

# Applied Linear Models with SAS

DANIEL ZELTERMAN

CAMBRIDGE

CAMBRIDGE

[www.cambridge.org/9780521761598](http://www.cambridge.org/9780521761598)

This page intentionally left blank

# Applied Linear Models with SAS

---

This textbook for a second course in basic statistics for undergraduates or first-year graduate students introduces linear regression models and describes other linear models including Poisson regression, logistic regression, proportional hazards regression, and nonparametric regression. Numerous examples drawn from the news and current events with an emphasis on health issues illustrate these concepts.

Assuming only a pre-calculus background, the author keeps equations to a minimum and demonstrates all computations using SAS. Most of the programs and output are displayed in a self-contained way, with an emphasis on the interpretation of the output in terms of how it relates to the motivating example. Plenty of exercises conclude every chapter. All of the datasets and SAS programs are available from the book's Web site, along with other ancillary material.

**Dr. Daniel Zelterman** is Professor of Epidemiology and Public Health in the Division of Biostatistics at Yale University. His application areas include work in genetics, HIV, and cancer. Before moving to Yale in 1995, he was on the faculty of the University of Minnesota and at the State University of New York at Albany. He is an elected Fellow of the American Statistical Association. He serves as associate editor of *Biometrics* and other statistical journals. He is the author of *Models for Discrete Data* (1999), *Advanced Log-Linear Models Using SAS* (2002), *Discrete Distributions: Application in the Health Sciences* (2004), and *Models for Discrete Data: 2nd Edition* (2006). In his spare time he plays the bassoon in orchestral groups and has backpacked hundreds of miles of the Appalachian Trail.





# Applied Linear Models with SAS

Daniel Zelterman

Yale University



**CAMBRIDGE**  
UNIVERSITY PRESS

CAMBRIDGE UNIVERSITY PRESS  
Cambridge, New York, Melbourne, Madrid, Cape Town, Singapore,  
São Paulo, Delhi, Dubai, Tokyo

Cambridge University Press  
The Edinburgh Building, Cambridge CB2 8RU, UK

Published in the United States of America by Cambridge University Press, New York

[www.cambridge.org](http://www.cambridge.org)

Information on this title: [www.cambridge.org/9780521761598](http://www.cambridge.org/9780521761598)

© Daniel Zelterman 2010

This publication is in copyright. Subject to statutory exception and to the provision of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published in print format 2010

ISBN-13 978-0-511-77109-5 eBook (Adobe Reader)

ISBN-13 978-0-521-76159-8 Hardback

Cambridge University Press has no responsibility for the persistence or accuracy of urls for external or third-party internet websites referred to in this publication, and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

# Contents

<i>Preface</i>	<i>page</i> ix
<i>Acknowledgments</i>	xiii
1 Introduction	1
1.1 What Is Statistics?	1
1.2 Statistics in the News: The Weather Map	4
1.3 Mathematical Background	6
1.4 Calculus	7
1.5 Calculus in the News: New Home Sales	9
1.6 Statistics in the News: IMF Loans and Tuberculosis	11
1.7 Exercises	13
2 Principles of Statistics	21
2.1 Binomial Distribution	21
2.2 Confidence Intervals and the Hubble Constant	25
2.3 Normal Distribution	26
2.4 Hypothesis Tests	30
2.5 The Student t-Test	34
2.6 The Chi-Squared Test and $2 \times 2$ Tables	42
2.7 What Are Degrees of Freedom?	47
2.8 SAS, in a Nutshell	49
2.9 Survey of the Rest of the Book	51
2.10 Exercises	52
3 Introduction to Linear Regression	58
3.1 Low-Birth-Weight Infants	58
3.2 The Least Squares Regression Line	59
3.3 Regression in SAS	63
3.4 Statistics in the News: Future Health Care Costs	65
3.5 Exercises	66

4	Assessing the Regression	75
4.1	Correlation	75
4.2	Statistics in the News: Correlations of the Global Economy	77
4.3	Analysis of Variance	78
4.4	Model Assumptions and Residual Plots	81
4.5	Exercises	84
5	Multiple Linear Regression	90
5.1	Introductory Example: Maximum January Temperatures	90
5.2	Graphical Displays of Multivariate Data	94
5.3	Leverage and the Hat Matrix Diagonal	96
5.4	Jackknife Diagnostics	99
5.5	Partial Regression Plots and Correlations	102
5.6	Model-Building Strategies	105
5.7	Exercises	110
6	Indicators, Interactions, and Transformations	120
6.1	Indicator Variables	120
6.2	Synergy in the News: Airline Mergers	127
6.3	Interactions of Explanatory Variables	128
6.4	Transformations	132
6.5	Additional Topics: Longitudinal Data	137
6.6	Exercises	138
7	Nonparametric Statistics	150
7.1	A Test for Medians	150
7.2	Statistics in the News: Math Achievement Scores	153
7.3	Rank Sum Test	155
7.4	Nonparametric Methods in SAS	156
7.5	Ranking and the Healthiest State	157
7.6	Nonparametric Regression: LOESS	160
7.7	Exercises	163
8	Logistic Regression	169
8.1	Example	169
8.2	The Logit Transformation	170
8.3	Logistic Regression in SAS	173
8.4	Statistics in the News: The New York Mets	177
8.5	Key Points	178
8.6	Exercises	179
9	Diagnostics for Logistic Regression	187
9.1	Some Syntax for <code>proc logistic</code>	188
9.2	Residuals for Logistic Regression	190

