

Contents

Preface	xv
1 Introductory Probability Concepts	1
1.1 Definitions	1
1.2 Kolmogorov axioms	2
1.3 Counting methods	5
1.3.1 Permutations and combinations	6
1.3.2 Partitioning sets	7
1.3.3 Generating functions	9
1.4 Conditional probability	9
1.4.1 Law of total probability	11
1.4.2 Bayes rule	11
1.5 Independent events	12
1.5.1 Repeated trials and mutual independence	13
1.6 Laboratory problems	13
1.6.1 Laboratory: Introductory concepts	13
1.6.2 Additional problem notebooks	13
2 Discrete Probability Distributions	15
2.1 Definitions	15
2.1.1 PDF and CDF for discrete distributions	16
2.2 Univariate distributions	17
2.2.1 Example: Discrete uniform distribution	17
2.2.2 Example: Hypergeometric distribution	17
2.2.3 Distributions related to Bernoulli experiments	18
2.2.4 Simple random samples	20
2.2.5 Example: Poisson distribution	21
2.3 Joint distributions	22
2.3.1 Bivariate distributions; marginal distributions	22
2.3.2 Conditional distributions; independence	23
2.3.3 Example: Bivariate hypergeometric distribution	24
2.3.4 Example: Trinomial distribution	25
2.3.5 Survey analysis	25
2.3.6 Discrete multivariate distributions	26
2.3.7 Probability generating functions	26

2.4	Laboratory problems	27
2.4.1	Laboratory: Discrete models	27
2.4.2	Additional problem notebooks	27
3	Continuous Probability Distributions	29
3.1	Definitions	29
3.1.1	PDF and CDF for continuous random variables	29
3.1.2	Quantiles; percentiles	31
3.2	Univariate distributions	31
3.2.1	Example: Uniform distribution	31
3.2.2	Example: Exponential distribution	32
3.2.3	Euler gamma function	33
3.2.4	Example: Gamma distribution	33
3.2.5	Distributions related to Poisson processes	34
3.2.6	Example: Cauchy distribution	35
3.2.7	Example: Normal or Gaussian distribution	35
3.2.8	Example: Laplace distribution	36
3.2.9	Transforming continuous random variables	36
3.3	Joint distributions	38
3.3.1	Bivariate distributions; marginal distributions	38
3.3.2	Conditional distributions; independence	40
3.3.3	Example: Bivariate uniform distribution	41
3.3.4	Example: Bivariate normal distribution	41
3.3.5	Transforming continuous random variables	42
3.3.6	Continuous multivariate distributions	43
3.4	Laboratory problems	44
3.4.1	Laboratory: Continuous models	44
3.4.2	Additional problem notebooks	44
4	Mathematical Expectation	45
4.1	Definitions and properties	45
4.1.1	Discrete distributions	45
4.1.2	Continuous distributions	46
4.1.3	Properties	47
4.2	Mean, variance, standard deviation	48
4.2.1	Properties	48
4.2.2	Chebyshev inequality	49
4.2.3	Markov inequality	50
4.3	Functions of two or more random variables	50
4.3.1	Properties	51
4.3.2	Covariance, correlation	51
4.3.3	Sample summaries	54
4.3.4	Conditional expectation; regression	55
4.4	Linear functions of random variables	56
4.4.1	Independent normal random variables	57

Contents	ix
4.5 Laboratory problems	58
4.5.1 Laboratory: Mathematical expectation	58
4.5.2 Additional problem notebooks	58
5 Limit Theorems	59
5.1 Definitions	59
5.2 Law of large numbers	60
5.2.1 Example: Monte Carlo evaluation of integrals	60
5.3 Central limit theorem	61
5.3.1 Continuity correction	62
5.3.2 Special cases	62
5.4 Moment generating functions	63
5.4.1 Method of moment generating functions	65
5.4.2 Relationship to the central limit theorem	65
5.5 Laboratory problems	66
5.5.1 Laboratory: Sums and averages	66
5.5.2 Additional problem notebooks	66
6 Transition to Statistics	69
6.1 Distributions related to the normal distribution	69
6.1.1 Chi-square distribution	69
6.1.2 Student t distribution	70
6.1.3 F ratio distribution	71
6.2 Random samples from normal distributions	71
6.2.1 Sample mean, sample variance	72
6.2.2 Approximate standardization of the sample mean	73
6.2.3 Ratio of sample variances	74
6.3 Multinomial experiments	75
6.3.1 Multinomial distribution	75
6.3.2 Goodness-of-fit: Known model	75
6.3.3 Goodness-of-fit: Estimated model	77
6.4 Laboratory problems	79
6.4.1 Laboratory: Transition to statistics	79
6.4.2 Additional problem notebooks	79
7 Estimation Theory	81
7.1 Definitions	81
7.2 Properties of point estimators	82
7.2.1 Bias; unbiased estimator	82
7.2.2 Efficiency for unbiased estimators	82
7.2.3 Mean squared error	83
7.2.4 Consistency	83
7.3 Interval estimation	84
7.3.1 Example: Normal distribution	84
7.3.2 Approximate intervals for means	86
7.4 Method of moments estimation	86
7.4.1 Single parameter estimation	86

7.4.2	Multiple parameter estimation	87
7.5	Maximum likelihood estimation	87
7.5.1	Single parameter estimation	87
7.5.2	Cramer–Rao lower bound	89
7.5.3	Approximate sampling distribution	90
7.5.4	Multiple parameter estimation	93
7.6	Laboratory problems	94
7.6.1	Laboratory: Estimation theory	94
7.6.2	Additional problem notebooks	94
8	Hypothesis Testing Theory	97
8.1	Definitions	97
8.1.1	Neyman–Pearson framework	97
8.1.2	Equivalent tests	99
8.2	Properties of tests	100
8.2.1	Errors, size, significance level	100
8.2.2	Power, power function	101
8.3	Example: Normal distribution	103
8.3.1	Tests of $\mu = \mu_o$	103
8.3.2	Tests of $\sigma^2 = \sigma_o^2$	104
8.4	Example: Bernoulli/binomial distribution	105
8.5	Example: Poisson distribution	106
8.6	Approximate tests of $\mu = \mu_o$	107
8.7	Likelihood ratio tests	107
8.7.1	Likelihood ratio statistic; Neyman–Pearson lemma	107
8.7.2	Generalized likelihood ratio tests	109
8.7.3	Approximate sampling distribution	111
8.8	Relationship with confidence intervals	114
8.9	Laboratory problems	115
8.9.1	Laboratory: Hypothesis testing	115
8.9.2	Additional problem notebooks	115
9	Order Statistics and Quantiles	117
9.1	Order statistics	117
9.1.1	Approximate mean and variance	120
9.2	Confidence intervals for quantiles	121
9.2.1	Approximate distribution of the sample median	121
9.2.2	Exact confidence interval procedure	122
9.3	Sample quantiles	123
9.3.1	Sample quartiles, sample IQR	123
9.3.2	Box plots	124
9.4	Laboratory problems	125
9.4.1	Laboratory: Order statistics and quantiles	125
9.4.2	Additional problem notebooks	125
10	Two Sample Analysis	127
10.1	Normal distributions: Difference in means	127

Contents	xi
10.1.1	Known variances 128
10.1.2	Pooled t methods 129
10.1.3	Welch t methods 130
10.2	Normal distributions: Ratio of variances 131
10.3	Large sample: Difference in means 134
10.4	Rank sum test 135
10.4.1	Rank sum statistic 136
10.4.2	Tied observations; midranks 138
10.4.3	Mann–Whitney U statistic 139
10.4.4	Shift models 140
10.5	Sampling models 143
10.5.1	Population model 144
10.5.2	Randomization model 144
10.6	Laboratory problems 145
10.6.1	Laboratory: Two sample analysis 145
10.6.2	Additional problem notebooks 145
11	Permutation Analysis 147
11.1	Introduction 147
11.1.1	Permutation tests 148
11.1.2	Example: Difference in means test 149
11.1.3	Example: Smirnov two sample test 151
11.2	Paired sample analysis 152
11.2.1	Example: Signed rank test 153
11.2.2	Shift models 156
11.2.3	Example: Fisher symmetry test 157
11.3	Correlation analysis 159
11.3.1	Example: Correlation test 159
11.3.2	Example: Rank correlation test 161
11.4	Additional tests and extensions 162
11.4.1	Example: One sample trend test 162
11.4.2	Example: Two sample scale test 164
11.4.3	Stratified analyses 165
11.5	Laboratory problems 166
11.5.1	Laboratory: Permutation analysis 166
11.5.2	Additional problem notebooks 167
12	Bootstrap Analysis 169
12.1	Introduction 169
12.1.1	Approximate conditional estimation 171
12.2	Bootstrap estimation 172
12.2.1	Error distribution 173
12.2.2	Simple approximate confidence interval procedures . . . 173
12.2.3	Improved intervals: Nonparametric case 175
12.3	Applications of bootstrap estimation 176
12.3.1	Single random sample 176

