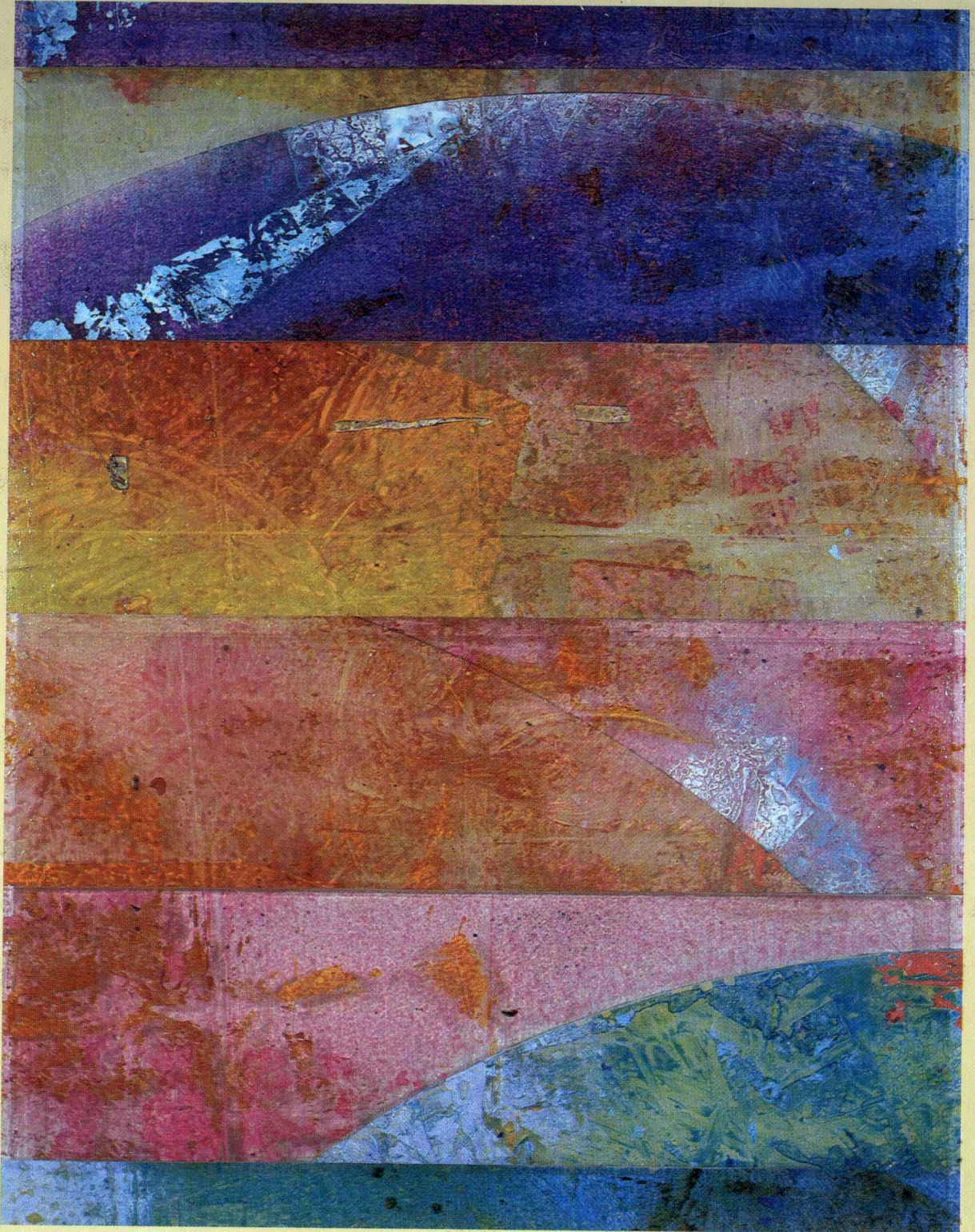SECOND EDITION

# STATISTICAL METHODS FOR THE SOCIAL SCIENCES

ALAN AGRESTI AND BARBARA FINLAY

# SECOND EDITION

# Statistical Methods for the Social Sciences

**Alan Agresti**
Department of Statistics
University of Florida

**Barbara Finlay**
Department of Sociology
Texas A & M University

*On the cover:* The cover, *Floating World*, was executed by San Francisco artist Gregg Renfrow in 1984. The work is on polymer plexiglass and measures 67¼" × 53".

*Chapter opening pages:* The chapter opening pages present a series of drawings by Sam Richardson. The drawings refer to the wedge shape as a vehicle for surface textures and sculptural elements such as pries and ropes.

