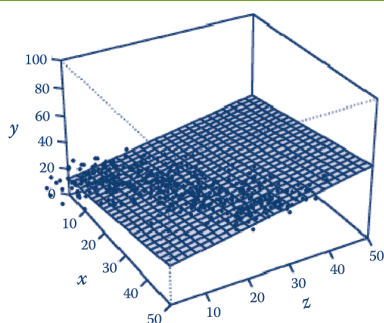
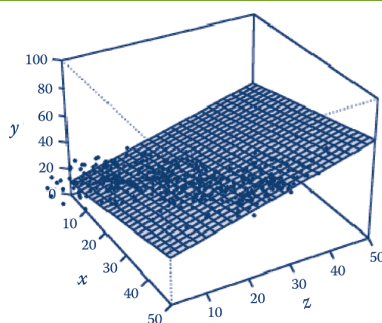


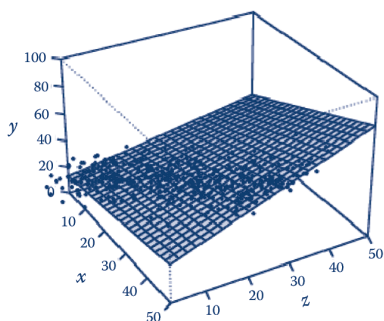
STATISTICAL THINKING IN EPIDEMIOLOGY



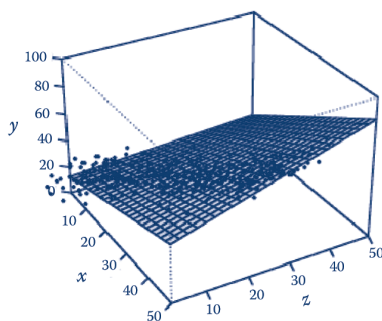
Model 1



Model 2



Model 3



Model 4

YU-KANG TU
MARK S. GILTHORPE



CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

STATISTICAL THINKING IN EPIDEMIOLOGY

STATISTICAL THINKING IN EPIDEMIOLOGY

YU-KANG TU
MARK S. GILTHORPE



CRC Press

Taylor & Francis Group

Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group an **informa** business
A CHAPMAN & HALL BOOK

CRC Press
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2012 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works
Version Date: 20110617

International Standard Book Number-13: 978-1-4200-9992-8 (eBook - PDF)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the CRC Press Web site at
<http://www.crcpress.com>

Contents

Preface.....	xi
1. Introduction.....	1
1.1 Uses of Statistics in Medicine and Epidemiology	1
1.2 Structure and Objectives of This Book	2
1.3 Nomenclature in This Book.....	5
1.4 Glossary	5
2. Vector Geometry of Linear Models for Epidemiologists	7
2.1 Introduction	7
2.2 Basic Concepts of Vector Geometry in Statistics	7
2.3 Correlation and Simple Regression in Vector Geometry	9
2.4 Linear Multiple Regression in Vector Geometry.....	11
2.5 Significance Testing of Correlation and Simple Regression in Vector Geometry.....	12
2.6 Significance Testing of Multiple Regression in Vector Geometry.....	14
2.7 Summary.....	15
3. Path Diagrams and Directed Acyclic Graphs	17
3.1 Introduction	17
3.1 Path Diagrams.....	17
3.1.1 The Path Diagram for Simple Linear Regression.....	18
3.1.1.1 Regression Weights, Path Coefficients and Factor Loadings	19
3.1.1.2 Exogenous and Endogenous Variables	19
3.1.2 The Path Diagram for Multiple Linear Regression.....	20
3.2 Directed Acyclic Graphs	21
3.2.1 Identification of Confounders	22
3.2.2 Backdoor Paths and Colliders	23
3.2.3 Example of a Complex DAG.....	24
3.3 Direct and Indirect Effects.....	25
3.4 Summary.....	26
4. Mathematical Coupling and Regression to the Mean in the Relation between Change and Initial Value.....	27
4.1 Introduction	27
4.2 Historical Background	29
4.3 Why Should Change Not Be Regressed on Initial Value? A Review of the Problem	29