12.3.2	Independent random samples	178
12.4	Bootstrap hypothesis testing	179
12.5	Laboratory problems	181
12.5.1	Laboratory: Bootstrap analysis	181
12.5.2	Additional problem notebooks	181
13	Multiple Sample Analysis	183
13.1	One-way layout	183
13.1.1	Example: Analysis of variance	183
13.1.2	Example: Kruskal–Wallis test	187
13.1.3	Example: Permutation f test	189
13.2	Blocked design	190
13.2.1	Example: Analysis of variance	190
13.2.2	Example: Friedman test	194
13.3	Balanced two-way layout	196
13.3.1	Example: Analysis of variance	196
13.3.2	Example: Permutation f tests	202
13.4	Laboratory problems	202
13.4.1	Laboratory: Multiple sample analysis	203
13.4.2	Additional problem notebooks	203
14	Linear Least Squares Analysis	205
14.1	Simple linear model	205
14.1.1	Least squares estimation	206
14.1.2	Permutation confidence interval for slope	208
14.2	Simple linear regression	208
14.2.1	Confidence interval procedures	209
14.2.2	Predicted responses and residuals	211
14.2.3	Goodness-of-fit	212
14.3	Multiple linear regression	214
14.3.1	Least squares estimation	214
14.3.2	Analysis of variance	219
14.3.3	Confidence interval procedures	220
14.3.4	Regression diagnostics	221
14.4	Bootstrap methods	223
14.5	Laboratory problems	225
14.5.1	Laboratory: Linear least squares analysis	225
14.5.2	Additional problem notebooks	225
15	Contingency Table Analysis	227
15.1	Independence analysis	227
15.1.1	Example: Pearson’s chi-square test	227
15.1.2	Example: Rank correlation test	229
15.2	Homogeneity analysis	230
15.2.1	Example: Pearson’s chi-square test	230
15.2.2	Example: Kruskal–Wallis test	232
15.3	Permutation chi-square tests	233

Contents	xiii
<hr/>	
15.4 Fourfold tables	235
15.4.1 Odds ratio analysis	235
15.4.2 Small sample analyses	238
15.5 Laboratory problems	239
15.5.1 Laboratory: Contingency table analysis	240
15.5.2 Additional problem notebooks	240
Bibliography	241
Index	251

Copyright ©2005 by the Society for Industrial and Applied Mathematics

This electronic version is for personal use and may not be duplicated or distributed.

From "Mathematica Laboratories for Mathematical Statistics: Emphasizing Simulation and Computer Intensive Methods" by
Jenny Baglivo.

Buy this book from SIAM at www.ec-securehost.com/SIAM/SA14.html

Preface

There is no doubt that the computer has revolutionized the practice of statistics in recent years. Computers allow us to analyze data more quickly using classical techniques, to analyze much larger data sets, to replace classical data analytic methods—whose assumptions may not be met—with more flexible computer intensive approaches, and to solve problems with no satisfactory classical solution.

Nor is there doubt that undergraduate mathematics and statistics courses could benefit from the integration of computer technology. Computer laboratories can be used to illustrate and reinforce important concepts; allow students to simulate experiments and visualize their results; and allow them to compare the results of classical methods of data analysis with those using alternative techniques. The problem is how best to introduce these techniques in the curriculum.

This book introduces an approach to incorporating technology in the mathematical statistics sequence, with an emphasis on simulation and computer intensive methods. The printed book is a concise introduction to the concepts of probability theory and mathematical statistics. The accompanying electronic materials are a series of in-class and take-home computer laboratory problems designed to reinforce the concepts and to apply the techniques in real and realistic settings.

The laboratory materials are written as *Mathematica* Version 5 notebooks [112] and are designed so that students with little or no experience in *Mathematica* will be able to complete the work. *Mathematica* notebooks contain text, data, computations, and graphics; they are particularly well suited for presenting concepts and problems and for writing solutions.

Laboratory problems, custom tools designed to enhance the capabilities of *Mathematica*, an introduction to using *Mathematica* for probability and statistics, and additional materials are included in an accompanying CD. An instructor's CD is available to those who adopt the book. The instructor's CD contains complete solutions to all laboratory problems, instructor guides, and hints on developing additional tools and laboratory problems.

The materials are written to be used in the mathematical statistics sequence given at most colleges and universities (two courses of four semester hours each or three courses of three semester hours each). Multivariable calculus and familiarity with the basics of set theory, vectors and matrices, and problem solving using a computer are assumed. The order of topics generally follows that of a standard sequence. Chapters 1 through 5 cover concepts in probability. Chapters 6 through 10 cover introductory mathematical statistics. Chapters 11 and 12 are on permutation

and bootstrap methods. In each case, problems are designed to expand on ideas from previous chapters so that instructors could choose to use some of the problems earlier in the course. Permutation and bootstrap methods also appear in the later chapters. Chapters 13, 14, and 15 are on multiple sample analysis, linear least squares, and analysis of contingency tables, respectively. References for specialized topics in Chapters 10 through 15 are given at the beginning of each chapter.

The materials can also be used profitably by statistical practitioners or consultants interested in a computer-based introduction to mathematical statistics, especially to computer intensive methods.

Laboratory problems

Each chapter has a main laboratory notebook, containing between five and seven problems, and a series of additional problem notebooks. The problems in the main laboratory notebook are for basic understanding and can be used for in-class work or assigned for homework. The additional problem notebooks reinforce and/or expand the ideas from the main laboratory notebook and are generally longer and more involved.

There are a total of 238 laboratory problems. Each main laboratory notebook and many of the problem notebooks contain examples for students to work before starting the assigned problems. One hundred twenty-three examples and problems use simulation, permutation, and bootstrap methods. One hundred twenty-five problems use real data.

Many problems are based on recent research reports or ongoing research—for example, analyses of the spread of an infectious disease in the cultured oyster population in the northeastern United States [18], [42], [100]; analyses of the ecological effects of the introduction of the Asian shore crab to the eastern United States [19], [20]; comparison of modeling strategies for occurrences of earthquakes in southern California [35]; comparison of spatial distributions of earthquakes [60] and of animal species [105]; comparison of treatments for multiple sclerosis [63], [8]; and analyses of associations between cellular telephone use and car accidents [88], between genetics and longevity [114], and between incidence of childhood leukemia and distance to a hazardous waste site [111]. Whimsical examples include comparisons of world-class sprinters [108] and of winning baseball players and teams [98].

Note to the student

Concepts from probability and statistics are used routinely in fields as diverse as actuarial science, ecology, economics, engineering, genetics, health sciences, marketing, and quality management. The ideas discussed in each chapter of the text will give you a basic understanding of the important concepts. The last section in each chapter outlines the laboratory problems.

Although formal proofs are not emphasized, the logical progression of the ideas in a proof is given whenever possible. Comments, including reminders about topics from calculus and pointers to where concepts will be applied, are enclosed in boxes throughout the text.

The accompanying CD contains two folders:

1. The `PDFFiles` folder contains documents in Acrobat PDF format. You will need a current copy of Adobe Acrobat Reader to open and print these files. Adobe Acrobat Reader is available for free from `adobe.com`.
2. The `MMAFiles` folder contains *Mathematica* files. You will need a copy of *Mathematica* Version 5 to work with these files.

The `PDFFiles` folder includes two appendices to the printed text and 15 laboratory workbooks. Appendix A is an introduction to the *Mathematica* commands used in the laboratory problems. Print Appendix A and keep it for reference. Appendix B contains tables of probabilities and quantiles suitable for solving problems when you are not using the computer. Print Appendix B and keep it for reference. There is one laboratory workbook for each chapter of the text. Print the ones you need for your course.

The `MMAFiles` folder includes 15 folders of laboratory problems and a folder of customized tools (`StatTools`). The `StatTools` folder should be placed in the user base directory or other appropriate directory on your system. Consult the online help within the *Mathematica* system for details, or speak to your instructor.

Note to the instructor

The material in the text is sufficient to support a problem-oriented mathematical statistics sequence, where the computer is used throughout the sequence. In fact, the first lab can be scheduled after three or four class meetings. Students are introduced to parametric, nonparametric, permutation, and bootstrap methods and will learn about data analysis, including diagnostic methods. (See the chapter outlines below.)

The text does not include exercises intended to be done by hand. You will need to supplement the text with by-hand exercises from other books or with ones that you design yourself. Suggestions for by-hand exercises that complement certain laboratory problems are given in the instructor's CD.

In addition, the printed text does not include *Mathematica* commands. Step-by-step instructions for using *Mathematica* commands are given in examples in the electronic materials. Online help is available, and Appendix A on the CD can be used as a reference.

Chapter outlines

Chapter 1 covers counting methods, axioms of probability, conditional probability, and independence. The first laboratory session is intended to be scheduled early in the term, as soon as the counting methods, axioms, and first examples are discussed. Students become familiar with using *Mathematica* commands to compute and graph binomial coefficients and hypergeometric probabilities (called "urn probabilities" in the lab) and get an informal introduction to maximum likelihood and likelihood ratio methods using custom tools. The additional problem notebooks reinforce these ideas

and include problems on frequency generating functions, conditional probability, and independence.

Chapters 2 and 3 are on discrete and continuous families of probability distributions, respectively. In the laboratory sessions, students become familiar with using *Mathematica* commands for computing probabilities and pseudorandom samples from univariate distributions, and with using custom tools for graphing models and samples. The additional problem notebooks reinforce these ideas, give students an informal introduction to goodness-of-fit, and include problems on probability generating functions, bivariate distributions, and transformations.

Chapter 4 is on mathematical expectation. In the laboratory and additional problem notebooks, students work with *Mathematica* commands for model and sample summaries, use sample summaries to estimate unknown parameters, apply the Chebyshev and Markov inequalities, and work with conditional expectations.

Chapter 5 is on limit theorems. In the laboratory session, students use custom tools to study sequences of running sums and averages, and answer a variety of questions on exact and approximate distributions of sums. The additional problem notebooks reinforce and expand on these ideas, and include several problems on probability and moment generating functions.

Chapter 6 serves as a transition from probability to statistics. The chi-square, Student *t*, and *f* ratio distributions are defined, and several applications are introduced, including the relationship of the chi-square distribution to the sampling distribution of the sample variance of a random sample from a normal distribution and the application of the chi-square distribution to the multinomial goodness-of-fit problem. In the laboratory session, students become familiar with chi-square and multinomial distributions, and use a custom tool for carrying out a goodness-of-fit analysis using Pearson's test (including analysis of standardized residuals). The additional problem notebooks contain simulation studies and applications of Pearson's goodness-of-fit test, and introduce students to minimum chi-square and method of moments estimates. The chapter is intended to precede formal statistical inference.

