

DOCUMENT RESUME

ED 103 480

TN 004 326

TITLE Expected Value of a Sample Estimate.
INSTITUTION Statistical Reporting Service (DOA), Washington, D.C.
REPORT NO SRS-19
PUB DATE Sep 74
NOTE 146p.

EDRS PRICE MF-\$0.76 HC-\$6.97 PLUS POSTAGE
DESCRIPTORS *Agriculture; Algebra; Classification; Data Collection; Design; *Mathematical Applications; Mathematics; Matrices; Probability; *Probability Theory; *Sampling; Set Theory; *Statistical Analysis; Statistical Bias; Statistical Surveys
IDENTIFIERS Random Variables; Variance (Statistical)

ABSTRACT

Intended as a reference for the convenience of students in sampling, this monograph attempts to express relevant, introductory mathematics and probability in the context of sample surveys. Although some proofs are presented, the emphasis is more on exposition of mathematical language and concepts than on the mathematics per se and rigorous proofs. Many problems are given as exercises so a student may test his interpretation or understanding of the concepts. Most of the mathematics is elementary. Each chapter begins with simple explanations and ends at a much more advanced level. Students with only high school algebra should have no difficulty with the first parts of each chapter. Chapters 1 and 2 were added as background for Chapter 3 which discusses expected values of random variables. Chapter 4 focuses attention on the distribution of an estimate which is the basis for comparing the accuracy of alternative sampling plans as well as a basis for statements about the accuracy of an estimate from a sample. The content of chapter 4 is included in books on sampling, but it is important that students hear or read more than one discussion of the distribution of an estimate, especially with reference to estimates from actual sample surveys. (Author/RC)

BEST COPY AVAILABLE

Expected Value of a Sample Estimate

4 - OCT 1 6

Copy 1974

Statistical Reporting Service • U.S. Department of Agriculture • SRS No. 19

AE20	1F1F1D1C	2B2D1917	21241D1D	2D2F1718	2524221F	2F2D1918	21211C1C	24261
AE40	1F1F1A18	2A2A1617	1F241C23	2D2D1817	27272525	2F2F1818	211F1C18	31311
AE60	24271C25	2F2F1917	2724281D	332D1818	211F1D18	31331A1C	1F21181C	312F1
AE80	211E1A17	2B2B1618	1E1E1717	2B331918	1C1E1214	352B1E16	1E1E1414	1717C
AEA0	1F1F1212	11C006C2	1F1F1211	0C0C02D2	1F1F1111	0A0B0202	1E1E1111	0B0B0
AEC0	1E1F1215	001704C8	21241D1D	2B2B1517	211F1D18	282B1715	1C1E110F	190B0
AEB0	1E1E0F0F	0A0A0101	1E1F0F14	0A0A0202	1E1F0E0F	0A0A0202	1F1C110F	0B0B0
AFC0	1F21141C	14260814	241F1D1C	2A2A1615	1E24181E	2B2F1618	2521221C	2F2D1
AE20	211F1C1A	2B2A1715	1E1F1414	2A2D1918	1E1F171A	2D2B1917	1F241A1D	2D2A1
AE40	1E1F1214	1316C7D8	1E1C1114	2B2D1519	1C1F1717	2B2B1616	1E151717	2B2A1
AE60	21211F22	2F2A1715	21211D1D	2B26121C	21212325	26261213	1F211D22	1F26C
AE80	27292825	31291615	1F1F1C	2A2D1417	252C132A	333A1A1D	2C252A22	3A33
AEA0	21251C23	31331919	2721231C	2F2A1715	1F21181C	262B1413	211F1F1C	2A2A1
AFC0	2427222A	2A2F1517	2524231D	2F2B1615	21241C2	2B2D1617	2A25211D	31351
AEB0	2F303C30	38381A19	30303030	3031A1	2C272A22	35311919	2122315	31351
BO20	24272223	2F311118	2724252	2F2B1715	211F1D17	272A1514	1C24171F	2A2B1
BO40	1F1C1814	2D2B1714	1E1F111C	2A2B1515	24291F2C	30331619	2C272A22	2F2A1
BO60	21251D25	2A311519	2524231F	212F1A1C	24242223	312E1919	21211F1C	2F2A1
BO80	24241C1F	2D2F1716	21241D22	2D2B1416	211F1F1C	212F1719	212F1C1C	2D2D1
BOA0	21211C1C	2B271515	211F1C1D	2F2B1513	1211D1D	212A1113	211F1C1A	2A2B1
BOC0	1F1F1817	2F211918	1E1E1815	2A211713	11C1212	212A1317	1E1C1214	2B2B1
BOE0	1C1C1412	26241415	1C1F1215	24261714	31B1412	212B1416	1E1C1214	2622
BOF0	1C1F171A	242B1318	2121111D	2D2D1918	1F211C1D	212F1919	1F211D1C	2F2D1
BL20	1F211C1C	2F2D191A	211F1C1C	2D2F1817	211F1C1A	1F2F1519	1F211C1D	2F2F1
BL40	21241F1F	2D2F1119	27241F1E	2F2B1E18	1F211D1C	2F2F1618	1F1E1515	2F2D1
BL60	21241C1D	2D2B1716	211F1D18	24261215	1F1F1517	2D2B1A19	1F1F181C	2D2D1
BL80	211F1C1B	31311A11	1E1E1718	2F2D1H19	1F211C1D	2D2D1817	211F1C1A	2D2B1
BLA0	1E1C1715	2B261615	1C1F1418	2B2B1516	21211122	2F2D1716	2524251F	2F2D1
BLC0	24271F25	2D331619	21212D2D	3B2D1C17	21211F4F	262B1313	25242322	2B2D1
BLE0	21211F1F	2D2D1817	24271F1E	2B2B1615	21212222	2A2B1314	21211F22	2A2A1
BLF0	2521221D	2B2B1414	211F1F1D	2D2A0417	21241F1F	2B2B1513	211F1F1D	2F1B1
BLG0	1F1F181E	1B260514	21241C1F	2F311919	25242222	31311A1A	24242323	3331
BMA0	25252522	33311A17	21241C1D	2B2F1618	24241D1D	2F2F1818	21211D1C	2F2D1
BMC0	21211C1C	2B2D1616	21241D1F	2F2F1718	21241D1D	312F1819	21241D1D	2F2D1
BME0	24242222	2D2F1817	21211F1C	2D2D1718	24211D1D	2D2F1717	21211C1C	2D2F1
BMG0	21211D1D	31311A19	242C1F30	313F1A15	2C242F22	3D2A1D16	21241F22	2B2F1
BMA0	21211D1F	2D2F1818	2521221D	2F2F1818	24291D28	312F1817	2721231C	31311
BMC0	21211D1C	3331191C	21211A1D	3D331C1A	211F1A1A	2F2F1818	21211C1D	2D2F1
BME0	21211F1C	2D2F1818	1F1F1A1C	2A2B1514	1F1F1817	2D2F1918	1F1F1818	2F2D1
BMA0	1F1F171B	2F2F1A19	1F1F1814	332F1B1A	1C1F171C	2F261813	1F1F191C	2D2D1
BMC0	1F1F181B	22241112	1F1F1817	24241312	1F1F1817	26221312	1F1F1718	221F1
BME0	1F1F181A	1E1F1D0F	1F1F1C1C	1D1B0C0D	1F1F1A1C	1D1B0C0D	1F1F1A1A	1D1B1

FOREWORD

The Statistical Reporting Service (SRS) has been engaged for many years in the training of agricultural statisticians from around the world. Most of these participants come under the support of the Agency for International Development (AID) training programs; however, many also come under sponsorship of the Food and Agriculture Organization into the International Statistical Programs Center of the Bureau of the Census, with which SRS is cooperating.

This treatise was developed by the SRS with the cooperation of AID and the Center, in an effort to provide improved materials for teaching and reference in the area of agricultural statistics, not only for foreign students but also for development of staff working for these agencies.

HARRY C. TRELOGAN
Administrator
Statistical Reporting Service

Washington, D. C.

September 1974

PREFACE

The author has felt that applied courses in sampling should give more attention to elementary theory of expected values of a random variable. The theory pertaining to a random variable and to functions of random variables is the foundation for probability sampling. Interpretations of the accuracy of estimates from probability sample surveys are predicated on, among other things, the theory of expected values.

There are many students with career interests in surveys and the application of probability sampling who have very limited backgrounds in mathematics and statistics. Training in sampling should go beyond simply learning about sample designs in a descriptive manner. The foundations in mathematics and probability should be included. It can (1) add much to the breadth of understanding of bias, random sampling error, components of error, and other technical concepts; (2) enhance one's ability to make practical adaptations of sampling principals and correct use of formulas; and (3) make communication with mathematics' statisticians easier and more meaningful.

This monograph is intended as a reference for the convenience of students in sampling. It attempts to express relevant, introductory mathematics and probability in the context of sample surveys. Although some proofs are presented, the emphasis is more on exposition of mathematical language and concepts than on the mathematics per se and rigorous proofs. Many problems are given as exercises so a student may test his interpretation or understanding of the concepts. Most of the mathematics is elementary. If a formula looks involved, it is probably because it represents a long sequence of arithmetic operations.

Each chapter begins with very simple explanations and ends at a much more advanced level. Most students with only high school algebra should have no difficulty with the first parts of each chapter. Students with a few courses in college mathematics and statistics might review the first parts of each chapter and spend considerable time studying the latter parts. In fact, some students might prefer to start with Chapter III and refer to Chapters I and II only as needed.

Discussion of expected values of random variables, as in Chapter III, was the original purpose of this monograph. Chapters I and II were added as background for Chapter III. Chapter IV focuses attention on the distribution of an estimate which is the basis for comparing the accuracy of alternative sampling plans as well as a basis for statements about the accuracy of an estimate from a sample. The content of Chapter IV is included in books on sampling, but it is important that students hear or read more than one discussion of the distribution of an estimate, especially with reference to estimates from actual sample surveys.

The author's interest and experience in training has been primarily with persons who had begun careers in agricultural surveys. I appreciate the opportunity, which the Statistical Reporting Service has provided, to prepare this monograph.

Earl E. Houseman
Statistician

CONTENTS

	<u>Page</u>
Chapter I. Notation and Summation	1
1.1 Introduction	1
1.2 Notation and the Symbol for Summation	1
1.3 Frequency Distributions	9
1.4 Algebra	10
1.5 Double Indexes and Summation	14
1.5.1 Cross Classification	15
1.5.2 Hierarchical or Nested Classification	22
1.6 The Square of a Sum	26
1.7 Sums of Squares	29
1.7.1 Nested Classification	29
1.7.2 Cross Classification	32
 Chapter II. Random Variables and Probability	 33
2.1 Random Variables	33
2.2 Addition of Probabilities	35
2.3 Multiplication of Probabilities	41
2.4 Sampling With Replacement	42
2.5 Sampling Without Replacement	45
2.6 Simple Random Samples	47
2.7 Some Examples of Restricted Random Sampling	50
2.8 Two-Stage Sampling	57
 Chapter III. Expected Values of Random Variables	 63
3.1 Introduction	63
3.2 Expected Value of the Sum of Two Random Variables	67
3.3 Expected Value of an Estimate	72
3.4 Variance of a Random Variable	77
3.4.1 Variance of the Sum of Two Independent Random Variables	77
3.4.2 Variance of the Sum of Two Dependent Random Variables	79

CONTENTS (Continued)

	<u>Page</u>
3.5 Variance of an Estimate	81
3.5.1 Equal Probability of Selection	80
3.5.2 Unequal Probability of Selection	92
3.6 Variance of a Linear Combination	84
3.7 Estimation of Variance	89
3.7.1 Simple Random Sampling	89
3.7.2 Unequal Probability of Selection	92
3.8 Ratio of Two Random Variables	94
3.9 Conditional Expectation	97
3.10 Conditional Variance	103
 Chapter IV. The Distribution of an Estimate	 113
4.1 Properties of Simple Random Samples	113
4.2 Shape of the Sampling Distribution	117
4.3 Sample Design	119
4.4 Response Error	125
4.5 Bias and Standard Error	135

CHAPTER I. NOTATION AND SUMMATION

1.1 INTRODUCTION

To work with large amounts of data, an appropriate system of notation is needed. The notation must identify data by individual elements, and provide meaningful mathematical expressions for a wide variety of summaries from individual data. This chapter describes notation and introduces summation algebra, primarily with reference to data from censuses and sample surveys. The purpose is to acquaint students with notation and summation rather than to present statistical concepts. Initially some of the expressions might seem complex or abstract, but nothing more than sequences of operations involving addition, subtraction, multiplication, and division is involved. Exercises are included so a student may test his interpretation of different mathematical expressions. Algebraic manipulations are also discussed and some algebraic exercises are included. To a considerable degree, this chapter could be regarded as a manual of exercises for students who are interested in sampling but are not fully familiar with the summation symbol, Σ . Familiarity with the mathematical language will make the study of sampling much easier.

1.2 NOTATION AND THE SYMBOL FOR SUMMATION

"Element" will be used in this monograph as a general expression for a unit that a measurement pertains to. An element might be a farm, a person, a school, a stalk of corn, or an animal. Such units are sometimes called units of observation or reporting units. Generally, there are several characteristics or items of information about an element that one might be interested in.

"Measurement" or "value" will be used as general terms for the numerical value of a specified characteristic for an element. This includes assigned values. For example, the element might be a farm and the characteristic could be whether wheat is being grown or is not being grown on a farm. A value of "1" could be assigned to a farm growing wheat and a value of "0" to a farm not growing wheat. Thus, the "measurement" or "value" for a farm growing wheat would be "1" and for a farm not growing wheat the value would be "0."

Typically, a set of measurements of N elements will be expressed as follows: X_1, X_2, \dots, X_N where X refers to the characteristic that is measured and the index (subscript) to the various elements of the population (or set). For example, if there are N persons and the characteristic X is a person's height, then X_1 is the height of the first person, etc. To refer to any one of elements, not a specific element, a subscript "i" is used. Thus, X_i (read X sub i) means the value of X for any one of the N elements. A common expression would be " X_i is the value of X for the i^{th} element."

The Greek letter Σ (capital sigma) is generally used to indicate a sum. When found in an equation, it means "the sum of." For example,

$\sum_{i=1}^N X_i$ represents the sum of all values of X from X_1 to X_N ; that is,

$\sum_{i=1}^N X_i = X_1 + X_2 + \dots + X_N$. The lower and upper limits of the index of

summation are shown below and above the summation sign. For example, to

specify the sum of X for elements 11 thru 20 one would write $\sum_{i=11}^{20} X_i$.

You might also see notation such as " $\sum_{i=1}^N X_i$ " where $i = 1, 2, \dots, N$ " which indicates there are N elements (or values) in the set indexed by serial numbers 1 thru N , or for part of a set you might see " $\sum_{i=11}^{20} X_i$ " where $i = 11, 12, \dots, 20$." Generally the index of summation starts with 1; so you will

often see a summation written as $\sum_{i=1}^N X_i$. That is, only the upper limit of the summation is shown and it is understood that the summation begins with $i=1$. Alternatively, when the set of values being summed is clearly understood, the lower and upper limits might not be shown. Thus, it is understood that $\sum_{i=1} X_i$ or $\sum X_i$ is the sum of X over all values of the set under consideration. Sometimes a writer will even drop the subscript and use $\sum X$ for the sum of all values of X . Usually the simplest notation that is adequate for the purpose is adopted. In this monograph, there will be some deliberate variation in notation to familiarize students with various representations of data.

An average is usually indicated by a "bar" over the symbol. For example, \bar{X} (read "X bar," or sometimes "bar X") means the average value of

X . Thus, $\bar{X} = \frac{\sum_{i=1}^N X_i}{N}$. In this case, showing the upper limit, N , of the summation makes it clear that the sum is being divided by the number of elements and \bar{X} is the average of all elements. However, $\frac{\sum X_i}{N}$ would also be interpreted as the average of all values of X unless there is an indication to the contrary.

Do not try to study mathematics without pencil and paper. Whenever the shorthand is not clear, try writing it out in long form. This will often reduce any ambiguity and save time.

Here are some examples of mathematical shorthand:

- (1) Sum of the reciprocals of X

$$\sum_{i=1}^N \frac{1}{X_i} = \frac{1}{X_1} + \frac{1}{X_2} + \dots + \frac{1}{X_N}$$
- (2) Sum of the differences between X_i and a constant, C

$$\sum_{i=1}^N (X_i - C) = (X_1 - C) + (X_2 - C) + \dots + (X_N - C)$$
- (3) Sum of the deviations of X_i from the average of X

$$\sum_{i=1}^N (X_i - \bar{X}) = (X_1 - \bar{X}) + (X_2 - \bar{X}) + \dots + (X_N - \bar{X})$$
- (4) Sum of the absolute values of the differences between X_i and \bar{X} . (Absolute value, indicated by the vertical lines, means the positive value of the difference)

$$\sum |X_i - \bar{X}| = |X_1 - \bar{X}| + |X_2 - \bar{X}| + \dots + |X_N - \bar{X}|$$
- (5) Sum of the squares of X_i

$$\sum X_i^2 = X_1^2 + X_2^2 + X_3^2 + \dots + X_N^2$$
- (6) Sum of squares of the deviations of X from \bar{X}

$$\sum (X_i - \bar{X})^2 = (X_1 - \bar{X})^2 + \dots + (X_N - \bar{X})^2$$
- (7) Average of the squares of the deviations of X from \bar{X}

$$\frac{\sum_{i=1}^N (X_i - \bar{X})^2}{N} = \frac{(X_1 - \bar{X})^2 + \dots + (X_N - \bar{X})^2}{N}$$
- (8) Sum of products of X and Y

$$\sum_{i=1}^N X_i Y_i = X_1 Y_1 + X_2 Y_2 + \dots + X_N Y_N$$
- (9) Sum of quotients of X divided by Y

$$\sum \frac{X_i}{Y_i} = \frac{X_1}{Y_1} + \frac{X_2}{Y_2} + \dots + \frac{X_N}{Y_N}$$
- (10) Sum of X divided by the sum of Y

$$\frac{\sum X_i}{\sum Y_i} = \frac{X_1 + X_2 + \dots + X_N}{Y_1 + Y_2 + \dots + Y_N}$$
- (11) Sum of the first N digits

$$\sum_{i=1}^N i = 1 + 2 + 3 + \dots + N$$
- (12)

$$\sum_{i=1}^N i X_i = X_1 + 2X_2 + 3X_3 + \dots + NX_N$$
- (13)

$$\sum_{i=1}^6 (-1)^i X_i = -X_1 + X_2 - X_3 + X_4 - X_5 + X_6$$

Exercise 1.1. You are given a set of four elements having the following values of X : $X_1 = 2$, $X_2 = 0$, $X_3 = 5$, $X_4 = 7$. To test your understanding of the summation notation, compute the values of the following algebraic expressions:

<u>Expression</u>	<u>Answer</u>
(1) $\sum_{i=1}^4 (X_i + 4)$	30
(2) $\sum 2(X_i - 1)$	20
(3) $2\sum (X_i - 1)$	20
(4) $\sum 2X_i - 1$	27
(5) $\bar{X} = \frac{\sum X_i}{N}$	3.5
(6) $\sum X_i^2$	78
(7) $\sum (-X_i)^2$	78
(8) $[\sum X_i]^2$	196
(9) $\sum (X_i^2 - X_i)$	64
(10) $\sum (X_i^2) - \sum X_i$	64
(11) $\sum 1(X_i)$	45
(12) $\sum (-1)^1 (X_i)$	0
(13) $\sum_{i=1}^4 (X_i^2 - 3)$	66
(14) $\sum_{i=1}^4 X_i^2 - \sum_{i=1}^4 (3)$	66

Note: $\sum_{i=1}^4 (3)$ means find the sum of four 3's

<u>Expression (Continued)</u>	<u>Answer</u>
(15) $\Sigma (X_1 - \bar{X})$	0
(16) $\frac{\Sigma (X_1 - \bar{X})^2}{N-1}$	$\frac{29}{3}$
(17) $\frac{\Sigma [X_1^2 - 2X_1\bar{X} + \bar{X}^2]}{N-1}$	$\frac{29}{3}$
(18) $\frac{\Sigma X_1^2 - N\bar{X}^2}{N-1}$	$\frac{29}{3}$

Definition 1.1. The variance of X where $X = X_1, X_2, \dots, X_N$, is defined in one of two ways:

$$\sigma^2 = \frac{\sum_{i=1}^N (X_i - \bar{X})^2}{N}$$

or

$$s^2 = \frac{\sum_{i=1}^N (X_i - \bar{X})^2}{N-1}$$

The reason for the two definitions will be explained in Chapter III. The variance formulas provide measures of how much the values of X vary (deviate) from the average. The square root of the variance of X is called the standard deviation of X . The central role that the above definitions of variance and standard deviation play in sampling theory will become apparent as you study sampling. The variance of an estimate from a sample is one of the measures needed to judge the accuracy of the estimate and to evaluate alternative sampling designs. Much of the algebra and notation in this chapter is related to computation of variance. For

complex sampling plans, variance formulas are complex. This chapter should help make the mathematics used in sampling more readable and more meaningful when it is encountered.

