

Introduction

There are several types of professional economists working in the world today. Academic economists in universities often derive and test theoretical models of various aspects of the economy. Economists in the civil service often study the merits and demerits of policies under consideration by government. Economists employed by a central bank often give advice on whether or not interest rates should be raised, while in the private sector, economists often predict future variables such as exchange rate movements and their effect on company exports.

For all of these economists, the ability to work with data is an important skill. To decide between competing theories, to predict the effect of policy changes, or to forecast what may happen in the future, it is necessary to appeal to facts. In economics, we are fortunate in having at our disposal an enormous amount of facts (in the form of “data”) that we can analyze in various ways to shed light on many economic issues.

The purpose of this book is to present the basics of data analysis in a simple, non-mathematical way, emphasizing graphical and verbal intuition. It focusses on the tools that economists apply in practice (primarily regression) and develops computer skills that are necessary in virtually any career path that the economics student may choose to follow.

To explain further what this book does, it is perhaps useful to begin by discussing what it does **not** do. **Econometrics** is the name given to the study of quantitative tools for analyzing economic data. The field of econometrics is based on probability and statistical theory; it is a fairly mathematical field. This book does not attempt to teach probability and statistical theory. Neither does it contain much mathematical content. In both these respects, it represents a clear departure from traditional econometrics textbooks. Yet, it aims to teach most of the practical tools used by applied econometricians today.

Books that merely teach the student which buttons to press on a computer without providing an understanding of what the computer is doing, are commonly referred to as “cookbooks”. The present book is **not** a cookbook. Some econometricians may interject at this point: “But how can a book teach the student to use the tools of econometrics, without teaching the basics of probability and statistics?” My answer is that much of what the econometrician does in practice can be understood intuitively, without resorting to probability and statistical theory. Indeed, it is a contention of this book that most of the tools econometricians use can be mastered simply through a thorough understanding of the concept of correlation, and its generalization, regression. If a student understands correlation and regression well, then he/she can understand most of what econometricians do. In the vast majority of cases, it can be argued that regression will reveal most of the information in a data set. Furthermore, correlation and regression are fairly simple concepts that can be understood through verbal intuition or graphical methods. They provide the basis of explanation for more difficult concepts, and can be used to analyze many types of economic data.

This book focusses on the **analysis** of economic data; **it is not a book about collecting economic data**. With some exceptions, it treats the data as given, and does not explain how the data is collected or constructed. For instance, it does not explain how national accounts are created or how labor surveys are designed. It simply teaches the reader to make sense out of the data that has been gathered.

Statistical theory usually proceeds from the formal definition of general concepts, followed by a discussion of how these concepts are relevant to particular examples. The present book attempts to do the opposite. That is, **it attempts to motivate general concepts through particular examples**. In some cases formal definitions are not even provided. For instance, P-values and confidence intervals are important statistical concepts, providing measures relating to the accuracy of a fitted regression line (see Chapter 5). The chapter uses examples, graphs and verbal intuition to demonstrate how they might be used in practice. But no formal definition of a P-value nor derivation of a confidence interval is ever given. This would require the introduction of probability and statistical theory, which is not necessary for using these techniques sensibly in practice. For the reader wishing to learn more about the statistical theory underlying the techniques, many books are available; for instance *Introductory Statistics for Business and Economics* by Thomas Wonnacott and Ronald Wonnacott (Fourth edition, John Wiley & Sons, 1990). For those interested in how statistical theory is applied in econometric modeling, *Undergraduate Econometrics* by R. Carter Hill, William E. Griffiths and George G. Judge (Second edition, John Wiley & Sons, 2000) provides a useful introduction.

This book reflects my belief that the use of concrete examples is the best way to teach data analysis. Appropriately, each chapter presents several examples as a means of illustrating key concepts. One risk with such a strategy is that some students might

interpret the presence of so many examples to mean that myriad concepts must be mastered before they can ever hope to become adept at the practice of econometrics. This is not the case. At the heart of this book are only a few basic concepts, and they appear repeatedly in a variety of different problems and data sets. The best approach for teaching introductory econometrics, in other words, is to illustrate its specific concepts over and over again in a variety of contexts.