# PREFACE

This book is designed as a text for introductory courses in statistics at the under-graduate or beginning graduate level. It is intended especially for students in social science disciplines. The social science orientation determined both the choice of techniques to be discussed and the selection of examples and problems. The many examples and problems included in the book are based on research studies in sociology, political science, education, geography, anthropology, psychology, his-tory, journalism, and speech. A small amount of elementary algebra is the only mathematical prerequisite for understanding the material presented and working the problems.

The book is suited for use in either a single-term or a two-term sequence. Chapters 1–9 form a basis for a one-term introductory course. If the instructor has only one term and wishes to go further than Chapter 9 or wishes to cover some material in greater depth, some sections could easily be omitted without disturbing the con-tinuity of presentation. We suggest, for example, the following sections for possible omission: 6.5, 6.7, 7.3, 7.4, 8.3–8.5, 9.4, and 9.5. Also, Chapters 7–9 and Sections 13.1–13.2 are self-contained, and the instructor could move directly into any of these after covering the fundamentals in Chapters 1–6. Four possible paths for a one-term course are as follows:

1. Chapters 1–9 (possibly omitting sections noted above): Standard cross-section of methods, including basic descriptive and inferential statistics, two-sample procedures, categorical data, and linear regression

2. Chapters 1–6, 9–12: Emphasis on bivariate and multiple regression

3. Chapters 1–7, 13: Emphasis on group comparisons and analysis of variance

4. Chapters 1–8, 10, 15: Emphasis on nonparametrics and categorical data

Chapters 10–16 are primarily concerned with regression modeling, and they could be used as the basis of a second course. These chapters can also be naturally linked to an introduction to a computer package such as SPSS$^X$, the SAS System, BMDP, or OSIRIS, using data sets and printouts provided in the text.

This edition contains some changes and additions in content, compared to the first edition. The technical level has been lowered somewhat in the first nine chapters, to make the book more easily accessible to undergraduate students. The chapters on categorical data have been reorganized, with a separate chapter given to the impor-tant concepts of multivariate relationships and the use of statistical control. A new chapter (15) has been added on advanced methods for categorical data, with emphasis on logistic regression and loglinear models—methods that have recently become very popular in social science research. Many homework problems have been added to each chapter, with particular emphasis on elementary problems.

The book is organized first by type of procedure (e.g., comparing two groups in Chapter 7, bivariate association in Chapters 8–9). For each type of procedure, the sections and chapters are also classified according to the levels of measurement of the variables analyzed. Throughout the book, we have attempted to provide understanding through intuition and example, rather than through mathematical derivations. We believe it is of primary importance that the student develop an understanding of the purposes of the various procedures and their interpretations. We feel that this is best accomplished by exposure to realistic but simple examples and to numerous homework problems. Thus, many problems appear at the end of each chapter.

In view of the increasing reliance of users of statistics on sophisticated pocket calculators, home computers, and computer software packages, we have omitted many traditional shortcut computational formulas and formulas for coded and grouped data. In later chapters, computationally complex procedures such as multiple regression and analysis of variance are taught by explaining how to interpret output from a computer package program. Since it is unlikely that students will ever do such analyses by hand, the matrix-based formulas for obtaining parameter estimates, sums of squares, and standard errors have been omitted from this book.

Unlike many other statistics textbooks, we have integrated the presentation of descriptive and inferential methods from a very early point. We believe that artificially dividing these topics into separate presentations leads to confusion for many students. Hence, we present the foundations of inference early in the text (Chapters 5 and 6). In later chapters, we present new descriptive measures together with the corresponding inferential procedures so that the student becomes familiar with a unified process of describing and making inferences for a wide variety of problems.

Some material appears in the text in end-of-chapter Notes or in paragraphs or subsections marked by asterisks. This material is optional, being of lesser importance to beginning students, although it often provides a broader coverage of topics for the interested student. Likewise, nearly every chapter contains a number of optional problems identified by asterisks. These problems either refer to the optional text material, introduce additional material related to the chapter, or pose questions of a more difficult or theoretical nature than the unstarred ones.

Finally, we have not attempted to provide a catalog of every technique for every situation. This text is meant to be primarily a teaching tool, not an encyclopedic cookbook. We do feel that we have covered the most important procedures for social science research, however, and we have included some methods that are not usually described in introductory statistics books, which are useful to social scientists, for example:

1. Procedures that are more powerful than chi-square when categories in a cross-classification table are ordered
2. Controlling for variables, and testing for causal relationships
3. Models for nonlinear relationships and for statistical interaction
4. Analysis of variance and covariance using dummy variables
5. Loglinear and logit models for categorical data
6. Introductions to path analysis, factor analysis, and LISREL

We hope that the combination of our respective fields of expertise has led to a book that is statistically sound as well as relevant to social science problems. We believe that the student who works through this book successfully will have acquired a good foundation in statistical methodology.