- 4.4 Proposed Solutions in the Literature30
 - 4.4.1 Blomqvist’s Formula30
 - 4.4.2 Oldham’s Method: Testing Change and Average30
 - 4.4.3 Geometrical Presentation of Oldham’s Method32
 - 4.4.4 Variance Ratio Test32
 - 4.4.5 Structural Regression34
 - 4.4.6 Multilevel Modelling35
 - 4.4.7 Latent Growth Curve Modelling36
- 4.5 Comparison between Oldham’s Method and Blomqvist’s Formula37
- 4.6 Oldham’s Method and Blomqvist’s Formula Answer Two Different Questions38
- 4.7 What Is Galton’s Regression to the Mean?39
- 4.8 Testing the Correct Null Hypothesis40
 - 4.8.1 The Distribution of the Correlation Coefficient between Change and Initial Value41
 - 4.8.2 Null Hypothesis for the Baseline Effect on Treatment43
 - 4.8.3 Fisher’s Z-Transformation43
 - 4.8.4 A Numerical Example44
 - 4.8.5 Comparison with Alternative Methods44
- 4.9 Evaluation of the Categorisation Approach45
- 4.10 Testing the Relation between Changes and Initial Values When There Are More than Two Occasions47
- 4.11 Discussion48
- 5. Analysis of Change in Pre-/Post-Test Studies51**
 - 5.1 Introduction51
 - 5.2 Analysis of Change in Randomised Controlled Trials51
 - 5.3 Comparison of Six Methods53
 - 5.3.1 Univariate Methods54
 - 5.3.1.1 Test the Post-Treatment Scores Only54
 - 5.3.1.2 Test the Change Scores54
 - 5.3.1.3 Test the Percentage Change Scores54
 - 5.3.1.4 Analysis of Covariance54
 - 5.3.2 Multivariate Statistical Methods55
 - 5.3.2.1 Random Effects Model55
 - 5.3.2.2 Multivariate Analysis of Variance56
 - 5.3.3 Simulation Results59
 - 5.6 Analysis of Change in Non-Experimental Studies: Lord’s Paradox60
 - 5.6.1 Controversy around Lord’s Paradox62
 - 5.6.1.1 Imprecise Statement of Null Hypothesis62
 - 5.6.1.2 Causal Inference in Non-Experimental Study Design63

- 5.6.2 Variable Geometry of Lord’s Paradox64
 - 5.6.3 Illustration of the Difference between ANCOVA and Change Scores in RCTs Using Variable Space Geometry65
- 5.7 ANCOVA and *t*-Test for Change Scores Have Different Assumptions65
 - 5.7.1 Scenario One: Analysis of Change for Randomised Controlled Trials67
 - 5.7.2 Scenario Two: The Analysis of Change for Observational Studies69
- 5.8 Conclusion.....71
- 6. Collinearity and Multicollinearity73**
 - 6.1 Introduction: Problems of Collinearity in Linear Regression73
 - 6.2 Collinearity76
 - 6.3 Multicollinearity.....77
 - 6.4 Mathematical Coupling and Collinearity79
 - 6.5 Vector Geometry of Collinearity.....79
 - 6.5.1 Reversed Relation between the Outcome and Covariate due to Collinearity.....80
 - 6.5.2 Unstable Regression Models due to Collinearity81
 - 6.5.3 The Relation between the Outcome–Explanatory Variable’s Correlations and Collinearity82
 - 6.6 Geometrical Illustration of Principal Components Analysis as a Solution to Multicollinearity.....84
 - 6.7 Example: Mineral Loss in Patients Receiving Parenteral Nutrition.....85
 - 6.8 Solutions to Collinearity89
 - 6.8.1 Removal of Redundant Explanatory Variables89
 - 6.8.2 Centring.....89
 - 6.8.3 Principal Component Analysis.....90
 - 6.8.4 Ridge Regression.....93
 - 6.9 Conclusion.....94
- 7. Is ‘Reversal Paradox’ a Paradox?.....97**
 - 7.1 A Plethora of Paradoxes: The Reversal Paradox97
 - 7.2 Background: The Foetal Origins of Adult Disease Hypothesis (Barker’s Hypothesis)98
 - 7.2.1 Epidemiological Evidence on the Foetal Origins Hypothesis.....99
 - 7.2.2 Criticisms of the Foetal Origins Hypothesis100
 - 7.2.3 Reversal Paradox and Suppression in Epidemiological Studies on the Foetal Origins Hypothesis102
 - 7.2.4 Catch-Up Growth and the Foetal Origins Hypothesis.....104

