It is fully conceivable that practitioners start with Part Three (Differentials: the practice) and, dependent on their predilections, carry on to Parts Five or Six, which deal with applications. Those who want a full understanding of the underlying theory should read the whole book, although even then they could go through the necessary matrix algebra only when the specific need arises.

Matrix differential calculus as presented in this book is based on differentials, and this sets the book apart from other books in this area. The approach via differentials is, in our opinion, superior to any other existing approach. Our principal idea is that differentials are more congenial to multivariable functions as they crop up in econometrics, mathematical statistics, or psychometrics than derivatives, although from a theoretical point of view the two concepts are equivalent.

The book falls into six parts. Part One deals with matrix algebra. It lists, and also often proves, items like the Schur, Jordan, and singular-value decompositions; concepts like the Hadamard and Kronecker products; the vec operator; the commutation and duplication matrices; and the Moore-Penrose

inverse. Results on bordered matrices (and their determinants) and (linearly restricted) quadratic forms are also presented here.

Part Two, which forms the theoretical heart of the book, is entirely devoted to a thorough treatment of the theory of differentials, and presents the essentials of calculus but geared to and phrased in terms of differentials. First and second differentials are defined, ‘identification’ rules for Jacobian and Hessian matrices are given, and chain rules derived. A separate chapter on the theory of (constrained) optimization in terms of differentials concludes this part.

Part Three is the practical core of the book. It contains the rules for working with differentials, lists the differentials of important scalar, vector, and matrix functions (*inter alia* eigenvalues, eigenvectors, and the Moore-Penrose inverse) and supplies ‘identification’ tables for Jacobian and Hessian matrices.

Part Four, treating inequalities, owes its existence to our feeling that econometricians should be conversant with inequalities, such as the Cauchy-Schwarz and Minkowski inequalities (and extensions thereof), and that they should also master a powerful result like Poincaré’s separation theorem. This part is to some extent also the case history of a disappointment. When we started writing this book we had the ambition to derive all inequalities by means of matrix differential calculus. After all, every inequality can be rephrased as the solution of an optimization problem. This proved to be an illusion, due to the fact that the Hessian matrix in most cases is singular at the optimum point.

Part Five is entirely devoted to applications of matrix differential calculus to the linear regression model. There is an exhaustive treatment of estimation problems related to the fixed part of the model under various assumptions concerning ranks and (other) constraints. Moreover, it contains topics relating to the stochastic part of the model, viz. estimation of the error variance and prediction of the error term. There is also a small section on sensitivity analysis. An introductory chapter deals with the necessary statistical preliminaries.

Part Six deals with maximum likelihood estimation, which is of course an ideal source for demonstrating the power of the propagated techniques. In the first of three chapters, several models are analysed, *inter alia* the multivariate normal distribution, the errors-in-variables model, and the nonlinear regression model. There is a discussion on how to deal with symmetry and positive definiteness, and special attention is given to the information matrix. The second chapter in this part deals with simultaneous equations under normality conditions. It investigates both identification and estimation problems, subject to various (non)linear constraints on the parameters. This part also discusses full-information maximum likelihood (FIML) and limited-information maximum likelihood (LIML), with special attention to the derivation of asymptotic variance matrices. The final chapter addresses itself to various psychometric problems, *inter alia* principal components, multimode component analysis, factor analysis, and canonical correlation.

All chapters contain many exercises. These are frequently meant to be complementary to the main text.

A large number of books and papers have been published on the theory and applications of matrix differential calculus. Without attempting to describe their relative virtues and particularities, the interested reader may wish to consult Dwyer and Macphail (1948), Bodewig (1959), Wilkinson (1965), Dwyer (1967), Neudecker (1967, 1969), Tracy and Dwyer (1969), Tracy and Singh (1972), McDonald and Swaminathan (1973), MacRae (1974), Balestra (1976), Bentler and Lee (1978), Henderson and Searle (1979), Wong and Wong (1979, 1980), Nel (1980), Rogers (1980), Wong (1980, 1985), Graham (1981), McCulloch (1982), Schönemann (1985), Magnus and Neudecker (1985), Pollock (1985), Don (1986), and Kollo (1991). The papers by Henderson and Searle (1979) and Nel (1980), and Rogers' (1980) book contain extensive bibliographies.

The two authors share the responsibility for Parts One, Three, Five, and Six, although any new results in Part One are due to Magnus. Parts Two and Four are due to Magnus, although Neudecker contributed some results to Part Four. Magnus is also responsible for the writing and organization of the final text.

We wish to thank our colleagues F. J. H. Don, R. D. H. Heijmans, D. S. G. Pollock, and R. Ramer for their critical remarks and contributions. The greatest obligation is owed to Sue Kirkbride at the London School of Economics who patiently and cheerfully typed and retyped the various versions of the book. Partial financial support was provided by the Netherlands Organization for the Advancement of Pure Research (Z. W. O.) and the Suntory Toyota International Centre for Economics and Related Disciplines at the London School of Economics.

London/Amsterdam
April 1987

Jan R. Magnus
Heinz Neudecker

Preface to the first revised printing

Since this book first appeared — now almost four years ago — many of our colleagues, students, and other readers have pointed out typographical errors and have made suggestions for improving the text. We are particularly grateful to R. D. H. Heijmans, J. F. Kiviet, I. J. Steyn, and G. Trenkler. We owe the greatest debt to F. Gerrish, formerly of the School of Mathematics in the Polytechnic, Kingston-upon-Thames, who read Chapters 1–11 with awesome precision and care and made numerous insightful suggestions and constructive remarks. We hope that this printing will continue to trigger comments from our readers.

London/Tilburg/Amsterdam
February 1991

Jan R. Magnus
Heinz Neudecker

Preface to the second edition

A further seven years have passed since our first revision in 1991. We are happy to see that our book is still being used by colleagues and students. In this revision we attempted to reach three goals. First, we made a serious attempt to keep the book up-to-date by adding many recent references and new exercises. Second, we made numerous small changes throughout the text, improving the clarity of exposition. Finally, we corrected a number of typographical and other errors.

The structure of the book and its philosophy are unchanged. Apart from a large number of small changes, there are two major changes. First, we interchanged Sections 12 and 13 of Chapter 1, since complex numbers need to be discussed before eigenvalues and eigenvectors, and we corrected an error in Theorem 1.7. Second, in Chapter 17 on psychometrics, we rewrote Sections 8–10 relating to the Eckart-Young theorem.

We are grateful to Karim Abadir, Paul Bekker, Hamparsum Bozdogan, Michael Browne, Frank Gerrish, Kaddour Hadri, Tõnu Kollo, Shuangzhe Liu, Daan Nel, Albert Satorra, Kazuo Shigemasu, Jos ten Berge, Peter ter Berg, Götz Trenkler, Haruo Yanai, and many others for their thoughtful and constructive comments. Of course, we welcome further comments from our readers.

Tilburg/Amsterdam
March 1998

Jan R. Magnus
Heinz Neudecker

Preface to the third edition

Twenty years have passed since the appearance of the second edition and thirty years since the book first appeared. This is a long time, but the book still lives. Unfortunately, my coauthor Heinz Neudecker does not; he died in December 2017. Heinz was my teacher at the University of Amsterdam and I was fortunate to learn the subject of matrix calculus through differentials (then in its infancy) from his lectures and personal guidance. This technique is still a remarkably powerful tool, and Heinz Neudecker must be regarded as its founding father.

The original text of the book was written on a typewriter and then handed over to the publisher for typesetting and printing. When it came to the second edition, the typeset material could no longer be found, which is why the second edition had to be produced in an *ad hoc* manner which was not satisfactory. Many people complained about this, to me and to the publisher, and the publisher offered us to produce a new edition, freshly typeset, which would look good. In the mean time, my Russian colleagues had proposed to translate the book into Russian, and I realized that this would only be feasible if they had a good English L^AT_EX text. So, my secretary Josette Janssen at Tilburg University and I produced a L^AT_EX text with expert advice from Jozef Pijenburg. In the process of retyping the manuscript, many small changes

were made to improve the readability and consistency of the text, but the structure of the book was not changed. The English L^AT_EX version was then used as the basis for the Russian edition,

*Matrichnoe Differentsial'noe Ischislenie s Prilozhenijami
k Statistike i Ekonometrike,*

translated by my friends Anatoly Peresetsky and Pavel Katyshev, and published by Fizmatlit Publishing House, Moscow, 2002. The current third edition is based on this English L^AT_EX version, although I have taken the opportunity to make many improvements to the presentation of the material.

Of course, this was not the only reason for producing a third edition. It was time to take a fresh look at the material and to update the references. I felt it was appropriate to stay close to the original text, because this is the book that Heinz and I conceived and the current text is a new edition, not a new book. The main changes relative to the second edition are as follows:

- Some subjects were treated insufficiently (some of my friends would say ‘incorrectly’) and I have attempted to repair these omissions. This applies in particular to the discussion on matrix functions (Section 1.21), complex differentiation (Section 5.17), and Jacobians of transformations (Section 9.17).
- The text on differentiating eigenvalues and eigenvectors and associated continuity issues has been rewritten, see Sections 8.7–8.11.
- Chapter 10 has been completely rewritten, because I am now convinced that it is not useful to define Hessian matrices for vector or matrix functions. So I now define Hessian matrices only for scalar functions and for individual components of vector functions and individual elements of matrix functions. This makes life much easier.
- I have added two additional sections at the end of Chapter 17 on psychometrics, relating to correspondence analysis and linear discriminant analysis.
- Chapter 18 is new. It can be read without consulting the other chapters and provides a summary of the whole book. It can therefore be used as an introduction to matrix calculus for advanced undergraduates or Master’s and PhD students in economics, statistics, mathematics, and engineering who want to know how to apply matrix calculus without going into all the theoretical details.