Chapters 7 and 8 are on estimation theory and hypothesis testing theory, respectively. In the first laboratory session, students become familiar with *Mathematica* commands for constructing confidence intervals for normal means and variances, and use custom tools to study the concepts of confidence interval and maximum likelihood estimation. In the second laboratory session, students become familiar with *Mathematica* commands for carrying out tests for normal means and variances, construct power curves, use a custom tool to construct tests and compute power at fixed alternatives, and compute sample sizes. The additional problem notebooks reinforce and expand on these ideas, contain simulation studies, introduce the idea of inverting tests to produce confidence intervals, and include applications of the likelihood ratio goodness-of-fit test.

Chapter 9 is on order statistics and quantiles. In the laboratory session, students apply custom tools for visualizing order-statistic distributions, for quantile estimation, and for constructing box plots in a variety of problems. The additional problem notebooks reinforce and expand on these ideas, introduce probability plots, study order statistics for uniform models, and contain simulation studies.

Chapter 10 is on parametric and nonparametric two sample analysis. In the laboratory session, students apply *Mathematica* commands for analyzing independent random samples from normal distributions and custom tools for the Wilcoxon rank sum test in a variety of problems. Normal probability plots of standardized observations are used to determine whether parametric methods should be used. The additional problem notebooks reinforce and expand on these ideas, contain simulation studies, introduce custom tools for quantile-quantile plots and inverting the Wilcoxon rank sum test under the shift model, and consider the randomization model for two sample analysis.

Chapter 11 is an introduction to permutation analysis, using nonparametric analyses of two samples and paired samples as first examples. In the laboratory session, students apply the rank sum, Smirnov, correlation, and signed rank tests in a variety of problems. The additional problem notebooks introduce a variety of different applications of permutation methods (using a variety of different test statistics) and use frequency generating functions to construct certain permutation distributions. Custom tools are used throughout, including tools for signed rank analyses, for constructing random reorderings of data, and for visualizing random reorderings of data.

Chapter 12 is an introduction to parametric and nonparametric bootstrap analysis. In the laboratory and additional problem notebooks, students consider the performance of the bootstrap and apply bootstrap estimation and testing methods in a variety of situations. Custom tools are used to construct random resamples, to visualize random resamples, to summarize the results of bootstrap analyses, and to construct approximate bootstrap confidence intervals using Efron's BC_α method in the nonparametric setting.

Chapter 13 is on parametric, nonparametric, and permutation methods for analysis of multiple samples. In the laboratory session, students use simulation to study analysis of variance for one-way layouts and blocked designs and to study Kruskal–Wallis and Friedman tests and apply these techniques in a variety of situations. Normal probability plots of standardized residuals are used to check analysis of variance assumptions. The additional problem notebooks reinforce these ideas and contain simulation studies and problems on analysis of variance in the balanced two-way layout setting. Custom tools are used throughout, including tools for analysis of variance, Bonferroni analysis, and Kruskal–Wallis and Friedman tests.

Chapter 14 is on linear least squares, including simple and multiple linear regression, permutation and bootstrap methods, and regression diagnostics. In the laboratory session, students use simulation to study the components of a linear regression analysis and apply the techniques in a variety of situations. The additional problem notebooks reinforce these ideas and contain problems on goodness-of-fit for simple linear models, analysis of covariance, model building, and locally weighted regression. Custom tools are provided for permutation analysis of slope in the simple linear setting, locally weighted regression, and diagnostic plots.

Chapter 15 is on large sample and small sample analyses of contingency tables, including diagnostic methods. In the laboratory session, students apply custom tools for large sample analyses of I -by- J tables and for constructing large sample confidence intervals for odds ratios to data from four studies. The additional problem

notebooks reinforce these ideas, consider the relationship between odds ratios and risk ratios, introduce McNemar's test for paired samples, and contain problems on permutation methods for fourfold and I -by- J tables.

Acknowledgments

This work was supported by Boston College through its faculty research programs and by the National Science Foundation through its Division of Undergraduate Education (NSF DUE 9555178). Boston College provided generous released time over several years while the materials were in development. NSF provided summer support for me, stipends for six additional faculty members, and generous support for an assistant.

Boston College Professors Dan Chambers and Charlie Landraitis, College of St. Catherine Professor Adele Rothan, C.S.J., and Stetson University Professor Erich Freedman used earlier versions of the laboratory materials in their classes and provided helpful comments and suggestions. Mt. Holyoke College Professor George Cobb and Harvard University Professor Marcello Pagano provided guidance on project design. I consulted with Boston College Professor Peg Kenney on assessment issues. University of Ottawa Professor John Nash and Boston College Professor Rob Gross provided interesting problem ideas and expert \LaTeX advice. Ms. Sarah Quebec worked with me as an undergraduate and masters student at Boston College and then as a research assistant on this project; her thoughtful comments helped shape the final product. The comments provided by students in my classes were uniformly helpful in improving the laboratory materials. I extend a warm thank you to SIAM's editorial team, especially Linda Thiel, and to the reviewers of the text and laboratory materials.

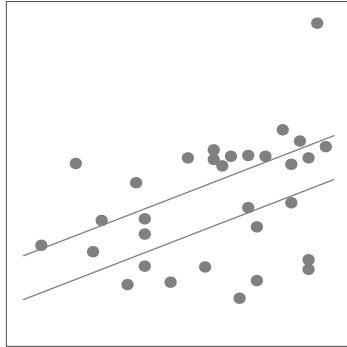
Data sets were kindly provided by Fairfield University Biology Professor Diane Brousseau, Boston College Geophysics Professors John Ebel and Alan Kafka, Dr. Susan Ford (Haskin Shellfish Laboratory, Rutgers University), and Dr. Roxana Smolowitz (Marine Biological Laboratories, University of Pennsylvania).

Copyright ©2005 by the Society for Industrial and Applied Mathematics

This electronic version is for personal use and may not be duplicated or distributed.

From "Mathematica Laboratories for Mathematical Statistics: Emphasizing Simulation and Computer Intensive Methods" by
Jenny Baglivo.

Buy this book from SIAM at www.ec-securehost.com/SIAM/SA14.html



Chapter 14

Linear Least Squares Analysis

Linear least squares methods allow researchers to study how variables are related. For example, a researcher might be interested in determining the relationship between the weight of an individual and such variables as height, age, sex, and general body dimensions.

Sections 1 and 2 introduce methods used to analyze how one variable can be used to predict another (for example, how height can be used to predict weight). Section 3 introduces methods to analyze how several variables can be used to predict another (for example, how the combination of height, age, sex, and general body dimensions can be used to predict weight). Bootstrap applications are given in Section 4. Section 5 outlines the laboratory problems. References for regression diagnostic methods are [12], [28], [49].

14.1 Simple linear model

A *simple linear model* is a model of the form

$$Y = \alpha + \beta X + \epsilon,$$

where X and ϵ are independent random variables, and the distribution of ϵ has mean 0 and standard deviation σ . Y is called the *response* variable, and X is called the *predictor* variable. ϵ represents the measurement error.

The response variable Y can be written as a linear function of the predictor variable X plus an error term. The linear prediction function has slope β and intercept α .

The objective is to estimate the parameters in the conditional mean formula

$$E(Y|X = x) = \alpha + \beta x$$

using a list of paired observations. The observed pairs are assumed to be either the values of a random sample from the joint (X, Y) distribution or a collection of

independent responses made at predetermined levels of the predictor. Analysis is done conditional on the observed values of the predictor variable.

14.1.1 Least squares estimation

Assume that

$$Y_i = \alpha + \beta x_i + \epsilon_i \text{ for } i = 1, 2, \dots, N$$

are independent random variables with means $E(Y_i) = \alpha + \beta x_i$, that the collection $\{\epsilon_i\}$ is a random sample from a distribution with mean 0 and standard deviation σ , and that all parameters (α , β , and σ) are unknown.

Least squares is a general estimation method introduced by A. Legendre in the early 1800's. In the simple linear case, the *least squares* (LS) estimators of α and β are obtained by minimizing the following sum of squared deviations of observed from expected responses:

$$S(\alpha, \beta) = \sum_{i=1}^N (Y_i - (\alpha + \beta x_i))^2.$$

Multivariable calculus can be used to demonstrate that the LS estimators of slope and intercept can be written in the form

$$\hat{\beta} = \sum_{i=1}^N \left[\frac{(x_i - \bar{x})}{S_{xx}} \right] Y_i \quad \text{and} \quad \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{x} = \sum_{i=1}^N \left[\frac{1}{N} - \frac{(x_i - \bar{x})\bar{x}}{S_{xx}} \right] Y_i,$$

where \bar{x} and \bar{Y} are the mean values of predictor and response, respectively, and S_{xx} is the sum of squared deviations of observed predictors from their sample mean:

$$S_{xx} = \sum_{i=1}^N (x_i - \bar{x})^2.$$

Formulas for $\hat{\alpha}$ and $\hat{\beta}$ can be written in many different ways. The method used here emphasizes that each estimator is a linear combination of the response variables.