9.3	Influence in Logistic Regression	193
9.4	Exercises	197
10	Poisson Regression	204
10.1	Statistics in the News: Lottery Winners	204
10.2	Poisson Distribution Basics	204
10.3	Regression Models for Poisson Data	206
10.4	Statistics in the News: Attacks in Iraq	208
10.5	Poisson Regression in SAS	209
10.6	Exercises	215
11	Survival Analysis	225
11.1	Censoring	225
11.2	The Survival Curve and Its Estimate	227
11.3	The Log-Rank Test and SAS Program	232
11.4	Exercises	235
12	Proportional Hazards Regression	237
12.1	The Hazard Function	237
12.2	The Model of Proportional Hazards Regression	239
12.3	Proportional Hazards Regression in SAS	241
12.4	Exercises	243
13	Review of Methods	247
13.1	The Appropriate Method	247
13.2	Other Review Questions	249
	<i>Appendix: Statistical Tables</i>	255
A.1	Normal Distribution	255
A.2	Chi-squared Tables	257
	<i>References</i>	259
	<i>Selected Solutions and Hints</i>	263
	<i>Index</i>	269



# Preface

Linear models are a powerful and useful set of methods in a large number of settings. Very briefly, there is some outcome measurement that is very important to us and we want to explain variations in its values in terms of other measurements in the data. The heights of several trees can be explained in terms of the trees' ages, for example. It is not a straight line relationship, of course, but knowledge of a tree's age offers us a large amount of explanatory value. We might also want to take into account the effects of measurements on the amount of light, water, nutrients, and weather conditions experienced by each tree. Some of these measurements will have greater explanatory value than others and we may want to quantify the relative usefulness of these different measures. Even after we are given all of this information, some trees will appear to thrive and others will remain stunted, when all are subjected to identical conditions. This variability is the whole reason for statistics existing as a scientific discipline. We usually try to avoid the use of the word "prediction" because this assumes that there is a cause-and-effect relationship. A tree's age does not directly cause it to grow, for example, but rather, a cumulative process associated with many environmental factors results in increasing height and continued survival. The best estimate we can make is a statement about the behavior of the average tree under identical conditions.

Many of my students go on to work in the pharmaceutical or health-care industry after graduating with a masters degree. Consequently, the choice of examples has a decidedly health/medical bias. We expect our students to be useful to their employers the day they leave our program so there is not a lot of time to spend on advanced theory that is not directly applicable. Not all of the examples are from the health sciences. Diverse examples such as the number of lottery winners and temperatures in various US cities are part of our common knowledge. Such examples do not need a lengthy explanation for the reader to appreciate many of the aspects of the data being presented.

How is this book different? The mathematical content and notation are kept to an absolute minimum. To paraphrase the noted physicist Steven Hawking, who

has written extensively for the popular audience, every equation loses half of your audience. There is really no need for formulas and their derivations in a book of this type if we rely on the computer to calculate quantities of interest. Long gone are the days of doing statistics with calculators or on the back of an envelope. Students of mathematical statistics should be able to provide the derivations of the formulas but they represent a very different audience. All of the important formulas are programmed in software so there is no need for the general user to know these.

The three important skills needed by a well-educated student of applied statistics are

1. Recognize the appropriate method needed in a given setting.
2. Have the necessary computer skills to perform the analysis.
3. Be able to interpret the output and draw conclusions in terms of the original data.

This book gives examples to introduce the reader to a variety of commonly encountered settings and provides guidance through these to complete these three goals. Not all possible situations can be described, of course, but the chosen settings include a broad survey of the type of problems the student of applied statistics is likely to run into.

What do I ask of my readers? We still need to use a lot of mathematical concepts such as the connection between a linear equation and drawing the line on  $X - Y$  coordinates. There will be algebra and special functions such as square roots and logarithms. Logarithms, while we are on the subject, are always to the base  $e (=2.718)$  and not base 10.

We will also need a nodding acquaintance with the concepts of calculus. Many of us may have taken calculus in college, a long time ago, and not had much need to use it in the years since then. Perhaps we intentionally chose a course of study that avoided abstract mathematics. Even so, calculus represents an important and useful tool. The definition of the derivative of a function (What does this new function represent?) and integral (What does *this* new function represent?) are needed although we will never need to actually find a derivative or an integral. The necessary refresher to these important concepts is given in Section 1.4.

Also helpful is a previous course in statistics. The reader should be familiar with the mean and standard deviation, normal and binomial distributions, and hypothesis tests in general and the chi-squared and t-tests specifically. These important concepts are reviewed in Chapter 2 but an appreciation of these important ideas is almost a full course in itself. There is a large reliance on p-values in scientific research so it is important to know exactly what these represent.

There are a number of excellent general-purpose statistical packages available. We have chosen to illustrate our examples using SAS because of its wide acceptance and use in many industries but especially health care and pharmaceutical. Most of the examples given here are small, to emphasize interpretation and encourage practice. These datasets could be examined by most software packages. SAS, however, is

capable of handling huge datasets so the skills learned here can easily be used if and when much larger projects are encountered later.

The reader should already have some familiarity with running SAS on a computer. This would include using the editor to change the program, submitting the program, and retrieving and then printing the output. There are also popular point-and-click approaches to data analysis. While these are quick and acceptable, their ease of use comes with the price of not always being able to repeat the analysis because of the lack of a printed record of the steps that were taken. Data analysis, then, should be reproducible.

We will review some of the basics of SAS but a little hand-holding will prevent some of the agonizing frustrations that can occur when first starting out. Running the computer and, more generally, doing the exercises in this book are a very necessary part of learning statistics. Just as you cannot learn to play the piano simply by reading a book, statistical expertise, and the accompanying computer skills, can only be obtained through hours of active participation in the relevant act. Again, much like the piano, the instrument is not damaged by playing a wrong note. Nobody will laugh at you if you try something truly outlandish on the computer either. Perhaps something better will come of a new look at a familiar setting. Similarly, the reader is encouraged to look at the data and try a variety of different ways of looking, plotting, modeling, transforming, and manipulating. Unlike a mathematical problem with only one correct solution (contrary to many of our preconceived notions), there is often a lot of flexibility in the way statistics can be applied to summarize a set of data. As with yet another analogy to music, there are many ways to play the same song.



# Acknowledgments

Thanks to the many students and teaching assistants who have provided useful comments and suggestions to the exposition as well as the computer assignments. Also to Chang Yu, Steven Schwager, and Amelia Dziengeleski for their careful readings of early drafts of the manuscript. Lauren Cowles and her staff at Cambridge University Press provided innumerable improvements and links to useful Web sites.