In addition, many small changes have been made, references have been updated, and exercises have been added. Over the past 30 years, I received many queries, problems, and requests from readers, about once every 2 weeks, which amounts to about 750 queries in 30 years. I responded to all of them and a number of these problems appear in the current text as exercises.

I am grateful to Don Andrews, Manuel Arellano, Richard Baillie, Luc Bauwens, Andrew Chesher, Gerda Claeskens, Russell Davidson, Jean-Marie

Dufour, Ronald Gallant, Eric Ghysels, Bruce Hansen, Grant Hillier, Cheng Hsiao, Guido Imbens, Guido Kuersteiner, Offer Lieberman, Esfandiar Maasoumi, Whitney Newey, Kazuhiro Ohtani, Enrique Sentana, Cezary Sielużycki, Richard Smith, Götz Trenkler, and Farshid Vahid for general encouragement and specific suggestions; to Henk Pijls for answering my questions on complex differentiation and Michel van de Velden for help on psychometric issues; to Jan Brinkhuis, Chris Muris, Franco Peracchi, Andrey Vasnev, Wendun Wang, and Yuan Yue on commenting on the new Chapter 18; to Ang Li for exceptional research assistance in updating the literature; and to Ilka van de Werve for expertly redrawing the figures. No blame attaches to any of these people in case there are remaining errors, ambiguities, or omissions; these are entirely my own responsibility, especially since I have not always followed their advice.

Cross-References. The numbering of theorems, propositions, corollaries, figures, tables, assumptions, examples, and definitions is with two digits, so that Theorem 3.5 refers to Theorem 5 in Chapter 3. Sections are numbered 1, 2, . . . within each chapter but always referenced with two digits so that Section 5 in Chapter 3 is referred to as Section 3.5. Equations are numbered (1), (2), . . . within each chapter, and referred to with one digit if it refers to the same chapter; if it refers to another chapter we write, for example, see Equation (16) in Chapter 5. Exercises are numbered 1, 2, . . . after a section.

Notation. Special symbols are used to denote the derivative (matrix) D and the Hessian (matrix) H . The differential operator is denoted by d . The third edition follows the notation of earlier editions with the following exceptions. First, the symbol for the vector $(1, 1, \dots, 1)'$ has been altered from a calligraphic s to ι (dotless i); second, the symbol i for imaginary root has been replaced by the more common i ; third, $v(A)$, the vector indicating the essentially distinct components of a symmetric matrix A , has been replaced by $\text{vech}(A)$; fourth, the symbols for expectation, variance, and covariance (previously \mathcal{E} , \mathcal{V} , and \mathcal{C}) have been replaced by E , var , and cov , respectively; and fifth, we now denote the normal distribution by N (previously \mathcal{N}). A list of all symbols is presented in the Index of Symbols at the end of the book.

Brackets are used sparingly. We write $\text{tr } A$ instead of $\text{tr}(A)$, while $\text{tr } AB$ denotes $\text{tr}(AB)$, not $(\text{tr } A)B$. Similarly, $\text{vec } AB$ means $\text{vec}(AB)$ and dXY means $d(XY)$. In general, we only place brackets when there is a possibility of ambiguity.

I worked on the third edition between April and November 2018. I hope the book will continue to be useful for a few more years, and of course I welcome comments from my readers.

Amsterdam/Wapserveen
November 2018

Jan R. Magnus

Note to the reader:

This preview of the the third edition of *Matrix Differential Calculus* contains only Chapter 18, the final chapter. This chapter is different from the other chapters in the book and can be read independently. It attempts to summarize matrix calculus for the user who does not want to go into all the details.

Chapter 18 is available free of charge. When quoting results from this chapter, please refer to:

Magnus, J. R. and H. Neudecker (2019). *Matrix Differential Calculus with Applications in Statistics and Econometrics*, Third Edition, John Wiley, New York.

CHAPTER 18

Matrix calculus: the essentials

1 INTRODUCTION

This chapter differs from the other chapters in this book. It attempts to summarize the theory and the practical applications of matrix calculus in a few pages, leaving out all the subtleties that the typical user will not need. It also serves as an introduction for (advanced) undergraduates or Master's and PhD students in economics, statistics, mathematics, and engineering, who want to know how to apply matrix calculus without going into all the theoretical details. The chapter can be read independently of the rest of the book.

We begin by introducing the concept of a differential, which lies at the heart of matrix calculus. The key advantage of the differential over the more common derivative is the following. Consider the linear vector function $f(x) = Ax$ where A is an $m \times n$ matrix of constants. Then, $f(x)$ is an $m \times 1$ vector function of an $n \times 1$ vector x , and the derivative $Df(x)$ is an $m \times n$ matrix (in this case, the matrix A). But the differential df remains an $m \times 1$ vector. In general, the differential df of a vector function $f = f(x)$ has the same dimension as f , irrespective of the dimension of the vector x , in contrast to the derivative $Df(x)$.

The advantage is even larger for matrices. The differential dF of a matrix function $F(X)$ has the same dimension as F , irrespective of the dimension of the matrix X . The practical importance of working with differentials is huge and will be demonstrated through many examples.

We next discuss vector calculus and optimization, with and without constraints. We emphasize the importance of a correct definition and notation for the derivative, present the 'first identification theorem', which links the first differential with the first derivative, and apply these results to least squares. Then we extend the theory from vector calculus to matrix calculus and obtain the differentials of the determinant and inverse.

Matrix Differential Calculus with Applications in Statistics and Econometrics,

Third Edition. Jan R. Magnus and Heinz Neudecker.

© 2019 John Wiley & Sons Ltd. Published 2019 by John Wiley & Sons Ltd.

A brief interlude on quadratic forms follows, the primary purpose of which is to show that if $x'Ax = 0$ for all x , then this does not imply that A is zero, but only that $A' = -A$. We then define the second differential and the Hessian matrix, prove the ‘second identification theorem’, which links the second differential with the Hessian matrix, and discuss the chain rule for second differentials. The first part of this chapter ends with four examples.

In the second (more advanced) part, we introduce the vec operator and the Kronecker product, and discuss symmetry (commutation and duplication matrices). Many examples are provided to clarify the technique. The chapter ends with an application to maximum likelihood estimation, where all elements discussed in the chapter come together.

The following notation is used. Unless specified otherwise, ϕ denotes a scalar function, f a vector function, and F a matrix function. Also, x denotes a scalar or vector argument, and X a matrix argument. All functions and variables in this chapter are real. Parentheses are used sparingly. We write dX , $\text{tr } X$, and $\text{vec } X$ without parentheses, and also dXY , $\text{tr } XY$, and $\text{vec } XY$ instead of $d(XY)$, $\text{tr}(XY)$, and $\text{vec}(XY)$. However, we write $\text{vech}(X)$ with parentheses for historical reasons.

2 DIFFERENTIALS

We assume that the reader is familiar with high-school calculus. This includes not only simple derivatives, such as

$$\frac{dx^2}{dx} = 2x, \quad \frac{de^x}{dx} = e^x, \quad \frac{d \sin x}{dx} = \cos x, \quad (1)$$

but also the chain rule, for example:

$$\frac{d(\sin x)^2}{dx} = \frac{d(\sin x)^2}{d \sin x} \frac{d \sin x}{dx} = 2 \sin x \cos x = \sin(2x).$$

We now introduce the concept of a *differential*, by expressing (1) as

$$dx^2 = 2x dx, \quad de^x = e^x dx, \quad d \sin x = \cos x dx, \quad (2)$$

where we write d rather than d to emphasize that this is a differential rather than a derivative. The two concepts are closely related, but they are not the same.

The concept of differential may be confusing for students who remember their mathematics teacher explain to them that it is wrong to view dx^2/dx as a fraction. They might wonder what dx and dx^2 really are. What does $dx^2 = 2x dx$ mean? From a geometric point of view, it means that if we replace the graph of the function $\phi(x) = x^2$ at some value x by its linear approximation, that is, by the tangent line at the point (x, x^2) , then an increment dx in x leads to an increment $dx^2 = 2x dx$ in x^2 in linear approximation. From an algebraic point of view, if we replace x by $x + dx$ (‘increment dx ’), then $\phi(x)$ is replaced by

$$\phi(x + dx) = (x + dx)^2 = x^2 + 2x dx + (dx)^2.$$

For small dx , the term $(dx)^2$ will be *very* small and, if we ignore it, we obtain the linear approximation $x^2 + 2x dx$. The differential dx^2 is, for a given value of x , just a function of the real variable dx , given by the formula $dx^2 = 2x dx$.

