

MATHEMATICAL STATISTICS:

Basic Ideas and Selected Topics

MATHEMATICAL STATISTICS

Basic Ideas and Selected Topics

PETER J. BICKEL

*University of California,
Berkeley*

KJELL A. DOKSUM

*University of California,
Berkeley*

HOLDEN-DAY, INC.

San Francisco

*Düsseldorf Johannesburg London
Panama Singapore Sydney*

MATHEMATICAL STATISTICS

Copyright © 1977 by Holden-Day, Inc.,
500 Sansome Street, San Francisco, Ca.
All rights reserved. No part of this
book may be reproduced by mimeograph,
or any other means, without the permission
in writing of the publisher.

Library of Congress Catalog Card Number: 76-8724
ISBN: 0-8162-0784-4

Printed in the United States of America

1234567890 09876

PREFACE

This book presents our view of what an introduction to mathematical statistics for students with a good mathematics background should be. By a good mathematics background we mean linear algebra and matrix theory, and advanced calculus (but no measure theory). Since the book is an introduction to statistics we need probability theory and expect readers to have had a course at the level of, for instance, P. Hoel, S. Port and C. Stone's *Introduction to Probability Theory*. Our appendix does give all the probability that is needed. However, the treatment is abridged with few proofs and no examples or problems.

We feel such an introduction should at least do the following:

- (1) Describe the basic concepts of mathematical statistics indicating the relation of theory to practice.

- (2) Give careful proofs of the major "elementary" results such as the Neyman-Pearson lemma, the Lehmann-Scheffé theorem, the information inequality and the Gauss-Markoff theorem.

- (3) Give heuristic discussions of more advanced results such as the large sample theory of maximum likelihood estimates, and the structure of both Bayes and admissible solutions in decision theory. The extent to which holes in the discussion can be patched and where patches can be found should be clearly indicated.

(4) Show how the ideas and results apply in a variety of important subfields such as Gaussian linear models, multinomial models, and nonparametric models.

Although there are several good books available for this purpose we feel that none has quite the mix of coverage and depth desirable at this level. The work of Rao, *Linear Statistical Inference and Its Applications*, 2nd ed., covers most of the material we do and much more but at a more abstract level employing measure theory. At the other end of the scale of difficulty for books at this level is the work of Hogg and Craig, *Introduction to Mathematical Statistics*, 3rd ed. These authors also discuss most of the topics we deal with but in many instances do not include detailed discussion of topics we consider essential such as existence and computation of procedures and large sample behavior.

Our book contains more material than can be covered in two quarters. In the two quarter courses for graduate students in mathematics, statistics, the physical sciences and engineering that we have taught we cover the core Chapters 2 to 7 which go from modelling through estimation and testing to linear models. In addition we feel Chapter 10 on decision theory is essential and cover at least the first two sections. Finally we select topics from Chapter 8 on discrete data and Chapter 9 on nonparametric models.

Chapter 1 covers probability theory rather than statistics. Much of this material unfortunately does not appear in basic probability texts but we need to draw on it for the rest of the book. It may be integrated with the material of Chapters 2–7 as the course proceeds rather than being given at the start; or it may be included at the end of an introductory probability course which precedes the statistics course.

A special feature of the book is its many problems. They range from trivial numerical exercises and elementary problems intended to familiarize the students with the concepts to material more difficult than that worked out in the text. They are included both as a check on the student's mastery of the material and as pointers to the wealth of ideas and results that for obvious reasons of space could not be put into the body of the text.

Conventions: (i) In order to minimize the number of footnotes we have added a section of comments at the end of each chapter preceding the problem section. These comments are ordered by the section to which they pertain. Within each section of the text the presence of comments at the end of the chapter is signalled by one or more numbers, 1 for the first, 2 for the second, etc. The comments contain digressions, reservations and additional references. They need to be read only as the reader's curiosity is piqued.