Example: Olympic winning times

To illustrate the computations, consider the following 20 data pairs, where x is the time in years since 1900 and y is the Olympic winning time in seconds for men in the final round of the 100-meter event [50, p. 248]:

x	0	4	8	12	20	24	28	32	36	48
y	10.8	11.0	10.8	10.8	10.8	10.6	10.8	10.3	10.3	10.3
x	52	56	60	64	68	72	76	80	84	88
y	10.4	10.5	10.2	10.0	9.95	10.14	10.06	10.25	9.99	9.92

The data set covers all Olympic events held between 1900 and 1988. (Olympic games were not held in 1916, 1940, and 1944.) For these data, $\bar{x} = 45.6$, $\bar{y} = 10.396$, and

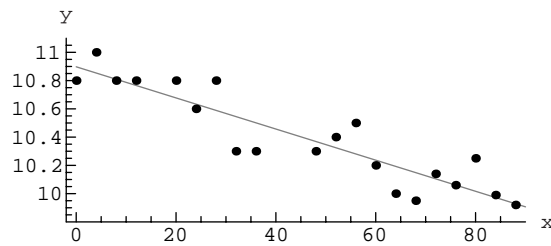


Figure 14.1. Olympic winning time in seconds for men's 100-meter finals (vertical axis) versus year since 1900 (horizontal axis). The gray line is the linear least squares fit, $y = 10.898 - 0.011x$.

the least squares estimates of slope and intercept are $\hat{\beta} = -0.011$ and $\hat{\alpha} = 10.898$, respectively. Figure 14.1 shows a scatter plot of the Olympic winning times data pairs superimposed on the least squares fitted line. The results suggest that the winning times have decreased at the rate of about 0.011 seconds per year during the 88 years of the study.

Properties of LS estimators

Theorem 4.4 can be used to demonstrate the following:

1. $E(\hat{\beta}) = \beta$ and $Var(\hat{\beta}) = \sigma^2/S_{xx}$.
2. $E(\hat{\alpha}) = \alpha$ and $Var(\hat{\alpha}) = (\sum_i x_i^2) \sigma^2 / (N S_{xx})$.

In addition, the following theorem, proven by Gauss and Markov, states that LS estimators are best (minimum variance) among all linear unbiased estimators of intercept and slope.

Theorem 14.1 (Gauss–Markov Theorem). *Under the assumptions of this section, the least squares (LS) estimators are the best linear unbiased estimators of α and β .*

For example, consider estimating β using a linear function of the response variables, say $W = c + \sum_i d_i Y_i$ for some constants c and d_1, d_2, \dots, d_N . If W is an unbiased estimator of β , then

$$Var(W) = Var\left(c + \sum_i d_i Y_i\right) = \sum_i d_i^2 Var(Y_i) = \left(\sum_i d_i^2\right) \sigma^2$$

is minimized when $d_i = (x_i - \bar{x})/S_{xx}$ and $c = 0$. That is, the variance is minimized when W is the LS estimator of β .

Although LS estimators are best among linear unbiased estimators, they may not be ML estimators. Thus, there may be other more efficient methods of estimation.

14.1.2 Permutation confidence interval for slope

Permutation methods can be used to construct confidence intervals for the slope parameter β in the simple linear model. Let

$$(x_i, y_i) \text{ for } i = 1, 2, \dots, N$$

be the observed pairs and π be a permutation of the indices $1, 2, \dots, N$ other than the identity. Then the quantity

$$b(\pi) = \frac{\sum_i (x_i - \bar{x})(y_{\pi(i)} - y_i)}{\sum_i (x_i - \bar{x})(x_{\pi(i)} - x_i)}$$

is an estimate of β , and the collection

$$\{b(\pi) : \pi \text{ is a permutation other than the identity}\}$$

is a list of $N! - 1$ estimates. The ordered estimates

$$b_{(1)} < b_{(2)} < b_{(3)} < \dots < b_{(N!-1)}$$

are used in constructing confidence intervals.

Theorem 14.2 (Slope Confidence Intervals). *Under the assumptions of this section, the interval*

$$[b_{(k)}, b_{(N!-k)}]$$

is a $100(1 - 2k/N!)%$ confidence interval for β .

The procedure given in Theorem 14.2 is an example of *inverting* a hypothesis test: A value β_o is in a $100(1 - \gamma)%$ confidence interval if the two sided permutation test of

$$H_o : \text{The correlation between } Y - \beta_o X \text{ and } X \text{ is zero}$$

is accepted at the γ significance level. For a proof, see [74, p. 120].

Since the number of permutations can be quite large, Monte Carlo analysis is used to estimate endpoints. For example, assume the Olympic times data (page 206) are the values of random variables satisfying the assumptions of this section. An approximate 95% confidence interval for the slope parameter (based on 5000 random permutations) is $[-0.014, -0.008]$.

14.2 Simple linear regression

In simple *linear regression*, the error distribution is assumed to be normal, and, as above, analyses are done conditional on the observed values of the predictor variable. Specifically, assume that

$$Y_i = \alpha + \beta x_i + \epsilon_i \text{ for } i = 1, 2, \dots, N$$

are independent random variables with means $E(Y_i) = \alpha + \beta x_i$, that the collection $\{\epsilon_i\}$ is a random sample from a normal distribution with mean 0 and standard deviation σ , and that all parameters are unknown.

In this setting, LS estimators are ML estimators.

Theorem 14.3 (Parameter Estimation). *Given the assumptions and definitions above, the LS estimators of α and β given on page 206 are ML estimators, and the statistics*

$$\begin{aligned}\hat{\epsilon}_i &= Y_i - (\hat{\alpha} + \hat{\beta}x_i) = Y_i - (\bar{Y} + \hat{\beta}(x_i - \bar{x})) \\ &= Y_i - \sum_j \left[\frac{1}{N} + \frac{(x_j - \bar{x})(x_i - \bar{x})}{S_{xx}} \right] Y_j\end{aligned}$$

are ML estimators of the error terms for $i = 1, 2, \dots, N$. Each estimator is a normal random variable, and each is unbiased. Further, the statistic

$$S^2 = \frac{1}{N-2} \sum_i (Y_i - (\hat{\alpha} + \hat{\beta}x_i))^2$$

is an unbiased estimator of the common variance σ^2 .

14.2.1 Confidence interval procedures

This section develops confidence interval procedures for the slope and intercept parameters, and for the mean response at a fixed value of the predictor variable.

Hypothesis tests can also be developed. Most computer programs automatically include both types of analyses.

Confidence intervals for β

Since the LS estimator $\hat{\beta}$ is a normal random variable with mean β and variance σ^2/S_{xx} , Theorem 6.2 can be used to demonstrate that

$$\hat{\beta} \pm t_{N-2}(\gamma/2) \sqrt{\frac{S^2}{S_{xx}}}$$

is a $100(1 - \gamma)\%$ confidence interval for β , where S^2 is the estimate of the common variance given in Theorem 14.3 and $t_{N-2}(\gamma/2)$ is the $100(1 - \gamma/2)\%$ point on the Student t distribution with $(N - 2)$ degrees of freedom.

Confidence intervals for α

Since the LS estimator $\hat{\alpha}$ is a normal random variable with mean α and variance $\sigma^2 (\sum_i x_i^2) / (N S_{xx})$, Theorem 6.2 can be used to demonstrate that

$$\hat{\alpha} \pm t_{N-2}(\gamma/2) \sqrt{\frac{S^2 (\sum_i x_i^2)}{N S_{xx}}}$$

is a $100(1 - \gamma)\%$ confidence interval for α , where S^2 is the estimate of the common variance given in Theorem 14.3 and $t_{N-2}(\gamma/2)$ is the $100(1 - \gamma/2)\%$ point on the Student t distribution with $(N - 2)$ degrees of freedom.

For example, if the Olympic times data (page 206) are the values of random variables satisfying the assumptions of this section, then a 95% confidence interval for the slope parameter is $[-0.013, -0.009]$, and a 95% confidence interval for the intercept parameter is $[10.765, 11.030]$.

Confidence intervals for mean response

The mean response $E(Y_o) = \alpha + \beta x_o$ at a new predictor-response pair, (x_o, Y_o) , can be estimated using the statistic

$$\hat{\alpha} + \hat{\beta}x_o = \bar{Y} + \hat{\beta}(x_o - \bar{x}) = \sum_{i=1}^N \left[\frac{1}{N} + \frac{(x_o - \bar{x})(x_i - \bar{x})}{S_{xx}} \right] Y_i.$$

This estimator is a normal random variable (by Theorem 4.6) with mean $\alpha + \beta x_o$ and

$$\text{Var}(\hat{\alpha} + \hat{\beta}x_o) = \sigma^2 \left(\frac{1}{N} + \frac{(x_o - \bar{x})^2}{S_{xx}} \right).$$

Thus, Theorem 6.2 can be used to demonstrate that

$$(\hat{\alpha} + \hat{\beta}x_o) \pm t_{N-2}(\gamma/2) \sqrt{S^2 \left(\frac{1}{N} + \frac{(x_o - \bar{x})^2}{S_{xx}} \right)}$$

is a $100(1 - \gamma)\%$ confidence interval for $\alpha + \beta x_o$, where S^2 is the estimate of the common variance given in Theorem 14.3 and $t_{N-2}(\gamma/2)$ is the $100(1 - \gamma/2)\%$ point on the Student t distribution with $(N - 2)$ degrees of freedom.

Example: Percentage of dead or damaged spruce trees

For example, as part of a study on the relationship between environmental stresses and the decline of red spruce tree forests in the Appalachian Mountains, data were collected on the percentage of dead or damaged trees at various altitudes in forests in the northeast. The paired data were of interest because concentrations of airborne pollutants tend to be higher at higher altitudes [49, p. 102].

Figure 14.2 is based on information gathered in 53 areas. For these data, the least squares fitted line is $y = 8.24x - 33.66$, suggesting that the percentage of damaged or dead trees increases at the rate of 8.24 percentage points per 100 meters elevation.