This may sound complicated, but working with differentials is easy. The passage from (1) to (2) holds generally for any (differentiable) real-valued function ϕ , and the differential $d\phi$ is thus given by the formula

$$d\phi = \frac{d\phi(x)}{dx} dx.$$

Put differently,

$$d\phi = \alpha(x) dx \iff \frac{d\phi(x)}{dx} = \alpha(x), \quad (3)$$

where α may depend on x , but not on dx . Equation (3) is a special case of the *first identification theorem* (Theorem 18.1) in the next section. It shows that we can *identify* the derivative from the differential (and vice versa), and it shows that the new concept differential is equivalent to the familiar concept derivative. We will always work with the differential, as this has great practical advantages.

The differential is an operator, in fact a linear operator, and we have

$$da = 0, \quad d(ax) = a dx,$$

for any scalar constant a , and

$$d(x + y) = dx + dy, \quad d(x - y) = dx - dy.$$

For the product and the ratio, we have

$$d(xy) = (dx)y + x dy, \quad d\left(\frac{1}{x}\right) = -\frac{dx}{x^2} \quad (x \neq 0),$$

and, in addition to the differential of the exponential function $de^x = e^x dx$,

$$d \log x = \frac{dx}{x} \quad (x > 0).$$

The chain rule, well-known for derivatives, also applies to differentials and is then called *Cauchy's rule of invariance*. For example,

$$d(\sin x)^2 = 2 \sin x d \sin x = 2 \sin x \cos x dx = \sin(2x) dx, \quad (4)$$

or

$$de^{x^2} = e^{x^2} dx^2 = 2e^{x^2} x dx,$$

or, combining the two previous examples,

$$\begin{aligned} de^{\sin x^2} &= e^{\sin x^2} d \sin x^2 = e^{\sin x^2} \cos x^2 dx^2 \\ &= 2x e^{\sin x^2} \cos x^2 dx. \end{aligned}$$

The chain rule is a good example of the general principle that things are easier — sometimes a bit, sometimes a lot — in terms of differentials than in terms of derivatives. The chain rule in terms of differentials states that taking differentials of functions preserves composition of functions. This is easier than the chain rule in terms of derivatives. Consider, for example, the function $z = h(x) = (\sin x)^2$ as the composition of the functions

$$y = g(x) = \sin x, \quad z = f(y) = y^2,$$

so that $h(x) = f(g(x))$. Then $dy = \cos x \, dx$, and hence

$$dz = 2y \, dy = 2y \cos x \, dx = 2 \sin x \cos x \, dx,$$

as expected.

The chain rule is, of course, a key instrument in differential calculus. Suppose we realize that x in (4) depends on t , say $x = t^2$. Then, we do not need to compute the differential of $(\sin t^2)^2$ all over again. We can use (4) and simply write

$$d(\sin t^2)^2 = \sin(2t^2) \, dt^2 = 2t \sin(2t^2) \, dt.$$

The chain rule thus allows us to apply the rules of calculus sequentially, one after another.

In this section, we have only concerned ourselves with scalar functions of a scalar argument, and the reader may wonder why we bother to introduce differentials. They do not seem to have a great advantage over the more familiar derivatives. This is true, but when we come to vector functions of vector arguments, then the advantage will become clear.

3 VECTOR CALCULUS

Let x ($n \times 1$) and y ($m \times 1$) be two vectors and let y be a function of x , say $y = f(x)$. What is the derivative of y with respect to x ? To help us answer this question, we first consider the linear equation

$$y = f(x) = Ax,$$

where A is an $m \times n$ matrix of constants. The derivative is A and we write

$$\frac{\partial f(x)}{\partial x'} = A. \tag{5}$$

The notation $\partial f(x)/\partial x'$ is just notation, nothing else. We sometimes write the derivative as $Df(x)$ or as Df , but we avoid the notation $f'(x)$ because this may cause confusion with the transpose. The proposed notation emphasizes that we differentiate an $m \times 1$ column vector f with respect to a $1 \times n$ row vector x' , resulting in an $m \times n$ derivative matrix.

More generally, the derivative of $f(x)$ is an $m \times n$ matrix containing all partial derivatives $\partial f_i(x)/\partial x_j$, but in a specific ordering, namely

$$\frac{\partial f(x)}{\partial x'} = \begin{pmatrix} \partial f_1(x)/\partial x_1 & \partial f_1(x)/\partial x_2 & \dots & \partial f_1(x)/\partial x_n \\ \partial f_2(x)/\partial x_1 & \partial f_2(x)/\partial x_2 & \dots & \partial f_2(x)/\partial x_n \\ \vdots & \vdots & & \vdots \\ \partial f_m(x)/\partial x_1 & \partial f_m(x)/\partial x_2 & \dots & \partial f_m(x)/\partial x_n \end{pmatrix}. \quad (6)$$

There is only one definition of a vector derivative, and this is it. Of course, one can organize the mn partial derivatives in different ways, but these other combinations of the partial derivatives are not derivatives, have no practical use, and should be avoided.

Notice that each row of the derivative in (6) contains the partial derivatives of *one* element of f with respect to *all* elements of x , and that each column contains the partial derivatives of *all* elements of f with respect to *one* element of x . This is an essential characteristic of the derivative. As a consequence, the derivative of a scalar function $y = \phi(x)$, such as $y = a'x$ (where a is a vector of constants), is a row vector; in this case, a' . So the derivative of $a'x$ is a' , not a .

The rules in the previous section imply that the following rules apply to vector differentials, where x and y are vectors and a is a vector of real constants, all of the same order:

$$da = 0, \quad d(x') = (dx)', \quad d(a'x) = a'dx,$$

$$d(x + y) = dx + dy, \quad d(x - y) = dx - dy,$$

and

$$d(x'y) = (dx)'y + x'dy.$$

Now we can see the advantage of working with differentials rather than with derivatives. When we have an $m \times 1$ vector y , which is a function of an $n \times 1$ vector of variables x , say $y = f(x)$, then the derivative is an $m \times n$ matrix, but the differential dy or df remains an $m \times 1$ vector. This is relevant for vector functions, and even more relevant for matrix functions and for second-order derivatives, as we shall see later. The practical advantage of working with differentials is therefore that the order does not increase but always stays the same.

Corresponding to the identification result (3), we have the following relationship between the differential and the derivative.

Theorem 18.1 (first identification theorem):

$$df = A(x) dx \iff \frac{\partial f(x)}{\partial x'} = A(x).$$

This theorem shows that there is a one-to-one correspondence between first-order differentials and first-order derivatives. In other words, the differential identifies the derivative.

Example 18.1: Consider the linear function $\phi(x) = a'x$, where a is a vector of constants. This gives

$$d\phi = a'dx,$$

so that the derivative is a' , as we have seen before.

Example 18.2a: Next, consider the quadratic function $\phi(x) = x'Ax$, where A is a matrix of constants. Here, we have

$$d\phi = (dx)'Ax + x'A dx = x'A'dx + x'A dx = x'(A + A') dx.$$

The derivative is $x'(A + A')$, and in the special case where A is symmetric, the derivative is $2x'A$.

Now suppose that $z = f(y)$ and that $y = g(x)$, so that $z = f(g(x))$. Then,

$$\frac{\partial z}{\partial x'} = \frac{\partial z}{\partial y'} \frac{\partial y}{\partial x'}.$$

This is the chain rule for vector functions. The corresponding result for differentials is the following.

Theorem 18.2 (chain rule for first differentials): Let $z = f(y)$ and $y = g(x)$, so that $z = f(g(x))$. Then,

$$dz = A(y)B(x) dx,$$

where $A(y)$ and $B(x)$ are defined through

$$dz = A(y) dy, \quad dy = B(x) dx.$$

Example 18.3: Let $x = (x_1, x_2, x_3)'$ and

$$f(x) = \begin{pmatrix} x_1^2 - x_2^2 \\ x_1x_2x_3 \end{pmatrix}.$$

Then, the differential is

$$\begin{aligned} df &= \begin{pmatrix} d(x_1^2) - d(x_2^2) \\ d(x_1x_2x_3) \end{pmatrix} = \begin{pmatrix} 2x_1 dx_1 - 2x_2 dx_2 \\ (dx_1)x_2x_3 + x_1(dx_2)x_3 + x_1x_2 dx_3 \end{pmatrix} \\ &= \begin{pmatrix} 2x_1 & -2x_2 & 0 \\ x_2x_3 & x_1x_3 & x_1x_2 \end{pmatrix} \begin{pmatrix} dx_1 \\ dx_2 \\ dx_3 \end{pmatrix}, \end{aligned}$$

which identifies the derivative as

$$\frac{\partial f(x)}{\partial x'} = \begin{pmatrix} 2x_1 & -2x_2 & 0 \\ x_2x_3 & x_1x_3 & x_1x_2 \end{pmatrix}.$$

Example 18.4a: Let $x = (x_1, x_2)'$, $y = (y_1, y_2)'$, and

$$\phi(y) = e^{y_1} \sin y_2, \quad y_1 = x_1x_2^2, \quad y_2 = x_1^2x_2.$$

Then,

$$d\phi = (de^{y_1}) \sin y_2 + e^{y_1} d \sin y_2 = a(y)' dy,$$

where

$$a(y) = e^{y_1} \begin{pmatrix} \sin y_2 \\ \cos y_2 \end{pmatrix}, \quad dy = \begin{pmatrix} dy_1 \\ dy_2 \end{pmatrix}.$$

Also,

$$dy = \begin{pmatrix} x_2^2 & 2x_1x_2 \\ 2x_1x_2 & x_1^2 \end{pmatrix} \begin{pmatrix} dx_1 \\ dx_2 \end{pmatrix} = B(x) dx.$$

Hence,

$$d\phi = a(y)' dy = a(y)' B(x) dx = c_1 dx_1 + c_2 dx_2,$$

where

$$\begin{aligned} c_1 &= x_2 e^{y_1} (x_2 \sin y_2 + 2x_1 \cos y_2), \\ c_2 &= x_1 e^{y_1} (x_1 \cos y_2 + 2x_2 \sin y_2), \end{aligned}$$

so that the derivative is $\partial\phi(x)/\partial x' = (c_1, c_2)$.