The DASL (pronounced “dazzle”) StatLib library maintained at Carnegie Mellon University is a great resource and provided data for many examples and exercises contained here. Ed Tufte’s books on graphics have taught me to look at data more carefully. His books are highly recommended.

I am grateful to *The New York Times* for their permission to use many graphic illustrations.

Finally, thanks to my wife Linda who provided important doses of encouragement and kept me on task. This work is dedicated to her memory.

The Pennsylvania State University Department of Meteorology supplied the graphics for the weather map in Fig. 1.1.

DANIEL ZELTERMAN  
Hamden, CT  
August 25, 2009

The following figures are copyright *The New York Times* and used with permission: Figure 1.4 (June 23, 2007); Figure 1.6 (August 15, 2008); Figure 1.7 (August 4, 2008); Figure 1.8 (August 23, 2008); Figure 1.9 (January 8, 2009); Figure 1.10 (October 14, 2008); Figure 4.4 (April 17, 2008); Figure 6.1 (January 17, 2008); Figure 7.1 (June 13, 2007); Figure 10.3 (May 30, 2007). All rights are reserved by *The New York Times*.

Table 2.1: From J. P. Frisby and J. L. Clatworthy, "Learning to see complex random-dot stereograms," *Perception* 4(2), pp. 173–78. Copyright 1975 Pion Limited, London.

Figure 2.1: Courtesy of John Huchra.

Table 2.8: From Marcello Pagano and Kimberlee Gauvreau. *Principles of Biostatistics*, 2E. Copyright 2000 Brooks/Cole, a part of Cengage Learning, Inc. Reproduced by permission. [www.cengage.com/permissions](http://www.cengage.com/permissions)

Table 2.10: From N. Teasdale, C. Bard, J. Larue, et al., "On the cognitive penetrability of posture control," *Experimental Aging Research*. Copyright 1993 Taylor & Francis, reprinted by permission of the publisher (Taylor & Francis Group, <http://www.informaworld.com>).

Tables 5.1 and 6.15: From Frederick Mosteller and John W. Tukey, *Data Analysis and Regression: A Second Course in Statistics*, Table 5.1 pp. 73–74 and Table 6.15 pp. 549–51. Copyright 1977 Addison-Wesley Publishing Company, Inc. Reproduced by permission of Pearson Education, Inc.

Table 6.10: From Douglas Bates and Donald Watts, *Nonlinear Regression Analysis and Its Applications*. Copyright 2007 John Wiley and Sons, Inc. Reprinted with permission of John Wiley and Sons, Inc.

Table 6.11: From James A. Koziol and Donna A. Maxwell, "A distribution-free test for tumor-growth curve analysis with application to an animal tumor immunotherapy experiment," *Biometrics* 37, pp. 383–90. Reprinted with permission of Oxford University Press.

Table 7.1: From A. J. Dobson, *Introduction to Generalized Linear Models*, 2E. Copyright 2001 Taylor & Francis Group LLC – Books. Reproduced with permission of Taylor & Francis Group LLC – Books in the format textbook via Copyright Clearance Center.

Table 9.1: From R. G. Miller et al., *Biostatistics Casebook 318*. Copyright 1980 John Wiley and Sons, Inc. Reprinted with permission of John Wiley and Sons, Inc.

Table 10.7: From P. J. Antsaklis and K. M. Passino, *Introduction to Intelligent and Autonomous Control*, Chapter 2, pp. 27–56; Figure 2 of James S. Albus, "A Reference Model Architecture for Intelligent Systems Design." Copyright 1993 Kluwer Academic Publishers with kind permission of Springer Science and Business Media.

Table 10.11: From Pieter Joost van Watum et al., "Patterns of response to acute Naxolone infusion in Tourette's Syndrome," *Movement Disorders* 15 (2000) pp. 1252–54. Reprinted with permission of Oxford University Press.

Table 11.2: From Jennifer Wheler et al., "Survival of patients in a phase 1 clinic," *Cancer* 115 (5), pp. 1091–99. Copyright 2009 American Cancer Society.

Table 11.4: From Thomas R. Fleming et al., "Modified Kolmogorov-Smirnov test procedures with application to arbitrarily right-censored data," *Biometrics* 36 (1980), pp. 607–25. Reprinted with permission of Oxford University Press.

Table 12.2: From Nicholas Lange et al., *Case Studies in Biometry*. Copyright 1994 John Wiley and Sons, Inc. Reprinted with permission of John Wiley and Sons, Inc.

Table 12.4: From J. M. Krall, V. A. Uthoff, and J. B. Harley. "A step-up procedure for selecting variables associated with survival," *Biometrics* 31 (1975), pp. 49–57. Reprinted with permission of Oxford University Press.

# Introduction

We are surrounded by data. With a tap at a computer keyboard, we have access to more than we could possibly absorb in a lifetime. But is this data the same as information? How do we get from numbers to understanding? How do we identify simplifying trends – but also find exceptions to the rule? The computers that provide access to the data also provide the tools to answer these questions. Unfortunately, owning a hammer does not enable us to build a fine house. It takes experience using the tools, knowing when they are appropriate, and also knowing their limitations.

The study of statistics provides the tools to create understanding out of raw data. Expertise comes with experience, of course. We need equal amounts of theory (in the form of statistical tools), technical skills (at the computer), and critical analysis (identifying the limitations of various methods for each setting). A lack of one of these cannot be made up by the other two.

This chapter provides a review of statistics in general, along with the mathematical and statistical prerequisites that will be used in subsequent chapters. Even more broadly, the reader will be reminded of the larger picture. It is very easy to learn many statistical methods only to lose sight of the point of it all.

## 1.1 What Is Statistics?

In an effort to present a lot of mathematical formulas, we sometimes lose track of the central idea of the discipline. It is important to remember the big picture when we get too close to the subject.

