

INTRODUCTION TO



STATISTICS

Robert G.D. Steel / James H. Torrie

0212
S813
(2)

9062531

f-26
贈閱

Introduction to **STATISTICS**

Robert G. D. Steel
Professor of Experimental Statistics
School of Physical Sciences
North Carolina State University

James H. Torrie
Professor of Agronomy
College of Agriculture
University of Wisconsin, Madison



The Foundation for Books to China

美国友好书刊基金会

贈 书



McGraw-Hill Book Company

*New York St. Louis San Francisco Auckland Düsseldorf Johannesburg Kuala Lumpur
London Mexico Montreal New Delhi Panama Paris São Paulo Singapore Sydney
Tokyo Toronto*

INTRODUCTION TO STATISTICS

Copyright © 1976 by McGraw-Hill, Inc. All rights reserved.
Printed in the United States of America. No part of this publication
may be reproduced, stored in a retrieval system, or transmitted, in any
form or by any means, electronic, mechanical, photocopying, recording, or
otherwise, without the prior written permission of the publisher.

1234567890DODO79876

This book was set in Times Roman.
The editors were A. Anthony Arthur, Alice Macnow, and Shelly Levine Langman;
the designer was Scott Chelius;
the production supervisor was Angela Kardovich.
The drawings were done by J & R Services, Inc.
R. R. Donnelley & Sons Company was printer and binder.

Library of Congress Cataloging in Publication Data

Steel, Robert George Douglas, date
Introduction to statistics.
(McGraw-Hill series in probability and statistics)
Bibliography: p.
Includes index.
1. Mathematical statistics. I. Torrie, James
Hiram, date joint author. II. Title.
QA276.S79 519.5 75-22235
ISBN 0-07-060918-7

Introduction to
STATISTICS

McGraw-Hill Series in Probability and Statistics

David Blackwell and Herbert Solomon, Consulting Editors

- BHARUCHA-REID** Elements of the Theory of Markov Processes and Their Applications
- DE GROOT** Optimal Statistical Decisions
- DRAKE** Fundamentals of Applied Probability Theory
- EHRENFELD AND LITTAUER** Introduction to Statistical Methods
- GIBBONS** Nonparametric Statistical Inference
- GRAYBILL** Introduction to Linear Statistical Models
- HODGES, KRECH, AND CRUTCHFIELD** StatLab: An Empirical Introduction to Statistics
- LI** Introduction to Experimental Statistics
- MOOD, GRAYBILL, AND BOES** Introduction to the Theory of Statistics
- MORRISON** Multivariate Statistical Methods
- RAJ** The Design of Sample Surveys
- RAJ** Sampling Theory
- STEEL AND TORRIE** Introduction to Statistics
- THOMASIAN** The Structure of Probability Theory with Applications
- WADSWORTH AND BRYAN** Applications of Probability and Random Variables
- WASAN** Parametric Estimation
- WOLF** Elements of Probability and Statistics

In Memoriam
Jonathan Chester Steel
1954–1975

Preface

This is an introductory text, presenting some of the basic concepts and techniques of statistical inference. These concepts and techniques have played an important role in a long list of very diverse fields, in fact, wherever one finds observational or experimental data, whether it be counts or measurements. It is because of the chance element in the origin and nature of data that statistics has become such an important tool, since statistics depends upon probability, which is concerned with chance.

The first four chapters of the book deal primarily with the presentation and summarization of large amounts of data; they indicate the need to distinguish between samples and populations and to consider the possibilities of drawing inferences about the latter from the former. At the same time, some of the procedures given are regularly applied to small samples.

The introductory chapters are followed by probability theory. It is here that the notions of sample space, random variable, and probability function are introduced. Several frequently encountered distributions are considered at some length.

The general concepts of statistical inference, namely the estimation of parameters and testing of hypotheses, are next discussed, with special reference to and examples for the binomial and normal distributions. However, each aspect of inference once applied to a specific distribution is then looked at from a non-parametric point of view. Consequently, the students' understanding of the topic as originally introduced is reinforced.