4 OPTIMIZATION

Let $\phi(x)$ be a scalar differentiable function that we wish to optimize with respect to an $n \times 1$ vector x . Then we obtain the differential $d\phi = a(x)' dx$, and set $a(x) = 0$. Suppose, for example, that we wish to minimize the function

$$\phi(x) = \frac{1}{2} x' Ax - b' x, \tag{7}$$

where the matrix A is positive definite. The differential is

$$d\phi = x' A dx - b' dx = (Ax - b)' dx.$$

(Recall that a positive definite matrix is symmetric, by definition.) The solution \hat{x} needs to satisfy $A\hat{x} - b = 0$, and hence $\hat{x} = A^{-1}b$. The function ϕ has an absolute minimum at \hat{x} , which can be seen by defining $y = x - \hat{x}$ and writing

$$y' Ay = (x - A^{-1}b)' A(x - A^{-1}b) = 2\phi(x) + b' A^{-1}b.$$

Since A is positive definite, $y' Ay$ has a minimum at $y = 0$ and hence $\phi(x)$ has a minimum at $x = \hat{x}$. This holds for the specific linear-quadratic function (7) and it holds more generally for any (strictly) convex function. Such functions attain a (strict) absolute minimum.

Next suppose there is a restriction, say $g(x) = 0$. Then we need to optimize subject to the restriction, and we need Lagrangian theory. This works as follows. First define the Lagrangian function, usually referred to as the Lagrangian,

$$\psi(x) = \phi(x) - \lambda g(x),$$

where λ is the Lagrange multiplier. Then we obtain the differential of ψ with respect to x ,

$$d\psi = d\phi - \lambda dg,$$

and set it equal to zero. The equations

$$\frac{\partial\phi(x)}{\partial x'} = \lambda \frac{\partial g(x)}{\partial x'}, \quad g(x) = 0$$

are the *first-order conditions*. From these $n + 1$ equations in $n + 1$ unknowns (x and λ), we solve x and λ .

If the constraint g is a vector rather than a scalar, then we have not one but several (say, m) constraints. In that case we need m multipliers and it works like this. First, define the Lagrangian

$$\psi(x) = \phi(x) - l'g(x),$$

where $l = (\lambda_1, \lambda_2, \dots, \lambda_m)'$ is a vector of Lagrange multipliers. Then, we obtain the differential of ψ with respect to x :

$$d\psi = d\phi - l' dg$$

and set it equal to zero. The equations

$$\frac{\partial\phi(x)}{\partial x'} = l' \frac{\partial g(x)}{\partial x'}, \quad g(x) = 0$$

constitute $n + m$ equations (the first-order conditions). If we can solve these equations, then we obtain the solutions, say \hat{x} and \hat{l} .

The Lagrangian method gives necessary conditions for a local constrained extremum to occur at a given point \hat{x} . But how do we know that this point is in fact a maximum or a minimum? Sufficient conditions are available but they may be difficult to verify. However, in the special case where ϕ is linear-quadratic (or more generally, convex) and g is linear, ϕ attains an absolute minimum at the solution \hat{x} under the constraint $g(x) = 0$.

5 LEAST SQUARES

Suppose we are given an $n \times 1$ vector y and an $n \times k$ matrix X with linearly independent columns, so that $r(X) = k$. We wish to find a $k \times 1$ vector β , such that $X\beta$ is 'as close as possible' to y in the sense that the 'error' vector $e = y - X\beta$ is minimized. A convenient scalar measure of the 'error' would be $e'e$ and our objective is to minimize

$$\phi(\beta) = \frac{e'e}{2} = \frac{(y - X\beta)'(y - X\beta)}{2}, \quad (8)$$

where we note that we write $e'e/2$ rather than $e'e$. This makes no difference, since any β which minimizes $e'e$ will also minimize $e'e/2$, but it is a common trick, useful because we know that we are minimizing a quadratic function, so that a '2' will appear in the derivative. The $1/2$ neutralizes this 2.

Differentiating ϕ in (8) gives

$$d\phi = e' de = e' d(y - X\beta) = -e' X d\beta.$$

Hence, the optimum is obtained when $X'e = 0$, that is, when $X'X\hat{\beta} = X'y$, from which we obtain

$$\hat{\beta} = (X'X)^{-1}X'y, \quad (9)$$

the least-squares solution.

If there are constraints on β , say $R\beta = r$, then we need to solve

$$\begin{array}{ll} \text{minimize} & \phi(\beta) \\ \text{subject to} & R\beta = r. \end{array}$$

We assume that the m rows of R are linearly independent, and define the Lagrangian

$$\psi(\beta) = (y - X\beta)'(y - X\beta)/2 - l'(R\beta - r),$$

where l is a vector of Lagrange multipliers. The differential is

$$\begin{aligned} d\psi &= d(y - X\beta)'(y - X\beta)/2 - l' d(R\beta - r) \\ &= (y - X\beta)' d(y - X\beta) - l' R d\beta \\ &= -(y - X\beta)' X d\beta - l' R d\beta. \end{aligned}$$

Setting the differential equal to zero and denoting the restricted estimators by $\tilde{\beta}$ and \tilde{l} , we obtain the first-order conditions

$$(y - X\tilde{\beta})' X + \tilde{l}' R = 0, \quad R\tilde{\beta} = r,$$

or, written differently,

$$X'X\tilde{\beta} - X'y = R'\tilde{l}, \quad R\tilde{\beta} = r.$$

We do not know $\tilde{\beta}$ but we know $R\tilde{\beta}$. Hence, we premultiply by $R(X'X)^{-1}$. Letting $\hat{\beta} = (X'X)^{-1}X'y$ as in (9), this gives

$$r - R\hat{\beta} = R(X'X)^{-1}R'\tilde{l}.$$

Since R has full row rank, we can solve for l :

$$\tilde{l} = (R(X'X)^{-1}R')^{-1}(r - R\hat{\beta}),$$

and hence for β :

$$\tilde{\beta} = \hat{\beta} + (X'X)^{-1}R'\tilde{l} = \hat{\beta} + (X'X)^{-1}R'(R(X'X)^{-1}R')^{-1}(r - R\hat{\beta}).$$

Since the constraint is linear and the function ϕ is linear-quadratic as in (7), it follows that the solution $\tilde{\beta}$ indeed minimizes $\phi(\beta) = e'e/2$ under the constraint $R\beta = r$.

6 MATRIX CALCULUS

We have moved from scalar calculus to vector calculus, now we move from vector calculus to matrix calculus. When discussing matrices we assume that the reader is familiar with matrix addition and multiplication, and also knows the concepts of a determinant $|A|$ and an inverse A^{-1} . An important function of a square matrix $A = (a_{ij})$ is its *trace*, which is defined as the sum of the diagonal elements of A : $\text{tr } A = \sum_i a_{ii}$. We have

$$\text{tr } A = \text{tr } A',$$

which is obvious because a matrix and its transpose have the same diagonal elements. Less obvious is

$$\text{tr } A'B = \text{tr } BA'$$

for any two matrices A and B of the same order (but not necessarily square). This follows because

$$\begin{aligned} \text{tr } A'B &= \sum_j (A'B)_{jj} = \sum_j \sum_i a_{ij}b_{ij} \\ &= \sum_i \sum_j b_{ij}a_{ij} = \sum_i (BA')_{ii} = \text{tr } BA'. \end{aligned} \quad (10)$$

The rules for vector differentials in Section 3 carry over to matrix differentials. Let A be a matrix of constants and let α be a scalar. Then, for any X ,

$$dA = 0, \quad d(\alpha X) = \alpha dX, \quad d(X') = (dX)',$$

and, for square X ,

$$d \text{tr } X = \text{tr } dX.$$

If X and Y are of the same order, then

$$d(X + Y) = dX + dY, \quad d(X - Y) = dX - dY,$$

and, if the matrix product XY is defined,

$$d(XY) = (dX)Y + XdY.$$

Two less trivial differentials are the determinant and the inverse. For nonsingular X we have

$$d|X| = |X| \operatorname{tr} X^{-1} dX, \tag{11}$$

and in particular, when $|X| > 0$,

$$d \log |X| = \frac{d|X|}{|X|} = \operatorname{tr} X^{-1} dX.$$

The proof of (11) is a little tricky and is omitted (in this chapter, but not in Chapter 8).

The differential of the inverse is, for nonsingular X ,

$$dX^{-1} = -X^{-1}(dX)X^{-1}. \tag{12}$$

This we can prove easily by considering the equation $X^{-1}X = I$. Differentiating both sides gives

$$(dX^{-1})X + X^{-1}dX = 0$$

and the result then follows by postmultiplying with X^{-1} .