## Acknowledgments

# CONTENTS

**Chapter 14**

**Comparison of Several Groups While Controlling for a Covariate** **441**

**Chapter 15**

**Models for Categorical Variables** **481**

**Chapter 16**

**An Introduction to Advanced Methodology** **509**

## Tables 525

# PART I

# INTRODUCTION TO UNIVARIATE STATISTICAL ANALYSIS

# CHAPTER 1

# Introduction

## Contents

## 1.1
## Introduction to
## Statistical
## Methodology

Within the past 20 to 30 years, almost all social science disciplines have seen a rapid increase in the use of statistics. The increase is evident in the changes in the content of articles published in major journals, in the styles of textbooks, and in the increasingly common requirement of academic departments that their majors take courses in statistics. A quick glance through recent issues of *American Political Science Review, American Sociological Review, Social Forces*, or other leading social science journals is enough to convince one of the central place of statistics in social science research. Job advertisements for social scientists commonly list a knowledge of statistics as an important work tool. Almost any student preparing for a career as a social scientist must become familiar with basic statistical methodology.

The study of statistics is also useful to those who do not expect to become professional social scientists. In fact, statistics is an important part of a general education for living in today's world. Almost daily we are confronted with advertising, news reporting, political campaigning, and other communications containing statistical arguments. The study of statistics helps us to understand and evaluate these arguments.

### Science and Observation

Much of science involves collecting and organizing information about the world around us. Although different sciences have different subject matters and therefore different methods, all sciences engage in some form of information-gathering through observation. The social sciences use a wide variety of information-gathering techniques. Among these are questionnaire surveys, telephone surveys, content analysis of literary materials, planned experiments, and direct observation of behavior in natural settings. In addition, social scientists often utilize information already observed and recorded for other purposes, such as police records, census materials, and hospital files. The information gathered through such processes is collectively called *data*. These data consist of measurements taken on the various characteristics being observed. The measurements are obtained by classifying observations according to specific rules. Thus, we can measure sex by classifying individuals according to gender; we can measure age by classifying people according to the number of years since birth; and we can measure the populations of cities by obtaining census counts for their numbers of residents.

### What Is Statistics?

The general field of statistics involves methods for (a) designing and carrying out research studies; (b) describing collected data; and (c) making decisions, predictions, or inferences about phenomena represented by the data. In this book, we deal primarily with the latter two aspects. This is not to imply that the other aspects of the research process are unimportant. If a study is poorly designed or if the data are improperly collected or recorded, then the conclusions may be worthless or misleading, no matter how good a statistical analysis is performed. Methods for designing and carrying out research studies are covered in detail in textbooks on research methods (e.g., Bailey, 1982).

The descriptive aspect of statistics allows researchers to summarize large quantities of data using measures that are easily understood by an observer. It would always be possible, of course, simply to present a long list of measurements for each characteristic being observed. In a study of ages at first marriage, for example, we might just present the reader with a listing of the ages of all persons marrying for the first time within the past year in a particular county. This kind of detail, however, is not easy to assess—the reader simply gets bogged down in numbers. Instead of presenting *all* observations, we could use one of several statistical measures that would summarize the *typical* age at marriage in the collection of data. This would be much more meaningful to most people than the complete listing.

The other aspect of statistics that we will study is making decisions or inferences about characteristics by interpreting data patterns. This process often takes the form of noting whether an *expected* or *predicted* pattern of data values is actually found in the observations. We might expect, for example, that Catholics would be more likely to believe in life after death than would Unitarians. In order to decide whether this expectation is true, we might survey a number of Catholics and Unitarians, questioning them about their views concerning an afterlife. We could then use the statistical methodology described in this textbook to make comparisons between the two groups and to make a decision about whether the prediction is true for Catholics and Unitarians in general. In addition, we could estimate how great the difference is between the beliefs of the two groups.

Inherent in both the descriptive and decision-making aspects of statistics is the development of explanations for complex social processes. Much of social science research deals with sorting out relationships among such factors as political party preference, income, degree of education, religion, race, sex, and so forth. *Statistical models*, mathematical representations of such relationships, provide a mechanism for doing so. We shall see that many relatively simple statistical procedures are based on the assumption of a particular model for a real-world situation.

## 1.2 Description and Inference

A statistical procedure is usually classified as *descriptive* or *inferential*. In order to see the distinction between these purposes, it is necessary to understand the terms *population* and *sample*.

---

### Population and Sample

The *population* is the total set of individual objects or persons of interest in a study.

A *sample* is a subset of the population that is actually observed.

---

Although research is directed toward learning about populations, it is often necessary to study only samples from those populations. For example, the Gallup and