Definition 1.2. "Population" is a statistical term that refers to a set of elements from which a sample is selected ("Universe" is often used instead of "Population").

Some examples of populations are farms, retail stores, students, households, manufacturers, and hospitals. A complete definition of a population is a detailed specification of the elements that compose it. Data to be collected also need to be defined. Problems of defining populations to be surveyed should receive much attention in courses on sampling. From a defined population a sample of elements is selected, information for each element in the sample is collected, and inferences from the sample are made about the population. Nearly all populations for sample surveys are finite so the mathematics and discussion in this monograph are limited to finite populations.

In the theory of sampling, it is important to distinguish between data for elements in a sample and data for elements in the entire population. Many writers use uppercase letters when referring to the population and lowercase letters when referring to a sample. Thus X_1, \dots, X_N would represent the values of some characteristic X for the N elements of the population; and x_1, \dots, x_n would represent the values of X in a sample of n elements. The subscripts in x_1, \dots, x_n simply index the different elements in a sample and do not correspond to the subscripts in X_1, \dots, X_N which index the elements of the population. In other words, x_1 could be any one of the X_i 's. Thus,

$$\frac{\sum_{i=1}^N X_i}{N} = \bar{X} \quad \text{represents the population mean, and}$$

$$\frac{\sum_{i=1}^n x_i}{n} = \bar{x} \quad \text{represents a sample mean}$$

In this chapter we will be using only uppercase letters, except for constants and subscripts, because the major emphasis is on symbolic representation of data for a set of elements and on algebra. For this purpose, it is sufficient to start with data for a set of elements and not be concerned with whether the data are for a sample of elements or for all elements in a population.

The letters X, Y, and Z are often used to represent different characteristics (variables) whereas the first letters of the alphabet are commonly used as constants. There are no fixed rules regarding notation. For example, four different variables or characteristics might be called X_1 , X_2 , X_3 , and X_4 . In that case X_{11} might be used to represent the 1th value of the variable X_1 . Typically, writers adopt notation that is convenient for their problems. It is not practical to completely standardize notation.

Exercise 1.2. In the list of expressions in Exercise 1.1 find the variance of X, that is, find S^2 . Suppose that X_4 is 15 instead of 7. How much is the variance of X changed? Answer: From $9\frac{2}{3}$ to $44\frac{1}{3}$.

Exercise 1.3. You are given four elements having the following values of X and Y

$X_1 = 2$	$X_2 = 0$	$X_3 = 5$	$X_4 = 7$
$Y_1 = 2$	$Y_2 = 3$	$Y_3 = 1$	$Y_4 = 14$

Find the value of the following expressions:

<u>Expression</u>	<u>Answer</u>	<u>Expression</u>	<u>Answer</u>
(1) $\sum X_i Y_i$	107	(7) $\sum X_i - \sum Y_i$	-6
(2) $(\sum X_i)(\sum Y_i)$	280	(8) $\sum (X_i - Y_i)^2$	74
(3) $\sum (X_i - \bar{X})(Y_i - \bar{Y})$	37	(9) $\sum (X_i^2 - Y_i^2)$	-132
(4) $\sum X_i Y_i - N\bar{X}\bar{Y}$	37	(10) $\sum X_i^2 - \sum Y_i^2$	-132
(5) $\frac{1}{N} \sum \frac{X_i}{Y_i}$	1.625	(11) $[\sum (X_i - Y_i)]^2$	36
(6) $\sum (X_i - Y_i)$	-6	(12) $[\sum X_i]^2 - [\sum Y_i]^2$	-204

1.3 FREQUENCY DISTRIBUTIONS

Several elements in a set of N might have the same value for some characteristic X . For example, many people have the same age. Let X_j be a particular age and let N_j be the number of people in a population (set) of N people who have the age X_j . Then $\sum_{j=1}^K N_j = N$ where K is the number of different ages found in the population. Also $\sum N_j X_j$ is the sum

of the ages of the N people in the population and $\frac{\sum N_j X_j}{\sum N_j}$ represents the average age of the N people. A listing of X_j and N_j is called the frequency distribution of X , since N_j is the number of times (frequency) that the age X_j is found in the population.

On the other hand, one could let X_i represent the age of the i^{th} individual in a population of N people. Notice that j was an index of age. We are now using i as an index of individuals, and the average age would

be written as $\frac{\sum X_i}{N}$. Note that $\sum N_j X_j = \sum X_i$ and that $\frac{\sum N_j X_j}{\sum N_j} = \frac{\sum X_i}{N}$. The

choice between these two symbolic representations of the age of people in the population is a matter of convenience and purpose.

Exercise 1.4. Suppose there are 20 elements in a set (that is, $N = 20$) and that the values of X for the 20 elements are: 4, 8, 3, 7, 3, 8, 3, 3, 7, 2, 8, 4, 8, 8, 3, 7, 8, 10, 3, 8.

- (1) List the values of X_j and N_j , where j is an index of the values 2, 3, 4, 7, 8, and 10. This is the frequency distribution of X .
- (2) What is K equal to?

Interpret and verify the following by making the calculations indicated:

$$(3) \quad \sum_{i=1}^N X_i = \sum_{j=1}^K N_j X_j$$

$$(4) \quad \frac{\sum X_i}{N} = \frac{\sum N_j X_j}{\sum N_j} = \bar{X}$$

$$(5) \quad \frac{\sum (X_i - \bar{X})^2}{N} = \frac{\sum N_j (X_j - \bar{X})^2}{\sum N_j}$$

1.4 ALGEBRA

In arithmetic and elementary algebra, the order of the numbers when addition or multiplication is performed does not affect the results. The familiar arithmetic laws when extended to algebra involving the summation symbol lead to the following important rules or theorems:

$$\text{Rule 1.1} \quad \sum (X_i - Y_i + Z_i) = \sum X_i - \sum Y_i + \sum Z_i$$

$$\text{or } \sum (X_{1i} + X_{2i} + \dots + X_{Ki}) = \sum X_{1i} + \sum X_{2i} + \dots + \sum X_{Ki}$$

$$\text{Rule 1.2} \quad \sum aX_i = a\sum X_i \text{ where } \underline{a} \text{ is a constant}$$

$$\text{Rule 1.3} \quad \sum (X_i + b) = \sum X_i + Nb \text{ where } \underline{b} \text{ is constant}$$

If it is not obvious that the above equations are correct, write both sides of each equation as series and note that the difference between the two sides is a matter of the order in which the summation (arithmetic) is performed. Note that the use of parentheses in Rule 1.3 means that b is contained in the series N times. That is,

$$\begin{aligned}\sum_{i=1}^N (X_i + b) &= (X_1 + b) + (X_2 + b) + \dots + (X_N + b) \\ &= (X_1 + X_2 + \dots + X_N) + Nb\end{aligned}$$

On the basis of Rule 1.1, we can write

$$\sum_{i=1}^N (X_i + b) = \sum_{i=1}^N X_i + \sum_{i=1}^N b$$

The expression $\sum_{i=1}^N b$ means "sum the value of b , which occurs N times." Therefore,

$$\sum_{i=1}^N b = Nb.$$

Notice that if the expression had been $\sum_{i=1}^N X_i + b$, then b is an amount to add

to the sum, $\sum_{i=1}^N X_i$.

In many equations \bar{X} will appear; for example, $\sum_{i=1}^N \bar{X} X_i$ or $\sum_{i=1}^N (X_i - \bar{X})$.

Since \bar{X} is constant with regard to the summation, $\sum_{i=1}^N \bar{X} X_i = \bar{X} \sum_{i=1}^N X_i$. Thus,

$$\sum_{i=1}^N (X_i - \bar{X}) = \sum_{i=1}^N X_i - \sum_{i=1}^N \bar{X} = \sum_{i=1}^N X_i - N\bar{X}. \quad \text{By definition, } \bar{X} = \frac{\sum_{i=1}^N X_i}{N}. \quad \text{Therefore,}$$

$$N\bar{X} = \sum_{i=1}^N X_i \quad \text{and} \quad \sum_{i=1}^N (X_i - \bar{X}) = 0.$$

To work with an expression like $\sum_{i=1}^N (X_i + b)^2$ we must square the quantity in parentheses before summing. Thus,

$$\begin{aligned}
 \sum_1 (X_1 + b)^2 &= \sum_1 (X_1^2 + 2bX_1 + b^2) \\
 &= \sum_1 X_1^2 + \sum_1 2bX_1 + \sum_1 b^2 \quad \text{Rule 1} \\
 &= \sum_1 X_1^2 + 2b\sum_1 X_1 + Nb^2 \quad \text{Rules 2 and 3}
 \end{aligned}$$

Verify this result by using series notation. Start with $(X_1+b)^2 + \dots + (X_N+b)^2$.

It is very important that the ordinary rules of algebra pertaining to the use of parentheses be observed. Students frequently make errors because inadequate attention is given to the placement of parentheses or to the interpretation of parentheses. Until you become familiar with the above rules, practice translating shorthand to series and series to shorthand. Study the following examples carefully:

$$(1) \quad \sum_1 (X_1)^2 \neq (\sum_1 X_1)^2$$

The left-hand side is the sum of the squares of X_1 . The right-hand side is the square of the sum of X_1 . On the right the parentheses are necessary. The left side could have been written $\sum_1 X_1^2$.

Rule 1.2 applies.

$$(2) \quad \sum \left[\frac{X_1}{N} \right]^2 = \frac{\sum X_1^2}{N^2}$$

$$(3) \quad \sum (X_1 + Y_1)^2 \neq \sum X_1^2 + \sum Y_1^2$$

A quantity in parentheses must be squared before taking a sum.

$$(4) \quad \sum (X_1^2 + Y_1^2) = \sum X_1^2 + \sum Y_1^2$$

Rule 1.1 applies

$$(5) \quad \sum_1 X_1 Y_1 \neq (\sum_1 X_1)(\sum_1 Y_1)$$

The left side is the sum of products. The right side is the product of sums.

$$(6) \quad \sum (X_1 - Y_1)^2 = \sum X_1^2 - 2\sum X_1 Y_1 + \sum Y_1^2$$

$$(7) \quad \sum_1^N a(X_1 - b) \neq a\sum_1^N X_1 - ab$$

$$(8) \quad \sum_{i=1}^N a(X_i - b) = a \sum_{i=1}^N X_i - Nab$$

$$(9) \quad a[\sum_{i=1}^N X_i - b] = a \sum_{i=1}^N X_i - ab$$

$$(10) \quad \sum_{i=1}^N X_i(X_i - Y_i) = \sum_{i=1}^N X_i^2 - \sum_{i=1}^N X_i Y_i$$

Exercise 1.5. Prove the following:

In all cases, assume $i = 1, 2, \dots, N$.

$$(1) \quad \sum (X_i - \bar{X}) = 0$$

$$(2) \quad \sum \frac{X_i Y_i}{X_i^2} = \sum \frac{Y_i}{X_i}$$

$$(3) \quad N\bar{X}^2 = \frac{(\sum X_i)^2}{N}$$

$$(4) \quad \sum_{i=1}^N (aX_i + bY_i + C) = a \sum X_i + b \sum Y_i + NC$$

Note: Equations (5) and (6) should be (or become) very familiar equations.

$$(5) \quad \sum (X_i - \bar{X})^2 = \sum X_i^2 - N\bar{X}^2$$

$$(6) \quad \sum (X_i - \bar{X})(Y_i - \bar{Y}) = \sum X_i Y_i - N\bar{X}\bar{Y}$$

$$(7) \quad \sum \left(\frac{X_i}{a} + Y_i \right)^2 = \frac{1}{a^2} \sum (X_i + aY_i)^2$$

$$(8) \quad \text{Let } Y_i = a + bX_i, \text{ show that } \bar{Y} = a + b\bar{X}$$

$$\text{and } \sum Y_i^2 = Na(a + 2b\bar{X}) + b^2 \sum X_i^2$$

(9) Assume that $X_i = 1$ for N_1 elements of a set and that $X_i = 0$ for N_0 of the elements. The total number of elements in the set is $N = N_1 + N_0$. Let $\frac{N_1}{N} = P$ and $\frac{N_0}{N} = Q$. Prove that

$$\frac{\sum (X_i - \bar{X})^2}{N} = PQ.$$

(10) $\sum (X_i - d)^2 = \sum (X_i - \bar{X})^2 + N(\bar{X} - d)^2$. Hint: Rewrite $(X_i - d)^2$ as $[(X_i - \bar{X}) + (\bar{X} - d)]^2$. Recall from elementary algebra that $(a+b)^2 = a^2 + 2ab + b^2$ and think of $(X_i - \bar{X})$ as a and of $(\bar{X} - d)$ as b . For what value of d is $\sum (X_i - d)^2$ a minimum?

1.5 DOUBLE INDEXES AND SUMMATION

When there is more than one characteristic for a set of elements, the different characteristics might be distinguished by using a different letter for each or by an index. For example, X_i and Y_i might represent the number of acres of wheat planted and the number of acres of wheat harvested on the i^{th} farm. Or, X_{ij} might be used where i is the index for the characteristics and j is the index for elements; that is, X_{ij} would be the value of characteristic X_i for the j^{th} element. However, when data on each of several characteristics for a set of elements are to be processed in the same way, it might not be necessary to use notation that distinguishes the characteristics. Thus, one might say

calculate $\frac{\sum (X_i - \bar{X})^2}{N-1}$ for all characteristics.

More than one index is needed when the elements are classified according to more than one criterion. For example, X_{ij} might represent the value of characteristic X for the j^{th} farm in the i^{th} county; or X_{ijk} might be the value of X for the k^{th} household in the j^{th} block in the i^{th} city. As another example, suppose the processing of data for farms involves classification of farms by size and type. We might let X_{ijk} represent the value of characteristic X for the k^{th} farm in the subset of farms classified as type j and size i . If N_{ij} is the number of farms classified

as type j and size i , then $\frac{\sum_{k=1}^N X_{ijk}}{N_{ij}} = \bar{X}_{ij}$. is the average value of X for

the subset of farms classified as type j and size i .

There are two general kinds of classification--cross classification and hierarchal or nested classification. Both kinds are often involved in the same problem. However, we will discuss each separately. An example of nested classification is farms within counties, counties within States, and States within regions. Cross classification means that the data can be arranged in two or more dimensions as illustrated in the next section.

1.5.1 CROSS CLASSIFICATION

As a specific illustration of cross classification and summation with two indexes, suppose we are working with the acreages of K crops on a set of N farms. Let X_{ij} represent the acreage of the i^{th} crop on the j^{th} farm where $i = 1, 2, \dots, K$ and $j = 1, 2, \dots, N$. In this case, the data could be arranged in a K by N matrix as follows:

Row (i)	Column (j)					Row total
	1		j		N	
1	X_{11}	...	X_{1j}	...	X_{1N}	$\sum_j X_{1j}$
.
i	X_{i1}	...	X_{ij}	...	X_{iN}	$\sum_j X_{ij}$
.
K	X_{K1}	...	X_{Kj}	...	X_{KN}	$\sum_j X_{Kj}$
Column total	$\sum_i X_{i1}$		$\sum_i X_{ij}$		$\sum_i X_{iN}$	$\sum_i \sum_j X_{ij}$

The expression $\sum_j^N X_{ij}$ (or $\sum_j X_{ij}$) means the sum of the values of X_{ij} for a fixed value of i . Thus, with reference to the matrix, $\sum_j X_{ij}$ is the total of the values of X in the i^{th} row; or, with reference to the example about farms and crop acreages, $\sum_j X_{ij}$ would be the total acreage on all farms of whatever the i^{th} crop is. Similarly, $\sum_i^K X_{ij}$ (or $\sum_i X_{ij}$) is the column total for the j^{th} column, which in the example is the total for the j^{th} farm of the acreages of the K crops under consideration. The sum of all values of X could be written as $\sum_{ij}^{KN} X_{ij}$ or $\sum_{ij} X_{ij}$.

Double summation means the sum of sums. Breaking a double sum into parts can be an important aid to understanding it. Here are two examples:

$$(1) \quad \sum_{ij}^{KN} X_{ij} = \sum_j^N X_{1j} + \sum_j^N X_{2j} + \dots + \sum_j^N X_{Kj} \quad (1.1)$$

With reference to the above matrix, Equation (1.1) expresses the grand total as the sum of row totals.

$$(2) \quad \sum_{ij}^{KN} X_{ij} (Y_{ij} + a) = \underbrace{\sum_j^N X_{1j} (Y_{1j} + a) + \dots + \sum_j^N X_{Kj} (Y_{Kj} + a)}_{\sum_j^N X_{1j} (Y_{1j} + a) = X_{11} (Y_{11} + a) + \dots + X_{1N} (Y_{1N} + a)} \quad (1.2)$$

In Equations (1.1) and (1.2) the double sum is written as the sum of K partial sums, that is, one partial sum for each value of i .

Exercise 1.6. (a) Write an equation similar to Equation (1.1) that expresses the grand total as the sum of column totals. (b) Involved in Equation (1.2) are KN terms, $X_{ij} (Y_{ij} + a)$. Write these terms in the form of a matrix.

The rules given in Section 1.4 also apply to double summation.

Thus,

$$\sum_{ij}^{KN} X_{ij} (Y_{ij} + a) = \sum_{ij}^{KN} X_{ij} Y_{ij} + a \sum_{ij}^{KN} X_{ij} \quad (1.3)$$

Study Equation (1.3) with reference to the matrix called for in Exercise 1.6(b). To fully understand Equation (1.3), you might need to write out intermediate steps for getting from the left-hand side to the right-hand side of the equation.