The chain rule also applies to matrix functions. More precisely, if $Z = F(Y)$ and $Y = G(X)$, so that $Z = F(G(X))$, then

$$dZ = A(Y)B(X) dX,$$

where $A(Y)$ and $B(X)$ are defined through

$$dZ = A(Y) dY, \quad dY = B(X) dX,$$

as in Theorem 18.2.

Regarding constrained optimization, treated for vector functions in Section 18.4, we note that this can be easily and elegantly extended to matrix constraints. If we have a matrix G (rather than a vector g) of constraints and a matrix X (rather than a vector x) of variables, then we define a matrix of multipliers $L = (\lambda_{ij})$ of the same dimension as $G = (g_{ij})$. The Lagrangian then becomes

$$\psi(X) = \phi(X) - \operatorname{tr} L'G(X),$$

where we have used the fact, also used in (10) above, that

$$\operatorname{tr} L'G = \sum_i \sum_j \lambda_{ij} g_{ij}.$$

7 INTERLUDE ON LINEAR AND QUADRATIC FORMS

Before we turn from first to second differentials, that is, from linear forms to quadratic forms, we investigate under what conditions a linear or quadratic form vanishes. The sole purpose of this section is to help the reader appreciate Theorem 18.3 in the next section.

A *linear* form is an expression such as Ax . When $Ax = 0$, this does not imply that either A or x is zero. For example, if

$$A = \begin{pmatrix} 1 & -1 \\ -2 & 2 \\ 3 & -3 \end{pmatrix}, \quad x = \begin{pmatrix} 1 \\ 1 \end{pmatrix},$$

then $Ax = 0$, but neither $A = 0$ nor $x = 0$.

However, when $Ax = 0$ for every x , then A must be zero, which can be seen by taking x to be each elementary vector e_i in turn. (The i th elementary vector is the vector with one in the i th position and zeros elsewhere.)

A *quadratic* form is an expression such as $x'Ax$. When $x'Ax = 0$, this does not imply that $A = 0$ or $x = 0$ or $Ax = 0$. Even when $x'Ax = 0$ for every x , it does not follow that $A = 0$, as the example

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

demonstrates. This matrix A is skew-symmetric, that is, it satisfies $A' = -A$. In fact, when $x'Ax = 0$ for every x then it follows that A must be skew-symmetric. This can be seen by taking $x = e_i$ which implies that $a_{ii} = 0$, and then $x = e_i + e_j$ which implies that $a_{ij} + a_{ji} = 0$.

In the special case where $x'Ax = 0$ for every x and A is symmetric, then A is both symmetric ($A' = A$) and skew-symmetric ($A' = -A$), and hence $A = 0$.

8 THE SECOND DIFFERENTIAL

The second differential is simply the differential of the first differential:

$$d^2f = d(df).$$

Higher-order differentials are similarly defined, but they are seldom needed.

Example 18.2b: Let $\phi(x) = x'Ax$. Then, $d\phi = x'(A + A')dx$ and

$$\begin{aligned} d^2\phi &= d(x'(A + A')dx) = (dx)'(A + A')dx + x'(A + A')d^2x \\ &= (dx)'(A + A')dx, \end{aligned}$$

since $d^2x = 0$.

The first differential leads to the first derivative (sometimes called the *Jacobian matrix*) and the second differential leads to the second derivative (called the *Hessian matrix*). We emphasize that the concept of Hessian matrix is only useful for scalar functions, not for vector or matrix functions. When we have a vector function f we shall consider the Hessian matrix of each element of f separately, and when we have a matrix function F we shall consider the Hessian matrix of each element of F separately.

Thus, let ϕ be a scalar function and let

$$d\phi = a(x)' dx, \quad da = (H\phi) dx, \quad (13)$$

where

$$a(x)' = \frac{\partial\phi(x)}{\partial x'}, \quad H\phi = \frac{\partial a(x)}{\partial x'} = \frac{\partial}{\partial x'} \left(\frac{\partial\phi(x)}{\partial x'} \right)'.$$

The ij th element of the Hessian matrix $H\phi$ is thus obtained by first calculating $a_j(x) = \partial\phi(x)/\partial x_j$ and then $(H\phi)_{ij} = \partial a_j(x)/\partial x_i$. The Hessian matrix contains all second-order partial derivatives $\partial^2\phi(x)/\partial x_i \partial x_j$, and it is *symmetric* if ϕ is twice differentiable.

The Hessian matrix is often written as

$$H\phi = \frac{\partial^2\phi(x)}{\partial x \partial x'}, \quad (14)$$

where the expression on the right-hand side is a notation, the precise meaning of which is given by

$$\frac{\partial^2\phi(x)}{\partial x \partial x'} = \frac{\partial}{\partial x'} \left(\frac{\partial\phi(x)}{\partial x'} \right)'. \quad (15)$$

Given (13) and using the symmetry of $H\phi$, we obtain the second differential as

$$d^2\phi = (da)' dx = (dx)'(H\phi) dx,$$

which shows that the second differential of ϕ is a quadratic form in dx .

Now, suppose that we have obtained, after some calculations, that $d^2\phi = (dx)'B(x) dx$. Then,

$$(dx)'(H\phi - B(x)) dx = 0$$

for all dx . Does this imply that $H\phi = B(x)$? No, it does not, as we have seen in the previous section. It does, however, imply that

$$(H\phi - B(x))' + (H\phi - B(x)) = 0,$$

and hence that $H\phi = (B(x) + B(x)')/2$, using the symmetry of $H\phi$. This proves the following result.

Theorem 18.3 (second identification theorem):

$$d^2\phi = (dx)'B(x) dx \iff H\phi = \frac{B(x) + B(x)'}{2}.$$

The second identification theorem shows that there is a one-to-one correspondence between second-order differentials and second-order derivatives, but only if we make the matrix $B(x)$ in the quadratic form symmetric. Hence, the second differential identifies the second derivative.

Example 18.2c: Consider again the quadratic function $\phi(x) = x'Ax$. Then we can start with $d\phi = x'(A + A')dx$, as in Example 18.2b, and obtain $d^2\phi = (dx)'(A + A')dx$. The matrix in the quadratic form is already symmetric, so we obtain directly $H\phi = A + A'$.

Alternatively — and this is often quicker — we differentiate ϕ twice without writing out the first differential in its final form. From

$$d\phi = (dx)'Ax + x'A dx,$$

we thus obtain

$$d^2\phi = 2(dx)'A dx, \quad (16)$$

which identifies the Hessian matrix as $H\phi = A + A'$. (Notice that the matrix A in (16) is not necessarily symmetric.)

Even with such a simple function as $\phi(x) = x'Ax$, the advantage and elegance of using differentials is clear. Without differentials we would need to prove first that $\partial a'x/\partial x' = a'$ and $\partial x'Ax/\partial x' = x'(A + A')$, and then use (15) to obtain

$$\frac{\partial^2 x'Ax}{\partial x \partial x'} = \frac{\partial (x'(A + A'))'}{\partial x'} = \frac{\partial (A + A')x}{\partial x'} = A + A',$$

which is cumbersome in this simple case and not practical in more complex situations.

9 CHAIN RULE FOR SECOND DIFFERENTIALS

Let us now return to Example 18.2b. The function ϕ in this example is a function of x , and x is the argument of interest. This is why $d^2x = 0$. But if ϕ is a function of x , which in turn is a function of t , then it is no longer true that d^2x equals zero. More generally, suppose that $z = f(y)$ and that $y = g(x)$, so that $z = f(g(x))$. Then,

$$dz = A(y) dy$$

and

$$d^2z = (dA) dy + A(y) d^2y. \quad (17)$$

This is true whether or not y depends on some other variables. If we think of z as a function of y , then $d^2y = 0$, but if y depends on x then d^2y is not zero; in fact,

$$dy = B(x) dx, \quad d^2y = (dB) dx.$$

This gives us the following result.

Theorem 18.4 (chain rule for second differentials): Let $z = f(y)$ and $y = g(x)$, so that $z = f(g(x))$. Then,

$$d^2z = (dA)B(x) dx + A(y)(dB) dx,$$

where $A(y)$ and $B(x)$ are defined through

$$dz = A(y) dy, \quad dy = B(x) dx.$$

In practice, one usually avoids Theorem 18.4 by going back to the first differential $dz = A(y) dy$ and differentiating again. This gives (17), from which we obtain the result step by step.