(ii) Various notational conventions and abbreviations are used in the text. A list of the most frequently occurring ones indicating where they are introduced is given at the end of the text.

(iii) Basic notation for probabilistic objects such as random variables and vectors, densities, distribution functions and moments is established in the Appendix.

We would like to acknowledge our indebtedness to colleagues, students, and friends who helped us during the various stages (notes, preliminary edition, final draft), through which this book passed. E. L. Lehmann's wise advice has played a decisive role at many points. R. Pyke's careful reading of a next-to-final version caught a number of infelicities of style and content. Many careless mistakes and typographical errors in an earlier version were caught by D. Minassian who sent us an exhaustive and helpful listing. W. Carmichael, in proofreading the final version, caught more mistakes than both authors together. A serious error in Problem 2.2.5 was discovered by F. Scholz. Among many others who helped in the same way we would like to mention C. Chen, S. J. Chou, G. Drew, C. Gray, U. Gupta, P. X. Quang, and A. Samulon. Without Winston Chow's lovely plots Section 9.6 would probably not have been written and without Julia Rubalcava's impeccable typing and tolerance this text would never have seen the light of day.

We would also like to thank the colleagues and friends who inspired and helped us to enter the field of statistics. The foundation of our statistical knowledge was obtained in the lucid, enthusiastic and stimulating lectures of Joe Hodges and Chuck Bell, respectively. Later, we were both very much influenced by Erich Lehmann whose ideas are strongly reflected in this book.

Last and most important we would like to thank our wives Nancy Kramer Bickel and Maria Delasalas Doksum and our families for support, encouragement, and active participation in an enterprise which at times seemed endless.

Peter J. Bickel
Kjell Doksum

Berkeley
1976

CONTENTS

1	SOME TOPICS IN PROBABILITY	1
1.1	Conditioning by a random variable or vector, 1	
	A The discrete case	
	B Conditional expectation for discrete variables	
	C Continuous variables	
	D Comments on the general case	
1.2	Distribution theory for transformations of random vectors, 9	
1.3	Distribution theory for samples from a normal population, 15	
	A The χ^2 , F and t distributions	
	B Orthogonal Transformations	
1.4	The bivariate normal distribution, 22	
1.5	Approximations to distributions and moments, 28	
	A Some examples	
	B Approximations to moments	
	C Variance stabilizing transforms	
	D Edgeworth and other approximations	
1.6	Prediction, 34	
1.7	Comments, 41	
1.8	Problems and complements, 42	
1.9	References, 55	
2	STATISTICAL MODELS	56
2.1	Formulations of statistical models, 56	
2.2	Sufficiency, 63	
2.3	Exponential families, 67	
	A The one parameter case	
	B The k parameter case	
2.4	Bayesian models, 73	
2.5	Comments, 79	
2.6	Problems and complements, 80	
2.7	References, 87	
3	METHODS OF ESTIMATION	89
3.1	Substitution principles, 90	
	A Frequency substitution	
	B Method of moments	

- 3.2 The method of least squares, 94
 - A General and linear regression models
 - B Weighted least squares
- 3.3 Maximum likelihood estimates, 99
 - A One parameter families
 - B Maximum likelihood in multiparameter families
 - C Maximum likelihood and other methods
- 3.4 Comments, 107
- 3.5 Problems and complements, 108
- 3.6 References, 114

4 COMPARISON OF ESTIMATES—OPTIMALITY THEORY

116

- 4.1 Criteria of estimation, 116
- 4.2 Uniformly minimum variance unbiased estimates, 120
- 4.3 The information inequality, 126
- 4.4 Large sample theory, 132
 - A Consistency
 - B Asymptotic normality and related properties
 - C Asymptotic efficiency and optimality
- 4.5 Unbiased and maximum likelihood estimates. A comparison, 141
- 4.6 Comments, 142
- 4.7 Problems and complements, 142
- 4.8 References, 151