To simplify notation, a system of dot notation is commonly used, for example:

$$\sum_j X_{ij} = X_{i.}$$

$$\sum_i X_{ij} = X_{.j}$$

$$\sum_{ij} X_{ij} = X_{..}$$

The dot in $X_{i.}$ indicates that an index in addition to i is involved and $X_{i.}$ is interpreted as the sum of the values of X for a fixed value of i . Similarly, $X_{.j}$ is the sum of X for any fixed value of j , and $X_{..}$ represents a sum over both indexes. As stated above, averages are indicated by use of a bar. Thus $\bar{X}_{i.}$ is the average of X_{ij} for a fixed value of i , namely

$$\frac{\sum_{j=1}^N X_{ij}}{N} = \bar{X}_{i.} \text{ and } \bar{X}_{..} \text{ would represent the average of all values of } X_{ij},$$

$$\text{namely } \frac{\sum_{ij} X_{ij}}{NK}.$$

Here is an example of how the dot notation can simplify an algebraic expression. Suppose one wishes to refer to the sum of the squares of the row totals in the above matrix. This would be written as $\sum_i (X_{i.})^2$. The sum

of squares of the row means would be $\sum_1 (\bar{X}_{1.})^2$. Without the dot notation the

corresponding expressions would be $\sum_1 \sum_j (\sum_j X_{1j})^2$ and $\sum_1 \left[\frac{\sum_j X_{1j}}{N} \right]^2$. It is very

important that the parentheses be used correctly. For example, $\sum_1 \sum_j (\sum_j X_{1j})^2$ is

not the same as $\sum_1 \sum_j X_{1j}^2$. Incidentally, what is the difference between the

last two expressions?

Using the dot notation, the variance of the row means could be written as follows:

$$V(\bar{X}_{1.}) = \frac{\sum_1 (\bar{X}_{1.} - \bar{X}_{..})^2}{K-1} \quad (1.4)$$

where V stands for variance and $V(\bar{X}_{1.})$ is an expression for the variance of $\bar{X}_{1.}$. Without the dot notation, or something equivalent to it, a formula for the variance of the row means would look much more complicated.

Exercise 1.7. Write an equation, like Equation (1.4), for the variance of the column means.

Exercise 1.8. Given the following values of X_{1j}

i	j			
	1	2	3	4
1	8	11	9	14
2	10	13	11	14
3	12	15	10	17

Find the value of the following algebraic expressions:

<u>Expression</u>	<u>Answer</u>	<u>Expression</u>	<u>Answer</u>
(1) $\sum_j^N X_{1j}$	42	(9) $\sum_j^N (\bar{X}_{1j} - \bar{X}_{..})^2$	54
(2) $\frac{\sum_j^N X_{2j}}{N}$	12	(10) $\sum_{ij}^{KN} (X_{ij} - \bar{X}_{.j} - \bar{X}_{i.} + \bar{X}_{..})^2$	6
(3) $\bar{X}_{3.}$	13.5	(11) $\sum_{ij}^{KN} X_{ij}^2 - \frac{\left[\sum_{ij}^{KN} X_{ij} \right]^2}{KN}$	78
(4) $\sum X_{14}$	45	(12) $\frac{\sum_{i.}^K X_{i.}^2}{N} - \frac{\left[\sum_{ij}^{KN} X_{ij} \right]^2}{KN}$	18
(5) $\sum_{ij}^{KN} X_{ij}$	144	(13) $\sum_j^N (X_{1j} - \bar{X}_{1.})^2$	21
(6) $\bar{X}_{..}$	12	(14) $\sum_{ij}^{KN} (X_{ij} - \bar{X}_{i.})^2$	60
(7) $\sum_{ij}^{KN} (X_{ij} - \bar{X}_{..})^2$	78		
(8) $\sum_i^K N \sum_{1.} (\bar{X}_{i.} - \bar{X}_{..})^2$	18		

Illustration 1.1. To introduce another aspect of notation, refer to the matrix on Page 15 and suppose that the values of X in row one are to be multiplied by a_1 , the values of X in row two by a_2 , etc. The matrix would then be

$$\begin{array}{cccc}
 a_1 X_{11} & \dots & a_1 X_{1j} & \dots & a_1 X_{1N} \\
 \vdots & & \vdots & & \vdots \\
 \vdots & & \vdots & & \vdots \\
 a_i X_{i1} & \dots & a_i X_{ij} & \dots & a_i X_{iN} \\
 \vdots & & \vdots & & \vdots \\
 \vdots & & \vdots & & \vdots \\
 a_K X_{K1} & \dots & a_K X_{Kj} & \dots & a_K X_{KN}
 \end{array}$$

The general term can be written as $a_i X_{ij}$ because the index of a and the

index i in X_{ij} are the same. The total of all KN values of $a_i X_{ij}$ is

$\sum_{ij}^{KN} a_i X_{ij}$. Since a_i is constant with respect to summation involving j ,

we can place a_i ahead of the summation symbol \sum_j . That is, $\sum_{ij}^{KN} a_i X_{ij} = \sum_i a_i \sum_j X_{ij}$.

Exercise 1.9. Refer to the matrix of values of X_{ij} in Exercise 1.8.

Assume that $a_1 = -1$, $a_2 = 0$, and $a_3 = 1$.

Calculate:

$$(1) \sum_{ij} a_i X_{ij}$$

$$(2) \sum_{ij} \frac{a_i X_{ij}}{N}$$

$$(3) \sum_{ij} a_i X_{ij}^2 \quad \text{Answer: } -296$$

Show algebraically that:

$$(4) \sum_{ij} a_i X_{ij} = \sum_j X_{3j} - \sum_j X_{1j}$$

$$(5) \sum_{ij} \frac{a_i X_{ij}}{N} = \bar{X}_3 - \bar{X}_1$$

$$(6) \sum_{ij} a_i X_{ij}^2 = \sum_j X_{3j}^2 - \sum_j X_{1j}^2$$

Exercise 1.10. Study the following equation and if necessary write

the summations as series to be satisfied that the equation is correct:

$$\sum_{ij}^{KN} (aX_{ij} + bY_{ij}) = a \sum_{ij} X_{ij} + b \sum_{ij} Y_{ij}$$

Illustration 1.2. Suppose

$$Y_{ij} = X_{ij} + a_i + b_j + c \quad \text{where } i = 1, 2, \dots, K \text{ and } j = 1, 2, \dots, N$$

The values of Y_{ij} can be arranged in matrix format as follows:

$$\begin{array}{ccccccc}
 Y_{11} = X_{11} + a_1 + b_1 + c & . & . & . & . & . & Y_{1N} = X_{1N} + a_1 + b_N + c \\
 & & & & & & \vdots \\
 & & Y_{ij} = X_{ij} + a_i + b_j + c & & & & \vdots \\
 & & & & & & \vdots \\
 Y_{K1} = X_{K1} + a_K + b_1 + c & . & . & . & . & . & Y_{KN} = X_{KN} + a_K + b_N + c
 \end{array}$$

Notice that a_i is a quantity that varies from row to row but is constant within a row and that b_j varies from column to column but is constant within a column. Applying the rules regarding the summation symbols we have

$$\begin{aligned}
 \sum_j Y_{ij} &= \sum_j (X_{ij} + a_i + b_j + c) \\
 &= \sum_j X_{ij} + Na_i + \sum_j b_j + Nc \\
 \sum_i Y_{ij} &= \sum_i (X_{ij} + a_i + b_j + c) \\
 &= \sum_i X_{ij} + \sum_i a_i + Kb_j + Kc \\
 \sum_{ij} Y_{ij} &= \sum_{ij} (X_{ij} + a_i + b_j + c) \\
 &= \sum_{ij} X_{ij} + N \sum_i a_i + K \sum_j b_j + KNC
 \end{aligned}$$

Illustration 1.3. We have noted that $\sum (X_i Y_i)$ does not equal

$(\sum X_i)(\sum Y_i)$. (See (1) and (2) in Exercise 1.3, and (5) on Page 12). But,

$\sum_{ij} X_i Y_j = (\sum_i X_i)(\sum_j Y_j)$ where $i = 1, 2, \dots, K$ and $j = 1, 2, \dots, N$. This becomes

clear if we write the terms of $\sum_{ij} X_i Y_j$ in matrix format as follows:

	Row Totals
$X_1 Y_1 + X_1 Y_2 + \dots + X_1 Y_N$	$X_1 \sum_j Y_j$
$+ X_2 Y_1 + X_2 Y_2 + \dots + X_2 Y_N$	$X_2 \sum_j Y_j$
\vdots	
$+ X_K Y_1 + X_K Y_2 + \dots + X_K Y_N$	$X_K \sum_j Y_j$
$\sum_{ij} X_i Y_j$	$\sum_i X_i \sum_j Y_j$

The sum of the terms in each row is shown at the right. The sum of these row totals is $X_1 \Sigma Y_j + \dots + X_K \Sigma Y_j = (X_1 + \dots + X_K) \Sigma Y_j = \Sigma X_i \Sigma Y_j$. One could get the same final result by adding the columns first. Very often intermediate summations are of primary interest.

Exercise 1.11. Verify that $\sum_{ij} X_i Y_j = (\Sigma X_i)(\Sigma Y_j)$ using the values of X and Y in Exercise 1.3. In Exercise 1.3 the subscript of X and the subscript of Y were the same index. In the expression $\sum_{ij} X_i Y_j$ that is no longer the case.

Exercise 1.12. Prove the following:

$$(1) \quad \sum_{ij} (a_i X_{ij} + b_j)^2 = \sum_i a_i^2 \sum_j X_{ij}^2 + 2 \sum_i a_i \sum_j b_j X_{ij} + K \sum_j b_j^2$$

$$(2) \quad \sum_{ij} a_i (X_{ij} - \bar{X}_{i.})^2 = \sum_i a_i \sum_j X_{ij}^2 - N \sum_i a_i \bar{X}_{i.}^2$$

$$(3) \quad \sum_{ij} a_i (X_{ij} - \bar{X}_{i.})(Y_{ij} - \bar{Y}_{.j}) = \sum_i a_i \sum_j X_{ij} Y_{ij} - N \sum_i a_i \bar{X}_{i.} \bar{Y}_{.j}$$

1.5.2 HIERARCHAL OR NESTED CLASSIFICATION

A double index does not necessarily imply that a meaningful cross classification of the data can be made. For example, X_{ij} might represent the value of X for the j^{th} farm in the i^{th} county. In this case, j simply identifies a farm within a county. There is no correspondence, for example, between farm number 5 in one county and farm number 5 in another. In fact the total number of farms varies from county to county. Suppose there are K counties and N_i farms in the i^{th} county. The total of X for the i^{th} county could be expressed as $X_{i.} = \sum_j X_{ij}$. In the present case $\sum_i X_{ij}$ is meaningless. The total of all values of X is $\sum_{ij} X_{ij}$.

When the classification is nested, the order of the subscripts (indexes) and the order of the summation symbols from left to right should be from the highest to lowest order of classification. Thus in the above example the index for farms was on the right and the summation symbol

involving this index is also on the right. In the expression $\sum_{ij}^{KN} X_{ij}$, summation with respect to i cannot take place before summation with regard to j . On the other hand, when the classification is cross classification the summations can be performed in either order.

In the example of K counties and N_i farms in the i^{th} county, and in similar examples, you may think of the data as being arranged in rows (or columns):

$$\begin{array}{l} X_{11}, X_{12}, \dots, X_{1N_1} \\ X_{21}, X_{22}, \dots, X_{2N_2} \\ \vdots \\ X_{K1}, X_{K2}, \dots, X_{KN_K} \end{array}$$

Here are two double sums taken apart for inspection:

$$(1) \quad \sum_{ij}^{KN} (X_{ij} - \bar{X}_{..})^2 = \underbrace{\sum_j^{N_1} (X_{1j} - \bar{X}_{..})^2}_{\text{}} + \dots + \sum_j^{N_K} (X_{Kj} - \bar{X}_{..})^2 \quad (1.5)$$

$$\sum_j^{N_1} (X_{1j} - \bar{X}_{..})^2 = (X_{11} - \bar{X}_{..})^2 + \dots + (X_{1N_1} - \bar{X}_{..})^2$$

Equation (1.5) is the sum of squares of the deviations, $(X_{ij} - \bar{X}_{..})$, of all values of X_{ij} from the overall mean. There are $\sum_i^K N_i$ values of X_{ij} , and

$\bar{X}_{..} = \frac{\sum_{i=1}^K \sum_{j=1}^{N_i} X_{ij}}{\sum_{i=1}^K N_i}$. If there was no interest in identifying the data by counties,

a single index would be sufficient. Equation (1.5) would then be $\sum_{i=1}^N (X_i - \bar{X})^2$.

$$(2) \quad \sum_{i,j}^{KN} (X_{ij} - \bar{X}_{i.})^2 = \underbrace{\sum_{j=1}^{N_1} (X_{1j} - \bar{X}_{1.})^2 + \dots + \sum_{j=1}^{N_K} (X_{Kj} - \bar{X}_{K.})^2}_{(1.6)}$$

$$\sum_{j=1}^{N_1} (X_{1j} - \bar{X}_{1.})^2 = (X_{11} - \bar{X}_{1.})^2 + \dots + (X_{1N_1} - \bar{X}_{1.})^2$$

With reference to Equation (1.6) do you recognize $\sum_{j=1}^{N_1} (X_{1j} - \bar{X}_{1.})^2$? It involves only the subset of elements for which $i = 1$, namely $X_{11}, X_{12}, \dots, X_{1N_1}$. Note that $\bar{X}_{1.}$ is the average value of X in this subset. Hence, $\sum_{j=1}^{N_1} (X_{1j} - \bar{X}_{1.})^2$ is the sum of the squares of the deviations of the X 's in this subset from the subset mean. The double sum is the sum of K terms and each of the K terms is a sum of squares for a subset of X 's, the index for the subsets being i .

Exercise 1.13. Let X_{ij} represent the value of X for the j^{th} farm in the i^{th} county. Also, let K be the number of counties and N_i be the number of farms in the i^{th} county. Suppose the values of X are as follows:

$$\begin{array}{lll} X_{11} = 3 & X_{12} = 1 & X_{13} = 5 \\ X_{21} = 4 & X_{22} = 6 & \\ X_{31} = 0 & X_{32} = 5 & X_{33} = 1 \quad X_{34} = 2 \end{array}$$

Find the value of the following expressions:

Expression

Answer

(1) $\sum_{i=1}^K N_i$

9

<u>Expression (Continued)</u>	<u>Answer</u>
(2) $\sum_{ij}^{KN} x_{ij}$	27
(3) $x_{..}$ and $\bar{x}_{..}$	27 3
(4) $\sum_j^N x_{1j} = x_{1.}$	9
(5) $x_{2.}$ and $x_{3.}$	10 8
(6) $\bar{x}_{1.}$, $\bar{x}_{2.}$, and $\bar{x}_{3.}$	3 5 2
(7) $\frac{\sum_i^N \bar{x}_{i.}}{\sum_i^N 1}$	3
(8) $\sum_i^K \left(\sum_j^N x_{ij} \right)^2$ or $\sum_i^K x_{i.}^2$	245
(9) $\sum_{ij} (x_{ij} - \bar{x}_{..})^2$	36
(10) $\sum_j^N (x_{1j} - \bar{x}_{1.})^2$	8
(11) $\sum_j^N (x_{ij} - \bar{x}_{i.})^2$	8, 2, and 14 for $i = 1, 2,$ and 3 respectively
(12) $\sum_{ij}^{KN} (x_{ij} - \bar{x}_{1.})^2$	24
(13) $\sum_i^K (\bar{x}_{i.} - \bar{x}_{..})^2$	12
(14) $\sum_i^K \frac{\left[\sum_j^N x_{ij} \right]^2}{N_i} - \frac{\left[\sum_{ij}^{KN} x_{ij} \right]^2}{\sum_j^N N_j}$	12
(15) $\sum_i^K \bar{x}_{i.}^2 - N \bar{x}_{..}^2$	12

Expressions (14) and (15) in Exercise 1.13 are symbolic representations of the same thing. By definition

$$\sum_j^N x_{1j} = x_{1.}, \quad \sum_{ij}^{KN} x_{ij} = x_{..}, \quad \text{and} \quad \sum_i^K N_i = N$$

Substitution in (14) gives

$$\sum_i^K \frac{x_{1.}^2}{N_i} - \frac{x_{..}^2}{N} \quad (1.7)$$

Also by definition $\frac{x_{1.}}{N_i} = \bar{x}_{1.}$ and $\frac{x_{..}}{N} = \bar{x}_{..}$. Therefore $\frac{x_{1.}^2}{N_i} = N_i \bar{x}_{1.}^2$ and

$\frac{x_{..}^2}{N} = N \bar{x}_{..}^2$. Hence, by substitution, Equation (1.7) becomes $\sum_i^K N_i \bar{x}_{1.}^2 - N \bar{x}_{..}^2$.

Exercise 1.14. Prove the following:

- (1) $\sum_{ij}^{KN} x_{1.} x_{ij} = \sum_i^K x_{1.}^2$
- (2) $\sum_{ij}^{KN} \bar{x}_{1.} (x_{ij} - \bar{x}_{1.}) = 0$
- (3) $\sum_i^K N_i (\bar{x}_{1.} - \bar{x}_{..})^2 = \sum_i^K N_i \bar{x}_{1.}^2 - N \bar{x}_{..}^2$

Note that this equates (13) and (15) in Exercise 1.13.

The proof is similar to the proof called for in part (5) of Exercise 1.5.

$$(4) \quad \sum_{ij}^{KN} (a_i x_{ij} - b_i)^2 = \sum_i^K a_i^2 \sum_j^N x_{ij}^2 - 2 \sum_i^K a_i b_i x_{i.} + \sum_i^K N_i b_i^2$$

1.6 THE SQUARE OF A SUM

In statistics, it is often necessary to work algebraically with the square of a sum. For example,

$$(\sum x_i)^2 = (x_1 + x_2 + \dots + x_N)^2 = x_1^2 + x_1 x_2 + \dots + x_2^2 + x_2 x_1 + \dots + x_N^2 + x_N x_1 + \dots$$

The terms in the square of the sum can be written in matrix form as follows:

$$\begin{array}{cccccc}
 X_1 X_1 & X_1 X_2 & \dots & X_1 X_j & \dots & X_1 X_N \\
 X_2 X_1 & X_2 X_2 & \dots & X_2 X_j & \dots & X_2 X_N \\
 \vdots & \vdots & & \vdots & & \vdots \\
 \vdots & \vdots & & \vdots & & \vdots \\
 X_i X_1 & X_i X_2 & \dots & X_i X_j & \dots & X_i X_N \\
 \vdots & \vdots & & \vdots & & \vdots \\
 \vdots & \vdots & & \vdots & & \vdots \\
 X_N X_1 & X_N X_2 & \dots & X_N X_j & \dots & X_N X_N
 \end{array}$$

The general term in this matrix is $X_i X_j$ where X_i and X_j come from the same set of X 's, namely, X_1, \dots, X_N . Hence, i and j are indexes of the same set. Note that the terms along the main diagonal are the squares of the value of X and could be written as $\sum X_i^2$. That is, on the main diagonal $i = j$ and $X_i X_j = X_i X_i = X_i^2$. The remaining terms are all products of one value of X with some other value of X . For these terms the indexes are never equal. Therefore, the sum of all terms not on the main diagonal can be expressed as $\sum_{i \neq j} X_i X_j$ where $i \neq j$ is used to express the fact that the summation includes all terms where i is not equal to j , that is, all terms other than those on the main diagonal. Hence, we have shown that $(\sum X_i)^2 = \sum X_i^2 + \sum_{i \neq j} X_i X_j$.