The analysis of variance and regression analysis are included, as are nonparametric alternatives to these techniques.

At no time should one find a need for more than a working knowledge of high school algebra. Even when such algebra is used, as in finding the expectation of a binomial sum, one can generally go straight to the results. For those who enjoy algebra, there are a number of problems with algebraic content.

The book was written to provide sufficient material for a two-semester course. However, Chapters 1 to 13, with some topics deleted, together with topics selected from later chapters, perhaps 15 and 17 in particular, can serve as a one-semester

PREFACE

course. For those instructors with an interest in nonparametric statistics, Chapters 14 and 18 can substitute for 15 and 17. Section titles should help with one's choice of topics to be deleted or included.

We are indebted to the literary executor of the late Sir Ronald A. Fisher, F.R.S., to Dr. Frank Yates, F.R.S., and to Longman Group Ltd., London, for permission to reprint Table III from their book "Statistical Tables for Biological, Agricultural and Medical Research."

We are also indebted to E. S. Pearson and H. O. Hartley, editors of "Biometrika Tables for Statisticians," vol. I, and to the *Biometrika* trustees for permission to make use of Tables 8 and 18; to John W. Tukey and to Thomas E. Kurtz and Prentice-Hall, Inc., for permission to reprint Tables B-4 and B-5; to the late Frank Wilcoxon, Roberta A. Wilcox, and the American Cyanamid Company for permission to reprint Table B-6; and to Dr. James H. Goodnight for producing Table B-8 for us.

In addition, we are indebted to many authors, journal editors, and publishers for their generous permission to use the data that appear in the text and exercises. Their names are listed in the references.

Finally, we are grateful to a number of people for their contributions during the preparation of the manuscript, in particular to Mrs. Dorothy Green for her careful typing, and to Marcia and Jonathan Steel for their helpful suggestions toward a final draft.

ROBERT G. D. STEEL
JAMES H. TORRIE

Contents

Preface	xiii
1 INTRODUCTION	1
1-1 You and statistics	1
1-2 What is statistics?	2
1-3 What does a statistician do?	4
1-4 Where do statisticians work?	4
1-5 Aim of the text	5
2 COLLECTION, ORGANIZATION, AND PRESENTATION OF DATA	6
2-1 Introduction	6
2-2 Variables	7
<i>Exercise</i>	7
2-3 Frequency distributions	8
<i>Exercises</i>	14
2-4 Graphic presentation	15
<i>Exercises</i>	22
3 LOCATING A DISTRIBUTION	25
3-1 Introduction	25
3-2 Populations and samples; parameters and statistics	26
3-3 Notation and algebraic rules	27
<i>Exercises</i>	30
3-4 The arithmetic mean	31
3-5 The weighted mean	34
3-6 The median	35
<i>Exercises</i>	38
3-7 Quantiles and midrange	39
<i>Exercises</i>	41
3-8 The mode, geometric mean, and harmonic mean	41
<i>Exercises</i>	44
3-9 Selecting a measure of central tendency	44

CONTENTS

4	MEASURING VARIATION	46
	4-1 Introduction	46
	4-2 Measuring spread	47
	4-3 The standard deviation and the variance	47
	<i>Exercises</i>	57
	4-4 The range	58
	<i>Exercises</i>	59
	4-5 Other measures of variation	59
	<i>Exercises</i>	61
5	PROBABILITY	62
	5-1 Introduction	62
	5-2 Elementary probability	63
	5-3 Random experiments and sample spaces	65
	<i>Exercises</i>	67
	5-4 Probabilities	68
	<i>Exercises</i>	70
	5-5 Unions and intersections	71
	<i>Exercises</i>	76
6	INDEPENDENT AND DEPENDENT EVENTS	80
	6-1 Introduction	80
	6-2 Independent events	81
	<i>Exercises</i>	82
	6-3 Dependent events and conditional probabilities	83
	<i>Exercises</i>	85
7	PERMUTATIONS AND COMBINATIONS	88
	7-1 Introduction	88
	7-2 Permutations	88
	<i>Exercises</i>	91
	7-3 Combinations	92
	<i>Exercises</i>	93
8	RANDOM VARIABLES AND PROBABILITY DISTRIBUTIONS	95
	8-1 Introduction	95
	8-2 Random variables and probability functions for discrete variables	96
	<i>Exercises</i>	99
	8-3 The rectangular or uniform distribution	100
	<i>Exercises</i>	105