5 FROM ESTIMATION TO CONFIDENCE INTERVALS AND TESTING

153

- 5.1 Precision, confidence intervals, and bounds, 153
 - A The one dimensional case
 - B Confidence regions of higher dimension
 - C Other concepts of confidence regions
- 5.2 The elements of hypothesis testing, 163
 - A Introduction and the Neyman-Pearson framework
 - B The p value: the test statistic as evidence
 - C Power and sample size: indifference regions
- 5.3 Confidence procedures and hypothesis testing, 177
 - A The duality between tests and confidence regions
 - B Confidence intervals and power
 - C Applications of confidence intervals to comparisons and selections
- 5.4 Comments, 184
- 5.5 Problems and complements, 184
- 5.6 References, 191

**6 OPTIMAL TESTS AND CONFIDENCE INTERVALS:
LIKELIHOOD RATIO TESTS AND RELATED PROCEDURES**

192

- 6.1 The Neyman-Pearson lemma, 192
- 6.2 Uniformly most powerful tests, 198
- 6.3 Uniformly most accurate confidence bounds, 206
- 6.4 Likelihood ratio and related procedures, 209
 - A Tests for the mean of a normal distribution—matched pair experiments
 - B Tests and confidence intervals for the difference in means of two normal populations
 - C The two-sample problem with unequal variances
- 6.5 Likelihood ratio procedures for bivariate normal distributions, 219
 - A Testing independence, confidence intervals for p
 - B Tests for the bivariate mean vector
- 6.6 Large sample approximations in testing, 225
 - A Approximations to the distribution of test statistics under H
 - B Consistency and local power
- 6.7 Comments, 232
- 6.8 Problems and complements, 233
- 6.9 References, 247

7 LINEAR MODELS—REGRESSION AND ANALYSIS OF VARIANCE

- 7.1 Introduction to the general linear model, 248
 - A Some examples of linear models
 - B Statement and assumptions of the general linear model
 - C What does assuming a linear model mean?
 - D Matrix formulation of the linear model
 - E Related models
- 7.2 Estimation in linear models, 260
 - A The canonical form
 - B Estimation of linear functions of the means: relations to least squares and unbiasedness theory
 - C The variance of linear least squares estimates: the Gauss-Markoff theorem
 - D Estimation of the error variance
 - E Distribution theory: confidence intervals
- 7.3 Tests in linear models, 273
 - A General theory
 - B Linear regression
 - C Analysis of variance models
- 7.4 Simultaneous confidence intervals and multiple comparisons, 288
 - A The Tukey methods
 - B The Scheffé method

7.5	Comments, 296	
7.6	Problems and complements, 297	
7.7	References, 311	
8	ANALYSIS OF DISCRETE DATA	312
8.1	Goodness of fit to a single hypothesis, 312	
8.2	Goodness of fit to families of distributions: contingency tables, 317	
8.3	The p sample model and "regression" for binomial variables, 325	
	A The p sample model	
	B "Regression" (logit) model	
8.4	Comments, 332	
8.5	Problems and complements, 333	
8.6	References, 343	
9	NONPARAMETRIC MODELS	344
9.1	Rank methods for comparing two populations, 345	
	A The Wilcoxon statistic	
	B Confidence intervals and estimates for comparing two popula- tions	
	C Rank methods for tied observations	
9.2	The sign and Wilcoxon signed rank tests, 357	
	A The sign test	
	B The Wilcoxon signed rank test	
9.3	Rank tests for the one-way layout, 363	
9.4	Linear regression and independence, 365	
	A Linear regression	
	B Tests for independence	
9.5	Robust estimates and related procedures, 369	
9.6	Goodness of fit and model selection, 378	
	A The Kolmogorov test	
	B Studying distributional shape	
	C Testing goodness of fit to the normal shape	
	D A question	
9.7	Comments, 389	
9.8	Problems and complements, 390	
9.9	References, 404	
10	DECISION THEORY	407
10.1	The elements of decision theory, 409	
10.2	Comparison of decision procedures, 412	
10.3	Computation of Bayes procedures, 419	