Organization of the book

In organizing the book, I have attempted to adhere to the general philosophy outlined above. Each chapter covers a topic and includes a general discussion. However, most of the chapter is devoted to empirical examples that illustrate and, in some cases, introduce important concepts. Exercises, which further illustrate these concepts, are included in the text. Data required to work through the empirical examples and exercises can be found in the website which accompanies this book <http://www.wileyeurope.com/go/koopdata2ed>. By including real-world data, it is hoped that students will not only replicate the examples, but will feel comfortable extending and/or experimenting with the data in a variety of ways. Exposure to real-world data sets is essential if students are to master the conceptual material and apply the techniques covered in this book.

The empirical examples in this book are designed for use in conjunction with the computer package Excel. The website associated with this book contains Excel files. Excel is a simple and common software package. It is also one that students are likely to use in their economic careers. However, the data can be analyzed using many other computer software packages, not just Excel. Many of these packages recognize Excel files and the data sets can be imported directly into them. Alternatively, the website also contains all of the data files in ASCII text form. Appendix B at the end of the book provides more detail about the data.

Mathematical material has been kept to a minimum throughout this book. In some cases, a little bit of mathematics will provide additional intuition. For students familiar with mathematical techniques, appendices have been included at the end of some chapters. However, students can choose to omit this material without any detriment to their understanding of the basic concepts.

The content of the book breaks logically into two parts. Chapters 1–7 cover all the basic material relating to graphing, correlation and regression. A very short course would cover only this material. Chapters 8–12 emphasize time series topics and analyze some of the more sophisticated econometric models in use today. The focus on the underlying intuition behind regression means that this material should be easily accessible to students. Nevertheless, students will likely find that these latter chapters are more difficult than Chapters 1–7.

Useful background

As mentioned, this book assumes very little mathematical background beyond the pre-university level. Of particular relevance are:

1. Knowledge of simple equations. For instance, the equation of a straight line is used repeatedly in this book.
2. Knowledge of simple graphical techniques. For instance, this book is full of graphs that plot one variable against another (i.e. standard XY -graphs).
3. Familiarity with the summation operator is useful occasionally.
4. In a few cases, logarithms are used.

For the reader unfamiliar with these topics, the appendix at the end of this chapter provides a short introduction. In addition, these topics are discussed elsewhere, in many introductory mathematical textbooks.

This book also has a large computer component, and much of the computer material is explained in the text. There are myriad computer packages that could be used to implement the procedures described in this book. In the places where I talk directly about computer programs, I will use the language of spreadsheets and, particularly, that most common of spreadsheets, Excel. I do this largely because the average student is more likely to have knowledge of and access to a spreadsheet rather than a specialized statistics or econometrics package such as E-views, Stata or MicroFit.¹ I assume that the student knows the basics of Excel (or whatever computer software package he/she is using). In other words, students should understand the basics of spreadsheet terminology, be able to open data sets, cut, copy and paste data, etc. If this material is unfamiliar to the student, simple instructions can be found in Excel's on-line documentation. For computer novices (and those who simply want to learn more about the computing side of data analysis) *Computing Skills for Economists* by Guy Judge (John Wiley & Sons, 2000) is an excellent place to start.

Appendix 1.1: Mathematical concepts used in this book

This book uses very little mathematics, relying instead on intuition and graphs to develop an understanding of key concepts (including understanding how to interpret the numbers produced by computer programs such as Excel). For most students, previous study of mathematics at the pre-university level should give you all the background knowledge you need. However, here is a list of the concepts used in this book along with a brief description of each.

The equation of a straight line

Economists are often interested in the relationship between two (or more) variables. Examples of variables include house prices, gross domestic product (GDP), interest rates, etc. In our context a variable is something the economist is interested in and can collect data on. I use capital letters (e.g. Y or X) to denote variables. A very general way of denoting a relationship is through the concept of a function. A common mathematical notation for a function of X is $f(X)$. So, for instance, if the economist is interested in the factors that explain why some houses are worth more than others, he/she may think that the price of a house depends on the size of the house. In mathematical terms, he/she would then let Y denote the variable “price of the house” and X denote the variable “size of the house” and the fact that Y depends on X is written using the notation:

$$Y = f(X)$$

This notation should be read “ Y is a function of X ” and captures the idea that the value for Y depends on the value of X . There are many functions that one could use, but in this book I will usually focus on linear functions. Hence, I will not use this general “ $f(X)$ ” notation in this book.