7.2.5	Residual Current Body Weight: A Proposed Alternative Approach.....	107
7.2.6	Numerical Example	108
7.3	Vector Geometry of the Foetal Origins Hypothesis	109
7.4	Reversal Paradox and Adjustment for Current Body Size: Empirical Evidence from Meta-Analysis	111
7.5	Discussion	112
7.5.1	The Reversal Paradox and the Foetal Origins Hypothesis	112
7.5.2	Multiple Adjustments for Current Body Sizes	115
7.5.3	Catch-Up Growth and the Foetal Origins Hypothesis ...	116
7.6	Conclusion.....	117
8.	Testing Statistical Interaction	119
8.1	Introduction: Testing Interactions in Epidemiological Research	119
8.1	Testing Statistical Interaction between Categorical Variables....	121
8.2	Testing Statistical Interaction between Continuous Variables ...	124
8.3	Partial Regression Coefficient for Product Term in Regression Models	128
8.4	Categorization of Continuous Explanatory Variables	130
8.5	The Four-Model Principle in the Foetal Origins Hypothesis.....	131
8.6	Categorization of Continuous Covariates and Testing Interaction	132
8.6.1	Simulations	132
8.6.2	Numerical Example	133
8.7	Discussion	134
8.8	Conclusion.....	137
9.	Finding Growth Trajectories in Lifecourse Research	139
9.1	Introduction	139
9.1.1	Example: Catch-Up Growth and Impaired Glucose Tolerance	140
9.1.2	Galton and Regression to the Mean	141
9.1.3	Revisiting the Growth Trajectory of Men with Impaired Glucose Tolerance.....	144
9.2	Current Approaches to Identifying Postnatal Growth Trajectories in Lifecourse Research	146
9.2.1	The Lifecourse Plot.....	147
9.2.2	Regression with Changes Scores	150
9.2.3	Latent Growth Curve Models	151
9.2.4	Growth Mixture Models.....	159
9.3	Discussion	162

10. Partial Least Squares Regression for Lifecourse Research..... 165

10.1 Introduction 165

10.2 Data 166

10.3 OLS Regression..... 166

10.4 PLS Regression 167

10.4.1 History of PLS 167

10.4.2 PCA Regression..... 170

10.4.3 PLS Regression 171

10.4.4 PLS and Perfect Collinearity 172

10.4.5 Singular Value Decomposition, PCA and PLS
Regression 173

10.4.6 Selection of PLS Component 176

10.4.7 PLS Regression for Lifecourse Data Using Weight
z-Scores at Birth, Changes in z-Scores, and Current
Weight z-Scores 177

10.4.8 The Relationship between OLS Regression and PLS
Regression Coefficients 178

10.4.9 PLS Regression for Lifecourse Data Using Weight
z-Scores Measured at Six Different Ages..... 181

10.4.10 PLS Regression for Lifecourse Data Using Weight
z-Scores Measured at Six Different Ages and Five
Changes in z-Scores 182

10.5 Discussion 182

10.6 Conclusion..... 184

11. Concluding Remarks 187

References 189

Preface

Many books have been written on epidemiological methods. The main difference between them and this book is that we emphasise statistical thinking more than applications of specific statistical methods in epidemiological research, because we believe it is vital to appreciate context and for this to happen, one has to stop and reflect, rather than plough in. This book is therefore not a textbook for the statistical methods commonly used by biostatisticians and epidemiologists; instead, we assume readers have a basic understanding of generalised linear models, such as multiple regression and logistic regression. We do, however, discuss some basic methods in great detail, such as Pearson's correlation and the analysis of covariance, and we show that sometimes it requires a lot of careful thinking to use these simple methods correctly.