Let us consider a vast wall that separates our lives from the place where the information resides. It is impossible to see over or around this wall, but every now and then we have the good fortune of having some pieces of data thrown over to us. On the basis of this fragmentary sampled data, we are supposed to infer the composition of the remainder on the other side. This is the aim of *statistical inference*.

The population is usually vast and infinite, whereas the sample is just a handful of numbers.

In statistical inference we infer properties of the population from the sample.

There is an enormous possibility for error, of course. If all of the left-handed people I know also have artistic ability, am I allowed to generalize this to a statement that all left-handed people are artistic? I may not know very many left-handed people. In this case I do not have much data to make my claim, and my statement should reflect a large possibility of error. Maybe most of my friends are also artists. In this case we say that the sampled data is *biased* because it contains more artists than would be found in a representative sample of the population.

The population in this example is the totality of all left-handed people. Maybe the population should be *all* people, if we also want to show that artistic ability is greater in left-handed people than in right-handed people. We can't possibly measure such a large group. Instead, we must resign ourselves to the observed or *empirical* data made up of the people we know. This is called a *convenience sample* because it is not really random and may not be representative.

Consider next the separate concepts of sample and population for numerically valued data. The sample *average* is a number that we use to infer the value of the population *mean*. The average of several numbers is itself a number that we obtain. The population mean, however, is on the other side of the imaginary wall and is not observable. In fact, the population mean is almost an unknowable quantity that could not be observed even after a lifetime of study. Fortunately, statistical inference allows us to make statements about the population mean on the basis of the sample average. Sometimes we forget that this inference is taking place and will confuse the sample statistic with the population attribute.

Statistics are functions of the sampled data. Parameters are properties of the population.

Often the sampled data comes at great expense and through personal hardship, as in the case of clinical trials of new therapies for life-threatening diseases. In a clinical trial involving cancer, for example, costs are typically many thousands of dollars per patient enrolled. Innovative therapies can easily cost ten times that amount. Sometimes the most important data consists of a single number, such as how long the patient lived, recorded only after the patient loses the fight with his or her disease.

Sometimes we attempt to collect all of the data, as in the case of a *census*. The U.S. Constitution specifically mandates that a complete census of the population be performed every ten years.<sup>1</sup> The writers of the Constitution knew that in order to

<sup>1</sup> Article 1, Section 2 reads, in part: "Representatives and direct Taxes shall be apportioned among the several States which may be included within this Union, according to their respective Numbers, which shall be

have a representative democracy and a fair tax system, we also need to know where the people live and work. The composition of the House of Representatives is based on the decennial census. Locally, communities need to know about population shifts to plan for schools and roads. Despite the importance of the census data, there continues to be controversy on how to identify and count certain segments of the population, including the homeless, prison inmates, migrant workers, college students, and foreign persons living in the country without appropriate documentation.

Statistical inference is the process of generalizing from a sample of data to the larger population. The sample average is a simple statistic that immediately comes to mind. The Student t-test is the principal method used to make inferences about the population mean on the basis of the sample average. We review this method in Section 2.5. The sample *median* is the value at which half of the sample is above and half is below. The median is discussed in Chapter 7.

The standard deviation measures how far individual observations deviate from their average.

The sample *standard deviation* allows us to estimate the scale of variability in the population. On the basis of the normal distribution (Section 2.3), we usually expect about 68% of the population to appear within one standard deviation (above or below) of the mean. Similarly, about 95% of the population should occur within two standard deviations of the population mean.

The standard error measures the sampling variability of the mean.

A commonly used measure related to the standard deviation is the *standard error*, also called the *standard error of the mean* and often abbreviated SEM. These two similar-sounding quantities refer to very different measures. The standard error estimates the standard deviation associated with the sample average. As the sample size increases, the standard deviation (which refers to individuals in the population) should not appreciably change. On the other hand, a large sample size is associated with a precise estimate of the population mean as a consequence of a small standard error. This relationship provides the incentive for larger sample sizes, allowing us to estimate the population mean more accurately. The relationship is

$$\text{Standard error} = \frac{\text{Standard deviation}}{\sqrt{\text{Sample size}}}$$

determined by adding to the whole Number of free Persons, including those bound to Service for a Term of Years, and excluding Indians not taxed, three fifths of all other Persons. The actual Enumeration shall be made within three Years after the first Meeting of the Congress of the United States, and within every subsequent Term of ten Years, in such Manner as they shall by Law direct.”

Consider a simple example. We want to measure the heights of a group of people. There will always be tall people, and there will always be short people, so changing the sample size does not appreciably alter the standard deviation of the data. Individual variations will always be observed. If we were interested in estimating the average height, then the standard error will decrease with an increase in the sample size (at a rate of  $1/\sqrt{\text{sample size}}$ ), motivating the use of ever-larger samples. The average will be measured with greater precision, and this precision is described in terms of the standard error. Similarly, if we want to measure the average with twice the precision, then we will need a sample size four times larger.

Another commonly used term associated with the standard deviation is *variance*. The relationship between the variance and the standard deviation is

$$\text{Variance} = (\text{Standard deviation})^2$$

The standard deviation and variance are obtained in SAS using `proc univariate`, for example. The formula appears often, and the reader should be familiar with it, even though its value will be calculated using a computer.

Given observed sample values  $x_1, x_2, \dots, x_n$ , we compute the *sample variance* from

$$s^2 = \text{sample variance} = \frac{1}{n-1} \sum_i (x_i - \bar{x})^2, \quad (1.1)$$