Example 18.4b: Let

$$\phi(y_1, y_2) = e^{y_1} \sin y_2, \quad y_1 = x_1 x_2^2, \quad y_2 = x_1^2 x_2.$$

Then, by Theorem 18.4,

$$d^2\phi = (da)'B(x) dx + a(y)'(dB) dx,$$

where

$$a(y) = e^{y_1} \begin{pmatrix} \sin y_2 \\ \cos y_2 \end{pmatrix}, \quad B(x) = \begin{pmatrix} x_2^2 & 2x_1 x_2 \\ 2x_1 x_2 & x_1^2 \end{pmatrix}.$$

Now, letting

$$C(y) = e^{y_1} \begin{pmatrix} \sin y_2 & \cos y_2 \\ \cos y_2 & -\sin y_2 \end{pmatrix}$$

and

$$D_1(x) = 2 \begin{pmatrix} 0 & x_2 \\ x_2 & x_1 \end{pmatrix}, \quad D_2(x) = 2 \begin{pmatrix} x_2 & x_1 \\ x_1 & 0 \end{pmatrix},$$

we obtain

$$da = C(y) dy = C(y)B(x) dx$$

and

$$dB = (dx_1)D_1(x) + (dx_2)D_2(x).$$

It is convenient to write dx_1 and dx_2 in terms of dx , which can be done by defining $e_1 = (1, 0)'$ and $e_2 = (0, 1)'$. Then, $dx_1 = e_1' dx$ and $dx_2 = e_2' dx$, and hence

$$\begin{aligned} d^2\phi &= (da)'B(x) dx + a(y)'(dB) dx \\ &= (dx)'B(x)C(y)B(x) dx + a(y)'((dx_1)D_1(x) + (dx_2)D_2(x)) dx \\ &= (dx)'B(x)C(y)B(x) dx + (dx)'e_1 a(y)'D_1(x) dx + (dx)'e_2 a(y)'D_2(x) dx \\ &= (dx)' [B(x)C(y)B(x) + e_1 a(y)'D_1(x) + e_2 a(y)'D_2(x)] dx. \end{aligned}$$

Some care is required where to position the scalars $e_1'dx$ and $e_2'dx$ in the matrix product. A scalar can be positioned anywhere in a matrix product, but we wish to position the two scalars in such a way that the usual matrix multiplication rules still apply.

Having obtained the second differential in the desired form, Theorem 18.3 implies that the Hessian is equal to

$$\begin{aligned} \mathbf{H}\phi &= B(x)C(y)B(x) + \frac{1}{2} (e_1a(y)'D_1(x) + D_1(x)a(y)e_1') \\ &\quad + \frac{1}{2} (e_2a(y)'D_2(x) + D_2(x)a(y)e_2'). \end{aligned}$$

10 FOUR EXAMPLES

Let us provide four examples to show how the second differential can be obtained. The first three examples relate to scalar functions and the fourth example to a matrix function. The matrix X has order $n \times q$ in Examples 18.5a and 18.6a, and order $n \times n$ in Examples 18.7a and 18.8a.

Example 18.5a: Let $\phi(X) = \text{tr } X'AX$. Then,

$$\begin{aligned} d\phi &= d(\text{tr } X'AX) = \text{tr } d(X'AX) \\ &= \text{tr}(dX)'AX + \text{tr } X'A dX = \text{tr } X'(A + A') dX \end{aligned}$$

and

$$d^2\phi = d \text{tr } X'(A + A') dX = \text{tr}(dX)'(A + A') dX.$$

Example 18.6a: Let $\phi(X) = \log |X'X|$. Then,

$$\begin{aligned} d\phi &= d \log |X'X| = \text{tr}(X'X)^{-1} d(X'X) \\ &= \text{tr}(X'X)^{-1} (dX)'X + \text{tr}(X'X)^{-1} X' dX = 2 \text{tr}(X'X)^{-1} X' dX \end{aligned}$$

and

$$\begin{aligned} d^2\phi &= 2 d (\text{tr}(X'X)^{-1} X' dX) \\ &= 2 \text{tr}(d(X'X)^{-1}) X' dX + 2 \text{tr}(X'X)^{-1} (dX)' dX \\ &= -2 \text{tr}(X'X)^{-1} (dX'X) (X'X)^{-1} X' dX + 2 \text{tr}(X'X)^{-1} (dX)' dX \\ &= -2 \text{tr}(X'X)^{-1} (dX)' X (X'X)^{-1} X' dX \\ &\quad - 2 \text{tr}(X'X)^{-1} X' (dX) (X'X)^{-1} X' dX + 2 \text{tr}(X'X)^{-1} (dX)' dX \\ &= 2 \text{tr}(X'X)^{-1} (dX)' M dX - 2 \text{tr}(X'X)^{-1} X' (dX) (X'X)^{-1} X' dX, \end{aligned}$$

where $M = I_n - X(X'X)^{-1}X'$. Let us explain some of the steps in more detail. The second equality follows from considering $(X'X)^{-1}X'dX$ as a product

of three matrices: $(X'X)^{-1}$, X' , and dX (a matrix of constants), the third equality uses the differential of the inverse in (12), and the fourth equality separates $dX'X$ into $(dX)'X + X'dX$.

Example 18.7a: Let $\phi(X) = \text{tr } X^k$ for $k = 1, 2, \dots$. Then, for $k \geq 1$,

$$\begin{aligned} d\phi &= \text{tr}(dX)X^{k-1} + \text{tr } X(dX)X^{k-2} + \dots + \text{tr } X^{k-1}dX \\ &= k \text{tr } X^{k-1}dX, \end{aligned}$$

and for $k \geq 2$,

$$d^2\phi = k \text{tr}(dX^{k-1})dX = k \sum_{j=0}^{k-2} \text{tr } X^j(dX)X^{k-2-j}dX.$$

Example 18.8a: Let $F(X) = AX^{-1}B$. Then,

$$dF = A(dX^{-1})B = -AX^{-1}(dX)X^{-1}B$$

and

$$\begin{aligned} d^2F &= -A(dX^{-1})(dX)X^{-1}B - AX^{-1}(dX)(dX^{-1})B \\ &= 2AX^{-1}(dX)X^{-1}(dX)X^{-1}B. \end{aligned}$$

These four examples provide the second differential; they do not yet provide the Hessian matrix. In Section 18.14, we shall discuss the same four examples and obtain the Hessian matrices.

11 THE KRONECKER PRODUCT AND VEC OPERATOR

The theory and the four examples in the previous two sections demonstrate the elegance and simplicity of obtaining first and second differentials of scalar, vector, and matrix functions. But we also want to relate these first and second differentials to Jacobian matrices (first derivatives) and Hessian matrices (second derivatives). For this we need some more machinery, namely the vec operator and the Kronecker product.

First, the vec operator. Consider an $m \times n$ matrix A . This matrix has n columns, say a_1, \dots, a_n . Now define the $mn \times 1$ vector $\text{vec } A$ as the vector which stacks these columns one underneath the other:

$$\text{vec } A = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix}.$$

For example, if

$$A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{pmatrix},$$

then $\text{vec } A = (1, 4, 2, 5, 3, 6)'$. Of course, we have

$$d \text{vec } X = \text{vec } dX. \quad (18)$$

If A and B are matrices of the same order, then we know from (10) that $\text{tr } A'B = \sum_{ij} a_{ij}b_{ij}$. But $(\text{vec } A)'(\text{vec } B)$ is also equal to this double sum. Hence,

$$\text{tr } A'B = (\text{vec } A)'(\text{vec } B), \quad (19)$$

an important equality linking the vec operator to the trace.

We also need the Kronecker product. Let A be an $m \times n$ matrix and B a $p \times q$ matrix. The $mp \times nq$ matrix defined by

$$\begin{pmatrix} a_{11}B & \dots & a_{1n}B \\ \vdots & & \vdots \\ a_{m1}B & \dots & a_{mn}B \end{pmatrix}$$

is called the *Kronecker product* of A and B and is written as $A \otimes B$. The Kronecker product $A \otimes B$ is thus defined for any pair of matrices A and B , unlike the matrix product AB which exists only if the number of columns in A equals the number of rows in B or if either A or B is a scalar.

The following three properties justify the name *Kronecker product*:

$$\begin{aligned} A \otimes B \otimes C &= (A \otimes B) \otimes C = A \otimes (B \otimes C), \\ (A + B) \otimes (C + D) &= A \otimes C + A \otimes D + B \otimes C + B \otimes D, \end{aligned}$$

if A and B have the same order and C and D have the same order (not necessarily equal to the order of A and B), and

$$(A \otimes B)(C \otimes D) = AC \otimes BD,$$

if AC and BD exist.

The transpose of a Kronecker product is

$$(A \otimes B)' = A' \otimes B'.$$

If A and B are square matrices (not necessarily of the same order), then

$$\text{tr}(A \otimes B) = (\text{tr } A)(\text{tr } B),$$

and if A and B are nonsingular, then

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

The Kronecker product and the vec operator are related through the equality

$$\text{vec } ab' = b \otimes a,$$

where a and b are column vectors of arbitrary order. Using this inequality, we see that

$$\begin{aligned} \text{vec}(Abe'C) &= \text{vec}(Ab)(C'e)' = (C'e) \otimes (Ab) \\ &= (C' \otimes A)(e \otimes b) = (C' \otimes A) \text{vec } be' \end{aligned}$$

for any vectors b and e . Then, writing $B = \sum_j b_j e_j'$ where b_j and e_j denote the j th column of B and I , respectively, we obtain the following important relationship, which is used frequently.

Theorem 18.5: For any matrices A , B , and C for which the product ABC is defined, we have

$$\text{vec } ABC = (C' \otimes A) \text{vec } B.$$

12 IDENTIFICATION

When we move from vector calculus to matrix calculus, we need an ordering of the functions and of the variables. It does not matter *how* we order them (any ordering will do), but an ordering is essential. We want to define matrix derivatives within the established theory of vector derivatives in such a way that trivial changes such as relabeling functions or variables have only trivial consequences for the derivative: rows and columns are permuted, but the rank is unchanged and the determinant (in the case of a square matrix) is also unchanged, apart possibly from its sign. This is what we need to achieve. The arrangement of the partial derivatives matters, because a derivative is more than just a collection of partial derivatives. It is a mathematical concept, a mathematical unit.