Notice the symmetry of terms above and below the main diagonal:

$X_1 X_2 = X_2 X_1, X_1 X_3 = X_3 X_1$, etc. When symmetry like this occurs, instead of

$\sum_{i \neq j} X_i X_j$ you might see an equivalent expression $2 \sum_{i < j} X_i X_j$. The sum of all

terms above the main diagonal is $\sum_{i < j} X_i X_j$. Owing to the symmetry, the sum

of the terms below the main diagonal is the same. Therefore, $\sum_{i \neq j} X_i X_j = 2 \sum_{i < j} X_i X_j$.

Exercise 1.15. Express the terms of $[\sum_{i=1}^4 X_i]^2 = [X_1 + X_2 + X_3 + X_4]^2$ in matrix format. Let $X_1 = 2$, $X_2 = 0$, $X_3 = 5$, and $X_4 = 7$. Compute the values of $\sum X_i^2$, $2 \sum_{i < j} X_i X_j$, and $[\sum X_i]^2$. Show that $[\sum X_i]^2 = \sum X_i^2 + 2 \sum_{i < j} X_i X_j$.

An important result, which we will use in Chapter 3, follows from the fact that

$$[\sum X_i]^2 = \sum X_i^2 + \sum_{i \neq j} X_i X_j \quad (1.8)$$

Let $X_i = Y_i - \bar{Y}$. Substituting $(Y_i - \bar{Y})$ for X_i in Equation 1.8 we have

$$[\sum (Y_i - \bar{Y})]^2 = \sum (Y_i - \bar{Y})^2 + \sum_{i \neq j} (Y_i - \bar{Y})(Y_j - \bar{Y})$$

We know that $[\sum (Y_i - \bar{Y})]^2 = 0$ because $\sum (Y_i - \bar{Y}) = 0$. Therefore,

$$\sum (Y_i - \bar{Y})^2 + \sum_{i \neq j} (Y_i - \bar{Y})(Y_j - \bar{Y}) = 0$$

It follows that $\sum_{i \neq j} (Y_i - \bar{Y})(Y_j - \bar{Y}) = -\sum (Y_i - \bar{Y})^2$ (1.9)

Exercise 1.16. Consider

$$\begin{aligned} \sum_{i \neq j} (Y_i - \bar{Y})(Y_j - \bar{Y}) &= \sum_{i \neq j} (Y_i Y_j - \bar{Y} Y_i - \bar{Y} Y_j + \bar{Y}^2) \\ &= \sum_{i \neq j} Y_i Y_j - \bar{Y} \sum_{i \neq j} Y_i - \bar{Y} \sum_{i \neq j} Y_j + \sum_{i \neq j} \bar{Y}^2 \end{aligned}$$

Do you agree that $\sum_{i \neq j} \bar{Y}^2 = N(N-1)\bar{Y}^2$? With reference to the matrix layout,

\bar{Y}^2 appears N^2 times but the specification is $i \neq j$ so we do not want to count the N times that \bar{Y}^2 is on the main diagonal. Try finding the values of $\sum_{i \neq j} X_i$ and $\sum_{i \neq j} X_j$ and then show that

$$\sum_{i=j} (Y_i - \bar{Y})(Y_j - \bar{Y}) = \sum_{i \neq j} Y_i Y_j - N(N-1)\bar{Y}^2$$

Hint: Refer to a matrix layout. In $\sum_{i \neq j} Y_i$ how many times does Y_1 appear?

Does Y_2 appear the same number of times?

1.7 SUMS OF SQUARES

For various reasons statisticians are interested in components of variation, that is, measuring the amount of variation attributable to each of more than one source. This involves computing sums of squares that correspond to the different sources of variation that are of interest. We will discuss a simple example of nested classification and a simple example of cross classification.

1.7.1 NESTED CLASSIFICATION

To be somewhat specific, reference is made to the example of K counties and N_i farms in the i^{th} county. The sum of the squares of the deviations of X_{ij} and $\bar{X}_{i.}$ can be divided into two parts as shown by the following formula:

$$\sum_{i,j}^{KN} (X_{ij} - \bar{X}_{..})^2 = \sum_i^K N_i (\bar{X}_{i.} - \bar{X}_{..})^2 + \sum_{i,j}^{KN} (X_{ij} - \bar{X}_{i.})^2 \quad (1.10)$$

The quantity on the left-hand side of Equation (1.10) is called the total sum of squares. In Exercise 1.13, Part (9), the total sum of squares was 36.

The first quantity on the right-hand side of the equation involves the squares of $(\bar{X}_{i.} - \bar{X}_{..})$, which are deviations of the class means from the over-all mean. It is called the between class sum of squares or with reference to the example the between county sum of squares. In Exercise 1.13, Part (13), the between county sum of squares was computed. The answer was 12.

The last term is called the within sum of squares because it involves deviations within the classes from the class means. It was presented previously. See Equation (1.6) and the discussion pertaining to it. In Exercise 1.13, the within class sum of squares was 24, which was calculated in Part (12). Thus, from Exercise 1.13, we have the total sum of squares, 36, which equals the between, 12, plus the within, 24. This verifies Equation (1.10).

The proof of Equation 1.10 is easy if one gets started correctly. Write $X_{1j} - \bar{X}_{..} = (X_{1j} - \bar{X}_{1.}) + (\bar{X}_{1.} - \bar{X}_{..})$. This simple technique of adding and subtracting $\bar{X}_{1.}$ divides the deviation $(X_{1j} - \bar{X}_{..})$ into two parts. The proof proceeds as follows:

$$\begin{aligned} \sum_{1j}^{KN} (X_{1j} - \bar{X}_{..})^2 &= \sum_{1j} [(X_{1j} - \bar{X}_{1.}) + (\bar{X}_{1.} - \bar{X}_{..})]^2 \\ &= \sum_{1j} [(X_{1j} - \bar{X}_{1.})^2 + 2(X_{1j} - \bar{X}_{1.})(\bar{X}_{1.} - \bar{X}_{..}) + (\bar{X}_{1.} - \bar{X}_{..})^2] \\ &= \sum_{1j} (X_{1j} - \bar{X}_{1.})^2 + 2 \sum_{1j} (X_{1j} - \bar{X}_{1.})(\bar{X}_{1.} - \bar{X}_{..}) + \sum_{1j} (\bar{X}_{1.} - \bar{X}_{..})^2 \end{aligned}$$

Exercise 1.17. Show that $\sum_{1j}^{KN} (X_{1j} - \bar{X}_{1.})(\bar{X}_{1.} - \bar{X}_{..}) = 0$

$$\text{and that } \sum_{1j}^{KN} (\bar{X}_{1.} - \bar{X}_{..})^2 = \sum_i^K N_i (\bar{X}_{1.} - \bar{X}_{..})^2$$

Completion of Exercise 1.17 completes the proof.

Equation (1.10) is written in a form which displays its meaning rather than in a form that is most useful for computational purposes. For computation purposes, the following relationships are commonly used:

$$\text{Total} = \sum_{1j}^{KN} (X_{1j} - \bar{X}_{..})^2 = \sum_{1j} X_{1j}^2 - N \bar{X}_{..}^2$$

$$\text{Between} = \sum_1^K N_1 (\bar{X}_{1.} - \bar{X}_{..})^2 = \sum_1^K N_1 \bar{X}_{1.}^2 - N \bar{X}_{..}^2$$

$$\text{Within} = \sum_{1j}^{KN} (X_{1j} - \bar{X}_{1.})^2 = \sum_{1j} \sum X_{1j}^2 - \sum_1^K N_1 \bar{X}_{1.}^2$$

$$\text{where } N = \sum_1^K N_1, \quad \bar{X}_{1.} = \frac{\sum_j X_{1j}}{N_1}, \text{ and } \bar{X}_{..} = \frac{\sum_{1j} X_{1j}}{N}$$

Notice that the major part of arithmetic reduces to calculating $\sum_{1j}^{KN} X_{1j}^2$,

$\sum_1^K N_1 \bar{X}_{1.}^2$, and $N \bar{X}_{..}^2$. There are variations of this that one might use. For example, one could use $\sum_1^K \frac{X_{1.}^2}{N_1}$ instead of $\sum_1^K N_1 \bar{X}_{1.}^2$.

Exercise 1.18. Show that

$$\sum_{1j}^{KN} (X_{1j} - \bar{X}_{1.})^2 = \sum_{1j} \sum X_{1j}^2 - \sum_1^K N_1 \bar{X}_{1.}^2$$

A special case that is useful occurs when $N_1 = 2$. The within sum of squares becomes

$$\sum_{1j}^K 2 (X_{1j} - \bar{X}_{1.})^2 = \sum_1^K [(X_{11} - \bar{X}_{1.})^2 + (X_{12} - \bar{X}_{1.})^2]$$

Since $\bar{X}_{1.} = \frac{X_{11} + X_{12}}{2}$ it is easy to show that

$$(X_{11} - \bar{X}_{1.})^2 = \frac{1}{4} (X_{11} - X_{12})^2$$

$$\text{and } (X_{12} - \bar{X}_{1.})^2 = \frac{1}{4} (X_{11} - X_{12})^2$$

Therefore the within sum of squares is

$$\frac{1}{2} \sum_1^K (X_{11} - X_{12})^2$$

which is a convenient form for computation.

1.7.2 CROSS CLASSIFICATION

Reference is made to the matrix on Page 15 and to Exercise 1.8. The total sum of squares can be divided into three parts as shown by the following formula:

$$\sum_{ij}^{KN} (X_{ij} - \bar{X}_{..})^2 = \sum_i^K (\bar{X}_{i.} - \bar{X}_{..})^2 + \sum_j^N (\bar{X}_{.j} - \bar{X}_{..})^2 + \sum_{ij}^{KN} (X_{ij} - \bar{X}_{i.} - \bar{X}_{.j} + \bar{X}_{..})^2 \quad (1.11)$$

Turn to Exercise 1.8 and find the total sum of squares and the three parts. They are:

	<u>Sum of Squares</u>
Total	78
Rows	18
Columns	54
Remainder	6

The three parts add to the total which verifies Equation (1.11). In Exercise 1.8, the sum of squares called remainder was computed directly (see Part (10) of Exercise 1.8). In practice, the remainder sum of squares is usually obtained by subtracting the row and column sum of squares from the total.

Again, the proof of Equation (1.11) is not difficult if one makes the right start. In this case the deviation, $(X_{ij} - \bar{X}_{..})$, is divided into three parts by adding and subtracting $\bar{X}_{i.}$ and $\bar{X}_{.j}$ as follows:

$$(X_{ij} - \bar{X}_{..}) = (\bar{X}_{i.} - \bar{X}_{..}) + (\bar{X}_{.j} - \bar{X}_{..}) + (X_{ij} - \bar{X}_{i.} - \bar{X}_{.j} + \bar{X}_{..}) \quad (1.12)$$

Exercise 1.19. Prove Equation (1.11) by squaring both sides of Equation (1.12) and then doing the summation. The proof is mostly a matter of showing that the sums of the terms which are products (not squares) are zero.

For example, showing that $\sum_{ij}^{KN} (\bar{X}_{i.} - \bar{X}_{..})(X_{ij} - \bar{X}_{i.} - \bar{X}_{.j} + \bar{X}_{..}) = 0$.

CHAPTER II. RANDOM VARIABLES AND PROBABILITY

2.1 RANDOM VARIABLES

The word "random" has a wide variety of meanings. Its use in such terms as "random events," "random variable," or "random sample," however, implies a random process such that the probability of an event occurring is known a priori. To select a random sample of elements from a population, tables of random numbers are used. There are various ways of using such tables to make a random selection so any given element will have a specified probability of being selected.

The theory of probability sampling is founded on the concept of a random variable which is a variable that, by chance, might equal any one of a defined set of values. The value of a random variable on any particular occasion is determined by a random process in such a way that the chance (probability) of its being equal to any specified value in the set is known. This is in accord with the definition of a probability sample which states that every element of the population must have a known probability (greater than zero) of being selected. A primary purpose of this chapter is to present an elementary, minimum introduction or review of probability as background for the next chapter on expected values of a random variable. This leads to a theoretical basis for sampling and for evaluating the accuracy of estimates from a probability-sample survey.

In sampling theory, we usually start with an assumed population of N elements and a measurement for each element of some characteristic X . A typical mathematical representation of the N measurements or values is $X_1, \dots, X_1, \dots, X_N$ where X_i is the value of the characteristic X for the i^{th} element. Associated with the i^{th} element is a probability P_i , which is the probability of obtaining it when one element is selected at random from the

set of N . The P_i 's will be called selection probabilities. If each element has an equal chance of selection, $P_i = \frac{1}{N}$. The P_i 's need not be equal, but we will specify that each $P_i > 0$. When referring to the probability of X being equal to X_i we will use $P(X_i)$ instead of P_i .

We need to be aware of a distinction between selection probability and inclusion probability, the latter being the probability of an element being included in a sample. In this chapter, much of the discussion is oriented toward selection probabilities because of its relevance to finding expected values of estimates from samples of various kinds.

Definition 2.1. A random variable is a variable that can equal any value X_i , in a defined set, with a probability $P(X_i)$.

When an element is selected at random from a population and a measurement of a characteristic of it is made, the value obtained is a random variable. As we shall see later, if a sample of elements is selected at random from a population, the sample average and other quantities calculated from the sample are random variables.

Illustration 2.1. One of the most familiar examples of a random variable is the number of dots that happen to be on the top side of a die when it comes to rest after a toss. This also illustrates the concept of probability that we are interested in; namely, the relative frequency with which a particular outcome will occur in reference to a defined set of possible outcomes. With a die there are six possible outcomes and we expect each to occur with the same frequency, $1/6$, assuming the die is tossed a very large or infinite number of times. Implicit in a statement that each side of a die has a probability of $1/6$ of being the top side are some assumptions about the physical structure of the die and the "randomness" of the toss.

The additive and multiplicative laws of probability can be stated in several ways depending upon the context in which they are to be used. In sampling, our interest is primarily in the outcome of one random selection or of a series of random selections that yields a probability sample. Hence, the rules or theorems for the addition or multiplication of probabilities will be stated or discussed only in the context of probability sampling.

2.2 ADDITION OF PROBABILITIES

Assume a population of N elements and a variable X which has a value X_i for the i^{th} element. That is, we have a set of values of X , namely $X_1, \dots, X_i, \dots, X_N$. Let $P_1, \dots, P_i, \dots, P_N$ be a set of selection probabilities where P_i is the probability of selecting the i^{th} element when a random selection is made. We specify that each P_i must be greater than zero and

that $\sum_{i=1}^N P_i = 1$. When an element is selected at random, the probability that it is either the i^{th} element or the j^{th} element is $P_i + P_j$. This addition rule can be stated more generally. Let P_s be the sum of the selection probabilities for the elements in a subset of the N elements. When a random selection is made from the whole set, P_s is the probability that the element selected is from the subset and $1 - P_s$ is the probability that it is not from the subset. With reference to the variable X , let $P(X_i)$ represent the probability that X equals X_i . Then $P(X_i) + P(X_j)$ represents the probability that X equals either X_i or X_j ; and $P_s(X)$ could be used to represent the probability that X is equal to one of the values in the subset.

Before adding (or subtracting) probabilities one should determine whether the events are mutually exclusive and whether all possible events have been accounted for. Consider two subsets of elements, subset A and

subset B, of a population of N elements. Suppose one element is selected at random. What is the probability that the selected element is a member of either subset A or subset B? Let $P(A)$ be the probability that the selected element is from subset A; that is, $P(A)$ is the sum of the selection probabilities for elements in subset A. $P(B)$ is defined similarly. If the two subsets are mutually exclusive, which means that no element is in both subsets, the probability that the element selected is from either subset A or subset B is $P(A) + P(B)$. If some elements are in both subsets, see Figure 2.1, then event A (which is the selected element being a member of subset A) and event B (which is the selected element being a member of subset B) are not mutually exclusive events. Elements included in both subsets are counted once in $P(A)$ and once in $P(B)$. Therefore, we must subtract $P(A,B)$ from $P(A) + P(B)$ where $P(A,B)$ is the sum of the probabilities for the elements that belong to both subset A and subset B. Thus,

$$P(A \text{ or } B) = P(A) + P(B) - P(A,B)$$

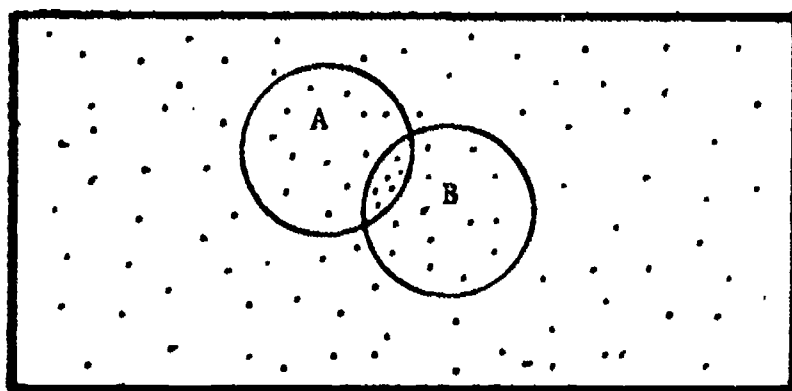


Figure 2.1

To summarize, the additive law of probability as used above could be stated as follows: If A and B are subsets of a set of all possible outcomes that could occur as a result of a random trial or selection, the probability

that the outcome is in subset A or in subset B is equal to the probability that the outcome is in A plus the probability that it is in B minus the probability that it is in both A and B.

The additive law of probability extends without difficulty to three or more subsets. Draw a figure like Figure 2.1 with three subsets so that some points are common to all three subsets.. Observe that the additive law extends to three subsets as follows:

$$P(A \text{ or } B \text{ or } C) = P(A) + P(B) + P(C) - P(A, B) - P(A, C) - P(B, C) + P(A, B, C)$$

As a case for further discussion purposes, assume a population of N elements and two criteria for classification. A two-way classification of the elements could be displayed in the format of Table 2.1.

Table 2.1--A two-way classification of N elements

Y class	X class			Total
	1	...	j ... s	
1	N_{11}, P_{11}	...	N_{1j}, P_{1j} ... N_{1s}, P_{1s}	$N_{1.}, P_{1.}$
.
.
i	N_{i1}, P_{i1}	...	N_{ij}, P_{ij} ... N_{is}, P_{is}	$N_{i.}, P_{i.}$
.
.
t	N_{t1}, P_{t1}	...	N_{tj}, P_{tj} ... N_{ts}, P_{ts}	$N_{t.}, P_{t.}$
Total	$N_{.1}$	$N_{.j}$	$N_{.s}$	$N, P=1$

The columns represent a classification of the elements in terms of criterion X; the rows represent a classification in terms of criterion Y; N_{ij} is the number of elements in X class j and Y class i; and P_{ij} is the sum of the

selection probabilities for the elements in X class j and Y class i . Any one of the N elements can be classified in one and only one of the t times s cells.

Suppose one element from the population of N is selected. According to the additive law of probability we can state that

$\sum_i P_{ij} = P_{.j}$ is the probability that the element selected is from X class j , and

$\sum_j P_{ij} = P_{i.}$ is the probability that the element selected is from Y class i , where

P_{ij} is the probability that the element selected is from (belongs to both) X class j and Y class i .

The probabilities $P_{.j}$ and $P_{i.}$ are called marginal probabilities.

The probability that one randomly selected element is from X class j or from Y-class i is $P_{.j} + P_{i.} - P_{ij}$. (The answer is not $P_{.j} + P_{i.}$ because in $P_{.j} + P_{i.}$ there are N_{ij} elements in X class j and Y class i that are counted twice.)

If the probabilities of selection are equal, $P_{ij} = \frac{N_{ij}}{N}$, $P_{.j} = \frac{N_{.j}}{N}$, and $P_{i.} = \frac{N_{i.}}{N}$.