CONTENTS

9 THE BINOMIAL AND OTHER DISCRETE DISTRIBUTIONS	106
9-1 Introduction	106
9-2 The binomial distribution	107
<i>Exercises</i>	111
9-3 The hypergeometric distribution	113
<i>Exercises</i>	115
9-4 The Poisson distribution	116
<i>Exercises</i>	119
9-5 Other discrete distributions	120
10 THE NORMAL DISTRIBUTION	121
10-1 Introduction	121
10-2 The normal distribution	122
10-3 The normal distribution with zero mean and unit variance	124
<i>Exercises</i>	128
10-4 The normal distribution with any mean and variance	129
<i>Exercises</i>	132
11 EXPECTATIONS FOR DISCRETE DISTRIBUTIONS	134
11-1 Introduction	134
11-2 Expectations of random variables	135
<i>Exercises</i>	137
11-3 Expectations and parameters	138
<i>Exercises</i>	139
11-4 The expected value of a binomial sum	141
<i>Exercises</i>	143
11-5 The population variance	143
<i>Exercises</i>	146
12 THE NORMAL DISTRIBUTION IN STATISTICS	148
12-1 Introduction	148
12-2 The mean and variance of a normal random variable	149
<i>Exercises</i>	151
12-3 The central limit theorem	152
<i>Exercises</i>	155
12-4 The normal approximation to the binomial distribution	155
<i>Exercises</i>	159
12-5 Sampling distributions	160
13 ESTIMATION OF PARAMETERS: KNOWN DISTRIBUTIONS	162
13-1 Introduction	162
13-2 Estimation and inference	163
<i>Exercise</i>	165

CONTENTS

13-3 Properties of estimators	165
13-4 Estimation of proportions	167
<i>Exercises</i>	171
13-5 Estimation of means in the normal distribution	172
<i>Exercises</i>	176
13-6 Sampling from small populations	176
13-7 Estimation of sample size	177
<i>Exercises</i>	179
13-8 Estimation of σ^2 in the normal population	180
<i>Exercises</i>	181
14 NONPARAMETRIC ESTIMATION OF LOCATION	182
14-1 Introduction	182
14-2 Estimation of the median; sign-test procedure	183
<i>Exercises</i>	187
14-3 Estimation of median: signed-rank-sum test procedure	188
<i>Exercises</i>	190
15 TESTS OF HYPOTHESES I: BINOMIAL DISTRIBUTION	192
15-1 Introduction	192
<i>Exercises</i>	193
15-2 Hypothesis testing; discrete data	194
<i>Exercises</i>	205
15-3 Hypothesis testing—Illustrations	206
<i>Exercises</i>	209
15-4 The normal approximation for hypothesis testing	210
<i>Exercises</i>	211
15-5 The χ^2 test criterion	212
<i>Exercises</i>	214
16 TESTS OF HYPOTHESES II: DISCRETE DATA IN TWO-WAY TABLES	216
16-1 Introduction	216
16-2 Testing the equality of binomial parameters: independent samples	217
<i>Exercises</i>	222
16-3 The completely random 2×2 table	225
<i>Exercises</i>	227
16-4 Testing the equality of binomial parameters: related samples	229
<i>Exercises</i>	231
16-5 The $r \times c$ contingency table	233
<i>Exercises</i>	234
17 TESTS OF HYPOTHESES III: NORMAL-DISTRIBUTION PARAMETERS	237
17-1 Introduction	237
17-2 Tests of a single mean	238
<i>Exercises</i>	241