- 10.4 Computing minimax procedures and establishing admissibility, 424
- 10.5 Comments, 429
- 10.6 Problems and complements, 430
- 10.7 References, 435

APPENDIX A REVIEW OF BASIC PROBABILITY THEORY

437

- A.1 The basic model, 437
- A.2 Elementary properties of probability models, 439
- A.3 Discrete probability models, 439
- A.4 Conditional probability and independence, 440
- A.5 Compound experiments, 441
- A.6 Binomial trials, sampling with and without replacement, 443
- A.7 Probabilities on Euclidean space, 444
- A.8 Random variables and vectors: transformations, 446
- A.9 Independence of random variables and vectors, 449
- A.10 The expectation of a random variable, 450
- A.11 Moments, 451
- A.12 Moment generating functions, 454
- A.13 Some classical discrete and continuous distributions, 455
- A.14 Modes of convergence of random variables and limit theorems, 460
- A.15 Further limit theorems, 462
- A.16 Poisson Process, 466
- A.17 References, 467

TABLES

469

- I Area under the normal curve, 469
- II(a) χ^2 upper tail probabilities for $k = 2, 3, 4, 5$ degrees of freedom, 470
- II(b) Quantiles $x(1 - \alpha)$ of χ^2 with k degrees of freedom, 471
- III Quantiles $t(1 - \alpha)$ of the t distribution, 472
- IV Quantiles $f(1 - \alpha)$ of the F distribution, 473
- V Wilcoxon hypothesis distribution, 476
- VI A short table of the $\mathcal{B}(n, \frac{1}{2})$ distribution function, 478
- VII Wilcoxon signed rank distribution, 479
- VIII Distribution of Spearman's statistic, 482
- IX Critical values k_α of Kolmogorov's statistic, 483

ACKNOWLEDGEMENTS AND SOURCES FOR TABLES AND FIGURES, 484

AUTHOR INDEX, 485

SUBJECT INDEX, 487

NOTATIONAL INDEX, 493

CHAPTER 1

SOME TOPICS IN PROBABILITY

In this chapter we will give some results in probability theory, which are essential in our treatment of statistics and which may not be treated in enough detail in some probability texts.

Measure theory will not be used. We make the blanket assumption that all sets and functions considered are measurable.

1.1. CONDITIONING BY A RANDOM VARIABLE OR VECTOR

The concept of conditioning is important in studying associations between random variables or vectors. In this section we present some results useful for prediction theory, estimation theory, and regression.

1.1.A. The Discrete Case

The reader is already familiar with the notion of the conditional probability of an event A given that another event B has occurred. If X and Y are discrete random

vectors possibly of different dimensions we want to study the conditional probability structure of \mathbf{X} given that \mathbf{Y} has taken on a particular value \mathbf{y} .

Define the *conditional frequency function* $p(\cdot | \mathbf{y})$ of \mathbf{X} given $\mathbf{Y} = \mathbf{y}$ by,

$$(1.1.1) \quad p(\mathbf{x} | \mathbf{y}) = P[\mathbf{X} = \mathbf{x} | \mathbf{Y} = \mathbf{y}] = \frac{p(\mathbf{x}, \mathbf{y})}{p_{\mathbf{Y}}(\mathbf{y})}$$

where p and $p_{\mathbf{Y}}$ are the frequency functions of (\mathbf{X}, \mathbf{Y}) and \mathbf{Y} . The conditional frequency function p is defined only for values of \mathbf{y} such that $p_{\mathbf{Y}}(\mathbf{y}) > 0$. With this definition it is clear that $p(\cdot | \mathbf{y})$ is the frequency function of a probability distribution, since

$$\sum_{\mathbf{x}} p(\mathbf{x} | \mathbf{y}) = \frac{\sum_{\mathbf{x}} p(\mathbf{x}, \mathbf{y})}{p_{\mathbf{Y}}(\mathbf{y})} = \frac{p_{\mathbf{Y}}(\mathbf{y})}{p_{\mathbf{Y}}(\mathbf{y})} = 1$$

by (A.8.11). This probability distribution is called the *conditional distribution of \mathbf{X} given that $\mathbf{Y} = \mathbf{y}$* .