An estimate of the mean response at 1000 meters ($x_o = 10$) is 48.76% damaged or dead. If these data are the values of independent random variables satisfying the assumptions of this section, then a 95% confidence interval for the mean response at 1000 meters is $[48.44, 49.07]$.

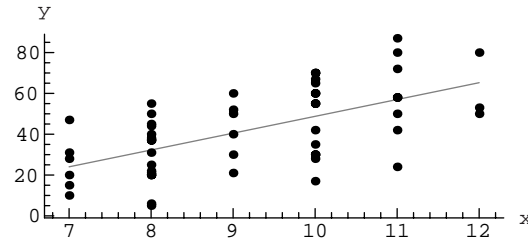


Figure 14.2. Percentage dead or damaged red spruce trees (vertical axis) versus elevation in 100 meters (horizontal axis) at 53 locations in the northeast. The gray line is the linear least squares fit, $y = 8.24x - 33.66$.

Comparison of procedures

The confidence interval procedure for β given in this section is valid when the error distribution is normal. When the error distribution is not normal, the permutation procedure given in Theorem 14.2 can be used.

The confidence interval procedures given in this section assume that the values of the predictor variable are known with certainty (the procedures are conditional on the observed values of the predictor) and assume that the error distributions are normal. Approximate bootstrap confidence interval procedures can also be developed under broader conditions; see Section 14.4.

14.2.2 Predicted responses and residuals

The i^{th} estimated mean (or *predicted response*) is the random variable

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}x_i = \sum_j \left[\frac{1}{N} + \frac{(x_j - \bar{x})(x_i - \bar{x})}{S_{xx}} \right] Y_j \quad \text{for } i = 1, 2, \dots, N,$$

and the i^{th} estimated error (or *residual*) is

$$\hat{\epsilon}_i = Y_i - \hat{Y}_i \quad \text{for } i = 1, 2, \dots, N.$$

Each random variable is a linear function of the response variables. Theorem 4.5 can be used to demonstrate that $\text{Cov}(\hat{Y}_i, \hat{\epsilon}_i) = 0$.

Although the error terms in the simple linear model have equal variances, the estimated errors do not. Specifically, the variance of the i^{th} residual is

$$\text{Var}(\hat{\epsilon}_i) = \sigma^2 \left[\left(1 - \frac{1}{N} - \frac{(x_i - \bar{x})^2}{S_{xx}} \right)^2 + \sum_{j \neq i} \left(\frac{1}{N} + \frac{(x_j - \bar{x})(x_i - \bar{x})}{S_{xx}} \right)^2 \right] = \sigma^2 c_i.$$

The i^{th} estimated *standardized residual* is defined as follows:

$$r_i = \hat{\epsilon}_i / \sqrt{S^2 c_i} \quad \text{for } i = 1, 2, \dots, N,$$

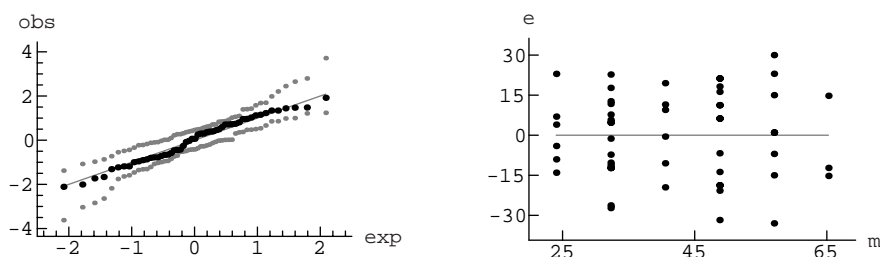


Figure 14.3. Enhanced normal probability plot of standardized residuals (left plot) and scatter plot of residuals versus estimated means (right plot) for the spruce trees example.

where S^2 is the estimate of the common variance given in Theorem 14.3 and c_i is the constant in brackets above.

Predicted responses, residuals, and estimated standardized residuals are used in diagnostic plots of model assumptions. For example, the left plot in Figure 14.3 is an enhanced normal probability of the estimated standardized residuals from the spruce trees example (page 210), and the right plot is a scatter plot of residuals (vertical axis) versus predicted responses (horizontal axis). The left plot suggests that the error distribution is approximately normally distributed; the right plot exhibits no relationship between the estimated errors and estimated means.

The scatter plot of residuals versus predicted responses should show no relationship between the variables. Of particular concern are the following:

1. If $\hat{\epsilon}_i \approx h(\hat{y}_i)$ for some function h , then the assumption that the conditional mean is a linear function of the predictor may be wrong.
2. If $SD(\hat{\epsilon}_i) \approx h(\hat{y}_i)$ for some function h , then the assumption of equal standard deviations may be wrong.

14.2.3 Goodness-of-fit

Suppose that the N predictor-response pairs can be written in the following form:

$$(x_i, Y_{i,j}) \text{ for } j = 1, 2, \dots, n_i, i = 1, 2, \dots, I.$$

(There are a total of n_i observed responses at the i^{th} level of the predictor variable for $i = 1, 2, \dots, I$, and $N = \sum_i n_i$.) Then it is possible to use an analysis of variance technique to test the goodness-of-fit of the simple linear model.

Assumptions

The responses are assumed to be the values of I independent random samples

$$Y_{i,1}, Y_{i,2}, \dots, Y_{i,n_i} \text{ for } i = 1, 2, \dots, I$$

from normal distributions with a common unknown standard deviation σ .

Let μ_i be the mean of responses in the i^{th} sample: $\mu_i = E(Y_{i,j})$ for all j . Of interest is to test the null hypothesis that $\mu_i = \alpha + \beta x_i$ for $i = 1, 2, \dots, I$.

Sources of variation

The formal goodness-of-fit analysis is based on writing the sum of squared deviations of the response variables from the predicted responses using the linear model (known as the *error* sum of squares),

$$SS_e = \sum_{i,j} (Y_{i,j} - (\hat{\alpha} + \hat{\beta}x_i))^2,$$

as the sum of squared deviations of the response variables from the estimated group means (known as the *pure error* sum of squares),

$$SS_p = \sum_{i,j} (Y_{i,j} - \bar{Y}_i)^2,$$

plus the weighted sum of squared deviations of the group means from the predicted responses (known as the *lack-of-fit* sum of squares),

$$SS_\ell = \sum_i n_i (\bar{Y}_i - (\hat{\alpha} + \hat{\beta}x_i))^2.$$

Pure error and lack-of-fit mean squares

The *pure error* mean square, MS_p , is defined as follows:

$$MS_p = \frac{1}{N-I} SS_p = \frac{1}{N-I} \sum_{i,j} (Y_{i,j} - \bar{Y}_i)^2.$$

MS_p is equal to the pooled estimate of the common variance. Theorem 6.1 can be used to demonstrate that $(N-I)MS_p/\sigma^2$ is a chi-square random variable with $(N-I)$ degrees of freedom.

The *lack-of-fit* mean square, MS_ℓ , is defined as follows:

$$MS_\ell = \frac{1}{I-2} SS_\ell = \frac{1}{I-2} \sum_i n_i (\bar{Y}_i - (\hat{\alpha} + \hat{\beta}x_i))^2.$$

Properties of expectation can be used to demonstrate that

$$E(MS_\ell) = \sigma^2 + \frac{1}{I-2} \sum_i n_i (\mu_i - (\alpha + \beta x_i))^2.$$

If the null hypothesis that the means follow a simple linear model is true, then the expected value of MS_ℓ is σ^2 ; otherwise, values of MS_ℓ will tend to be larger than σ^2 . The following theorem relates the pure error and lack-of-fit mean squares.

Theorem 14.4 (Distribution Theorem). *Under the general assumptions of this section and if the null hypothesis is true, then the ratio $F = MS_\ell/MS_p$ has an f ratio distribution with $(I-2)$ and $(N-I)$ degrees of freedom.*

Table 14.1. Goodness-of-fit analysis of the spruce tree data.

Source	df	SS	MS	F	p value
Lack-of-Fit	4	132.289	33.072	0.120	0.975
Pure Error	47	12964.3	275.835		
Error	51	13096.5			

Goodness-of-fit test: Observed significance level

Large values of $F = MS_{\ell}/MS_p$ support the alternative hypothesis that the simple linear model does not hold. For an observed ratio, f_{obs} , the p value is $P(F \geq f_{\text{obs}})$.

For example, assume the spruce trees data (page 210) satisfy the general assumptions of this section. Table 14.1 shows the results of the goodness-of-fit test. There were 6 observed predictor values. The observed ratio of the lack-of-fit mean square to the pure error mean square is 0.120. The observed significance level, based on the f ratio distribution with 4 and 47 degrees of freedom, is 0.975. The simple linear model fits the data quite well.

14.3 Multiple linear regression

A *linear model* is a model of the form

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_{p-1} X_{p-1} + \epsilon,$$

where each X_i is independent of ϵ , and the distribution of ϵ has mean 0 and standard deviation σ . Y is called the *response* variable, each X_i is a *predictor* variable, and ϵ represents the measurement error.

The response variable Y can be written as a linear function of the $(p-1)$ predictor variables plus an error term. The linear prediction function has p parameters.

In multiple *linear regression*, the error distribution is assumed to be normal, and analyses are done conditional on the observed values of the predictor variables. Observations are called *cases*.

14.3.1 Least squares estimation

Assume that

$$Y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \cdots + \beta_{p-1} x_{p-1,i} + \epsilon_i \text{ for } i = 1, 2, \dots, N$$

are independent random variables with means

$$E(Y_i) = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \cdots + \beta_{p-1} x_{p-1,i} \text{ for all } i,$$

that the collection of errors $\{\epsilon_i\}$ is a random sample from a normal distribution with mean 0 and standard deviation σ , and that all parameters (including σ) are unknown.