The equation of a straight line (what was called a “linear function” above) is used throughout this book. Any straight line can be written in terms of an equation:

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

where α and β are **coefficients**, which determine a particular line. So, for instance, setting $\alpha = 1$ and $\beta = 2$ defines one particular line while $\alpha = 4$ and $\beta = -5$ defines a different line.

It is probably easiest to understand straight lines by using a graph (and it might be worthwhile for you to sketch one at this stage). In terms of an XY graph (i.e. one which measures Y on the vertical axis and X on the horizontal axis) any line can be defined by its intercept and slope. In terms of the equation of a straight line, α is the intercept and β the slope. The intercept is the value of Y when $X = 0$ (i.e. point at which the line cuts the Y -axis). The slope is a measure of how much Y changes when X is changed. Formally, it is the amount Y changes when X changes by one unit. For the student with a knowledge of calculus, the slope is the first derivative, $\frac{dY}{dX}$.

Summation notation

At several points in this book, subscripts are used to denote different observations from a variable. For instance, a labor economist might be interested in the wage of every one of 100 people in a certain industry. If the economist uses Y to denote this variable, then he/she will have a value of Y for the first individual, a value of Y for the second individual, etc. A compact notation for this is to use subscripts so that Y_1

is the wage of the first individual, Y_2 the wage of the second individual, etc. In some contexts, it is useful to speak of a generic individual and refer to this individual as the i -th. We can then write, Y_i for $i = 1, \dots, 100$ to denote the set of wages for all individuals.

With the subscript notation established, summation notation can now be introduced. In many cases we want to add up observations (e.g. when calculating an average you add up all the observations and divide by the number of observations). The Greek symbol, Σ , is the summation (or “adding up”) operator and superscripts and subscripts on Σ indicate the observations that are being added up. So, for instance,

$$\sum_{i=1}^{100} Y_i = Y_1 + Y_2 + \dots + Y_{100}$$

adds up the wages for all of the 100 individuals. As other examples,

$$\sum_{i=1}^3 Y_i$$

adds up the wages for the first 3 individuals and

$$\sum_{i=47}^{48} Y_i$$

adds up the wages for the 47th and 48th individuals.

Sometimes, where it is obvious from the context (usually when summing over all individuals), the subscript and superscript will be dropped and I will simply write:

$$\sum Y_i.$$

Logarithms

For various reasons (which are explained later on), in some cases the researcher does not work directly with a variable but with a transformed version of this variable. Many such transformations are straightforward. For instance, in comparing the incomes of different countries the variable GDP per capita is used. This is a transformed version of the variable GDP. It is obtained by dividing GDP by population.

One particularly common transformation is the logarithmic one. The logarithm (to the base B) of a number, A , is the power to which B must be raised to give A . The notation for this is: $\log_B(A)$. So, for instance, if $B = 10$ and $A = 100$ then the logarithm is 2 and we write $\log_{10}(100) = 2$. This follows since $10^2 = 100$. In economics, it is common to work with the so-called natural logarithm which has $B = e$ where $e \approx 2.71828$. We will not explain where e comes from or why this rather unusual-looking base is chosen. The natural logarithm operator is denoted by \ln ; i.e. $\ln(A) = \log_e(A)$.

In this book, you do not really have to understand the material in the previous paragraph. The key thing to note is that the natural logarithmic operator is a common one (for reasons explained later on) and it is denoted by $\ln(A)$. In practice, it can be easily calculated in a spreadsheet such as Excel (or on a calculator).

Endnote

1. I expect that most readers of this book will have access to Excel (or a similar spreadsheet or statistics software package) through their university computing labs or on their home computers (note, however, that some of the methods in this book require the Excel Analysis ToolPak add-in which is not included in some basic installations of Microsoft Works). However, computer software can be expensive and, for the student who does not have access to Excel and is financially constrained, there is an increasing number of free statistics packages designed using open source software. R. Zelig, which is available at <http://gking.harvard.edu/zelig/>, is a good example of such a package.