We use a few real examples, some of which remain controversial in epidemiological research, to demonstrate our statistical thinking. However, we by no means feel that our thinking is the only approach to the problems we highlight, which are chosen because we believe they have an appeal to general readers. Postgraduate students in biostatistics and epidemiology may use this book as a supplementary reading to standard texts on statistical or epidemiological methods. Lecturers of postgraduate courses may use this book (we hope) as a good example for teaching statistical thinking in epidemiological research. Experienced researchers may find our book both intellectually entertaining and challenging, as some issues discussed are still controversial. We try to show how our statistical thinking of specific research questions develops and eventually leads us to a set of solutions, but our thinking is inevitably framed by our training, knowledge, and experience. For instance, we believe vector geometry is a very useful tool for intuitive understanding of the basic concepts and nuances of linear models, but we acknowledge that not all our readers will agree with us as some may not find thinking geometrically at all intuitive or helpful. Therefore, we welcome feedback from our readers to broaden our vision and improve our thinking, hopefully giving rise to better solutions to those problems discussed in this book.

Our students, collaborators, and colleagues have been very helpful in sharpening and improving our thinking, and the following is an incomplete list of those we would wish to acknowledge and thank for their help in this regard: our colleagues in the Division of Biostatistics and the Division of Epidemiology at the University of Leeds; Professor George Ellison at the London Metropolitan University; Professor David Gunnell, Professor Jonathan Sterne and Dr. Kate Tilling at the University of Bristol; Professor Vibeke Baelum at Aarhus University, Denmark; Dr. Samuel Manda at the

Biostatistics Unit, Medical Research Council, South Africa; and Professor Kuo-Liong Chien at the National Taiwan University, Taiwan. Dr. Chris Metcalfe at the University of Bristol and Dr. Jim Lewsey at Glasgow University kindly acted as external reviewers for our book and gave many useful comments and suggestions. Nevertheless, we take full responsibility for any weakness and errors in our thinking in this book. We would also like to thank Sarah Morris and Rob Calver at Chapman & Hall for their patience with this project.

Part of this book was written when the first author (YKT) was on study leave in Taiwan supported by an international joint project grant from the Royal Society and the National Science Council in Taiwan. In the last 5 years, the first author has been supported by a UK Research Council Fellowship jointly housed by the School of Medicine (Division of Biostatistics) and the Leeds Dental Institute, having enjoyed a large degree of academic freedom in research supported by both the second author (head of the Division of Biostatistics) and Dr. Margaret Kellett (dean of the Leeds Dental Institute). Their generous support is greatly appreciated.

Finally, we would like to thank our wives, Jorin and Amy, and our daughters, Emma and Zarana, for their unconditional love and support.

Introduction

1.1 Uses of Statistics in Medicine and Epidemiology

Correlation and regression analyses are among the most commonly used statistical methods in biomedical research. Correlation tests the linear relationship between two (usually continuous) variables, and linear regression tests the relationship between one outcome variable (also known as the dependent variable) and one or more explanatory variables (also known as independent variables or covariates).

As powerful personal computers and graphical user-interface (GUI) statistical software packages have become available and readily affordable in the last decade, complex statistical methods such as multivariable regression become more and more frequently used in medicine and epidemiology, despite involving complex numerical calculations. For instance, to obtain the regression coefficients for a multiple linear regression model with five covariates, one needs to invert a 5×5 matrix, which is very complex and might take days to do by hand. With the power of computers, it now only takes a few seconds to perform these calculations, and all the information relevant (or irrelevant) to the research questions can be shown on the computer screen.