The *least squares* (LS) estimators of the coefficients in the linear prediction function are obtained by minimizing the following sum of squared deviations of observed from expected responses:

$$\begin{aligned} S(\beta_0, \beta_1, \dots, \beta_{p-1}) &= \sum_{k=1}^N (Y_k - (\beta_0 + \beta_1 x_{1,k} + \beta_2 x_{2,k} + \dots + \beta_{p-1} x_{p-1,k}))^2 \\ &= \sum_{k=1}^N \left(Y_k - \sum_{j=0}^{p-1} \beta_j x_{j,k} \right)^2, \quad \text{where } x_{0,k} = 1 \text{ for all } k. \end{aligned}$$

The first step in the analysis is to compute the partial derivative with respect to β_i for each i . Partial derivatives have the following form:

$$\frac{\partial S}{\partial \beta_i} = -2 \left[\sum_{k=1}^N Y_k x_{i,k} - \sum_{j=0}^{p-1} \beta_j \left(\sum_{k=1}^N x_{j,k} x_{i,k} \right) \right].$$

The next step is to solve the p -by- p system of equations

$$\frac{\partial S}{\partial \beta_i} = 0, \quad \text{for } i = 0, 1, \dots, p-1,$$

or, equivalently,

$$\sum_{j=0}^{p-1} \left(\sum_{k=1}^N x_{j,k} x_{i,k} \right) \beta_j = \sum_{k=1}^N Y_k x_{i,k}, \quad \text{for } i = 0, 1, \dots, p-1.$$

In matrix notation, the system becomes

$$(\mathbf{X}^T \mathbf{X}) \underline{\beta} = \mathbf{X}^T \underline{Y},$$

where $\underline{\beta}$ is the p -by-1 vector of unknown parameters, \underline{Y} is the N -by-1 vector of response variables, \mathbf{X} is the N -by- p matrix whose (i, j) element is $x_{j,i}$, and \mathbf{X}^T is the transpose of the \mathbf{X} matrix. The \mathbf{X} matrix is often called the *design matrix*. Finally, the p -by-1 vector of LS estimators is

$$\hat{\underline{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underline{Y}.$$

Estimates exist as long as $(\mathbf{X}^T \mathbf{X})$ is invertible.

The rows of the design matrix correspond to the observations (or cases). The columns correspond to the predictors. The terms in the first column of the design matrix are identically equal to one.

For example, in the simple linear case, the matrix product

$$(\mathbf{X}^T \mathbf{X}) = \begin{bmatrix} N & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix}$$

has inverse

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{NS_{xx}} \begin{bmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & N \end{bmatrix}, \quad \text{where } S_{xx} = \sum_i (x_i - \bar{x})^2,$$

and the LS estimators are

$$\begin{bmatrix} \hat{\alpha} \\ \hat{\beta} \end{bmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underline{Y} = \begin{bmatrix} \sum_i \left(\frac{1}{N} - \frac{(x_i - \bar{x})\bar{x}}{S_{xx}} \right) Y_i \\ \sum_i \left(\frac{(x_i - \bar{x})}{S_{xx}} \right) Y_i \end{bmatrix}.$$

The estimators here correspond exactly to those given on page 206.

Model in matrix form

The model can be written as

$$\underline{Y} = \mathbf{X}\underline{\beta} + \underline{\epsilon},$$

where \underline{Y} and $\underline{\epsilon}$ are N -by-1 vectors of responses and errors, respectively, $\underline{\beta}$ is the p -by-1 coefficient vector, and \mathbf{X} is the N -by- p design matrix.

Theorem 14.5 (Parameter Estimation). *Given the assumptions and definitions above, the vector of LS estimators of $\underline{\beta}$ given on page 215 is a vector of ML estimators, and the vector*

$$\hat{\underline{\epsilon}} = \underline{Y} - \mathbf{X}\hat{\underline{\beta}} = (\mathbf{I} - \mathbf{H}) \underline{Y},$$

where $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ and \mathbf{I} is the N -by- N identity matrix, is a vector of ML estimators of the error terms. Each estimator is a normal random variable, and each is unbiased. Further, the statistic

$$S^2 = \frac{1}{N-p} \sum_{i=1}^N (Y_i - \hat{Y}_i)^2,$$

where \hat{Y}_i is the i^{th} estimated mean, is an unbiased estimator of σ^2 .

The i^{th} estimated mean (or *predicted response*) is the random variable

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i,1} + \cdots + \hat{\beta}_{p-1} x_{i,p-1} \text{ for } i = 1, 2, \dots, N.$$

Further, the matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is often called the *hat matrix* since it is the matrix that transforms the response vector to the predicted response vector

$$\hat{\underline{Y}} = \mathbf{X}\hat{\underline{\beta}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underline{Y} = \mathbf{H} \underline{Y}$$

(the vector of Y_i 's is transformed to the vector of Y_i hats).

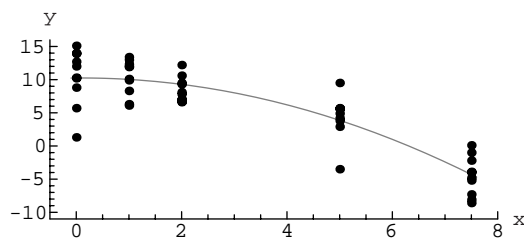


Figure 14.4. Change in weight in grams (vertical axis) versus dosage level in 100 mg/kg/day (horizontal axis) for data from the toxicology study. The gray curve is the linear least squares fit, $y = 10.2475 + 0.053421x - 0.2658x^2$.

Variability of LS estimators

If \underline{V} is an m -by-1 vector of random variables and \underline{W} is an n -by-1 vector of random variables, then $\Sigma(\underline{V}, \underline{W})$ is the m -by- n matrix whose (i, j) term is $Cov(V_i, W_j)$. The matrix $\Sigma(\underline{V}, \underline{W})$ is called a *covariance matrix*.

Theorem 14.6 (Covariance Matrices). Under the assumptions of this section, the following hold:

1. The covariance matrix of the coefficient estimators is

$$\sum (\hat{\underline{\beta}}, \hat{\underline{\beta}}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}.$$

2. The covariance matrix of the error estimators is

$$\sum (\hat{\underline{\epsilon}}, \hat{\underline{\epsilon}}) = \sigma^2 (\mathbf{I} - \mathbf{H}).$$

3. The covariance matrix of error estimators and predicted responses is

$$\sum (\hat{\underline{\epsilon}}, \hat{\underline{Y}}) = \mathbf{0},$$

where $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is the hat matrix, \mathbf{I} is the N -by- N identity matrix, and $\mathbf{0}$ is the N -by- N matrix of zeros.

The diagonal elements of $\sum (\hat{\underline{\beta}}, \hat{\underline{\beta}})$ and $\sum (\hat{\underline{\epsilon}}, \hat{\underline{\epsilon}})$ are the variances of the coefficient and error estimators, respectively. The last statement in the theorem says that error estimators and predicted responses are uncorrelated.

Example: Toxicology study

To illustrate some of the computations, consider the data pictured in Figure 14.4, collected as part of a study to assess the adverse effects of a proposed drug for the treatment of tuberculosis [40].

Ten female rats were given the drug for a period of 14 days at each of five dosage levels (in 100 milligrams per kilogram per day). The vertical axis in the plot

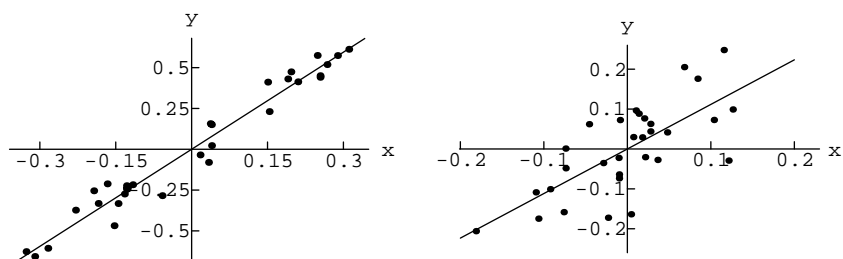


Figure 14.5. Partial regression plots for the data from the timber yield study. The left plot pictures residuals of log-volume (vertical axis) versus log-diameter (horizontal axis) with the effect of log-height removed. The right plot pictures residuals of log-volume (vertical axis) versus log-height (horizontal axis) with the effect of log-diameter removed. The gray lines are $y = 1.983x$ in the left plot and $y = 1.117x$ in the right plot.

shows the weight change in grams (WC), defined as the weight at the end of the period minus the weight at the beginning of the period; the horizontal axis shows the dose in 100 mg/kg/day. A linear model of the form

$$WC = \beta_0 + \beta_1 \text{dose} + \beta_2 \text{dose}^2 + \epsilon$$

was fit to the 50 (dose, WC) cases. (The model is linear in the unknown parameters and quadratic in the dosage level.) The LS prediction equation is shown in the plot.

Example: Timber yield study

As part of a study to find an estimate for the volume of a tree (and therefore its yield) given its diameter and height, data were collected on the volume (in cubic feet), diameter at 54 inches above the ground (in inches), and height (in feet) of 31 black cherry trees in the Allegheny National Forest [50, p. 159]. Since a multiplicative relationship is expected among these variables, a linear model of the form

$$\log\text{-volume} = \beta_0 + \beta_1 \log\text{-diameter} + \beta_2 \log\text{-height} + \epsilon$$

was fit to the 31 (log-diameter, log-height, log-volume) cases, using the natural logarithm function to compute log values.

The LS prediction equation is

$$\log\text{-volume} = -6.632 + 1.983 \log\text{-diameter} + 1.117 \log\text{-height}.$$

Figure 14.5 shows *partial regression* plots of the timber yield data.