Illustration 2.2. Suppose there are 5,000 students in a university. Assume there are 1,600 freshmen, 1,400 sophomores, and 500 students living in dormitory A. From a list of the 5,000 students, one student is selected at random. Assuming each student had an equal chance of selection, the probability that the selected student is a freshman is $\frac{1600}{5000}$, that he is a sophomore is $\frac{1400}{5000}$, and that he is either a freshman or a sophomore is $\frac{1600}{5000} + \frac{1400}{5000}$. Also, the probability that the selected student lives in dormitory A

is $\frac{500}{5000}$. But, what is the probability that the selected student is either a freshman or lives in dormitory A? The question involves two classifications: one pertaining to the student's class and the other to where the student lives. The information given about the 5000 students could be arranged as follows:

Dormitory	Class			Total
	Freshmen	Sophomores	Others	
A				500
Other				4500
Total	1600	1400	2000	5000

From the above format, one can readily observe that the answer to the question depends upon how many freshmen live in dormitory A. If the problem had stated that 200 freshmen live in dormitory A, the answer would have been $\frac{1600}{5000} + \frac{500}{5000} - \frac{200}{5000}$.

Statements about probability need to be made and interpreted with great care. For example, it is not correct to say that a student has a probability of 0.1 of living in dormitory A simply because 500 students out of 5000 live in A. Unless students are assigned to dormitories by a random process with known probabilities there is no basis for stating a student's probability of living in (being assigned to) dormitory A. We are considering the outcome of a random selection.

Exercise 2.1. Suppose one has the following information about a population of 1000 farms:

600 produce corn

500 produce soybeans

300 produce wheat

100 produce wheat and corn

200 have one or more cows

all farms that have cows also produce corn

200 farms do not produce any crops

One farm is selected at random with equal probability from the list of 1000. What is the probability that the selected farm,

- (a) produces corn? Answer: 0.6
- (b) does not produce wheat?
- (c) produces corn but no wheat? Answer: 0.5
- (d) produces corn or wheat but not both?
- (e) has no cows? Answer: 0.8
- (f) produces corn or soybeans?
- (g) produces corn and has no cows? Answer: 0.4
- (h) produces either corn, cows, or both?
- (i) does not produce corn or wheat?

One of the above questions cannot be answered.

Exercise 2.2. Assume a population of 10 elements and selection probabilities as follows:

<u>Element</u>	<u>X_1</u>	<u>P_1</u>	<u>Element</u>	<u>X_1</u>	<u>P_1</u>
1	2	.05	6	11	.15
2	7	.10	7	2	.20
3	12	.08	8	8	.05
4	0	.02	9	6	.05
5	8	.20	10	3	.10

One element is selected at random with probability P_i .

Find:

- (a) $P(X=2)$, the probability that $X = 2$.
- (b) $P(X>10)$, the probability that X is greater than 10.
- (c) $P(X\leq 2)$, the probability that X is equal to or less than 2.
- (d) $P(3<X<10)$, the probability that X is greater than 3 and less than 10
- (e) $P(X\leq 3 \text{ or } X\geq 10)$, the probability that X is either equal to or less than 3 or is equal to or greater than 10.

Note: The answer to (d) and the answer to (e) should add to 1.

So far, we have been discussing the probability of an event occurring as a result of a single random selection. When more than one random selection occurs simultaneously or in succession the multiplicative law of probability is useful.

2.3 MULTIPLICATION OF PROBABILITIES

Assume a population of N elements and selection probabilities $P_1, \dots, P_1, \dots, P_N$. Each P_i is greater than zero and $\sum_{i=1}^N P_i = 1$. Suppose two elements are selected but before the second selection is made the first element selected is returned to the population. In this case the outcome of the first selection does not change the selection probabilities for the second selection. The two selections (events) are independent. The probability of selecting the i^{th} element first and the j^{th} element second is, $P_i P_j$, the product of the selection probabilities P_i and P_j . If a selected element is not returned to the population before the next selection is made, the selection probabilities for the next selection are changed. The selections are dependent.

The multiplicative law of probability, for two independent events A and B, states that the joint probability of A and B happening in the order A,B is equal to the probability that A happens times the probability that B happens. In equation form, $P(AB) = P(A)P(B)$. For the order B,A, $P(BA) = P(B)P(A)$ and we note that $P(AB) = P(BA)$. Remember, independence means that the probability of B happening is not affected by the occurrence of A and vice versa. The multiplicative law extends to any number of independent events. Thus, $P(ABC) = P(A)P(B)P(C)$.

For two dependent events A and B, the multiplicative law states that the joint probability of A and B happening in the order A,B is equal to the probability of A happening times the probability that B happens under the condition that A has already happened. In equation form $P(AB) = P(A)P(B|A)$; or for the order B,A we have $P(BA) = P(B)P(A|B)$. The vertical bar can usually be translated as "given" or "given that." The notation on the left of the bar refers to the event under consideration and the notation on the right to a condition under which the event can take place. $P(B|A)$ is called conditional probability and could be read "the probability of B, given that A has already happened," or simply "the probability of B given A." When the events are independent, $P(B|A) = P(B)$; that is, the conditional probability of B occurring is the same as the unconditional probability of B. Extending the multiplication rule to a series of three events A,B,C occurring in that order, we have $P(ABC) = P(A)P(B|A)P(C|AB)$ where $P(C|AB)$ is the probability of C occurring, given that A and B have already occurred.

2.4 SAMPLING WITH REPLACEMENT

When a sample is drawn and each selected element is returned to the population before the next selection is made, the method of sampling is

called "sampling with replacement." In this case, the outcome of one selection does not change the selection probabilities for another selection.

Suppose a sample of n elements is selected with replacement. Let the values of X in the sample be x_1, x_2, \dots, x_n where x_1 is the value of X obtained on the first selection, x_2 the value obtained on the second selection, etc. Notice that x_1 is a random variable that could be equal to any value in the population set of values X_1, X_2, \dots, X_N , and the probability that x_1 equals X_1 is P_1 . The same statement applies to x_2 , etc. Since the selections are independent, the probability of getting a sample of n in a particular order is the product of the selection probabilities namely, $p(x_1)p(x_2)\dots p(x_n)$ where $p(x_1)$ is the P_1 for the element selected on the first draw, $p(x_2)$ is the P_1 for the element selected on the second draw, etc.

Illustration 2.3. As an illustration, consider a sample of two elements selected with equal probability and with replacement from a population of four elements. Suppose the values of some characteristic X for the four elements are X_1, X_2, X_3 , and X_4 . There are 16 possibilities:

X_1, X_1	X_2, X_1	X_3, X_1	X_4, X_1
X_1, X_2	X_2, X_2	X_3, X_2	X_4, X_2
X_1, X_3	X_2, X_3	X_3, X_3	X_4, X_3
X_1, X_4	X_2, X_4	X_3, X_4	X_4, X_4

In this illustration $p(x_1)$ is always equal to $\frac{1}{4}$ and $p(x_2)$ is always $\frac{1}{4}$. Hence each of the 16 possibilities has a probability of $(\frac{1}{4})(\frac{1}{4}) = \frac{1}{16}$.

Each of the 16 possibilities is a different permutation that could be regarded as a separate sample. However, in practice (as we are not concerned about which element was selected first or second) it is more logical to disregard the order of selection. Hence, as possible samples and the probability of each occurring, we have:

<u>Sample</u>	<u>Probability</u>	<u>Sample</u>	<u>Probability</u>
X_1, X_1	1/16	X_2, X_3	1/8
X_1, X_2	1/8	X_2, X_4	1/8
X_1, X_3	1/8	X_3, X_3	1/16
X_1, X_4	1/8	X_3, X_4	1/8
X_2, X_2	1/16	X_4, X_4	1/16

Note that the sum of the probabilities is 1. That must always be the case if all possible samples have been listed with the correct probabilities. Also note that, since the probability (relative frequency of occurrence) of each sample is known, the average for each sample is a random variable. In other words, there were 10 possible samples, and any one of 10 possible sample averages could have occurred with the probability indicated. This is a simple illustration of the fact that the sample average satisfies the definition of a random variable. As the theory of sampling unfolds, we will be examining the properties of a sample average that exist as a result of its being a random variable.

Exercise 2.3. With reference to Illustration 2.3, suppose the probabilities of selection were $P_1 = \frac{1}{4}$, $P_2 = \frac{1}{8}$, $P_3 = \frac{3}{8}$, and $P_4 = \frac{1}{4}$. Find the probability of each of the ten samples. Remember the sampling is with replacement. Check your results by adding the 10 probabilities.

The sum should be 1. Partial answer: For the sample composed of elements 2 and 4 the probability is $(\frac{1}{8})(\frac{1}{4}) + (\frac{1}{4})(\frac{1}{8}) = \frac{1}{16}$.

2.5 SAMPLING WITHOUT REPLACEMENT

When a selected element is not returned to the population before the next selection is made, the sampling method is called sampling without replacement. In this case, the selection probabilities change from one draw to the next; that is, the selections (events) are dependent.

As above, assume a population of N elements with values of some characteristic X equal to X_1, X_2, \dots, X_N . Let the selection probabilities for the first selection be $P_1, \dots, P_1, \dots, P_N$ where each $P_i > 0$ and $\sum P_i = 1$. Suppose three elements are selected without replacement. Let x_1, x_2 , and x_3 be the values of X obtained on the first, second, and third random draws, respectively. What is the probability that $x_1 = X_5, x_2 = X_6$, and $x_3 = X_7$? Let $P(X_5, X_6, X_7)$ represent this probability, which is the probability of selecting elements 5, 6, and 7 in that order.

According to the multiplicative probability law for dependent events,

$$P(X_5, X_6, X_7) = P(X_5)P(X_6|X_5)P(X_7|X_5, X_6)$$

It is clear that $P(X_5) = P_5$. For the second draw the selection probabilities (after element 5 is eliminated) must be adjusted so they add to 1. Hence, for the second draw the selection probabilities are

$$\frac{P_1}{1-P_5}, \frac{P_2}{1-P_5}, \frac{P_3}{1-P_5}, \frac{P_4}{1-P_5}, \frac{P_6}{1-P_5}, \dots, \frac{P_N}{1-P_5}. \quad \text{That is, } P(X_6|X_5) = \frac{P_6}{1-P_5}.$$

$$\text{Similarly, } P(X_7|X_5, X_6) = \frac{P_7}{1-P_5-P_6}.$$

$$\text{Therefore, } P(X_5, X_6, X_7) = (P_5) \left(\frac{P_6}{1-P_5} \right) \left(\frac{P_7}{1-P_5-P_6} \right) \quad (2.1)$$

Observe that $P(X_6, X_5, X_7) = (P_6) \left(\frac{P_5}{1-P_6} \right) \left(\frac{P_7}{1-P_6-P_5} \right)$. Hence, $P(X_5, X_6, X_7) \neq P(X_6, X_5, X_7)$ unless $P_5 = P_6$. In general, each permutation of n elements has a different probability of occurrence unless the P_i 's are all equal. To obtain the exact probability of selecting a sample composed of elements 5, 6, and 7, one would need to compute the probability for each of the six possible permutations and get the sum of the six probabilities.

Incidentally, in the actual process of selection, it is not necessary to compute a new set of selection probabilities after each selection is made. Make each selection in the same way that the first selection was made. If an element is selected which has already been drawn, ignore the random number and continue the same process of random selection until a new element is drawn.

As indicated by the very brief discussion in this section, the theory of sampling without replacement and with unequal probability of selection can be very complex. However, books on sampling present ways of circumventing the complex problems. In fact, it is practical and advantageous in many cases to use unequal probability of selection in sampling. The probability theory for sampling with equal probability of selection and without replacement is relatively simple and will be discussed in more detail.

Exercise 2.4. For a population of 4 elements there are six possible samples of two when sampling without replacement. Let $P_1 = \frac{1}{4}$, $P_2 = \frac{1}{8}$, $P_3 = \frac{3}{8}$, and $P_4 = \frac{1}{4}$. List the six possible samples and find the probability of getting each sample. Should the probabilities for the six samples add to 1? Check your results.

Exercise 2.5. Suppose two elements are selected with replacement and with equal probability from a population of 100 elements. Find the probability: (a) that element number 10 is not selected, (b) that element number 10 is selected only once, and (c) that element number 10 is selected twice? As a check, the three probabilities should add to 1. Why? Find the probability of selecting the combination of elements 10 and 20.

Exercise 2.6. Refer to Exercise 2.5 and change the specification "with replacement" to "without replacement." Answer the same questions. Why is the probability of getting the combination of elements 10 and 20 greater than it was in Exercise 2.5?

2.6 SIMPLE RANDOM SAMPLES

In practice, nearly all samples are selected without replacement. Selection of a random sample of n elements, with equal probability and without replacement, from a population of N elements is called simple random sampling (srs). One element must be selected at a time, that is, n separate random selections are required.

First, the probability of getting a particular combination of n elements will be discussed. Refer to Equation (2.1) and the discussion preceding it. The P_i 's are all equal to $\frac{1}{N}$ for simple random sampling. Therefore, Equation (2.1) becomes $P(X_5, X_6, X_7) = \left(\frac{1}{N}\right)\left(\frac{1}{N-1}\right)\left(\frac{1}{N-2}\right)$. All permutations of the three elements 5, 6, and 7 have the same probability of occurrence. There are $3! = 6$ possible permutations. Therefore, the probability that the sample is composed of the elements 5, 6, and 7 is $\frac{(1)(2)(3)}{N(N-1)(N-2)}$. Any other combination of three elements has the same probability of occurrence.

BEST COPY AVAILABLE

In general, all possible combinations of n elements have the same chance of selection and any particular combination of n has the following probability of being selected:

$$\frac{(1)(2)(3)\dots(n)}{N(N-1)(N-2)\dots(N-n+1)} = \frac{n!(N-n)!}{N!} \quad (2.2)$$

According to a theorem on number of combinations, there are $\frac{N!}{n!(N-n)!}$ possible combinations (samples) of n elements. If each combination of n elements has the same chance of being the sample selected, the probability of selecting a specified combination must be the reciprocal of the number of combinations. This checks with Equation (2.2).

An important feature of srs that will be needed in the chapter on expected values is the fact that the j^{th} element of the population is as likely to be selected at the i^{th} random draw as any other. A general expression for the probability that the j^{th} element of the population is selected at the i^{th} drawing is

$$\left(\frac{N-1}{N}\right)\left(\frac{N-2}{N-1}\right)\left(\frac{N-3}{N-2}\right)\dots\left(\frac{N-i+1}{N-i+2}\right)\left(\frac{1}{N-i+1}\right) = \frac{1}{N} \quad (2.3)$$

Let us check Equation 2.3 for $i = 3$. The equation becomes

$$\left(\frac{N-1}{N}\right)\left(\frac{N-2}{N-1}\right)\left(\frac{1}{N-2}\right) = \frac{1}{N}$$

The probability that the j^{th} element of the population is selected at the third draw is equal to the probability that it was not selected at either the first or second draw times the conditional probability of being selected at the third draw, given that it was not selected at the first or second draw. (Remember, the sampling is without replacement). Notice that $\frac{N-1}{N}$ is the probability that the j^{th} element is not selected at the first draw and $\frac{N-2}{N-1}$ is the conditional probability that it was not selected at the second draw. Therefore, $\left(\frac{N-1}{N}\right)\left(\frac{N-2}{N-1}\right)$ is the probability that the j^{th}

element has not been selected prior to the third draw. When the third draw is made, the conditional probability of selecting the j^{th} element is $\frac{1}{N-2}$. Hence the probability of selecting the j^{th} element at the third draw is $(\frac{N-1}{N})(\frac{N-2}{N-1})(\frac{1}{N-2}) = \frac{1}{N}$. This verifies Equation (2.3) for $i = 3$.

To summarize, the general result for any size of sample is that the j^{th} element in a population has a probability equal to $\frac{1}{N}$ of being selected at the i^{th} drawing. It means that x_i (the value of X obtained at the i^{th} draw) is a random variable that has a probability of $\frac{1}{N}$ of being equal to any value of the set X_1, \dots, X_N .

What probability does the j^{th} element have of being included in a sample of n ? We have just shown that it has a probability of $\frac{1}{N}$ of being selected at the i^{th} drawing. Therefore, any given element of the population has n chances, each equal to $\frac{1}{N}$, of being included in a sample. The element can be selected at the first draw, or the second draw, ..., or the n^{th} draw and it cannot be selected twice because the sampling is without replacement. Therefore the probabilities, $\frac{1}{N}$ for each of the n draws, can be added which gives $\frac{n}{N}$ as the probability of any given element being included in the sample.

Illustration 2.4. Suppose one has a list of 1,000 farms which includes some farms that are out-of-scope (not eligible) for a survey. There is no way of knowing in advance whether a farm on the list is out-of-scope. A simple random sample of 200 farms is selected from the list. All 200 farms are visited but only the ones found to be in scope are included in the sample. What probability does an in-scope farm have of being in the sample? Every farm on the list of 1000 farms has a probability equal to $\frac{1}{5}$

of being in the sample of 200. All in-scope farms in the sample of 200 are included in the final sample. Therefore, the answer is $\frac{1}{5}$.

Exercise 2.7. From the following set of 12 values of X a srs of three elements is to be selected: 2, 10, 5, 8, 1, 15, 7, 8, 13, 4, 6, and 2. Find $P(\bar{x} \geq 12)$ and $P(3 < \bar{x} < 12)$. Remember that the total possible number of samples of 3 can readily be obtained by formula. Since every possible sample of three is equally likely, you can determine which samples will have an $\bar{x} \leq 3$ or an $\bar{x} \geq 12$ without listing all of the numerous possible samples. Answer: $P(\bar{x} \geq 12) = \frac{3}{220}$; $P(\bar{x} \leq 3) = \frac{9}{220}$; $P(3 < \bar{x} < 12) = \frac{208}{220}$.

2.7 SOME EXAMPLES OF RESTRICTED RANDOM SAMPLING

There are many methods other than srs that will give every element an equal chance of being in the sample, but some combinations of n elements do not have a chance of being the sample selected unless srs is used. For example, one might take every k^{th} element beginning from a random starting point between 1 and k . This is called systematic sampling. For a five percent sample k would be 20. The first element for the sample would be a random number between 1 and 20. If it is 12, then elements 12, 32, 52, etc., compose the sample. Every element has an equal chance, $\frac{1}{20}$, of being in the sample, but there are only 20 combinations of elements that have a chance of being the sample selected. Simple random sampling could have given the same sample but it is the method of sampling that characterizes a sample and determines how error due to sampling is to be estimated. One may think of sample design as a matter of choosing a method of sampling; that is, choosing restrictions to place on the process of selecting a sample so the combinations which

have a chance of being the sample selected are generally "better" than many of the combinations that could occur with simple random sampling. At the same time, important properties that exist for simple random samples need to be retained. The key properties of srs will be developed in the next two chapters.

Another common method of sampling involves classification of all elements of a population into groups called strata. A sample is selected from each stratum. Suppose N_i elements of the population are in the i^{th} stratum and a simple random sample of n_i elements is selected from it. This is called stratified random sampling. It is clear that every element in the i^{th} stratum has a probability equal to $\frac{n_i}{N_i}$ of being in the sample. If the sampling fraction, $\frac{n_i}{N_i}$, is the same for all strata, every element of the population has an equal chance, namely $\frac{n_i}{N_i}$, of being in the sample. Again every element of the population has an equal chance of selection and of being in the sample selected, but some combinations that could occur when the method is srs cannot occur when stratified random sampling is used.