Example 1.1.1. Let $\mathbf{X} = (X_1, \dots, X_n)$, where the X_i are the indicators of a set of n binomial trials with probability of success p . Let $Y = \sum_{i=1}^n X_i$, the total number of successes. Then Y has a binomial, $\mathcal{B}(n, p)$, distribution and

$$(1.1.2) \quad p(\mathbf{x} | y) = \frac{P[\mathbf{X} = \mathbf{x}, Y = y]}{\binom{n}{y} p^y (1-p)^{n-y}} = \frac{p^y (1-p)^{n-y}}{\binom{n}{y} p^y (1-p)^{n-y}} = \frac{1}{\binom{n}{y}}$$

if the x_i are all 0 or 1 and $\sum x_i = y$.

Thus, if we are told we obtained k successes in n binomial trials, then these successes are as likely to occur on one set of trials as on any other. ■

Example 1.1.2. Let X and Y have the joint frequency function given by the table

TABLE 1.1.1

$x \backslash y$	0	10	20	$p_X(x)$
0	0.25	0.05	0.05	0.35
1	0.05	0.15	0.05	0.25
2	0.05	0.10	0.25	0.40
$p_Y(y)$	0.35	0.30	0.35	1

For instance, suppose Y is the number of cigarettes that a person picked at random from a certain population smokes per day (to the nearest 10), and X is a general

health rating for the same person with 0 corresponding to good, 2 to poor, and 1 to neither. We find for $y = 20$

x	0	1	2
$p(x 20)$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{2}$

These figures would indicate an association between heavy smoking and poor health, since $p(2|20)$ is almost twice as large as $p_X(2)$. ■

The conditional distribution of X given $Y = y$ is easy to calculate in two special cases.

- (i) If X and Y are independent $p(x|y) = p_X(x)$ and the conditional distribution coincides with the marginal distribution.
- (ii) If X is a function of Y , $h(Y)$, then the conditional distribution of X is degenerate, $X = h(y)$ with probability 1.

Both of these assertions follow immediately from Definition (1.1.1).

Two important formulae follow from (1.1.1) and (A.4.5). Let $q(y|x)$ denote the conditional frequency function of Y given $X = x$. Then,

$$(1.1.3) \quad p(x, y) = p(x|y)p_Y(y)$$

$$(1.1.4) \quad p(x|y) = \frac{q(y|x)p_X(x)}{\sum_z q(y|z)p_X(z)} \quad \text{Bayes' Rule}$$

whenever the denominator of the right-hand side is positive.

Equation (1.1.3) can be used for model construction. For instance, suppose that the number Y of defectives in a lot of N produced by a manufacturing process has a $\mathcal{B}(N, \theta)$ distribution. Suppose the lot is sampled n times without replacement and let X be the number of defectives found in the sample. We know that given $Y = y$, X has a hypergeometric, $\mathcal{H}(y, N, n)$, distribution. We can now use (1.1.3) to write down the joint distribution of X and Y

$$(1.1.5) \quad P[X = x, Y = y] = \binom{N}{y} \theta^y (1 - \theta)^{n-y} \frac{\binom{y}{x} \binom{N-y}{n-x}}{\binom{N}{n}}$$

where the combinatorial coefficients $\binom{a}{b}$ vanish unless a, b are integers with $b \leq a$.