- (i) In the left plot, the log-volume and log-diameter variables are adjusted to remove the effects of log-height. Specifically, the residuals from the simple linear regression of log-volume on log-height (vertical axis) are plotted against the residuals from the simple linear regression of log-diameter on log-height (horizontal axis). The relationship between the adjusted variables can be described using the linear equation $y = 1.983x$.

- (ii) In the right plot, the log-volume and log-height variables are adjusted to remove the effects of log-diameter. The relationship between the adjusted variables can be described using the linear equation $y = 1.117x$.

The slopes of the lines in the partial regression plots correspond to the LS estimates in the prediction equation above. The plots suggest that a linear relationship between the response variable and each of the predictors is reasonable.

14.3.2 Analysis of variance

The linear regression model can be reparametrized as follows:

$$Y_i = \mu + \sum_{j=1}^{p-1} \beta_j(x_{j,i} - \bar{x}_j) + \epsilon_i \text{ for } i = 1, 2, \dots, N,$$

where μ is the overall mean

$$\mu = \frac{1}{N} \sum_{i=1}^N E(Y_i) = \beta_0 + \sum_{j=1}^{p-1} \beta_j \bar{x}_j.$$

and \bar{x}_j is the mean of the j^{th} predictor for all j . The difference $E(Y_i) - \mu$ is called the i^{th} deviation (or the i^{th} regression effect). The sum of the regression effects is zero.

This section develops an analysis of variance F test for the null hypothesis that the regression effects are identically zero (equivalently, a test of the null hypothesis that $\beta_i = 0$ for $i = 1, 2, \dots, p - 1$).

If the null hypothesis is accepted, then the $(p - 1)$ predictor variables have no predictive ability; otherwise, they have some predictive ability.

Sources of variation; coefficient of determination

In the first step of the analysis, the sum of squared deviations of the response variables from the mean response (the *total* sum of squares),

$$SS_t = \sum_{i=1}^N (Y_i - \bar{Y})^2,$$

is written as the sum of squared deviations of the response variables from the predicted responses (the *error* sum of squares),

$$SS_e = \sum_{i=1}^N (Y_i - \hat{Y}_i)^2,$$

plus the sum of squared deviations of the predicted responses from the mean response (the *model* sum of squares),

$$SS_m = \sum_{i=1}^N (\hat{Y}_i - \bar{Y})^2 = \sum_{i=1}^N \left(\sum_{j=1}^{p-1} \hat{\beta}_j (x_{j,i} - \bar{x}_j) \right)^2.$$

The ratio of the model to the total sums of squares, $R^2 = SS_m/SS_t$, is called the *coefficient of determination*. R^2 is the proportion of the total variation in the response variable that is explained by the model.

In the simple linear case, R^2 is the same as the square of the sample correlation.

Analysis of variance f test

The *error* mean square is the ratio $MS_e = SS_e/(N - p)$, and the *model* mean square is the ratio $MS_m = SS_m/(p - 1)$. The following theorem relates these random variables.

Theorem 14.7 (Distribution Theorem). *Under the general assumptions of this section and if the null hypothesis is true, then the ratio $F = MS_m/MS_e$ has an f ratio distribution with $(p - 1)$ and $(N - p)$ degrees of freedom.*

Large values of $F = MS_m/MS_e$ support the hypothesis that the proposed predictor variables have some predictive ability. For an observed ratio, f_{obs} , the p value is $P(F \geq f_{\text{obs}})$.

For the toxicology study example (page 217), $f_{\text{obs}} = 82.3$ and the p value, based on the f ratio distribution with 2 and 47 degrees of freedom, is virtually zero. The coefficient of determination is 0.778; the estimated linear model explains about 77.8% of the variation in weight change.

For the timber yield example (page 218), $f_{\text{obs}} = 613.2$ and the p value, based on the f ratio distribution with 2 and 28 degrees of freedom, is virtually zero. The coefficient of determination is 0.978; the estimated linear model explains about 97.8% of the variation in log-volume.

It is possible for the f test to reject the null hypothesis and the value of R^2 to be close to zero. In this case, the potential predictors have some predictive ability, but additional (or different) predictor variables are needed to adequately model the response.

14.3.3 Confidence interval procedures

This section develops confidence interval procedures for the β parameters and for the mean response at a fixed combination of the predictor variables.

Hypothesis tests can also be developed. Most computer programs automatically include both types of analyses.

Confidence intervals for β_i

Let v_i be the element in the (i, i) position of $(\mathbf{X}^T \mathbf{X})^{-1}$, and let S^2 be the estimate of the common variance given in Theorem 14.5. Since the LS estimator $\hat{\beta}_i$ is a normal random variable with mean β_i and variance $\sigma^2 v_i$, Theorem 6.2 can be used to demonstrate that

$$\hat{\beta}_i \pm t_{N-p}(\gamma/2) \sqrt{S^2 v_i}$$

is a $100(1 - \gamma)\%$ confidence interval for β_i , where $t_{N-p}(\gamma/2)$ is the $100(1 - \gamma/2)\%$ point on the Student t distribution with $(N - p)$ degrees of freedom.

For example, if the data in the toxicology study (page 217) are the values of random variables satisfying the assumptions of this section, then a 95% confidence interval for β_2 is $[-0.43153, -0.10008]$. Since 0 is not in the confidence interval, the result suggests that the model with the dose² term is significantly better than a simple linear model relating dose to weight change.

Confidence intervals for mean response

The mean response $E(Y_o) = \sum_{i=0}^{p-1} \beta_i x_{i,0}$ at a new predictors-response case can be estimated using the statistic

$$\sum_{i=0}^{p-1} \hat{\beta}_i x_{i,0} = \underline{x}_0^T \hat{\underline{\beta}} = \underline{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underline{Y},$$

where $\underline{x}_0^T = (1, x_{1,0}, x_{2,0}, \dots, x_{p-1,0})$. This estimator is a normal random variable with mean $E(Y_o)$ and variance

$$\sigma^2 \left(\underline{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \underline{x}_0 \right) = \sigma^2 v_o.$$

Thus, Theorem 6.2 can be used to demonstrate that

$$\sum_{i=0}^{p-1} \hat{\beta}_i x_{i,0} \pm t_{N-p}(\gamma/2) \sqrt{S^2 v_o}$$

is a $100(1 - \gamma)\%$ confidence interval for $E(Y_o)$, where S^2 is the estimate of the common variance given in Theorem 14.3 and $t_{N-p}(\gamma/2)$ is the $100(1 - \gamma/2)\%$ point on the Student t distribution with $(N - p)$ degrees of freedom.

For example, an estimate of the mean log-volume of a tree with diameter 11.5 inches and height 80 inches is 3.106 log-cubic inches. If these data are the values of random variables satisfying the assumptions of this section, then a 95% confidence interval for the mean response at this combination of the predictors is $[3.05944, 3.1525]$.

14.3.4 Regression diagnostics

Recall that the hat matrix $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is the matrix that transforms the vector of observed responses \underline{Y} to the vector of predicted responses $\hat{\underline{Y}}$. Each predicted response is a linear combination of the observed responses:

$$\hat{Y}_i = \sum_{j=1}^N h_{i,j} Y_j \text{ for } i = 1, 2, \dots, N,$$

where $h_{i,j}$ is the (i, j) element of \mathbf{H} . In particular, the diagonal element $h_{i,i}$ is the coefficient of Y_i in the formula for \hat{Y}_i .

Leverage

The *leverage* of the i^{th} response is the value $h_i = h_{i,i}$. Leverages satisfy the following properties:

1. For each i , $0 \leq h_i \leq 1$.
2. $\sum_{i=1}^N h_i = p$, where p is the number of parameters.

Ideally, the leverages should be about p/N each (the average value).

Residuals and standardized residuals

Theorem 14.6 implies that the variance of the i^{th} estimated error (or *residual*) is $\text{Var}(\hat{\epsilon}_i) = \sigma^2(1 - h_i)$, where h_i is the leverage. The i^{th} estimated *standardized residual* is defined as follows:

$$r_i = \hat{\epsilon}_i / \sqrt{S^2(1 - h_i)} \text{ for } i = 1, 2, \dots, N,$$

where S^2 is the estimate of the common variance given in Theorem 14.5.

Residuals and standardized residuals are used in diagnostic plots of model assumptions. See Section 14.2.2 for examples in the simple linear case.

Standardized influences

The *influence* of the i^{th} observation is the change in prediction if the i^{th} observation is deleted from the data set.

Specifically, the influence is the difference $\hat{Y}_i - \hat{Y}_i(i)$, where \hat{Y}_i is the predicted response using all N cases to compute parameter estimates, and $\hat{Y}_i(i)$ is the prediction at a “new” predictor vector \underline{x}_i , where parameter estimates have been computed using the list of $N - 1$ cases with the i^{th} case removed.