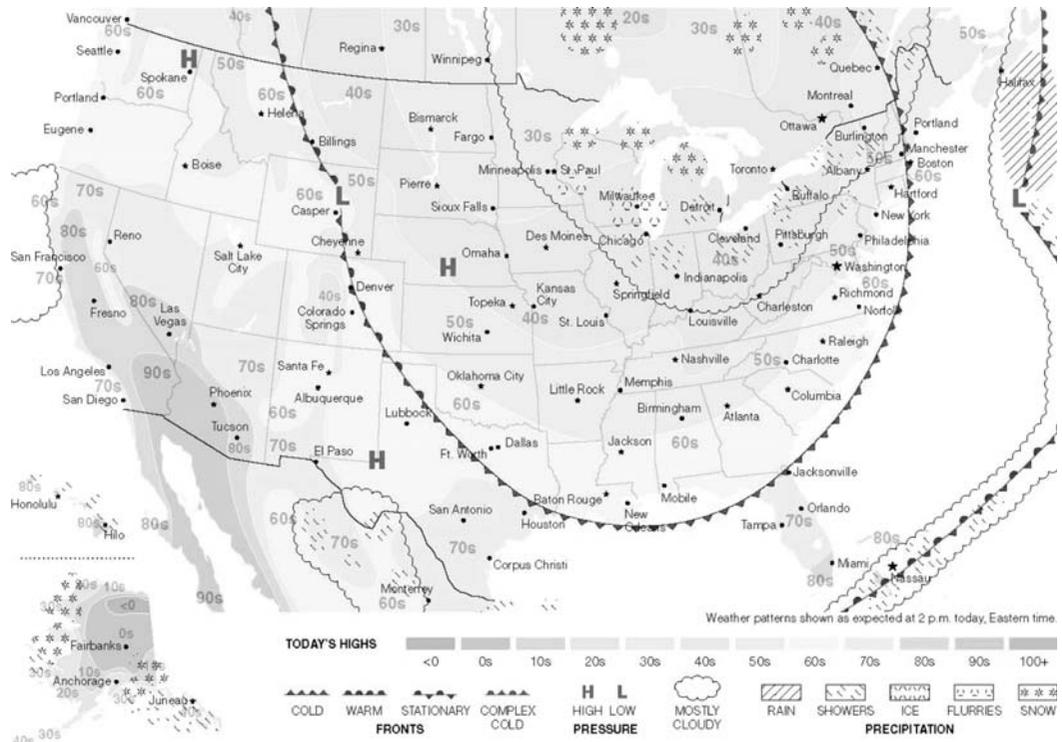
where  $\bar{x}$  is the average of the observed values.

This estimate is often denoted by the symbol  $s^2$ . Similarly, the estimated sample standard deviation  $s$  is the square root of this estimator. Intuitively, we see that (1.1) averages the squared difference between each observation and the sample average, except that the denominator is one less than the sample size. The “ $n - 1$ ” term is the degrees of freedom for this expression and is described in Sections 2.5 and 2.7.

## 1.2 Statistics in the News: The Weather Map

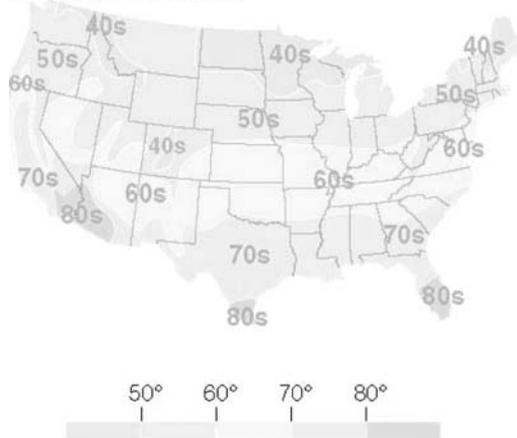
Sometimes it is possible to be overwhelmed with too much information. The business section of the newspaper is filled with stock prices, and the sports section has a wealth of scores and data on athletic endeavors. The business section frequently has several graphs and charts illustrating trends, rates, and prices. The sports pages have yet to catch up with the business section in terms of aids for the reader.

As an excellent way to summarize and display a huge amount of information, we reproduce the U.S. weather map from October 27, 2008, in Figure 1.1. There are several levels of information depicted here, all overlaid on top of one another. First we recognize the geographic-political map indicating the shorelines and state boundaries. The large map at the top provides the details of that day’s weather. The large Hs indicate the locations of high barometric pressure centers. Regions with



### Highlight: Temperature

Long-term normal highs today and tomorrow



Departure from normal highs today and tomorrow

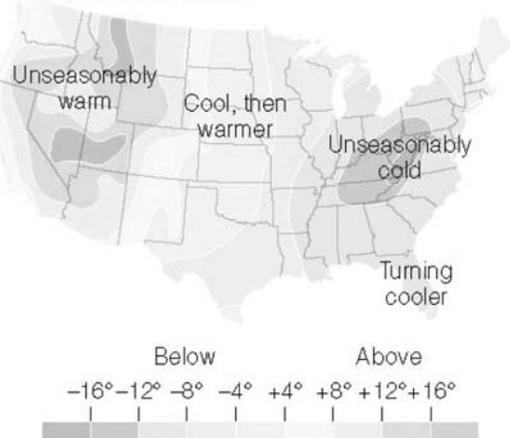


Figure 1.1 The U.S. weather map for October 27, 2008: Observed, expected, and residual data. Courtesy of Pennsylvania State University, Department of Meteorology.

similar temperatures are displayed in the same colors. The locations of rain and snow are indicated. An element of time and movement can also be inferred from this map: A large front has come across the country from the north, bringing cooler temperatures along with it. This figure represents the fine art of summarizing a huge amount of information.

The two smaller figures at the bottom provide a different kind of information. The lower-left map indicates the temperatures that we should expect to see on this date, based on previous years' experiences. The general pattern follows our preconception that southern states are warmer and northern states are cooler at this time of the year, with bands of constant temperature running east and west.

The figure on the lower right summarizes the differences between the normal pattern and the temperatures given in the large map at the top. Here we see that Florida is much cooler than what we would expect for late October. Similarly, Montana is cold at this time of year but is much warmer than typical.

The aim of statistics is to provide a similar reduction of a large amount of data into a succinct statement, generalizing, summarizing, and providing a clear message to your audience.

The goal of statistics is to start with the data and then prepare a concise summary of it.

### 1.3 Mathematical Background

---

We all need to start someplace. Let us decide on the common beginning point.

Many of us chose to study the health or social sciences and shunned engineering or physics in order to avoid the abstract rigor of mathematics. However, much of the research in the social and health fields is quantitative. We still need to demonstrate the benefit of any proposed intervention or social observation.