Thus motivated, we shall view the matrix function $F(X)$ as a vector function $f(x)$, where $f = \text{vec } F$ and $x = \text{vec } X$. We then obtain the following extension of the first identification theorem:

$$d \text{vec } F = A(X) d \text{vec } X \iff \frac{\partial \text{vec } F(X)}{\partial (\text{vec } X)'} = A(X),$$

and, similarly, for the second identification theorem:

$$d^2 \phi = (d \text{vec } X)' B(X) d \text{vec } X \iff H\phi = \frac{B(X) + B(X)'}{2},$$

where we notice, as in Section 18.8, that we only provide the Hessian matrix for scalar functions, not for vector or matrix functions.

13 THE COMMUTATION MATRIX

At this point, we need to introduce the commutation matrix. Let A be an $m \times n$ matrix. The vectors $\text{vec } A$ and $\text{vec } A'$ contain the same mn elements, but in a different order. Hence, there exists a unique $mn \times mn$ matrix, which transforms $\text{vec } A$ into $\text{vec } A'$. This matrix contains mn ones and $mn(mn - 1)$ zeros and is called the *commutation matrix*, denoted by K_{mn} . (If $m = n$, we write K_n instead of K_{nn} .) Thus,

$$K_{mn} \text{vec } A = \text{vec } A'. \quad (20)$$

It can be shown that K_{mn} is orthogonal, i.e. $K'_{mn} = K_{mn}^{-1}$. Also, premultiplying (20) by K_{nm} gives $K_{nm}K_{mn} \text{vec } A = \text{vec } A$, which shows that $K_{nm}K_{mn} = I_{mn}$. Hence,

$$K'_{mn} = K_{mn}^{-1} = K_{nm}.$$

The key property of the commutation matrix enables us to interchange (commute) the two matrices of a Kronecker product:

$$K_{pm}(A \otimes B) = (B \otimes A)K_{qn} \quad (21)$$

for any $m \times n$ matrix A and $p \times q$ matrix B . This is easiest shown, not by proving a matrix identity but by proving that the *effect* of the two matrices on an arbitrary vector is the same. Thus, let X be an arbitrary $q \times n$ matrix. Then, by repeated application of (20) and Theorem 18.5,

$$\begin{aligned} K_{pm}(A \otimes B) \text{vec } X &= K_{pm} \text{vec } BXA' = \text{vec } AX'B' \\ &= (B \otimes A) \text{vec } X' = (B \otimes A)K_{qn} \text{vec } X. \end{aligned}$$

Since X is arbitrary, (21) follows.

The commutation matrix has many applications in matrix theory. Its importance in matrix calculus stems from the fact that it transforms $d \text{vec } X'$ into $d \text{vec } X$. The simplest example is the matrix function $F(X) = X'$, where X is an $n \times q$ matrix. Then,

$$d \text{vec } F = \text{vec } dX' = K_{nq} \text{vec } dX,$$

so that the derivative is $D \text{vec } F = K_{nq}$.

The commutation matrix is also essential in identifying the Hessian matrix from the second differential. The second differential of a scalar function often takes the form of a trace, either $\text{tr } A(dX)'BdX$ or $\text{tr } A(dX)BdX$. We then have the following result, based on (19) and Theorem 18.5.

Theorem 18.6: Let ϕ be a twice differentiable real-valued function of an $n \times q$ matrix X . Then,

$$d^2\phi = \text{tr } A(dX)'BdX \iff H\phi = \frac{1}{2}(A' \otimes B + A \otimes B')$$

and

$$d^2\phi = \text{tr } A(dX)BdX \iff H\phi = \frac{1}{2}K_{qn}(A' \otimes B + B' \otimes A).$$

To identify the Hessian matrix from the first expression, we do not need the commutation matrix, but we do need the commutation matrix to identify the Hessian matrix from the second expression.

14 FROM SECOND DIFFERENTIAL TO HESSIAN

We continue with the same four examples as discussed in Section 18.10, showing how to obtain the Hessian matrices from the second differentials, using Theorem 18.6.

Example 18.5b: Let $\phi(X) = \text{tr } X'AX$, where X is an $n \times q$ matrix. Then,

$$d\phi = \text{tr } X'(A + A')dX = \text{tr } C'dX = (\text{vec } C)'d \text{vec } X,$$

where $C = (A + A')X$, and $d^2\phi = \text{tr}(dX)'(A + A')dX$. Hence, the derivative is $D\phi = (\text{vec } C)'$ and the Hessian is $H\phi = I_q \otimes (A + A')$.

Example 18.6b: Let $\phi(X) = \log |X'X|$, where X is an $n \times q$ matrix of full column rank. Then, letting $C = X(X'X)^{-1}$ and $M = I_n - X(X'X)^{-1}X'$,

$$d\phi = 2 \text{tr } C'dX = 2(\text{vec } C)'d \text{vec } X$$

and

$$d^2\phi = 2 \text{tr}(X'X)^{-1}(dX)'MdX - 2 \text{tr } C'(dX)C'dX.$$

This gives $D\phi = 2(\text{vec } C)'$ and

$$H\phi = 2(X'X)^{-1} \otimes M - 2K_{qn}(C \otimes C').$$

Example 18.7b: Let $\phi(X) = \text{tr } X^k$ for $k = 1, 2, \dots$, where X is a square $n \times n$ matrix. Then, for $k \geq 1$,

$$d\phi = k \text{tr } X^{k-1}dX = k(\text{vec } X'^{k-1})'d \text{vec } X,$$

and for $k \geq 2$,

$$d^2\phi = k \sum_{j=0}^{k-2} \text{tr } X^j(dX)X^{k-2-j}dX.$$

This gives $D\phi = k(\text{vec } X'^{k-1})'$ and

$$H\phi = (k/2) \sum_{j=0}^{k-2} K_n(X'^j \otimes X^{k-2-j} + X'^{k-2-j} \otimes X^j).$$

Example 18.8b: Let $F(X) = AX^{-1}B$, where X is a nonsingular $n \times n$ matrix. Then, using Theorem 18.5,

$$d \operatorname{vec} F = -((X^{-1}B)' \otimes (AX^{-1})) d \operatorname{vec} X,$$

and hence

$$D \operatorname{vec} F = -(X^{-1}B)' \otimes (AX^{-1}).$$

To obtain the Hessian matrix of the st th element of F , we let

$$C_{ts} = X^{-1}B e_t e_s' A X^{-1},$$

where e_s and e_t are elementary vectors with 1 in the s th (respectively, t th) position and zeros elsewhere. Then,

$$d^2 F_{st} = 2e_s' A X^{-1} (dX) X^{-1} (dX) X^{-1} B e_t = 2 \operatorname{tr} C_{ts} (dX) X^{-1} (dX)$$

and hence

$$H F_{st} = K_n (C_{ts}' \otimes X^{-1} + X'^{-1} \otimes C_{ts}).$$

15 SYMMETRY AND THE DUPLICATION MATRIX

Many matrices in statistics and econometrics are symmetric, for example variance matrices. When we differentiate with respect to symmetric matrices, we must take the symmetry into account and we need the duplication matrix.

Let A be a square $n \times n$ matrix. Then $\operatorname{vech}(A)$ will denote the $\frac{1}{2}n(n+1) \times 1$ vector that is obtained from $\operatorname{vec} A$ by eliminating all elements of A above the diagonal. For example, for $n = 3$,

$$\operatorname{vec} A = (a_{11}, a_{21}, a_{31}, a_{12}, a_{22}, a_{32}, a_{13}, a_{23}, a_{33})'$$

and

$$\operatorname{vech}(A) = (a_{11}, a_{21}, a_{31}, a_{22}, a_{32}, a_{33})'. \quad (22)$$

In this way, for symmetric A , $\operatorname{vech}(A)$ contains only the generically distinct elements of A . Since the elements of $\operatorname{vec} A$ are those of $\operatorname{vech}(A)$ with some repetitions, there exists a unique $n^2 \times \frac{1}{2}n(n+1)$ matrix which transforms, for symmetric A , $\operatorname{vech}(A)$ into $\operatorname{vec} A$. This matrix is called the *duplication matrix* and is denoted by D_n . Thus,

$$D_n \operatorname{vech}(A) = \operatorname{vec} A \quad (A = A'). \quad (23)$$

The matrix D_n has full column rank $\frac{1}{2}n(n+1)$, so that $D_n' D_n$ is nonsingular. This implies that $\operatorname{vech}(A)$ can be uniquely solved from (23), and we have

$$\operatorname{vech}(A) = (D_n' D_n)^{-1} D_n' \operatorname{vec} A \quad (A = A').$$

One can show (but we will not do so here) that the duplication matrix is connected to the commutation matrix by

$$K_n D_n = D_n, \quad D_n (D'_n D_n)^{-1} D'_n = \frac{1}{2} (I_{n^2} + K_n).$$

Much of the interest in the duplication matrix is due to the importance of the matrix $D'_n (A \otimes A) D_n$, where A is an $n \times n$ matrix. This matrix is important, because the scalar function $\phi(X) = \text{tr} AX'AX$ occurs frequently in statistics and econometrics, for example in the next section on maximum likelihood. When A and X are known to be symmetric we have

$$\begin{aligned} d^2 \phi &= 2 \text{tr} A(dX)'A dX = 2(d \text{vec } X)'(A \otimes A)d \text{vec } X \\ &= 2(d \text{vech}(X))'D'_n(A \otimes A)D_n d \text{vech}(X), \end{aligned}$$

and hence, $H\phi = 2D'_n(A \otimes A)D_n$.