So far, our discussion has referred to the selection of individual elements, which are the units that data pertain to. For sampling purposes a population must be divided into parts which are called sampling units. A sample of sampling units is then selected. Sampling units and elements could be identical. But very often, it is either not possible or not practical to use individual elements as sampling units. For example, suppose a sample of households is needed. A list of households does not exist but a list of blocks covering the area to be surveyed might be available. In this case, a sample of blocks might be selected and all households

within the selected blocks included in the sample. The blocks are the sampling units and the elements are households. Every element of the population should belong to one and only one sampling unit so the list of sampling units will account for all elements of the population without duplication or omission. Then, the probability of selecting any given element is the same as the probability of selecting the sampling unit that it belongs to.

Illustration 2.5. Suppose a population is composed of 1800 dwelling units located within 150 well-defined blocks. There are several possible sampling plans. A srs of 25 blocks could be selected and every dwelling unit in the selected blocks could be included in the sample. In this case, the sampling fraction is $\frac{1}{6}$ and every dwelling unit has a probability of $\frac{1}{6}$ of being in the sample. Is this a srs of dwelling units? No, but one could describe the sample as a random sample (or a probability sample) of dwelling units and state that every dwelling unit had an equal chance of being in the sample. That is, the term "simple random sample" would apply to blocks, not dwelling units. As an alternative sampling plan, if there were twelve dwelling units in each of the 150 blocks, a srs of two dwelling units could be selected from each block. This scheme, which is an example of stratified random sampling, would also give every dwelling unit a probability equal to $\frac{1}{6}$ of being in the sample.

Illustration 2.6. Suppose that a sample is desired of 100 adults living in a specified area. A list of adults does not exist, but a list of 4,000 dwelling units in the area is available. The proposed sampling plan is to select a srs of 100 dwelling units from the list. Then, the field staff is to visit the sample dwellings and list all adults living

in each. Suppose there are 220 adults living in the 100 dwelling units. A simple random sample of 100 adults is selected from the list of 220. Consider the probability that an adult in the population has of being in the sample of 100 adults.

Parenthetically, we should recognize that the discussion which follows overlooks important practical problems of definition such as the definition of a dwelling unit, the definition of an adult, and the definition of living in a dwelling unit. However, assume the definitions are clear, that the list of dwelling units is complete, that no dwelling is on the list more than once, and that no ambiguity exists about whether an adult lives or does not live in a particular dwelling unit. Incomplete definitions often lead to inexact probabilities or ambiguity that gives difficulty in analyzing or interpreting results. The many practical problems should be discussed in an applied course on sampling.

It is clear that the probability of a dwelling unit being in the sample is $\frac{1}{40}$. Therefore, every person on the list of 220 had a chance of $\frac{1}{40}$ of being on the list because, under the specifications, a person lives in one and only one dwelling unit, and an adult's chance of being on the list is the same as that of the dwelling unit he lives in.

The second phase of sampling involves selecting a simple random sample of 100 adults from the list of 220. The conditional probability of an adult being in the sample of 100 is $\frac{100}{220} = \frac{5}{11}$. That is, given the fact that an adult is on the list of 220, he now has a chance of $\frac{5}{11}$ of being in the sample of 100.

Keep in mind that the probability of an event happening is its relative frequency in repeated trials. If another sample were selected

following the above specifications, each dwelling unit in the population would again have a chance of $\frac{1}{40}$ of being in sample; but, the number of adults listed is not likely to be 220 so the conditional probability at the second phase depends upon the number of dwellings units in the sample blocks. Does every adult have the same chance of being in the sample? Examine the case carefully. An initial impression could be misleading. Every adult in the population has an equal chance of being listed in the first phase and every adult listed has an equal chance of being selected at the second phase. But, in terms of repetition of the whole sampling plan each person does not have exactly the same chance of being in the sample of 100. The following exercise will help clarify the situation and is a good exercise in probability.

Exercise 2.8. Assume a population of 5 d.u.'s (dwelling units) with the following numbers of adults:

<u>Dwelling Unit</u>	<u>No. of Adults</u>
1	2
2	4
3	1
4	2
5	3

A srs of two d.u.'s is selected. A srs of 2 adults is then selected from a list of all adults in the two d.u.'s. Find the probability that a specified adult in d.u. No. 1 has of being in the sample. Answer: 0.19. Find the probability that an adult in d.u. No. 2 has of being in the sample. Does the probability of an adult being in the sample appear to be related to the number of adults in his d.u.? In what way?

An alternative is to take a constant fraction of the adults listed instead of a constant number. For example, the specification might have been to select a random sample of $\frac{1}{2}$ of the adults listed in the first phase. In this case, under repeated application of the sampling specifications, the probability at the second phase does not depend on the outcome of the first phase and each adult in the population has an equal chance, $(\frac{1}{40})(\frac{1}{2}) = \frac{1}{80}$, of being selected in the sample. Notice that under this plan the number of adults in a sample will vary from sample to sample; in fact, the number of adults in the sample is a random variable.

For some surveys, interviewing more than one adult in a dwelling unit is inadvisable. Again, suppose the first phase of sampling is to select a srs of 100 dwelling units. For the second phase, consider the following: When an interviewer completes the listing of adults in a sample dwelling, he is to select one adult, from the list of those living in the dwelling, at random in accordance with a specified set of instructions. He then interviews the selected adult if available; otherwise, he returns at a time when the selected adult is available. What probability does an adult living in the area have of being in the sample? According to the multiplication theorem, the answer is $P'(D)P(A|D)$ where $P'(D)$ is the probability of the dwelling unit, in which the adult lives, being in the sample and $P(A|D)$ is the probability of the adult being selected given that his dwelling is in the sample. More specifically, $P'(D) = \frac{1}{40}$ and $P(A|D) = \frac{1}{k_i}$, where k_i is the number of adults in the i^{th} dwelling. Thus, an adult's chance, $(\frac{1}{40})(\frac{1}{k_i})$, of being in a sample is inversely proportional to the number of adults in his dwelling unit.

Exercise 2.9. Suppose there are five dwelling units and 12 persons living in the five dwelling units as follows:

<u>Dwelling Unit</u>	<u>Individuals</u>
1	1, 2
2	3, 4, 5, 6
3	7, 8
4	9
5	10, 11, 12

1. A sample of two dwelling units is selected with equal probability and without replacement. All individuals in the selected dwelling units are in the sample. What probability does individual number 4 have of being in the sample? Individual number 9?

2. Suppose from a list of the twelve individuals that one individual is selected with equal probability. From the selected individual two items of information are obtained: his age and the value of the dwelling in which he lives. Let X_1, X_2, \dots, X_{12} represent the ages of the 12 individuals and let Y_1, \dots, Y_5 represent the values of the five dwelling units. Clearly, the probability of selecting the i^{th} individual is $\frac{1}{12}$ and therefore $P(X_i) = \frac{1}{12}$. Find the five probabilities $P(Y_1), \dots, P(Y_5)$. Do you agree that $P(Y_j) = \frac{2}{12}$? As a check, $\sum P(Y_j)$ should equal one.

3. Suppose a sample of two individuals is selected with equal probability and without replacement. Let Y_{1j} be the value of Y_j obtained at the first draw and Y_{2j} be the value of Y_j obtained at the second draw. Does $P(Y_{1j}) = P(Y_{2j})$? That is, is the probability of getting Y_j on the second draw the same as it was on the first? If the answer is not evident, refer to Section 2.5.

Exercise 2.10. A small sample of third-grade students enrolled in public schools in a State is desired. The following plan is presented only

as an exercise and without consideration of whether it is a good one: A sample of 10 third-grade classes is to be selected. All students in the 10 classes will be included in the sample.

Step 1. Select a srs of 10 school districts.

Step 2. Within each of the 10 school districts, prepare a list of public schools having a third grade. Then select one school at random from the list.

Step 3. For each of the 10 schools resulting from Step 2, list the third-grade classes and select one class at random. (If there is only one third-grade class in the school, it is in the sample). This will give a sample of 10 classes.

Describe third-grade classes in the population which have relatively small chances of being selected. Define needed notation and write a mathematical expression representing the probability of a third-grade class being in the sample.

2.8 TWO-STAGE SAMPLING

For various reasons sampling plans often employ two or more stages of sampling. For example, a sample of counties might be selected, then within each sample county a sample of farms might be selected.

Units used at the first stage of sampling are usually called primary sampling units or psu's. The sampling units at the second stage of sampling could be called secondary sampling units. However, since there has been frequent reference earlier in this chapter to "elements of a population," the sampling units at the second stage will be called elements.

In the simple case of two-stage sampling, each element of the population is associated with one and only one primary sampling unit. Let i

be the index for psu's and let j be the index for elements within a psu. Thus X_{ij} represents the value of some characteristic X for the j^{th} element in the i^{th} psu. Also, let

M = the total number of psu's,

m = the number of psu's selected for a sample,

N_i = the total number of elements in the i^{th} psu, and

n_i = the number of elements in the sample from the i^{th} psu.

Then,

$\sum_{i=1}^M N_i = N$, the total number of elements in the population, and

$\sum_{i=1}^m n_i = n$, the total number of elements in the sample.

Now consider the probability of an element being selected by a two step process: (1) Select one psu, and (2) select one element within the selected psu. Let,

P_i = the probability of selecting the i^{th} psu,

$P_{j|i}$ = the conditional probability of selecting the j^{th} element in the i^{th} psu given that the i^{th} psu has already been selected, and

P_{ij} = the overall probability of selecting the j^{th} element in the i^{th} psu.

Then,

$$P_{ij} = P_i P_{j|i}$$

If the product of the two probabilities, P_i and $P_{j|i}$, is constant for every element, then every element of the population has an equal chance of

being selected. In other words, given a set of selection probabilities P_1, \dots, P_M for the psu's, one could specify that $P_{1j} = \frac{1}{N}$ and compute $P_{j|1}$, where $P_{j|1} = \frac{1}{NP_1}$, so every element of the population will have an equal chance of selection.

Exercise 2.11. Refer to Table 2.1. An element is to be selected by a three-step process as follows: (1) Select one of the Y classes (a row) with probability $\frac{N_1}{N}$, (2) within the selected row select an X class (a column) with probability $\frac{N_{1j}}{N_1}$, (3) within the selected cell select an element with equal probability. Does each element in the population of N elements have an equal probability of being drawn? What is the probability?

The probability of an element being included in a two-stage sample is given by

$$P'_{1j} = P'_1 P'_{j|1} \quad (2.4)$$

where

P'_1 = the probability that the i^{th} psu is in the sample of psu's, and

$P'_{j|1}$ = the conditional probability which the j element has of being in the sample, given that the i^{th} psu has been selected.

The inclusion probability P'_{ij} will be discussed very briefly for three important cases:

(1) Suppose a random sample of m psu's is selected with equal probability and without replacement. The probability, P'_1 , of the i^{th} psu being in the sample is $f_1 = \frac{m}{M}$ where f_1 is the sampling fraction for the first-stage units. In the second stage of sampling assume that, within each of the m psu's, a constant proportion, f_2 , of the elements is selected.

That is, in the i^{th} psu in the sample, a simple random sample of n_i elements out of N_i is selected, the condition being that $n_i = f_2 N_i$. Hence, the conditional probability of the j^{th} element in the i^{th} psu being in the sample is $P'_{j|i} = \frac{n_i}{N_i} = f_2$. Substituting in Equation 2.4, we have $P'_{ij} = f_1 f_2$ which shows that an element's probability of being in the sample is equal to the product of the sampling fractions at the two stages. In this case P'_{ij} is constant and is the overall sampling fraction.

Unless N_i is the same for all psu's, the size of the sample, $n_i = f_2 N_i$, varies from psu to psu. Also, since the psu's are selected at random the total size of the sample, $n = \sum_{i=1}^m n_i = f_2 \sum_{i=1}^m N_i$, is not constant with regard to repetition of the sampling plan. In practice variation in the size, n_i , of the sample from psu to psu might be very undesirable. If appropriate information is available, it is possible to select psu's with probabilities that will equalize the sample sizes n_i and also keep P'_{ij} constant.

(2) Suppose one psu is selected with probability $P_i = \frac{N_i}{N}$. This is commonly known as sampling with pps (probability proportional to size). Within the selected psu, assume that a simple random sample of k elements is selected. (If any N_i are less than k , consolidations could be made so all psu's have an N_i greater than k). Then,

$$P'_i = \frac{N_i}{N}, \quad P'_{j|i} = \frac{k}{N_i}, \quad \text{and} \quad P'_{ij} = \frac{N_i}{N} \frac{k}{N_i} = \frac{k}{N}$$

which means that every element of the population has an equal probability, $\frac{k}{N}$, of being included in a sample of k elements.

Extension of this sampling scheme to a sample of m psu's could encounter the complications indicated in Section 2.5. However, it was

stated that means exist for circumventing those complications. Sampling books 1/ discuss this matter quite fully so we will not include it in this monograph. The point is that one can select m psu's without replacement in such a way that $m \frac{N_i}{N}$ is the probability of including the i^{th} psu in the sample. That is, $P_i' = m \frac{N_i}{N}$. If a random sample of k elements is selected with equal probability from each of the selected psu's,

$$P_{j|i}' = \frac{k}{N_i} \quad \text{and}$$

$$P_{ij}' = \left(m \frac{N_i}{N}\right) \left(\frac{k}{N_i}\right) = \frac{mk}{N} = \frac{n}{N}$$

Thus, if the N_i are known exactly for all M psu's in the population, and if a list of elements in each psu is available, it is possible to select a two-stage sample of n elements so that k elements for the sample come from each of m psu's and every element of the population has an equal chance of being in the sample. In practice, however, one usually finds one of two situations: (a) there is no information on the number of elements in the psu's, or (b) the information that does exist is out-of-date. Nevertheless, out-of-date information on number of elements in the psu's can be very useful. It is also possible that a measure of size might exist which will serve, more efficiently, the purposes of sampling.

(3) Suppose that characteristic Y is used as a measure of size. Let Y_i be the value of Y for the i^{th} psu in the population and let $P_i = \frac{Y_i}{Y}$ where $Y = \sum_{i=1}^M Y_i$. A sample of m psu's is selected in such a way that $P_i' = m \frac{Y_i}{Y}$ is the probability that the i^{th} psu has of being in the sample.

1/ For example, Hansen, Hurwitz, and Madow. Sample Survey Methods and Theory. Volume I, Chapter 8. John Wiley and Sons. 1953.

With regard to the second stage of sampling, let f_{2i} be the sampling fraction for selecting a simple random sample within the i^{th} psu in the sample. That is, $P'_{j|i} = f_{2i}$. Then,

$$P'_{ij} = \left(m \frac{Y_i}{Y}\right) (f_{2i}) \quad (2.5)$$

In setting sampling specifications one would decide on a fixed value for P'_{ij} . In this context P'_{ij} is the overall sampling fraction or proportion of the population that is to be included in the sample. For example, if one wanted a 5 percent sample, P'_{ij} would be .05. Or, if one knew there were approximately 50,000 elements in the population and wanted a sample of about 2,000, he would set $P'_{ij} = .04$. Hence, we will let f be the overall sampling fraction and set P'_{ij} equal to f . Decisions are also made on the measure of size to be used and on the number, m , of psu's to be selected. In Equation 2.5, this leaves f_{2i} to be determined. Thus, f_{2i} is computed as follows for each psu in the sample:

$$f_{2i} = \frac{fY}{mY_i}$$

Use of the sampling fractions f_{2i} at the second stage of sampling will give every element of the population a probability equal to f of being in the sample. A sample wherein every element of the population has an equal chance of inclusion is often called a self-weighted sample.

CHAPTER III. EXPECTED VALUES OF RANDOM VARIABLES

3.1 INTRODUCTION

The theory of expected values of random variables is used extensively in the theory of sampling; in fact, it is the foundation for sampling theory. Interpretations of the accuracy of estimates from probability samples depend heavily on the theory of expected values.

The definition of a random variable was discussed in the previous chapter. It is a variable that can take (be equal to) any one of a defined set of values with known probability. Let X_i be the value of X for the i^{th} element in a set of N elements and let P_i be the probability that the i^{th} element has of being selected by some chance operation so that P_i is known a priori. What is the expected value of X ?

Definition 3.1. The expected value of a random variable X is

$\sum_{i=1}^N P_i X_i$ where $\sum_{i=1}^N P_i = 1$. The mathematical notation for the expected value of X is $E(X)$. Hence, by definition, $E(X) = \sum_{i=1}^N P_i X_i$.

Observe that $\sum P_i X_i$ is a weighted average of the values of X , the weights being the probabilities of selection. "Expected value" is a substitute expression for "average value." In other words, E means "the average value of" or "find the average value of" whatever follows E . For example, $E(X^2)$, read "the expected value of X^2 ," refers to the average value of the squares of the values that X can equal. That is, by definition,

$$E(X^2) = \sum_{i=1}^N P_i X_i^2.$$

If all of the N elements have an equal chance of being selected, all values of P_i must equal $\frac{1}{N}$ because of the requirement that $\sum P_i = 1$. In

this case, $E(X) = \sum_{i=1}^N \frac{1}{N} X_i = \frac{\sum X_i}{N} = \bar{X}$, which is the simple average of X for all N elements.

Illustration 3.1. Assume 12 elements having values of X as follows:

$X_1 = 3$	$X_5 = 5$	$X_9 = 10$
$X_2 = 9$	$X_6 = 3$	$X_{10} = 3$
$X_3 = 3$	$X_7 = 4$	$X_{11} = 8$
$X_4 = 5$	$X_8 = 3$	$X_{12} = 4$

For this set, $E(X) = \frac{3+9+\dots+4}{12} = 5$, assuming each element has the same chance of selection. Or, by counting the number of times that each unique value of X occurs, a frequency distribution of X can be obtained as follows:

X_j	N_j
3	5
4	2
5	2
8	1
9	1
10	1

where X_j is a unique value of X and N_j is the number of times X_j occurs.

We noted in Chapter I that $\sum N_j = N$, $\sum N_j X_j = \sum X_i$, and that $\frac{\sum N_j X_j}{\sum N_j} = \frac{\sum X_i}{N} = \bar{X}$.

Suppose one of the X_j values is selected at random with a probability equal

to P_j where $P_j = \frac{N_j}{N}$. What is the expected value of X_j ? By

definition $E(X_j) = \sum P_j X_j = \sum \frac{1}{N} X_j = \frac{\sum X_j}{N} = \bar{X}$. The student may verify that in this illustration $E(X_j) = 5$. Note that the selection specifications were equivalent to selecting one of the 12 elements at random with equal probability.

Incidentally, a frequency distribution and a probability distribution are very similar. The probability distribution with reference to X_j would be:

X_j	P_j
3	5/12
4	2/12
5	2/12
8	1/12
9	1/12
10	1/12

The 12 values, $P_j = \frac{1}{N}$, for the 12 elements are also a probability distribution. This illustration shows two ways of treating the set of 12 elements.

When finding expected values be sure that you understand the definition of the set of values that the random variable might equal and the probabilities involved.

Definition 3.2. When X is a random variable, by definition the expected value of a function of X is

$$E[f(X)] = \sum_{i=1}^N P_i [f(X_i)]$$

Some examples of simple functions of X are: $f(X) = aX$, $f(X) = X^2$, $f(X) = a + bX + cX^2$, and $f(X) = (X - \bar{X})^2$. For each value, X_i , in a defined set there is a corresponding value of $f(X_i)$.