For example, we all know the role that the ASPCA and other animal shelters perform in protecting homeless cats and dogs. It only takes a quick visit to their local facilities to assess the effectiveness of their efforts. We can easily count the number of charges under their care to quantify and measure what they do. In this example it is easy to separate the emotional appeal from the quantity of good such an organization supplies.

In contrast, we are shocked to see the brutality of whales being slaughtered. We are told about the majesty of their huge size and life under the sea. This is all fine and plays on our emotions. Before we send money to fund the appropriate charity, or decide to enforce global bans on whaling, we also should ask how many whales there are, and perhaps how this number has changed over the past decade. This information is much harder to get at and is outside our day-to-day experiences. We need to rely on estimates to quantify the problem. Perhaps we also need to question

who is providing these estimates and whether the estimates are biased to support a certain point of view. An objective estimate of the whale population may be difficult to obtain, yet it is crucial to quantifying the problem.

As a consequence, we need to use some level of mathematics. The computer will do most of the heavy lifting for us, but we will also need to understand what is going on behind the scenes. We need to use algebra and especially linear functions. So when we write

$$y = a + bx,$$

we recall that  $a$  is referred to as the *intercept* and  $b$  is called the *slope*. We need to recognize that this equation represents a straight-line relationship and be able to graph this relationship.

We will need to use logarithms. Logarithms, or logs for short, are always to the base  $e = 2.718\dots$  and never to base 10. The exponential function written as  $e^x$  or  $\exp(x)$  is the inverse process of the logarithm. That is,

$$\log(e^x) = x$$

and

$$e^{\log x} = \exp(\log x) = x.$$

Sometimes we will use the exponential notation when the argument is not a simple expression. It is awkward to write

$$e^{a+bw+cx+dy},$$

not to mention that it is difficult to read and that publishers hate this sort of expression.

It is easier on the reader to write this last expression as

$$\exp(a + bw + cx + dy).$$

## 1.4 Calculus

For those who took calculus a long time ago and have not used it since, the memories may be distant, fuzzy, and perhaps unpleasant. Calculus represents a collection of important mathematical tools that will be needed from time to time in our discussion later on in this book. We will need to use several useful results that require calculus.

Fortunately, there is no need to dig out and dust off long-forgotten textbooks. The actual mechanics of calculus will be reviewed here, but there will not be a need to actually perform the mathematics involved. The reader who is fluent in the relevant mathematics may be able to fill in the details that we will gloss over.

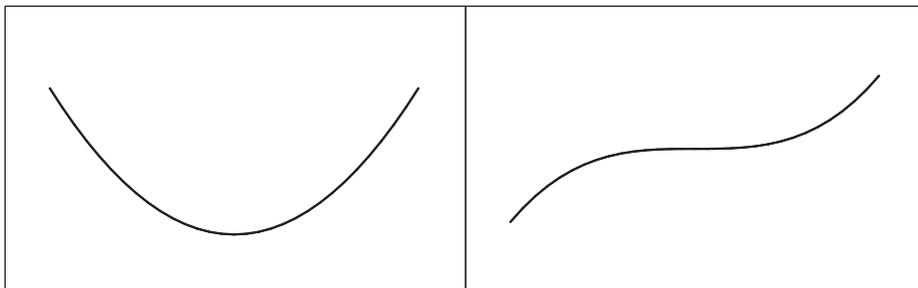


Figure 1.2 The slope is zero at the minimum of a function (left) and also at the saddle point of a function (right).

What is the point of calculus? If  $x$  and  $y$  have a straight-line relationship, we should be familiar with the concept of the *slope* of the line. When  $x$  changes by one unit, the slope is the amount of change in  $y$ .

For a nonlinear relationship, the concept of the slope remains the same, but it is a more local phenomenon. The idea of the slope depends on where in the  $x$ - $y$  relationship your interest lies. At any point in a curve, we can still talk about the slope, but we need to talk about the slope at each point of the curve. You might think of a curve as a lot of tiny linear segments all sewn together, end to end. In this case, the concept of slope is the ratio of a small change in  $y$  to the resulting small change in  $x$  at a given point on the curve. It still makes sense to talk about the ratio of these small amounts resulting in a definition of the slope of a curved line at every point  $x$ . In calculus, the *derivative* is a measure of the (local) slope at any given point in the function.

The derivative of a function provides its slope at each point.

The derivative is useful for identifying places where nonlinear functions achieve their minimums or maximums. Intuitively, we can see that a smooth function that decreases for a while and then increases has to pass through some point where the slope is zero. Solving for the places where the derivative is zero tells us where the original function is either maximized or minimized. See Figure 1.2 for an illustration of this concept.

Some functions also exhibit *saddle points* where the derivative is also zero. A saddle point is where an increasing function flattens out before resuming its increase. We will not concern ourselves with saddle points. Similarly, a zero value of the derivative may only indicate a local minimum or maximum (that is, there are either larger maximums or smaller minimums someplace else), but we will not be concerned with these topics either. A saddle point is illustrated in Figure 1.2.

Although we will not actually obtain derivatives in this book, on occasion we will need to minimize and maximize functions. When the need arises, we will recognize

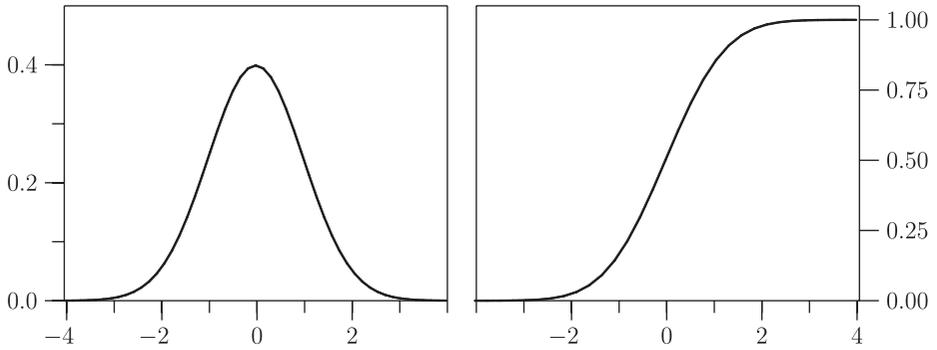


Figure 1.3 The normal density function (left) and its cumulative area (right).

the need to take a derivative and set it to zero in order to identify where the minimum occurs.