Unfortunately, whilst biomedical researchers may be able to follow instructions in the manuals accompanying the statistical software packages, they do not always have sufficient knowledge to choose the appropriate statistical methods and correctly interpret their results. Many biostatisticians have noticed that the misuses of statistical methods in clinical research are common and the quality of statistical reports in clinical journals needs to be improved (Altman 1991b, 1998, 2002; Andersen 1990). Besides, the increasing use of more advanced and complex statistical methods does not necessarily provide greater insight into the research questions and thereby obtain more reliable knowledge. On the contrary, careless use of these methods can sometimes generate confusing or even misleading results.

This book examines several common methodological and statistical problems in the use of correlation and regression in medical and epidemiological research: mathematical coupling, regression to the mean, collinearity, reversal paradox and statistical interaction. These problems are usually intertwined

with each other. For instance, in the analysis of the relation between change and initial value, mathematical coupling and regression to the mean are almost synonymous; mathematical coupling between explanatory variables usually gives rise to collinearity; and product interaction terms and their component variables are mathematically coupled and raise concerns over collinearity. Since the discussions of these five problems pervade all areas of medical and epidemiological research, it is necessary to study these problems in a framework that focuses on specific research themes. Therefore, this book aims to examine the problems of mathematical coupling and regression to the mean in the analysis of change in pre-test/post-test study designs. The evidence on the foetal origins of adult diseases hypothesis, which sometimes suffers problems of the reversal paradox caused by the adjustment of current body size as confounders, is used to illustrate the potential problems in selecting variables for statistical adjustment and testing statistical interaction.

The overall aim of this book is to explore these statistical problems, to critically evaluate the existing proposed solutions to these problems, and to develop new tools to overcome them. One specific feature of this book is that, wherever applicable, vector geometry is used as a mathematical tool, in novel ways, to illustrate and further develop the concepts behind these statistical methods. The advantage of a geometrical approach over an analytical (algebraic) approach is that geometry can be more readily visualised and may be more intuitive to non-statisticians. Some common uses and misuses of correlation and regression analyses can be more intuitively understood by using geometrical illustrations than using traditional algebraic formula. Moreover, geometry will not only enhance our understanding of statistical problems, but will also generate new insights into old problems, thereby paving the way to provide guidance as to how to avoid statistical errors or even to suggest solutions.

1.2 Structure and Objectives of This Book

This book begins in Chapter 2 with a concise introduction to the concepts underlying the uses of vector geometry in linear models. Although vector geometry has been shown to be a great tool for understanding statistical methods based on ordinary least squares techniques, discussions of vector geometry are, however, scattered in the statistical literature only, and none appear in the clinical literature. So the objective of Chapter 2 is to bring together the various elementary results in vector geometry and establish a set of basic geometric tools that will be utilised throughout this book.

In Chapter 3, we introduce the concepts of directed acyclic graphs (DAGs). DAGs have become popular in epidemiological research in recent years as it

provides a very useful tool for identifying potential confounders. It is very similar to path diagram for structural equation modelling, which is widely used in social sciences research for testing causal models.

Chapter 4 reviews a controversy in analysing the relation between change and initial value caused by mathematical coupling and regression to the mean. In the analyses of pre-test/post-test study design, it is difficult to distinguish between regression to the mean and mathematical coupling, and this difficulty has given rise to confusion amongst some statisticians surrounding how to solve the problem. The objectives of Chapter 4 are to demonstrate the inadequacy of current practice in clinical research, to correct a misconception amongst some statisticians and to review various statistical methods proposed to test the relation between change and initial value.

Chapter 5 explores the differences in statistical power and effect estimations between several statistical methods to analyse changes in the pre-test/post-test study design for both randomised and non-randomised studies. The objectives of Chapter 5 are first to illustrate, using vector geometry, the controversy about the adjustment of initial values, known as *Lord's paradox*, and to explain why and when this paradox arises. This is particularly pertinent to epidemiology since Lord's paradox arises where proper randomisation is not always feasible—as within most epidemiological research. Vector geometry and simulations are then used to show that the prevalent concept that analysis of covariance (ANCOVA) *always* has greater statistical power than other univariate methods is not strictly correct; only when sample sizes are large and/or the correlation between pre-test and post-test values are moderate or high does ANCOVA have greater statistical power than other univariate and multivariate methods.