We can also use this model to illustrate (1.1.4). Since we would usually only observe X , we may want to know what the conditional distribution of Y given $X = x$ is. By (1.1.4) this is,

$$P[Y = y|X = x] = \binom{N}{y} \theta^y (1 - \theta)^{n-y} \frac{\binom{y}{x} \binom{N-y}{n-x}}{\sum_z \binom{N}{z} \theta^z (1 - \theta)^{n-z} \binom{z}{x} \binom{N-z}{n-x}} / c(x)$$

where $c(x) = \sum_y \binom{N}{y} \theta^y (1 - \theta)^{N-y} \binom{N-x}{y-x}$. This formula simplifies to (see Problem 1.1.11) the binomial probability,

$$(1.1.6) \quad P[Y = y | X = x] = \binom{N-x}{y-x} \theta^{y-x} (1 - \theta)^{N-n-(y-x)}.$$

1.1.B. Conditional Expectation for Discrete Variables

Suppose that X is a random variable with $E(|X|) < \infty$. Define the *conditional expectation of X given $Y = y$* , written $E(X | Y = y)$, by

$$(1.1.7) \quad E(X | Y = y) = \sum_x x p(x | y).$$

Note that by (1.1.1), if $p_Y(y) > 0$,

$$(1.1.8) \quad \sum_x |x| p(x | y) \leq \sum_x |x| \frac{p_X(x)}{p_Y(y)} = \frac{E(|X|)}{p_Y(y)}.$$

Example 1.1.3. Suppose X and Y have the joint frequency function of Table 1.1.1 above. We find

$$E(X | Y = 20) = 0 \cdot \frac{1}{4} + 1 \cdot \frac{1}{4} + 2 \cdot \frac{5}{4} = \frac{11}{4} = 1.57.$$

Similarly, $E(X | Y = 10) = \frac{7}{6} = 1.17$ and $E(X | Y = 0) = \frac{3}{7} = 0.43$. Note that in the health versus smoking context, we can think of $E(X | Y = y)$ as the mean health rating for people who smoke y cigarettes a day. ■

Let $g(y) = E(X | Y = y)$. The random variable $g(Y)$ is written $E(X | Y)$ and is called the *conditional expectation of X given Y* .†

As an example we calculate $E(X_1 | Y)$ where X_1 and Y are given in Example 1.1.1. We have,

$$(1.1.9) \quad E(X_1 | Y = i) = P[X_1 = 1 | Y = i] = \frac{\binom{n-1}{i-1}}{\binom{n}{i}} = \frac{i}{n}.$$

The first of these equalities holds because X_1 is an indicator. The second follows from (1.1.2), since $\binom{n-1}{i-1}$ is just the number of ways i successes can occur in n trials with the first trial being a success. Therefore,

$$(1.1.10) \quad E(X_1 | Y) = \frac{Y}{n}.$$

The conditional distribution of a random vector X given $Y = y$ corresponds

†We shall follow the convention of also calling $E(X | Y)$ any variable which is equal to $g(Y)$ with probability 1.

to a single probability measure P_y on (Ω, \mathcal{A}) . Specifically, define for $A \in \mathcal{A}$,

$$(1.1.11) \quad P_y(A) = P(A | [Y = y]) \text{ if } p_Y(y) > 0.$$

This P_y is just the conditional probability measure on (Ω, \mathcal{A}) mentioned in (A.4.2). Now the conditional distribution of X given $Y = y$ is the same as the distribution of X if P_y is the probability measure on (Ω, \mathcal{A}) . Therefore, the conditional expectation is an ordinary expectation with respect to the probability measure P_y . It follows that all the properties of the expectation given in (A.10.3)–(A.10.8) hold for the conditional expectation given $Y = y$. For example,

$$(1.1.12) \quad E(\alpha X_1 + \beta X_2 | Y = y) = \alpha E(X_1 | Y = y) + \beta E(X_2 | Y = y)$$

identically in y for any X_1, X_2 such that $E(|X_1|), E(|X_2|)$ are finite. Since the identity holds for all y we have,

$$(1.1.13) \quad E(\alpha X_1 + \beta X_2 | Y) = \alpha E(X_1 | Y) + \beta E(X_2 | Y).$$

This process can be repeated for each of (A.10.3)–(A.10.8) to obtain analogous properties of the conditional expectation. An important example of this type of argument is given in Section 1.6.