Illustration 3.2. Suppose $f(X) = 2X+3$. With reference to the set of 12 elements discussed above, there are 12 values of $f(X_i)$ as follows:

$$f(X_1) = (2)(3) + 3 = 9$$

$$f(X_2) = (2)(9) + 3 = 21$$

.

$$f(X_{12}) = 2(4) + 3 = 11$$

Assuming $P_i = \frac{1}{N}$ the expected value of $f(X) = 2X+3$ would be

$$E(2X+3) = \sum_{i=1}^{12} \frac{1}{N}(2X_i+3) = \left(\frac{1}{12}\right)(9) + \left(\frac{1}{12}\right)(21) + \dots + \left(\frac{1}{12}\right)(11) = 13 \quad (3.1)$$

In algebraic terms, for $f(X) = aX+b$, we have

$$E(aX+b) = \sum_{i=1}^N P_i(aX_i+b) = \sum_{i=1}^N P_i(aX_i) + \sum_{i=1}^N P_i b$$

By definition $\sum_{i=1}^N P_i(aX_i) = E(aX)$, and $\sum_{i=1}^N P_i b = E(b)$. Therefore,

$$E(aX+b) = E(aX) + E(b) \quad (3.2)$$

Since b is constant and $\sum_{i=1}^N P_i = 1$, $\sum_{i=1}^N P_i b = b$, which leads to the first important theorem in expected values.

Theorem 3.1. The expected value of a constant is equal to the constant: $E(a) = a$.

By definition $E(aX) = \sum_{i=1}^N P_i(aX_i) = a \sum_{i=1}^N P_i X_i$. Since $\sum_{i=1}^N P_i X_i = E(X)$, we have another important theorem:

Theorem 3.2. The expected value of a constant times a variable equals the constant times the expected value of the variable: $E(aX) = aE(X)$.

Applying these two theorems to Equation (3.2) we have $E(aX+b) = aE(X) + b$. Therefore, with reference to Illustration 3.2, $E(2X+3) = 2E(X) + 3 = 2(5) + 3 = 13$, which is the same as the result found in

Equation (3.1).

Exercise 3.1. Suppose a random variable X can take any of the following four values with the probabilities indicated:

$$\begin{array}{llll} X_1 = 2 & X_2 = 5 & X_3 = 4 & X_4 = 6 \\ P_1 = 2/6 & P_2 = 2/6 & P_3 = 1/6 & P_4 = 1/6 \end{array}$$

- (a) Find $E(X)$ Answer: 4
 (b) Find $E(X^2)$ Answer: $18\frac{1}{3}$. Note that $E(X^2) \neq [E(X)]^2$
 (c) Find $E(X - \bar{X})$ Answer: 0 Note: By definition

$$E(X - \bar{X}) = \sum_{i=1}^4 P_i (X_i - \bar{X})$$

- (d) Find $E(X - \bar{X})^2$ Answer: $2\frac{1}{3}$. Note: By definition

$$E(X - \bar{X})^2 = \sum_{i=1}^4 P_i (X_i - \bar{X})^2$$

Exercise 3.2. From the following set of three values of Y_i one value is to be selected with a probability P'_i :

$$\begin{array}{lll} Y_1 = -2 & Y_2 = 2 & Y_3 = 4 \\ P'_1 = 1/4 & P'_2 = 2/4 & P'_3 = 1/4 \end{array}$$

- (a) Find $E(Y)$ Answer: $1\frac{1}{2}$
 (b) Find $E(\frac{1}{Y})$ Answer: $3/16$. Note: $\frac{1}{E(Y)} \neq E(\frac{1}{Y})$
 (c) Find $E(Y - \bar{Y})^2$ Answer: $4\frac{3}{4}$

3.2 EXPECTED VALUE OF THE SUM OF TWO RANDOM VARIABLES

The sum of two or more random variables is also a random variable.

If X and Y are two random variables, the expected value of $X + Y$ is equal to the expected value of X plus the expected value of Y : $E(X+Y) = E(X) + E(Y)$. Two numerical illustrations will help clarify the situation.

Illustration 3.3. Consider the two random variables X and Y in

Exercises 3.1 and 3.2:

$$\begin{array}{ll}
 X_1 = 2 & P_1 = \frac{2}{6} \\
 X_2 = 5 & P_2 = \frac{2}{6} \\
 X_3 = 4 & P_3 = \frac{1}{6} \\
 X_4 = 6 & P_4 = \frac{1}{6}
 \end{array}
 \qquad
 \begin{array}{ll}
 Y_1 = -2 & P'_1 = \frac{1}{4} \\
 Y_2 = 2 & P'_2 = \frac{2}{4} \\
 Y_3 = 4 & P'_3 = \frac{1}{4}
 \end{array}$$

Suppose one element of the first set and one element of the second set are selected with probabilities as listed above. What is the expected value of $X + Y$? The joint probability of getting X_i and Y_j is $P_i P'_j$ because the two selections are independent. Hence by definition

$$E(X + Y) = \sum_{i=1}^4 \sum_{j=1}^3 P_i P'_j (X_i + Y_j) \quad (3.3)$$

The possible values of $X + Y$ and the probability of each are as follows:

$X + Y$	$P_i P'_j$	$X + Y$	$P_i P'_j$
$X_1 + Y_1 = 0$	$P_1 P'_1 = \frac{2}{24}$	$X_3 + Y_1 = 2$	$P_3 P'_1 = \frac{1}{24}$
$X_1 + Y_2 = 4$	$P_1 P'_2 = \frac{4}{24}$	$X_3 + Y_2 = 6$	$P_3 P'_2 = \frac{2}{24}$
$X_1 + Y_3 = 6$	$P_1 P'_3 = \frac{2}{24}$	$X_3 + Y_3 = 8$	$P_3 P'_3 = \frac{1}{24}$
$X_2 + Y_1 = 3$	$P_2 P'_1 = \frac{2}{24}$	$X_4 + Y_1 = 4$	$P_4 P'_1 = \frac{1}{24}$
$X_2 + Y_2 = 7$	$P_2 P'_2 = \frac{4}{24}$	$X_4 + Y_2 = 8$	$P_4 P'_2 = \frac{2}{24}$
$X_2 + Y_3 = 9$	$P_2 P'_3 = \frac{2}{24}$	$X_4 + Y_3 = 10$	$P_4 P'_3 = \frac{1}{24}$

As a check the sum of the probabilities must be 1 if all possible sums have been listed and the probability of each has been correctly determined. Substituting the values of $X_i + Y_j$ and $P_i P'_j$ in Equation (3.3) we obtain 5.5 as follows for expected value of $X + Y$:

$$\left(\frac{2}{24}\right)(0) + \left(\frac{4}{24}\right)(4) + \dots + \left(\frac{1}{24}\right)(10) = 5.5$$

From Exercises 3.1 and 3.2 we have $E(X) = 4$ and $E(Y) = 1.5$. Therefore, $E(X) + E(Y) = 4 + 1.5 = 5.5$ which verifies the earlier statement that $E(X + Y) = E(X) + E(Y)$.

Illustration 3.4. Suppose a random sample of two is selected with replacement from the population of four elements used in Exercise 3.1. Let x_1 be the first value selected and let x_2 be the second. Then x_1 and x_2 are random variables and $x_1 + x_2$ is a random variable. The possible values of $x_1 + x_2$ and the probability of each, $P(x_1, x_2)$, are listed below. Notice that each possible order of selection is treated separately.

x_1	x_2	$P(x_1, x_2)$	$x_1 + x_2$	x_1	x_2	$P(x_1, x_2)$	$x_1 + x_2$
x_1	x_1	4/36	4	x_3	x_1	2/36	6
x_1	x_2	4/36	7	x_3	x_2	2/36	9
x_1	x_3	2/36	6	x_3	x_3	1/36	8
x_1	x_4	2/36	8	x_3	x_4	1/36	10
x_2	x_1	4/36	7	x_4	x_1	2/36	8
x_2	x_2	4/36	10	x_4	x_2	2/36	11
x_2	x_3	2/36	9	x_4	x_3	1/36	10
x_2	x_4	2/36	11	x_4	x_4	1/36	12

By definition $E(x_1 + x_2)$ is

$$\frac{4}{36}(4) + \frac{4}{36}(7) + \frac{2}{36}(6) + \dots + \frac{1}{36}(12) = 8$$

In Exercise 3.1 we found $E(X) = 4$. Since x_1 is the same random variable as X , $E(x_1) = 4$. Also, x_2 is the same random variable as X , and $E(x_2) = 4$. Therefore, $E(x_1) + E(x_2) = 8$, which verifies that $E(x_1 + x_2) = E(x_1) + E(x_2)$.

In general if X and Y are two random variables, where X might equal x_1, \dots, x_N and Y might equal y_1, \dots, y_M , then $E(X + Y) = E(X) + E(Y)$. The

proof is as follows: By definition $E(X+Y) = \sum_{ij}^{NM} P_{ij} (X_i + Y_j)$ where P_{ij} is the probability of getting the sum $X_i + Y_j$, and $\sum \sum P_{ij} = 1$. The double summation is over all possible values of $P_{ij} (X_i + Y_j)$. According to the rules for summation we may write

$$\sum_{ij}^{NM} P_{ij} (X_i + Y_j) = \sum_{ij}^{NM} P_{ij} X_i + \sum_{ij}^{NM} P_{ij} Y_j \quad (3.4)$$

In the first term on the right, X_i is constant with regard to the summation over j ; and in the second term on the right, Y_j is constant with regard to the summation over i . Therefore, the right-hand side of Equation (3.4) can be written as

$$\sum_i^N X_i \sum_j^M P_{ij} + \sum_j^M Y_j \sum_i^N P_{ij}$$

And, since $\sum_j^M P_{ij} = P_i$ and $\sum_i^N P_{ij} = P_j$, Equation (3.4) becomes

$$\sum_{ij}^{NM} P_{ij} (X_i + Y_j) = \sum_i^N X_i P_i + \sum_j^M Y_j P_j$$

By definition $\sum_i^N X_i P_i = E(X)$ and $\sum_j^M Y_j P_j = E(Y)$.

Therefore $E(X+Y) = E(X) + E(Y)$.

If the proof is not clear write the values of $P_{ij} (X_i + Y_j)$ in a matrix format. Then, follow the summation manipulations in the proof.

The above result extends to any number of random variables; that is, the expected value of a sum of random variables is the sum of the expected values of each. In fact, there is a very important theorem that applies to a linear combination of random variables.

Theorem 3.3. Let $u = a_1 u_1 + \dots + a_k u_k$, where u_1, \dots, u_k are random variables and a_1, \dots, a_k are constants. Then

$$E(u) = a_1 E(u_1) + \dots + a_k E(u_k)$$

or in summation notation

$$E(u) = E \sum_{i=1}^k a_i u_i = \sum_{i=1}^k a_i E(u_i)$$

The generality of Theorem 3.3 is impressive. For example, with reference to sampling from a population X_1, \dots, X_N , u_1 might be the value of X obtained at the first draw, u_2 the value obtained at the second draw, etc. The constants could be weights. Thus, in this case, u would be a weighted average of the sample measurements. Or, suppose $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_k$ are averages from a random sample for k different age groups. The averages are random variables and the theorem could be applied to any linear combination of the averages. In fact u_i could be any function of random variables. That is, the only condition on which the theorem is based is that u_i must be a random variable.

Illustration 3.5. Suppose we want to find the expected value of $(X + Y)^2$ where X and Y are random variables. Before Theorem 3.3 can be applied we must square $(X + Y)$. Thus $E(X + Y)^2 = E(X^2 + 2XY + Y^2)$.

The application of Theorem 3.3 gives $E(X + Y)^2 = E(X)^2 + 2E(XY) + E(Y)^2$.

Illustration 3.6. We will now show that

$$E(X - \bar{X})(Y - \bar{Y}) = E(XY) - \bar{X}\bar{Y} \quad \text{where} \quad E(X) = \bar{X} \quad \text{and} \quad E(Y) = \bar{Y}$$

Since $(X - \bar{X})(Y - \bar{Y}) = XY - \bar{X}Y - X\bar{Y} + \bar{X}\bar{Y}$ we have

$$E(X - \bar{X})(Y - \bar{Y}) = E(XY - \bar{X}Y - X\bar{Y} + \bar{X}\bar{Y})$$

and application of Theorem 3.3 gives

$$E(X - \bar{X})(Y - \bar{Y}) = E(XY) - E(\bar{X}Y) - E(X\bar{Y}) + E(\bar{X}\bar{Y})$$

Since \bar{X} and \bar{Y} are constant, $E(\bar{X}Y) = \bar{X} E(Y) = \bar{X}\bar{Y}$, $E(Y\bar{X}) = \bar{Y}\bar{X}$, and $E(\bar{X}\bar{Y}) = \bar{X}\bar{Y}$.

Therefore, $E(X-\bar{X})(Y-\bar{Y}) = E(XY) - \bar{X}\bar{Y}$

Exercise 3.3. Suppose $E(X) = 6$ and $E(Y) = 4$. Find

- (a) $E(2X+4Y)$ Answer: 28
- (b) $[E(2X)]^2$ Answer: 144
- (c) $\sqrt{E(Y)}$ Answer: 2
- (d) $E(5Y-X)$ Answer: 14

Exercise 3.4. Prove the following, assuming $E(X) = \bar{X}$ and $E(Y) = \bar{Y}$:

- (a) $E(X-\bar{X}) = 0$
- (b) $E(aX+bY) = a\bar{X} + b\bar{Y}$
- (c) $E[a(X-\bar{X}) + b(Y-\bar{Y})] = 0$
- (d) $E(X+a)^2 = E(X^2) + 2a\bar{X} + a^2$
- (e) $E(X-\bar{X})^2 = E(X^2) - \bar{X}^2$
- (f) $E(aX+bY) = 0$ for any values of a and b if $E(X) = 0$ and $E(Y) = 0$.

3.3 EXPECTED VALUE OF AN ESTIMATE

Theorem 3.3 will now be used to find the expected value of the mean of a simple random sample of n elements selected without replacement from a population of N elements. The term "simple random sample" implies equal probability of selection without replacement. The sample average is

$$\bar{x} = \frac{x_1 + \dots + x_n}{n}$$

where x_1 is the value of X for the 1th element in the sample. Without loss of generality, we can consider the subscript of x as corresponding to the 1th draw; i.e., x_1 is the value of X obtained on the first draw, x_2 the value on the second, etc. As each x_i is a random variable, \bar{x} is a linear combination of random variables. Therefore, Theorem 3.3 applies and

$$E(\bar{x}) = \frac{1}{n} [E(x_1) + \dots + E(x_n)]$$

In the previous chapter, Section 2.6, we found that any given element of the population had a chance of $\frac{1}{N}$ of being selected on the i^{th} draw.

This means that x_i is a random variable that has a probability equal to $\frac{1}{N}$ of being equal to any value of the population set X_1, \dots, X_N . Therefore,

$$E(x_1) = E(x_2) = \dots = E(x_n) = \bar{X}$$

Hence, $E(\bar{x}) = \frac{\bar{X} + \dots + \bar{X}}{n} = \bar{X}$. The fact that $E(\bar{x}) = \bar{X}$ is one of the very important properties of an average from a simple random sample. Incidentally, $E(\bar{x}) = \bar{X}$ whether the sampling is with or without replacement.

Definition 3.3. A parameter is a quantity computed from all values in a population set. The total of X , the average of X , the proportion of elements for which $X_i < A$, or any other quantity computed from measurements including all elements of the population is a parameter. The numerical value of a parameter is usually unknown but it exists by definition.

Definition 3.4. An estimator is a mathematical formula or rule for making an estimate from a sample. The formula for a sample average,

$\bar{x} = \frac{\sum x_i}{n}$, is a simple example of an estimator. It provides an estimate of

the parameter $\bar{X} = \frac{\sum X_i}{N}$.

Definition 3.5. An estimate is unbiased when its expected value equals the parameter that it is an estimate of. In the above example, \bar{x} is an unbiased estimate of \bar{X} because $E(\bar{x}) = \bar{X}$.

Exercise 3.5. Assume a population of only four elements having values of X as follows: $X_1 = 2$, $X_2 = 5$, $X_3 = 4$, $X_4 = 6$. For simple random samples of size 2 show that the estimator $N\bar{x}$ provides an unbiased estimate of the population total, $\sum X_i = 17$. List all six possible samples of two and

calculate $N\bar{x}$ for each. This will give the set of values that the random variable $N\bar{x}$ can be equal to. Consider the probability of each of the possible values of $N\bar{x}$ and show arithmetically that $E(N\bar{x}) = 17$.

A sample of elements from a population is not always selected by using equal probabilities of selection. Sampling with unequal probability is complicated when the sampling is without replacement, so we will limit our discussion to sampling with replacement.

Illustration 3.7. The set of four elements and the associated probabilities used in Exercise 3.1 will serve as an example of unbiased estimation when samples of two elements are selected with unequal probability and with replacement. Our estimator of the population total,

$2+5+4+6 = 17$, will be $x' = \frac{\sum_{i=1}^n \frac{x_i}{p_i}}{n}$. The estimate x' is a random variable.

Listed below are the set of values that x' can equal and the probability of each value occurring.

Possible Samples	x'	P
$x_1 x_1$	6	4/36
$x_1 x_2$	10.5	8/36
$x_1 x_3$	15	4/36
$x_1 x_4$	21	4/36
$x_2 x_2$	15	4/36
$x_2 x_3$	19.5	4/36
$x_2 x_4$	25.5	4/36
$x_3 x_3$	24	1/36
$x_3 x_4$	30	2/36
$x_4 x_4$	36	1/36

Exercise 3.6. Verify the above values of x'_j and P_j and find the expected value of x' . By definition $E(x') = \sum_j P_j x'_j$. Your answer should be 17 because x' is an unbiased estimate of the population total.

To put sampling with replacement and unequal probabilities in a general setting, assume the population is $X_1, \dots, X_j, \dots, X_N$ and the selection probabilities are $P_1, \dots, P_j, \dots, P_N$. Let x_i be the value of X for the i^{th} element in a sample of n elements and let p_i be the probability

which that element had of being selected. Then $x' = \frac{\sum_{i=1}^n \frac{x_i}{p_i}}{n}$ is an unbiased estimate of the population total. We will now show that $E(x') = \sum_{j=1}^N X_j$.

To facilitate comparison of x' with u in Theorem 3.3, x' may be written as follows:

$$x' = \frac{1}{n} \left(\frac{x_1}{p_1} \right) + \dots + \frac{1}{n} \left(\frac{x_n}{p_n} \right)$$

It is now clear that $a_i = \frac{1}{n}$ and $u_i = \frac{x_i}{p_i}$. Therefore,

$$E(x') = \frac{1}{n} \left[E\left(\frac{x_1}{p_1}\right) + \dots + E\left(\frac{x_n}{p_n}\right) \right] \quad (3.5)$$

The quantity $\frac{x_1}{p_1}$, which is the outcome of the first random selection from the population, is a random variable that might be equal to any one of the set of values $\frac{x_1}{p_1}, \dots, \frac{x_j}{p_j}, \dots, \frac{x_N}{p_N}$. The probability that $\frac{x_1}{p_1}$ equals $\frac{x_j}{p_j}$ is P_j .

Therefore, by definition

$$E\left(\frac{x_1}{p_1}\right) = \sum_{j=1}^N P_j \left(\frac{x_j}{p_j}\right) = \sum_{j=1}^N X_j$$

Since the sampling is with replacement it is clear that any $\frac{x_i}{p_i}$ is the same random variable as $\frac{x_1}{p_1}$.