The objective of Chapter 6 is to demonstrate that vector geometry can give new insights into the problems of collinearity and the detection of these problems in epidemiological research. Collinearity is especially important in observational studies, where the underlying premise of statistical adjustment in regression models is to control simultaneously for correlations amongst covariates and their correlations with the outcome. Vector geometry can elegantly show why principal component analysis, a commonly recommended solution to collinearity, is not always desirable.

Chapter 7 discusses the problem of the reversal paradox, which is perhaps better known in epidemiology as *Simpson's paradox* within categorical variable analysis. In regression analysis, or analysis of covariance, this phenomenon is known as Lord's paradox, and in general linear modelling, a widely discussed statistical paradox known as *suppression* or *enhancement* is another manifestation of the reversal paradox. The objectives of Chapter 7 are to use the foetal origins of adult diseases hypothesis as an example to show how and why the adjustment of current body size measures, such as current body

weight, can give rise to the reversal paradox in multiple linear regression. Vector geometry is used to illustrate the reversal paradox and question an alternative interpretation often proffered, namely, that catch-up growth has a greater impact than birth weight on health in later life. Computer simulations and evidence from empirical studies are then used to show that the negative relationships between birth weight and blood pressure are strengthened and positive relationships between birth weight and blood pressure are attenuated or even reversed after adjusting for one or more current body sizes.

Chapter 8 examines the role of statistical interaction within regression analyses by revisiting the four-model principle proposed in the foetal origins hypothesis literature. The objectives of this chapter are first to show that when three continuous variables, such as blood pressure, birth weight and current body weight, follow multivariate normality, the expected value of the partial regression coefficient for the product interaction term in the multiple regression model is zero. Computer simulations are used to show that categorising birth weight and/or current body weight, a common practice in epidemiological studies on the foetal origins hypothesis (and many epidemiological studies in general), gives rise to spurious interaction effects and can therefore potentially lead to seriously misleading conclusions.

Chapter 9 reviews the recent advances in identifying critical growth phases for health outcomes in later life. In recent years, some researchers' interests have shifted from birth size to the growth in body size during early childhood. In Chapter 9 we provide a concise introduction to those methods for identifying the critical growth phases and use a publicly available data set to illustrate and compare results. This chapter shows that whilst different statistical methods have been proposed to test the same hypothesis, the differences in those approaches seem to indicate that researchers may have different versions of the same hypothesis in mind. As statistical thinking is guided by the research hypothesis, differences in theorising and framing the research questions will inevitably lead to different thinkings.

After discussing the pros and cons of the proposed methods for identifying the critical growth phases within Chapter 9, we propose a new approach in Chapter 10 that uses partial least squares (PLS). PLS is widely used in chemometrics and bioinformatics, though it is rarely used in medical and epidemiological research. One apparent advantage of PLS is that it can deal with perfect collinearity that can arise due to the mathematical relationships amongst covariates. In Chapter 10 we use the same data set as in Chapter 9 to illustrate how PLS can be applied to lifecourse research.

Most of the data sets used in this book have been drawn from published literature that are easily obtainable. One data set was kindly made available by Dr. Lars Laurell.

1.3 Nomenclature in This Book

Italic capital letters (X) denote variables without measurement error.

Italic lower-case letters (x) denote observed variables measured with error.

Bold lower-case letters (\mathbf{x}) denote vectors.

Bold capital letters (\mathbf{X}) denote matrices.

Capital letters (X) with subscripts denote components of variables or vectors.

Subscripts of variables usually represent repeated observations on the same variables, though they may also represent levels within a multilevel framework.

1.4 Glossary

Due to our backgrounds and a number of extensive collaborations with dental researchers, we sometimes use examples from the dental literature to illustrate methodological problems. Consequently, non-dental readers might not be familiar with some of the terminology that frequently appears in some chapters. This section provides a non-technical explanation and illustration (Figure 1.1a and b) of the dental terminology.