CONTENTS

17-3 Comparison of two sample means: general	242
17-4 Comparison of two sample means: independent samples with equal variance	243
<i>Exercises</i>	247
17-5 Comparison of two sample means: paired observations	248
<i>Exercises</i>	251
17-6 The problem of independent samples with unequal variances	252
17-7 Equality of two variances	252
<i>Exercise</i>	254
18 NONPARAMETRIC TEST PROCEDURES	255
18-1 Introduction	255
18-2 The sign test with a single sample	256
<i>Exercises</i>	258
18-3 The Wilcoxon signed-rank test, with a single sample	259
<i>Exercises</i>	261
18-4 The sign test: comparison of two independent samples	261
<i>Exercises</i>	264
18-5 The Wilcoxon rank-sum test for comparing two populations	264
<i>Exercises</i>	266
18-6 The sign test: meaningfully paired observations	267
<i>Exercises</i>	268
18-7 The Wilcoxon signed-rank test: paired observations	269
<i>Exercises</i>	270
19 THE ANALYSIS OF VARIANCE	271
19-1 Introduction	271
19-2 The completely random design or one-way classification	272
19-3 Analysis of variance: completely random design	276
<i>Exercises</i>	284
19-4 The randomized block design or two-way classification	285
<i>Exercises</i>	288
19-5 The randomized block model	290
<i>Exercises</i>	295
19-6 Error-rate definitions and multiple comparisons	295
19-7 Generalizations of the completely random design	297
<i>Exercise</i>	298
19-8 Generalizations of the randomized block design	299
20 NONPARAMETRIC DATA ANALYSIS	300
20-1 Introduction	300
20-2 The sign test for the one-way classification	300
<i>Exercises</i>	304
20-3 Kruskal-Wallis t -sample test: one-way classification	304
<i>Exercises</i>	306

CONTENTS

20-4 The sign test for the two-way classification	307
<i>Exercises</i>	310
20-5 Friedman's rank test: randomized blocks	310
<i>Exercises</i>	311
20-6 Randomization tests	312
21 LINEAR REGRESSION	316
21-1 Introduction	316
21-2 Functional and statistical relations	317
21-3 The linear regression of Y on x	317
<i>Exercises</i>	323
21-4 Regression models	324
<i>Exercises</i>	325
21-5 Sources of variation in regression	325
<i>Exercise</i>	328
21-6 Confidence intervals and tests of hypotheses	328
<i>Exercises</i>	333
21-7 Confidence interval versus confidence band	333
21-8 Correlation and regression	334
<i>Exercises</i>	336
21-9 Nonparametric techniques for regression and correlation	336
<i>Exercises</i>	339
21-10 Multiple linear regression	339
APPENDIXES	340
A 1 Greek alphabet	341
2 Symbols and notation	342
B 1 Probability of a random value of $Z = (X - \mu)/\sigma$ being greater than the values tabulated in the margins	345
2 Values of Student's t	346
3 Values of χ^2	347
4 Values of $c_{\alpha,n}$, critical values for the sign-test count	348
5 Values of $s_{(\alpha,n)}$, critical values for the signed rank-sum test	348
6 Critical values for Wilcoxon's rank-sum test	349
7 Values of F	350
8 Suggested critical values of Spearman's r_s	354
References	355
Answers to Selected Problems	359
Index	375

Chapter 1

INTRODUCTION

1-1 YOU AND STATISTICS

Statistics plays a part in our daily lives in numerous ways—ways we may not always appreciate. This is a relatively new phenomenon, and is a consequence of the many problems that science and technology have chosen to study and deal with using statistical methods. Notice that the word “statistics” is used in the singular. The term no longer refers simply to the collection and presentation of data in tables, graphs, and descriptive figures. Today, statistics is a subject or discipline. As such, it is very much involved in the physical and social sciences, in business and industry, and indeed in most human activities that require decisions involving uncertain outcomes.

Statistics is used extensively in the production of many of our foodstuffs. The flour you use is most likely the result of breeding experiments with wheat, resulting in the selection of varieties or cultivars possessing high quality, and of baking tests requiring careful measurements, at various stages, of the appropriate

INTRODUCTION

properties. Both the breeding techniques and the baking tests are based on sound statistical procedures. The drinking water we use so casually for so many purposes is almost certainly under continuous statistical quality control. For much of our clothing, rigorous statistically analyzed trials have determined what mixtures of fibers and what insulating materials provide a reasonable compromise between maximum comfort and minimum care.