In two special cases we can calculate conditional expectations immediately. If X and Y are independent and $E(|X|) < \infty$, then

$$(1.1.14) \quad E(X | Y) = E(X).$$

This is clear by (i).

On the other hand by (ii)

$$(1.1.15) \quad E(h(Y) | Y) = h(Y).$$

The notion implicit in (1.1.15) is that given $Y = y$, Y acts as a constant. If we carry this further we have a relation which we shall call the *substitution theorem for conditional expectations*:

$$(1.1.16) \quad E(q(X, Y) | Y = y) = E(q(X, y) | Y = y),$$

for all y such that $p_Y(y) > 0$ and $q(X, Y)$ has a finite expectation. Assertion (1.1.16) is immediate, since

$$(1.1.17)$$

$$P[q(X, Y) = a | Y = y] = P[q(X, Y) = a, Y = y | Y = y] = P[q(X, y) = a | Y = y]$$

for any a .

If we put $q(X, Y) = r(X)h(Y)$, where h is bounded and $r(X)$ has a finite expectation, we obtain by (1.1.16),

$$(1.1.18) \quad E(r(X)h(Y) | Y = y) = E(r(X)h(y) | Y = y) = h(y)E(r(X) | Y = y).$$

Therefore,

$$(1.1.19) \quad E(r(X)h(Y)|Y) = h(Y)E(r(X)|Y).$$

Another intuitively reasonable result is that the mean of the conditional means is the mean:

$$(1.1.20) \quad E(E(X|Y)) = E(X),$$

whenever X has a finite expectation. We refer to this as the *double expectation theorem*.

To prove (1.1.20) we write, in view of (1.1.7) and (A.10.5),

$$\begin{aligned} E(E(X|Y)) &= \sum_y p_Y(y) [\sum_x x p(x|y)] \\ (1.1.21) \quad &= \sum_{x,y} x p(x|y) p_Y(y) \\ &= \sum_{x,y} x p(x, y) = E(X). \end{aligned}$$

The interchange of summation used is valid, since the finiteness of $E(|X|)$ implies that all sums converge absolutely.

As an illustration, we check (1.1.20) for $E(X_1|Y)$ given by (1.1.10). In this case,

$$(1.1.22) \quad E(E(X_1|Y)) = E\left(\frac{Y}{n}\right) = \frac{np}{n} = p = E(X_1).$$

If we apply (1.1.20) to $X = r(X)h(Y)$ and use (1.1.19), we obtain the *product expectation formula*:

Theorem 1.1.1. If $h(Y)$ is bounded and $E(|r(X)|) < \infty$, then

$$(1.1.23) \quad E(r(X)h(Y)) = E(h(Y)E(r(X)|Y)).$$

Note that we can express the conditional probability that $X \in A$ given $Y = y$ as

$$P[X \in A|Y = y] = E(I_A(X)|Y = y) = \sum_{x \in A} p(x|y).$$

Then by taking $r(X) = I_A(X)$, $h = 1$ in Theorem 1.1.1 we can express the (unconditional) probability that $X \in A$ as

$$(1.1.24) \quad P[X \in A] = E(E(r(X)|Y)) = \sum_y P[X \in A|Y = y] p_Y(y).$$

For example, if X and Y are as in (1.1.5),

$$P[X \leq x] = \sum_y \binom{N}{y} \theta^y (1 - \theta)^{n-y} H_y(x)$$

where H_y is the distribution function of a hypergeometric distribution with parameters (y, N, n) .