The function achieves a maximum or minimum where the derivative is zero.

Calculus is also concerned with *integrals of functions*. Briefly, an integral gives us the area between the function and the horizontal axis. As with the derivative, we will not actually need to derive one here. Many probabilities are determined according to the area under a curved function.

The integral of a function provides the area between the curve and the horizontal  $x$  axis.

Specifically, when we examine the normal distribution (Section 2.3), we will often draw the familiar bell-shaped curve. This curve is illustrated in Figure 1.3. For any value  $x$  on the horizontal axis, the curve on the right gives us the cumulative area under the left curve, up to  $x$ . The total area on the left is 1, and the cumulative area increases up to this value. The cumulative area under this curve is almost always of greater interest to us than the bell curve itself. Table A.1 in the appendix provides this area for us. It is very rare to see a table of the bell curve.

The area can be negative if the function is a negative number. Negative areas may seem unintuitive, but the example in the following section illustrates this concept.

## 1.5 Calculus in the News: New Home Sales

Home sales and building starts for new homes are both an important part of the economy. Builders will not start an expensive project unless they are reasonably sure that their investment will pay off. Home buyers will usually also purchase new furniture and carpets and will hire painters and carpenters to remodel as they move

**New-home starts plunge at fastest pace in decades**

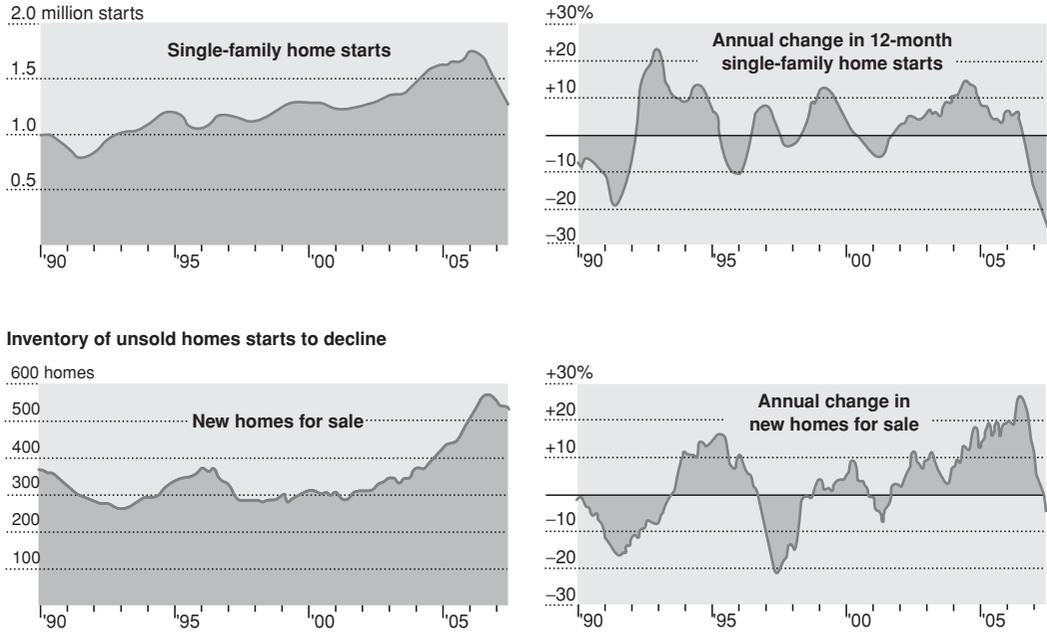


Figure 1.4 New home starts and sales. Source: New York Times.

in. Investors, economists, and government policy makers watch this data as a sign of the current state of the economy as well as future trends.

The graphs in Figure 1.4<sup>2</sup> depict new single-family home starts (upper left) and the number of new homes already on the market (lower left) over a period of a decade. There are always new homes being built and put up for sale, of course, but it is useful to know whether the trend is increasing or decreasing. The graphs on the right half of this figure show the trend more clearly in terms of the annual changes. More specifically, the graphs on the right show the slope of the line on the left at the corresponding point in time. When the figure on the left is increasing, then the figure on the right is positive. Decreasing rates on the left correspond to negative values on the right.

In words, the graphs on the right half of this figure are the derivatives of the graphs on the left half. Similarly, if we start at the values corresponding to the start of the year 1990, then the graphs on the left half are obtained by integrating the values on the right. Areas under the negative values on the right integrate to “negative areas” so that negative values on the right correspond to declining values on the left.

The times at which the derivatives on the right are zero correspond to turning points where maximums or minimums occur on the left. Remember that a zero slope is usually indicative of a change in direction. These maximums or minimums

<sup>2</sup> The graphs are available online at <http://www.nytimes.com/2007/06/23/business/23charts.html>.

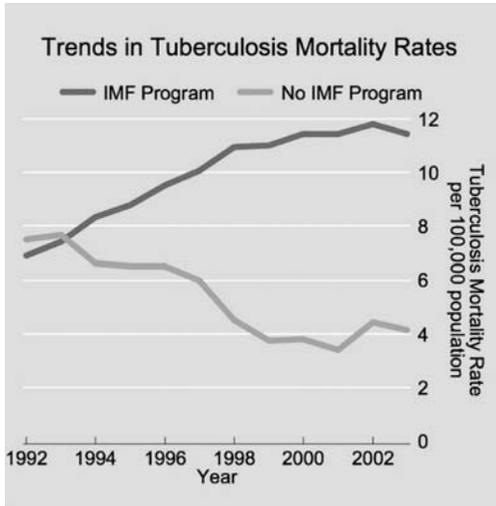


Figure 1.5 Cases of tuberculosis in Eastern Europe. *Source:* Stuckler *et al.* (2008).

may be short-lived, of course, and the underlying trend may continue after they end. The wide swings in the derivative often allow us to anticipate a change in direction of the underlying trend by a few months.