The paint on your house may have been selected for the market as a consequence of a weathering trial conducted according to statistical guidelines at a forest products or industrial chemistry research laboratory. The building code of your community may have been changed recently on the basis of a statistical comparison of a long-used building material with one developed since the code was adopted.

Consumer protective agencies evaluate statistical evidence that may lead to the removal from the market of clothing that could prove to be a fire hazard, food additives that could cause genetic damage, and drugs with undesirable side-effects.

Unfortunately, statistics in the form of numerical "facts" is sometimes used to confound us. For example, we may be told that a gasoline with a special additive has given up to 3 miles per gallon more than the average of a number of leading brands without the additive. The advertisers have not said that their average mileage is any better than the average given by the leading brands, though they almost certainly hope we will conclude so. Sometimes our would-be confounders compare wage rates or crime rates where definitions have changed over time or place. They may compare sales for consecutive time periods when the only significant comparison would involve comparing the same period in consecutive years, and they may quote percentages without saying whether the figures were based on a small or a large number of observations.

As a consequence of the ever-present role that statistics plays in our lives, it is a virtual requirement that we attain some degree of statistical literacy. Without such knowledge, we will be unable to interpret satisfactorily a reasonable number of the statements and situations with statistical content with which we come in contact.

1-2 WHAT IS STATISTICS?

Already we have implied that statistics is concerned with knowledge in the form of numbers relevant to various questions. The numbers initially collected are organized, processed arithmetically, interpreted, and finally presented in tabular or graphic form or as summary numbers such as averages or percentages. Usually interpretation involves a conclusion or inference.

The originally collected numbers or *observations* are called, collectively, *data*. Clearly they are almost invariably a fraction or *sample* of a larger set of possible observations. For example, only a small amount of flour will be used to provide an observation or *sample value* in a baking test. The same is true in measuring water quality, house-paint weathering, or gasoline mileage.

The set of observations which constitute all possible measurements is called a *population*.

In the cases mentioned, it is easy to visualize the population as consisting of more observations than we can number, that is, as an *infinite population*. However, if we are sampling United States income tax returns for 1974, say for the purpose of auditing them, we can number them all in sequence and so have a *finite population*. Note that the term population applies to the data rather than to the people or objects which are observed to provide the data.

If a census or complete enumeration of the observations in a finite population exists, one can find the true average for that population. This is a constant and not subject to any variation. Thus, if we have complete census data on ages in the American population, the computed average age will be the true average for the time of the census. In the case of an infinite population, we cannot have the numbers to do the arithmetic of finding a true average, but we can conceive of such a value which, like that for the finite population, is a constant.

Notice that for a population to exist at all, it must first be capable of definition. Definition of any population is a matter for careful thought. It involves providing a rule that will identify every unit on which an observation might be made. For example, what will be our rule to identify "farmer" for a study of farm incomes in a county? How do we identify each and every dwelling unit in a city block for a study of housing density? Answers, of course, depend on the aims of the investigation.

Any inference drawn from a sample must, of course, be one concerning the population sampled. Thus, if we want a sample of registered voters, we go to voter rolls rather than to telephone listings. Since we are reasoning from the part or sample to the whole or population, the end result is an *inductive* or *uncertain inference*. Because of sample-to-sample variation, even within the same population, it is evident that our inferences cannot be completely dependable; they must be qualified in some manner. To make sure that most of our inferences about descriptive measures like averages will be valid, we must have good samples, samples resulting from sound sampling schemes and experimental designs.

What is a good sample? We would certainly like any sample to be typical or representative of the population. However, there is no way to determine what typifies a population unless we can examine all of it, a course of action not available to us. Consequently, we need a technique to tell us, at least, how well our sampling and inferential procedures work. When the laws of probability theory apply, this theory provides such a technique.

Random sampling is done when each possible sample has the same probability of being selected or when its probability of being selected is known. As a consequence, the laws of probability can be applied.