Clinical attachment level (CAL): The distance between the cemento-enamel junction and the bottom of periodontal pocket. It is not common to measure CAL directly using a periodontal probe; instead, CAL is usually the sum of gingival recession (GR) and probing pocket depth (PPD).

Enamel matrix proteins derivatives (EMD): Extracts from developing teeth in young pigs; these mixtures of proteins are believed to be able to promote regeneration of lost periodontal tissue. A syringe is used to release EMD into periodontal infrabony defects.

Gingival recession (GR): The distance between the cemento-enamel junctions and the gingival margin. In a healthy condition, the level of gingival margin (i.e., gums) is assumed to be level with the cemento-enamel junction and, therefore, GR is zero. However, due to external trauma, periodontal diseases or aging, the gingival margins may recede; that is, move in the direction of the root apex, and GR

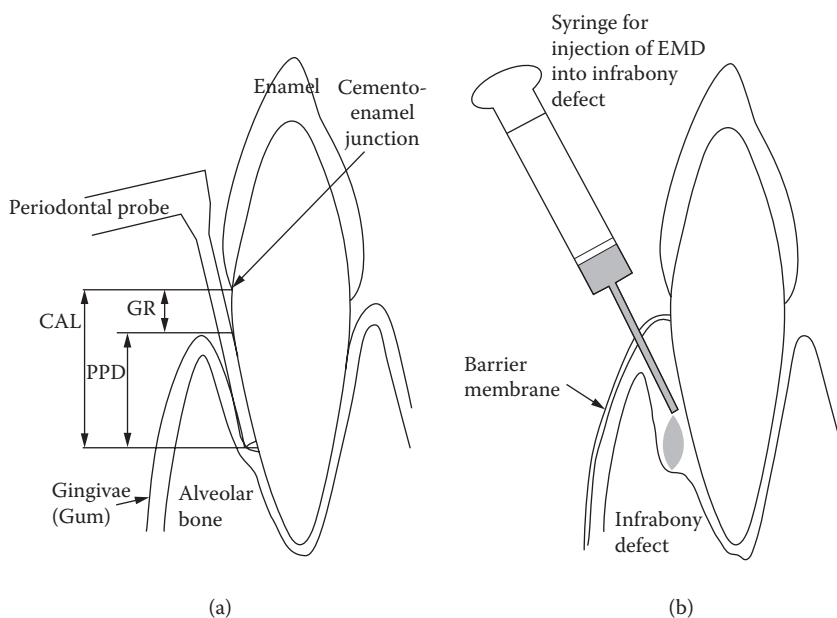
**FIGURE 1.1**

Illustration of the glossary. (a) Measurement of periodontal pocket depths; (b) application EMD in periodontal surgery.

is positive. In some circumstances, the gingivae (i.e., gums) might cover part of the crown of the tooth and the cemento-enamel junction becomes invisible, in which case GR is negative.

Guided tissue regeneration (GTR): A surgical procedure to regenerate lost periodontal tissue, especially in periodontal infrabony defects and furcation involvement in molars. A barrier membrane is placed around the tooth to cover the defect in order to allow progenitor cells from alveolar bone and/or periodontal membrane to repopulate on the root surface.

Probing (pocket) depth (PD): The distance between the gingival margin and the bottom of the periodontal pockets. Both PD and GR are usually measured with a periodontal probe with markings in millimetres. In this book, probing pocket depth, probing depth or pocket depth are all equivalent and used interchangeably.

2

Vector Geometry of Linear Models for Epidemiologists

2.1 Introduction

Vector geometry was first used in 1915 (Fisher 1915) by the great statistician Sir Ronald Fisher to address the problem of deriving the statistical distribution of the Pearson product-moment correlation coefficient. Fisher used his great geometric imagination to demonstrate that the correlation between two variables is the cosine function of the angle between two vectors in an n -dimensional space. However, as most statisticians do not have the same great geometrical insights as Fisher did, and most statisticians consider an algebraic approach more mathematically rigorous, the geometric approach has not thus far received as much attention as it perhaps deserves (Herr 1980). Nevertheless, some teachers of statistics have found a geometric approach to be more intuitive for people without a mathematical background and a very useful tool to develop the understanding of linear regression analyses. A few statistical textbooks (Saville and Wood 1991, 1996; Wickens 1995; Carroll et al. 1997) are written using mainly vector geometry, and some textbooks of statistics or econometrics on linear statistical models (Wonnacott and Wonnacott 1979, 1981; Fox 1997; Draper and Smith 1998) have chapters on vector geometry. However, vector geometry seems to have been used rarely to illustrate or explore specific methodological problems within biomedical research.