The upper headline says that there is a big decline in new home starts. But this decline is also down from the largest value in a decade. Similarly, the graph at the lower left shows a decline in the number of unsold homes on the market. Is that good for the economy because homes are selling more quickly, or bad for the economy because cautious sellers are waiting for a change in the market in order to anticipate a better price? Exercise 2.1 asks you to argue both cases: that, in the near term, the economy is improving, and also that the economy is getting worse.

## 1.6 Statistics in the News: IMF Loans and Tuberculosis

It is possible to learn all about statistics, computing, and data analysis and still come to an absurd conclusion. This sometimes leads to sensational headlines with often hilarious results as the story unfolds.

Figure 1.5 is reprinted from an article by Stuckler *et al.* (2008). A summary of the original article appeared in the *New York Times* on July 22, 2008.<sup>3</sup>

The article is about the relationship between International Monetary Fund (IMF) loans and the rate of tuberculosis (TB) in various eastern European countries. TB is an infectious disease that is spread through the air when an affected person sneezes or coughs. The disease is treated using antibiotics and is frequently fatal if left untreated. The elderly, those with diabetes, and those with a weakened immune system (such

<sup>3</sup> The original article is available online from *PLOS Medicine*: doi:10.1371/journal.pmed.0050143. A critique of this article appears at doi:10.1371/journal.pmed.0050162

as those with human immunodeficiency virus; or HIV) are at high risk for TB. We frequently see cases of TB in crowded living conditions with poor sanitation.

The IMF ([www.imf.org](http://www.imf.org)) is an international organization that oversees the global monetary system, including exchange rates and balance of payments, and sometimes making loans to foreign governments. Critics of the IMF claim that the conditions imposed on these loans will cause more harm than good to the population. These conditions have included forcing a nation to raise taxes or increase exports to the exclusion of food production needed at home. Critics will be able to point to this graph and claim that the IMF conditions result in cuts in preventative public health expenditures and reductions in the availability of necessary vaccines. Imposing these conditions on the recipient nations has resulted in crowded and unsanitary living conditions, thereby raising the incidence of TB.

In the original article, the authors performed a large number of statistical analyses that attempted to take into account differences in the various countries with respect to percent of the population living in urban settings, an index of democratization, differences in per capita incomes, whether or not the country was involved in a war, and population education levels. Most of the methods used in their article will be clear to the reader by the time we complete the material in Chapter 6.

Even so, not all countries are the same. There may be large differences between the countries that these analyses fail to correct for. Are there other factors that have not been taken into account? Could factors such as the age of the population or the rate of HIV infection result in the differences in TB rates, regardless of whether or not the nation received IMF loans?

How should we treat countries that applied for IMF loans but did not qualify? Should these be considered loan recipients? Similarly, some countries may have been offered loans but ultimately refused the money. Should these countries be considered as having received loans? What about the size of the loans: Would a large loan have the same effect as a small loan if few conditions were attached to it?

Even more importantly, this is an example of an *observational study*. Why did some countries receive loans while others did not? In what ways do these countries differ? We will never be able to know the effect on TB rates if a given country that did not receive a loan had been given one, or *vice versa*. Consider Exercise 1.1 for another possible interpretation of Figure 1.5. In an observational study, the subjects (in this case, individual countries) choose their causal treatment in some nonrandom fashion. In the present example, we do not know how countries were chosen to receive loans.

We could not randomly choose the countries that were to receive loans. A *randomized study*, in contrast to an observational study, allows us to randomly assign treatments to individuals. Differences in outcomes can then be attributed solely to the random assignment. In a medical study in which patients are randomly assigned to two different treatments, for example, any underlying imbalances in the two patient groups should be minimized by chance alone. Patients bearing a trait that is

unknown to us at the time of the randomization would be equally likely to appear in either of the two treatment groups, and then the trait would be averaged out when we examine the outcome. However, it is not possible to randomly give or withhold IMF loans to the various countries in the study.

One final comment on this example: Why was TB chosen to illustrate the effects of IMF loans? Of all the various disease rates that are reported, what is special about TB? Is it possible that the authors studied many different disease rates, but TB proved to be the most remarkable? We don't know how many diseases were compared between IMF loan and nonloan nations. We can intuit that if many comparisons were made, then it is virtually certain that some remarkable findings will be uncovered. One disease rate out of many must appear to have the largest difference between loan and nonloan nations.

This is the problem with *multiple comparisons*. If many comparisons are made, then the largest of these is not representative. We would need to make a correction for the number of different diseases that were studied. This topic and an appropriate adjustment for multiple comparisons are discussed again in Section 2.4.

There are many lessons that can be learned from this example. Ultimately, a study of statistical methods will provide you with a very useful set of tools. This book shows how these can be used to gain insight into underlying trends and patterns in your data. These tools, however, are only as good as the data that you provide. Of course, if you abuse the methods, it is possible to do more damage than good. You may be using the most sophisticated statistical methods available, but you are still responsible for the final conclusions that you draw.

Statistics is a useful tool, but it cannot think for you.

The same advice also holds for computers. The software can churn out numbers, but it cannot tell you whether the methods are appropriate or if the conditions for that method are valid. For other examples of statistics in action, consider the exercises at the end of this chapter. As with every chapter in this book, work out as many as you can.

## 1.7 Exercises

- 1.1 Argue that a country with a high rate of TB is more likely to receive an IMF loan. That is to say, use Figure 1.5 to claim that TB causes loans, rather than the other way around.
- 1.2 Do cell phones cause brain cancer? A prominent head of a cancer center sent an email to his staff urging limited use of cell phones because of a link to brain cancer. Does ownership and use of a cell phone constitute a randomized experiment or an observational study? Describe the person who is most likely to be a big user of cell phones. Is this a fair cross-section of the population?