For the model estimated using $N - 1$ cases only, linear algebra methods can be used to demonstrate that the predicted response is

$$\hat{Y}_i(i) = \hat{Y}_i - \hat{\epsilon}_i \frac{h_i}{1 - h_i}$$

and the estimated common variance is

$$S^2(i) = \frac{1}{N - p - 1} \left((N - p)S^2 - \frac{\hat{\epsilon}_i^2}{1 - h_i} \right).$$

The i^{th} standardized influence is the ratio of the influence to the standard deviation of the predicted response,

$$\frac{\hat{Y}_i - \hat{Y}_i(i)}{SD(\hat{Y}_i)} = \frac{\hat{\epsilon}_i h_i / (1 - h_i)}{\sqrt{\sigma^2 h_i}},$$

and the i^{th} estimated *standardized influence* is the value obtained by substituting $S^2(i)$ for σ^2 :

$$\delta_i = \frac{\hat{Y}_i - \hat{Y}_i(i)}{SD(\hat{Y}_i)} = \frac{\hat{\epsilon}_i h_i / (1 - h_i)}{\sqrt{S(i)^2 h_i}} = \frac{\hat{\epsilon}_i \sqrt{h_i}}{S(i)(1 - h_i)}.$$

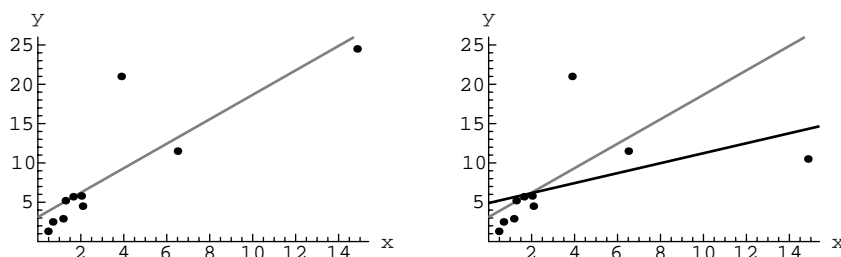


Figure 14.6. Scatter plots of example pairs (left plot) and altered example pairs (right plot). The gray line in both plots has equation $y = 3.11 + 1.55x$. The black line in the right plot has equation $y = 4.90 + 0.63x$.

Ideally, predicted responses should change very little if one case is removed from the list of N cases, and each δ_i should be close to zero. A general rule of thumb is that if $|\delta_i|$ is much greater than $2\sqrt{p/N}$, then the i^{th} case is highly influential.

Illustration

To illustrate the computations in the simple linear case, consider the following list of 10 (x, y) pairs:

x	0.47	0.69	1.17	1.28	1.64	2.02	2.08	3.88	6.50	14.86
y	1.30	2.50	2.90	5.20	5.70	5.80	4.50	21.00	11.50	24.50

The left plot in Figure 14.6 shows a scatter plot of the data pairs superimposed on the least squares fitted line, $y = 3.11 + 1.55x$. The following table gives the residuals, leverages, and standardized influences for each case:

i	1	2	3	4	5	6	7	8	9	10
$\hat{\epsilon}_i$	-2.54	-1.68	-2.03	0.10	0.04	-0.45	-1.85	11.86	-1.72	-1.72
h_i	0.15	0.14	0.13	0.13	0.12	0.11	0.11	0.10	0.15	0.85
δ_i	-0.25	-0.16	-0.18	0.01	0.00	-0.04	-0.14	4.15	-0.17	-2.33

Based on the rule of thumb above, cases 8 and 10 are highly influential. Case 8 has a very large residual, and case 10 has a very large leverage value.

The right plot in Figure 14.6 illustrates the concept of leverage. If the observed response in case 10 is changed from 24.5 to 10.5, then the predicted response changes from 26.2 to 14.32. The entire line has moved to accommodate the change in case 10.

Different definitions of δ_i appear in the literature, although most books use the definition above. The rule of thumb is from [12], where the notation DFFITS_i is used for δ_i .

14.4 Bootstrap methods

Bootstrap resampling methods can be applied to analyzing the relationship between one or more predictors and a response. This section introduces two examples.

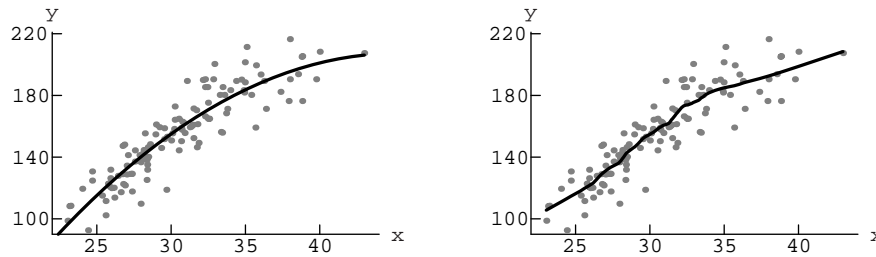


Figure 14.7. Scatter plots of weight in pounds (vertical axis) versus waist circumference in inches (horizontal axis) for 120 physically active young adults. In the left plot, the curve $y = -252.569 + 20.322x - 0.225x^2$ is superimposed. In the right plot, the 25% lowess smooth is superimposed.

Example: Unconditional analysis of linear models

If the observed cases are the values of a random sample from a joint distribution, then nonparametric bootstrap methods can be used to construct confidence intervals for parameters of interest without additional assumptions. (In particular, it is not necessary to condition on the observed values of the predictor variables.) Resampling is done from the list of N observed cases.

For example, the left plot in Figure 14.7 is a scatter plot of waist-weight measurements for 120 physically active young adults (derived from [53]) with a least squares fitted quadratic polynomial superimposed. An estimate of the mean weight for an individual with a 33-inch waist is 174.6 pounds. If the observed (x, y) pairs are the values of a random sample from a joint distribution satisfying a linear model of the form

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \epsilon,$$

then an approximate 95% confidence interval (based on 5000 resamples) for the mean weight when the waist size is 33 inches is [169.735, 176.944].

Example: Locally weighted regression

Locally weighted regression was introduced by W. Cleveland in the 1970's. Analysis is done conditional on the observed predictor values. In the single predictor case,

$$Y_i = g(x_i) + \epsilon_i \text{ for } i = 1, 2, \dots, N$$

are assumed to be independent random variables, the function g is assumed to be a differentiable function of *unknown* form, and the collection $\{\epsilon_i\}$ is assumed to be a random sample from a distribution with mean 0 and standard deviation σ .

The goal is to estimate the conditional mean function, $y = g(x)$. Since $g(x)$ is differentiable, and a differentiable function is approximately linear on a small x -interval, the curve can be estimated as follows:

- (i) For a given value of the predictor, say x_o , estimate the tangent line to $y = g(x)$ at $x = x_o$, and use the value predicted by the tangent line to estimate $g(x_o)$.
- (ii) Repeat this process for all observed predictor values.

For a given x_o , the tangent line is estimated using a method known as weighted linear least squares. Specifically, the intercept and slope of the tangent line are obtained by minimizing the weighted sum of squared deviations

$$S(\alpha, \beta) = \sum_{i=1}^N w_i (Y_i - (\alpha + \beta x_i))^2,$$

where the weights (w_i) are chosen so that pairs with x -coordinate near x_o have weight approximately 1; pairs with x -coordinate far from x_o have weight 0; and the weights decrease smoothly from 1 to 0 in a “window” centered at x_o .

The user chooses the proportion p of data pairs that will be in the “window” centered at x_o . When the process is repeated for each observed value of the predictor, the resulting estimated curve is called the $100p\%$ *lowess smooth*.

The right plot in Figure 14.7 shows the scatter plot of waist-weight measurements for the 120 physically active young adults with a 25% lowess smooth superimposed. The smoothed curve picks up the general pattern of the relationship between waist and weight measurements.

Lowess smooths allow researchers to approximate the relationship between predictor and response without specifying the function g . A bootstrap analysis can then be done, for example, to construct confidence intervals for the mean response at a fixed value of the predictor.

For the waist-weight pairs, a 25% smooth when $x = 33$ inches produced an estimated mean weight of 175.8 pounds. A bootstrap analysis (with 5000 random resamples) produced an approximate 95% confidence interval for mean weight when the waist size is 33 inches of [168.394, 182.650].

The lowess smooth algorithm implemented above uses tricube weights for smoothing and omits Cleveland’s robustness step. For details about the algorithm, see [25, p. 121].

14.5 Laboratory problems

Laboratory problems for this chapter introduce *Mathematica* commands for linear regression analysis, permutation analysis of slope in the simple linear case, locally weighted regression, and diagnostic plots. Problems are designed to reinforce the ideas of this chapter.

14.5.1 Laboratory: Linear least squares analysis

In the main laboratory notebook (Problems 1 to 5), you will use simulation and graphics to study the components of linear least squares analyses; solve a problem on correlated and uncorrelated factors in polynomial regression; and apply linear least squares methods to three data sets from a study of sleep in mammals [3], [30].

14.5.2 Additional problem notebooks

Problems 6, 7, and 8 are applications of simple linear least squares (and other) methods. Problem 6 uses several data sets from an ecology study [32], [77]. Problem 7

uses data from an arsenic study [103]. Problem 8 uses data from a study on ozone exposure in children [113].

Problems 9, 10, and 11 are applications of multiple linear regression (and other) methods. In each case, the *adjusted* coefficient of determination is used to help choose an appropriate prediction model. Problem 9 uses data from a hydrocarbon-emissions study [90]. Problem 10 uses data from a study of factors affecting plasma beta-carotene levels in women [104]. Problem 11 uses data from a study designed to find an empirical formula for predicting body fat in men using easily measured quantities only [59].

Problem 12 applies the goodness-of-fit analysis in simple linear regression to several data sets from a physical anthropology study [50].

Problems 13 and 14 introduce the use of “dummy” variables in linear regression problems. In Problem 13, the methods are applied to a study of the relationship between age and height in two groups of children [5]. In Problem 14, the methods are applied to a study of the pricing of diamonds [26]. Problem 13 also introduces a permutation method for the same problem.

Note that the use of dummy variables in Problem 13 is an example of a *covariance analysis* and the use of dummy variables in Problem 14 is an example of the analysis of an *unbalanced* two-way layout.

Problems 15 and 16 are applications of locally weighted regression and bootstrap methods. Problem 15 uses data from a study of ozone levels in the greater Los Angeles area [28]. Problem 16 uses data from a cholesterol-reduction study [36].