From the relationship (again not proved here)

$$D_n (D'_n D_n)^{-1} D'_n (A \otimes A) D_n = (A \otimes A) D_n,$$

which is valid for any $n \times n$ matrix A , not necessarily symmetric, we obtain the inverse

$$(D'_n (A \otimes A) D_n)^{-1} = (D'_n D_n)^{-1} D'_n (A^{-1} \otimes A^{-1}) D_n (D'_n D_n)^{-1}, \quad (24)$$

where A is nonsingular. Finally, we present the determinant:

$$|D'_n (A \otimes A) D_n| = 2^{\frac{1}{2}n(n-1)} |A|^{n+1}. \quad (25)$$

16 MAXIMUM LIKELIHOOD

This final section brings together most of the material that has been treated in this chapter: first and second differentials, the Hessian matrix, and the treatment of symmetry (duplication matrix).

We consider a sample of $m \times 1$ vectors y_1, y_2, \dots, y_n from the multivariate normal distribution with mean μ and variance Ω , where Ω is positive definite and $n \geq m + 1$. The density of y_i is

$$f(y_i) = (2\pi)^{-m/2} |\Omega|^{-1/2} \exp \left(-\frac{1}{2} (y_i - \mu)' \Omega^{-1} (y_i - \mu) \right),$$

and since the y_i are independent and identically distributed, the joint density of (y_1, \dots, y_n) is given by $\prod_i f(y_i)$. The 'likelihood' is equal to the joint density, but now thought of as a function of the parameters μ and Ω , rather than of the observations. Its logarithm is the 'loglikelihood', which here takes the form

$$\Lambda(\mu, \Omega) = -\frac{mn}{2} \log 2\pi - \frac{n}{2} \log |\Omega| - \frac{1}{2} \sum_{i=1}^n (y_i - \mu)' \Omega^{-1} (y_i - \mu).$$

The maximum likelihood estimators are obtained by maximizing the loglikelihood (which is the same, but usually easier, as maximizing the likelihood). Thus, we differentiate Λ and obtain

$$\begin{aligned}
d\Lambda &= -\frac{n}{2} d \log |\Omega| + \frac{1}{2} \sum_{i=1}^n (d\mu)' \Omega^{-1} (y_i - \mu) \\
&\quad - \frac{1}{2} \sum_{i=1}^n (y_i - \mu)' (d\Omega^{-1}) (y_i - \mu) + \frac{1}{2} \sum_{i=1}^n (y_i - \mu)' \Omega^{-1} d\mu \\
&= -\frac{n}{2} d \log |\Omega| - \frac{1}{2} \sum_{i=1}^n (y_i - \mu)' (d\Omega^{-1}) (y_i - \mu) + \sum_{i=1}^n (y_i - \mu)' \Omega^{-1} d\mu \\
&= -\frac{n}{2} \text{tr}(\Omega^{-1} d\Omega + S d\Omega^{-1}) + \sum_{i=1}^n (y_i - \mu)' \Omega^{-1} d\mu, \tag{26}
\end{aligned}$$

where

$$S = S(\mu) = \frac{1}{n} \sum_{i=1}^n (y_i - \mu)(y_i - \mu)'.$$

Denoting the maximum likelihood estimators by $\hat{\mu}$ and $\hat{\Omega}$, letting $\hat{S} = S(\hat{\mu})$, and setting $d\Lambda = 0$ then implies that

$$\text{tr} \left(\hat{\Omega}^{-1} - \hat{\Omega}^{-1} \hat{S} \hat{\Omega}^{-1} \right) d\Omega = 0$$

for all $d\Omega$ and

$$\sum_{i=1}^n (y_i - \hat{\mu})' \hat{\Omega}^{-1} d\mu = 0$$

for all $d\mu$. This, in turn, implies that

$$\hat{\Omega}^{-1} = \hat{\Omega}^{-1} \hat{S} \hat{\Omega}^{-1}, \quad \sum_{i=1}^n (y_i - \hat{\mu}) = 0.$$

Hence, the maximum likelihood estimators are given by

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y} \tag{27}$$

and

$$\hat{\Omega} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})(y_i - \bar{y})'. \tag{28}$$

We note that the condition that Ω is symmetric has not been imposed. But since the solution (28) is symmetric, imposing the condition would have made no difference.

The second differential is obtained by differentiating (26) again. This gives

$$\begin{aligned} d^2\Lambda &= -\frac{n}{2} \operatorname{tr} ((d\Omega^{-1})d\Omega + (dS)d\Omega^{-1} + Sd^2\Omega^{-1}) - n(d\mu)'\Omega^{-1}d\mu \\ &\quad + \sum_{i=1}^n (y_i - \mu)'(d\Omega^{-1})d\mu. \end{aligned} \quad (29)$$

We are usually not primarily interested in the Hessian matrix but in its expectation. Hence, we do not evaluate (29) further and first take expectations. Since $E(S) = \Omega$ and $E(dS) = 0$, we obtain

$$\begin{aligned} E d^2\Lambda &= -\frac{n}{2} \operatorname{tr} ((d\Omega^{-1})d\Omega + \Omega d^2\Omega^{-1}) - n(d\mu)'\Omega^{-1}d\mu \\ &= \frac{n}{2} \operatorname{tr} \Omega^{-1}(d\Omega)\Omega^{-1}d\Omega - n \operatorname{tr}(d\Omega)\Omega^{-1}(d\Omega)\Omega^{-1} - n(d\mu)'\Omega^{-1}d\mu \\ &= -\frac{n}{2} \operatorname{tr} \Omega^{-1}(d\Omega)\Omega^{-1}d\Omega - n(d\mu)'\Omega^{-1}d\mu, \end{aligned} \quad (30)$$

using the facts that $d\Omega^{-1} = -\Omega^{-1}(d\Omega)\Omega^{-1}$ and

$$\begin{aligned} d^2\Omega^{-1} &= -(d\Omega^{-1})(d\Omega)\Omega^{-1} - \Omega^{-1}(d\Omega)d\Omega^{-1} \\ &= 2\Omega^{-1}(d\Omega)\Omega^{-1}(d\Omega)\Omega^{-1}. \end{aligned}$$

To obtain the 'information matrix' we need to take the symmetry of Ω into account and this is where the duplication matrix appears. So far, we have avoided the vec operator and in practical situations one should work with differentials (rather than with derivatives) as long as possible. But we cannot go further than (30) without use of the vec operator. Thus, from (30),

$$\begin{aligned} -E d^2\Lambda &= \frac{n}{2} \operatorname{tr} \Omega^{-1}(d\Omega)\Omega^{-1}d\Omega + n(d\mu)'\Omega^{-1}d\mu \\ &= \frac{n}{2} (d \operatorname{vec} \Omega)'(\Omega^{-1} \otimes \Omega^{-1}) d \operatorname{vec} \Omega + n(d\mu)'\Omega^{-1}d\mu \\ &= \frac{n}{2} (d \operatorname{vech}(\Omega))' D'_m (\Omega^{-1} \otimes \Omega^{-1}) D_m d \operatorname{vech}(\Omega) + n(d\mu)'\Omega^{-1}d\mu. \end{aligned}$$

Hence, the information matrix for μ and $\operatorname{vech}(\Omega)$ is

$$\mathcal{F} = n \begin{pmatrix} \Omega^{-1} & 0 \\ 0 & \frac{1}{2} D'_m (\Omega^{-1} \otimes \Omega^{-1}) D_m \end{pmatrix}.$$

The results on the duplication matrix in Section 18.15 also allow us to obtain the inverse:

$$(\mathcal{F}/n)^{-1} = \begin{pmatrix} \Omega & 0 \\ 0 & 2(D'_m D_m)^{-1} D'_m (\Omega \otimes \Omega) D_m (D'_m D_m)^{-1} \end{pmatrix}$$

and the determinant:

$$|\mathcal{F}/n| = |\Omega| \cdot |2(D'_m D_m)^{-1} D'_m (\Omega \otimes \Omega) D_m (D'_m D_m)^{-1}| = 2^m |\Omega|^{m+2}.$$

FURTHER READING

§2–3. Chapter 5 discusses differentials in more detail, and contains the first identification theorem (Theorem 5.6) and the chain rule for first differentials (Theorem 5.9), officially called ‘Cauchy’s rule of invariance’.

§4. Optimization is discussed in Chapter 7.

§5. See Chapter 11, Sections 11.29–11.32 and Chapter 13, Sections 13.4 and 13.19.

§6. The trace is discussed in Section 1.10, the extension from vector calculus to matrix calculus in Section 5.15, and the differentials of the determinant and inverse in Sections 8.3 and 8.4.

§7. See Section 1.6 for more detailed results.

§8–9. Second differentials are introduced in Chapter 6. The second identification theorem is proved in Section 6.8 and the chain rule for second differentials in Section 6.11.

§11. See Chapter 2 for many more details on the vec operator and the Kronecker product. Theorem 2.2 is restated here as Theorem 18.5.

§12. See Sections 5.15 and 10.2.

§13 and §15. The commutation matrix and the duplication matrix are discussed in Chapter 3.

§16. Many aspects of maximum likelihood estimation are treated in Chapter 15.