The aim of this chapter is to provide an introduction to the various concepts of vector geometry within correlation and regression analyses. Vector geometry, as a mathematical tool, will be consistently and extensively explored in later chapters.

2.2 Basic Concepts of Vector Geometry in Statistics

The scatter plot, one of the most commonly used graphs to display the relationship between two continuous variables, can be viewed as a form of

geometry. In the scatter plot for two random variables, X and Y , each with n independent observations ($X_1 \dots X_n$) and ($Y_1 \dots Y_n$), there will be n points in *two-dimensional space* (i.e., on a plane). This is known as *variable space* geometry. The axes represent the variables X and Y , and the points are the observations made on each subject. In contrast, instead of using variables as axes, the same data can be displayed in what is termed *subject space* by using subjects as the axes; the variables X and Y are then two points in n -dimensional space. By drawing an arrow from the origin to each point, X and Y become two vectors with coordinates ($X_1 \dots X_n$) and ($Y_1 \dots Y_n$) in n -dimensional space (Wickens 1995; Fox 1997). Although it is impossible to visualise n -dimensional space, only a plane is required to visualise the relative relationship between the vectors representing X and Y in n -dimensional space by projecting onto a plane and effectively dropping the original axes. In general, the number of dimensions needed to draw the graph in subject space is no greater than the number of variables (Wickens 1995).

We use a numerical example to illustrate the difference in data presentation between variable-space (scatter plot) geometry and subject-space (vector) geometry. Suppose two men, Mr A and Mr B, have their body size measured: the heights of A and B are 180 cm and 170 cm, respectively, and their weights are 90 kg and 60 kg, respectively. When the data are presented in a scatter plot (i.e., variable-space), subjects A and B are two points in a two-dimensional plot, and the axes are *Height* and *Weight* (Figure 2.1a). When the same data are presented in vector geometry (subject space), *Height* and *Weight* become two vectors in a two-dimensional space, and subjects A and B are the axes (Figure 2.1b).

In Figure 2.1b, the two variables *Height* and *Weight* are presented in their raw data format. However, in general, it is very useful to represent the

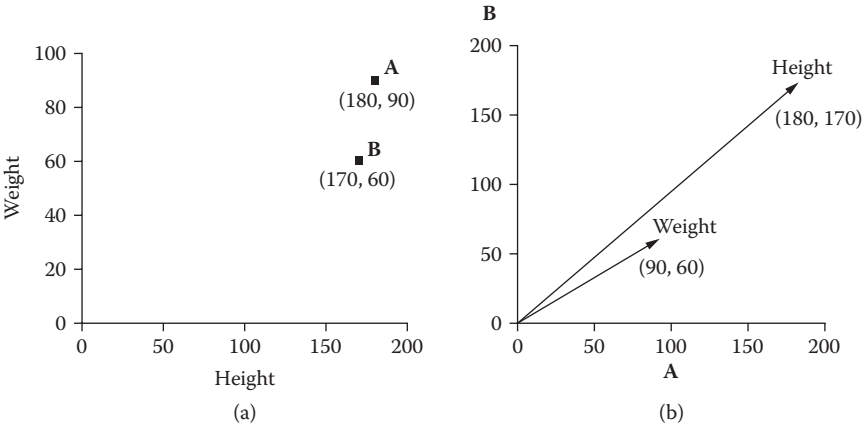


FIGURE 2.1
Illustrations of variable-space (a) and subject-space (b) geometry.