Therefore Equation (3.5) becomes

$$E(\bar{x}) = \frac{1}{n} \left[\sum_{j=1}^N x_j + \dots + \sum_{j=1}^N x_j \right]$$

Since there are n terms in the series it follows that

$$E(\bar{x}) = \sum_{j=1}^N x_j .$$

Exercise 3.7. As a corollary show that the expected value of $\frac{\bar{x}}{n}$ is equal to the population mean.

By this time, you should be getting familiar with the idea that an estimate from a probability sample is a random variable. Persons responsible for the design and selection of samples and for making estimates from samples are concerned about the set of values, and associated probabilities, that an estimate from a sample might be equal to.

Definition 3.6. The distribution of an estimate generated by probability sampling is the sampling distribution of the estimate.

The values of x_j and P_j in the numerical Illustration 3.7 are an example of a sampling distribution. Statisticians are primarily interested in three characteristics of a sampling distribution: (1) the mean (center) of the sampling distribution in relation to the value of the parameter being estimated, (2) a measure of the variation of possible values of an estimate from the mean of the sampling distribution, and (3) the shape of the sampling distribution. We have been discussing the first. When the expected value of an estimate equals the parameter being estimated, we know that the mean of the sampling distribution is equal to the parameter estimated. But, in practice, values of parameters are generally not known. To judge the accuracy of an estimate, we need

information on all three characteristics of the sampling distribution.

Let us turn now to the generally accepted measure of variation of a random variable.

3.4 VARIANCE OF A RANDOM VARIABLE

The variance of a random variable, X , is the average value of the squares of the deviation of X from its mean; that is, the average value of $(X - \bar{X})^2$. The square root of the variance is the standard deviation (error) of the variable.

Definition 3.7. In terms of expected values, the variance of a random variable, X , is $E(X - \bar{X})^2$ where $E(X) = \bar{X}$. Since X is a random variable, $(X - \bar{X})^2$ is a random variable and by definition of expected value,

$$E(X - \bar{X})^2 = \sum_{i=1}^N P_i (X_i - \bar{X})^2$$

In case $P_i = \frac{1}{N}$ we have the more familiar formula for variance, namely,

$$E(X - \bar{X})^2 = \frac{\sum_{i=1}^N (X_i - \bar{X})^2}{N} = \sigma_X^2$$

Commonly used symbols for variance include: σ^2 , σ_X^2 , V^2 , S^2 , $\text{Var}(X)$

and $V(X)$. Variance is often defined as $\frac{\sum_{i=1}^N (X_i - \bar{X})^2}{N-1}$. This will be discussed in Section 3.7.

3.4.1 VARIANCE OF THE SUM OF TWO INDEPENDENT RANDOM VARIABLES

Two random variables, X and Y , are independent if the joint probability, P_{ij} , of getting X_i and Y_j is equal to $(P_i)(P_j)$, where P_i is the probability of selecting X_i from the set of values of X , and P_j is the probability of selecting Y_j from the set of values of Y . The variance of the sum of two independent random variables is the sum of the variance of each. That is,

$$\sigma_{X+Y}^2 = \sigma_X^2 + \sigma_Y^2$$

Illustration 3.8. In Illustration 3.3, X and Y were independent. We had listed all possible values of $X_i + Y_j$ and the probability of each. From that listing we can readily compute the variance of $X+Y$. By definition

$$\sigma_{X+Y}^2 = E[(X+Y) - (\bar{X} + \bar{Y})]^2 = \sum_{i,j} p_i p_j [(X_i + Y_j) - (\bar{X} + \bar{Y})]^2 \quad (3.6)$$

Substituting in Equation (3.6) we have

$$\sigma_{X+Y}^2 = \frac{2}{24}(0-5.5)^2 + \frac{4}{24}(4-5.5)^2 + \dots + \frac{1}{24}(10-5.5)^2 = \frac{85}{12}$$

The variances of X and Y are computed as follows:

$$\sigma_X^2 = E(X - \bar{X})^2 = \frac{2}{3}(2-4)^2 + \frac{2}{6}(5-4)^2 + \frac{1}{6}(4-4)^2 + \frac{1}{6}(6-4)^2 = \frac{7}{3}$$

$$\sigma_Y^2 = E(Y - \bar{Y})^2 = \frac{1}{4}(-2-1.5)^2 + \frac{2}{4}(2-1.5)^2 + \frac{1}{4}(4-1.5)^2 = \frac{19}{4}$$

We now have $\sigma_X^2 + \sigma_Y^2 = \frac{7}{3} + \frac{19}{4} = \frac{85}{12}$ which verifies the above statement that the variance of the sum of two independent random variables is the sum of the variances.

Exercise 3.8. Prove that $E[(X+Y) - (\bar{X} + \bar{Y})]^2 = E(X+Y)^2 - (\bar{X} + \bar{Y})^2$. Then calculate the variance of $X+Y$ in Illustration 3.3 by using the formula $\sigma_{X+Y}^2 = E(X+Y)^2 - (\bar{X} + \bar{Y})^2$. The answer should agree with the result obtained in Illustration 3.8.

Exercise 3.9. Refer to Illustration 3.3 and the listing of possible values of $X + Y$ and the probability of each. Instead of $X_i + Y_j$ list the products $(X_i - \bar{X})(Y_j - \bar{Y})$ and show that $E(X_i - \bar{X})(Y_j - \bar{Y}) = 0$.

Exercise 3.10. Find $E(X - \bar{X})(Y - \bar{Y})$ for the numerical example used in Illustration 3.3 by the formula $E(XY) - \bar{X}\bar{Y}$ which was derived in Illustration 3.6.

3.4.2 VARIANCE OF THE SUM OF TWO DEPENDENT RANDOM VARIABLES

The variance of dependent random variables involves covariance which is defined as follows:

Definition 3.8. The covariance of two random variables, X and Y , is $E(X-\bar{X})(Y-\bar{Y})$ where $E(X) = \bar{X}$ and $E(Y) = \bar{Y}$. By definition of expected value

$$E(X-\bar{X})(Y-\bar{Y}) = \sum_{ij} P_{ij} (X_i - \bar{X})(Y_j - \bar{Y})$$

where the summation is over all possible values of X and Y .

Symbols commonly used for covariance are σ_{XY} , S_{XY} , and $\text{Cov}(X, Y)$.

Since $(X+Y) - (\bar{X}+\bar{Y}) = (X-\bar{X}) + (Y-\bar{Y})$ we can derive a formula for the variance of $X+Y$ as follows:

$$\begin{aligned} \sigma_{X+Y}^2 &= E[(X+Y) - (\bar{X}+\bar{Y})]^2 \\ &= E[(X-\bar{X}) + (Y-\bar{Y})]^2 \\ &= E[(X-\bar{X})^2 + (Y-\bar{Y})^2 + 2(X-\bar{X})(Y-\bar{Y})] \end{aligned}$$

Then, according to Theorem 3.3,

$$\sigma_{X+Y}^2 = E(X-\bar{X})^2 + E(Y-\bar{Y})^2 + 2E(X-\bar{X})(Y-\bar{Y})$$

and by definition we obtain,

$$\sigma_{X+Y}^2 = \sigma_X^2 + \sigma_Y^2 + 2\sigma_{XY}$$

Sometimes σ_{XX} is used instead of σ_X^2 to represent variance. Thus

$$\sigma_{X+Y}^2 = \sigma_{XX} + \sigma_{YY} + 2\sigma_{XY}$$

For two independent random variables, $P_{ij} = P_i P_j$. Therefore

$$E(X-\bar{X})(Y-\bar{Y}) = \sum_{ij} P_i P_j (X_i - \bar{X})(Y_j - \bar{Y})$$

Write out in longhand, if necessary, and be satisfied that the following is correct:

$$\sum_i \sum_j P_i P_j (X_i - \bar{X})(Y_j - \bar{Y}) = \sum_i P_i (X_i - \bar{X}) \sum_j P_j (Y_j - \bar{Y}) = 0 \quad (3.7)$$

which proves that the covariance σ_{X_i} is zero when X and Y are independent.

Notice that in Equation (3.7) $\sum_i P_i (X_i - \bar{X}) = E(X - \bar{X})$ and $\sum_j P_j (Y_j - \bar{Y}) = E(Y - \bar{Y})$

which, for independent random variables, proves that $E(X - \bar{X})(Y - \bar{Y}) = E(X - \bar{X}) E(Y - \bar{Y})$. When working with independent random variables the following important theorem is frequently very useful:

Theorem 3.4. The expected value of the product of independent random variables u_1, u_2, \dots, u_k is the product of their expected values:

$$E(u_1 u_2 \dots u_k) = E(u_1) E(u_2) \dots E(u_k)$$

3.5 VARIANCE OF AN ESTIMATE

The variance of an estimate from a probability sample depends upon the method of sampling. We will derive the formula for the variance of \bar{x} , the mean of a random sample selected with equal probability, with and without replacement. Then, the variance of an estimate of the population total will be derived for sampling with replacement and unequal probability of selection.

3.5.1 EQUAL PROBABILITY OF SELECTION

The variance of \bar{x} , the mean of a random sample of n elements selected with equal probabilities and with replacement from a population of N , is:

$$\text{Var}(\bar{x}) = \frac{\sigma_X^2}{n}, \quad \text{where } \sigma_X^2 = \frac{1}{N} \sum_{i=1}^N (X_i - \bar{X})^2$$

The proof follows:

By definition, $\text{Var}(\bar{x}) = E[\bar{x} - E(\bar{x})]^2$. We have shown that $E(\bar{x}) = \bar{X}$. Therefore, $\text{Var}(\bar{x}) = E(\bar{x} - \bar{X})^2$. By substitution and algebraic manipulation, we obtain

$$\begin{aligned}
 \text{Var}(\bar{x}) &= E\left[\frac{x_1 + \dots + x_n}{n} - \bar{X}\right]^2 \\
 &= E\left[\frac{(x_1 - \bar{X}) + \dots + (x_n - \bar{X})}{n}\right]^2 \\
 &= \frac{1}{n^2} E\left[\sum_{i=1}^n (x_i - \bar{X})^2 + \sum_{i \neq j} \sum (x_i - \bar{X})(x_j - \bar{X})\right].
 \end{aligned}$$

Applying Theorem 3.3 we now obtain

$$\text{Var}(\bar{x}) = \frac{1}{n^2} \left[\sum_{i=1}^n E(x_i - \bar{X})^2 + \sum_{i \neq j} \sum E(x_i - \bar{X})(x_j - \bar{X}) \right] \quad (3.8)$$

In series form, Equation (3.8) can be written as

$$\text{Var}(\bar{x}) = \frac{1}{n^2} [E(x_1 - \bar{X})^2 + E(x_2 - \bar{X})^2 + \dots + E(x_1 - \bar{X})(x_2 - \bar{X}) + E(x_1 - \bar{X})(x_3 - \bar{X}) + \dots]$$

Since the sampling is with replacement x_1 and x_j are independent and the expected value of all of the product terms is zero. For example, $E(x_1 - \bar{X})(x_2 - \bar{X}) = E(x_1 - \bar{X}) E(x_2 - \bar{X})$ and we know that $E(x_1 - \bar{X})$ and $E(x_2 - \bar{X})$ are zero. Next, consider $E(x_1 - \bar{X})^2$. We have already shown that x_1 is a random variable that can be equal to any one of the population set of values X_1, \dots, X_N with equal probability. Therefore

$$E(x_1 - \bar{X})^2 = \frac{\sum_{j=1}^N (X_j - \bar{X})^2}{N} = \sigma_X^2$$

The same argument applies to x_2, x_3 , etc. Therefore,

$$\sum_{i=1}^n E(x_i - \bar{X})^2 = \sigma_X^2 + \dots + \sigma_X^2 = n\sigma_X^2 \text{ and Equation (3.8) reduces to } \text{Var}(\bar{x}) = \frac{\sigma_X^2}{n}.$$

The mathematics for finding the variance of \bar{x} when the sampling is without replacement is the same as sampling with replacement down to and including Equation (3.8). The expected value of a product term in Equation (3.8) is not zero because x_1 and x_j are not independent. For example, on

the first draw an element has a probability of $\frac{1}{N}$ of being selected, but on the second draw the probability is conditioned by the fact that the element selected on the first draw was not replaced. Consider the first product term in Equation (3.8). To find $E(x_1 - \bar{X})(x_2 - \bar{X})$ we need to consider the set of values that $(x_1 - \bar{X})(x_2 - \bar{X})$ could be equal to. Reference to the following matrix is helpful:

$$\begin{array}{cccc} (x_1 - \bar{X})^2 & (x_1 - \bar{X})(x_2 - \bar{X}) & \dots & (x_1 - \bar{X})(x_N - \bar{X}) \\ (x_2 - \bar{X})(x_1 - \bar{X}) & (x_2 - \bar{X})^2 & \dots & (x_2 - \bar{X})(x_N - \bar{X}) \\ \vdots & \vdots & & \vdots \\ (x_N - \bar{X})(x_1 - \bar{X}) & (x_N - \bar{X})(x_2 - \bar{X}) & \dots & (x_N - \bar{X})^2 \end{array}$$

The random variable $(x_1 - \bar{X})(x_2 - \bar{X})$ has an equal probability of being any of the products in the above matrix, except for the squared terms on the main diagonal. There are $N(N-1)$ such products. Therefore,

$$E(x_1 - \bar{X})(x_2 - \bar{X}) = \frac{\sum_{i \neq j}^N \sum^N (x_i - \bar{X})(x_j - \bar{X})}{N(N-1)}$$

According to Equation (1.9) in Chapter 1,

$$\sum_{i \neq j}^N \sum^N (x_i - \bar{X})(x_j - \bar{X}) = - \sum_i^N (x_i - \bar{X})^2$$

Hence,

$$E(x_1 - \bar{X})(x_2 - \bar{X}) = - \frac{\sum_i^N (x_i - \bar{X})^2}{N(N-1)} = - \frac{\sigma_X^2}{N-1}$$

The same evaluation applies to all other product terms in Equation (3.8).

There are $n(n-1)$ product terms in Equation (3.8) and the expected value of

each is $-\frac{\sigma_X^2}{N-1}$. Thus, Equation (3.8) becomes

$$\text{Var}(\bar{x}) = \frac{1}{n^2} \left[\sum_1^n E(x_i - \bar{X})^2 - n(n-1) \frac{\sigma_X^2}{N-1} \right]$$

Recognizing that $E(x_i - \bar{X})^2 = \sigma_X^2$ and after some easy algebraic operations the answer as follows is obtained:

$$\text{Var}(\bar{x}) = \frac{N-n}{N-1} \frac{\sigma_X^2}{n} \quad (3.9)$$

The factor $\frac{N-n}{N-1}$ is called the correction for finite population because it does not appear when infinite populations are involved or when sampling with replacement which is equivalent to sampling from an infinite population.

For two characteristics, X and Y , of elements in the same simple random sample, the covariance of \bar{x} and \bar{y} is given by a formula analogous to Equation (3.9); namely,

$$\text{Cov}(\bar{x}, \bar{y}) = \frac{N-n}{N-1} \frac{\sigma_{XY}}{n} \quad (3.10)$$

3.5.2 UNEQUAL PROBABILITY OF SELECTION

In Section 3.3 we proved that $x' = \frac{\sum_1^n \frac{x_i}{p_i}}{n}$ is an unbiased estimate of the population total. This was for sampling with replacement and unequal probability of selection. We will now proceed to find the variance of x' .

By definition $\text{Var}(x') = E[x' - E(x')]^2$. Let $X = \sum_1^N X_i$. Then since $E(x') = X$, it follows that

$$\begin{aligned} \text{Var}(x') &= E \left[\frac{\frac{x_1}{p_1} + \dots + \frac{x_n}{p_n}}{n} - X \right]^2 = \frac{1}{n^2} E \left[\left(\frac{x_1}{p_1} - X \right) + \dots + \left(\frac{x_n}{p_n} - X \right) \right]^2 \\ &= \frac{1}{n^2} E \left[\sum_1^n \left(\frac{x_i}{p_i} - X \right)^2 + \sum_{i \neq k} \sum_1^n \left(\frac{x_i}{p_i} - X \right) \left(\frac{x_k}{p_k} - X \right) \right] \end{aligned}$$

(3.11)

Applying Theorem 3.3, $\text{Var}(x')$ becomes

$$\text{Var}(x') = \frac{1}{n^2} \left[\sum_i E\left(\frac{x_i}{p_i} - X\right)^2 + \sum_{i \neq j} \sum_i E\left(\frac{x_i}{p_i} - X\right)\left(\frac{x_j}{p_j} - X\right) \right] \quad (3.11)$$

Notice the similarity of Equations (3.8) and (3.11) and that the steps leading to these two equations were the same. Again, since the sampling is with replacement, the expected value of all product terms in Equation (3.11) is zero. Therefore Equation (3.11) becomes

$$\text{Var}(x') = \frac{1}{n^2} \left[\sum_i E\left(\frac{x_i}{p_i} - X\right)^2 \right]$$

By definition $E\left(\frac{x_i}{p_i} - X\right)^2 = \sum_i p_i \left(\frac{x_i}{p_i} - X\right)^2$

Therefore
$$\text{Var}(x') = \frac{\sum_i p_i \left(\frac{x_i}{p_i} - X\right)^2}{n} \quad (3.12)$$

Exercise 3.11. (a) Refer to Exercise 3.1 and compute the variance of x' for samples of two (that is, $n = 2$) using Equation (3.12). (b) Then turn to Illustration 3.7 and compute the variance of x' from the actual values of x' . Don't overlook the fact that the values of x' have unequal probabilities. According to Definition 3.7, the variance of x' is

$\sum_j p_j (x'_j - X)^2$ where $X = E(x')$, x'_j is one of the 10 possible values of x' , and p_j is the probability of x'_j .

3.6 VARIANCE OF A LINEAR COMBINATION

Before presenting a general theorem on the variance of a linear combination of random variables, a few key variance and covariance relationships will be given. In the following equations X and Y are random variables and a , b , c , and d are constants:

$$\text{Var}(X+a) = \text{Var}(X)$$

$$\text{Var}(aX) = a^2 \text{Var}(X)$$

$$\text{Var}(aX+b) = a^2 \text{Var}(X)$$

$$\text{Cov}(X+a, Y+b) = \text{Cov}(X, Y)$$

$$\text{Cov}(aX, bY) = ab \text{Cov}(X, Y)$$

$$\text{Cov}(aX+b, cY+d) = ac \text{Cov}(X, Y)$$

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

$$\text{Var}(X+Y+a) = \text{Var}(X+Y)$$

$$\text{Var}(aX+bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y) + 2ab \text{Cov}(X, Y)$$

Illustration 3.9. The above relationships are easily verified by using the theory of expected values. For example,

$$\begin{aligned} \text{Var}(aX+b) &= E[aX+b-E(aX+b)]^2 \\ &= E[aX+b-E(aX)-E(b)]^2 \\ &= E[aX-aE(X)]^2 \\ &= E[a(X-\bar{X})]^2 \\ &= a^2 E(X-\bar{X})^2 = a^2 \text{Var}(X) \end{aligned}$$

Exercise 3.12. As in Illustration 3.9 use the theory of expected values to prove that

$$\text{Cov}(aX+b, cY+d) = ac \text{Cov}(X, Y)$$

As in Theorem 3.3, let $u = a_1 u_1 + \dots + a_k u_k$ where a_1, \dots, a_k are constants and u_1, \dots, u_k are random variables. By definition the variance of u is

$$\text{Var}(u) = E[u-E(u)]^2$$

By substitution

$$\begin{aligned} \text{Var}(u) &= E[a_1 u_1 + \dots + a_k u_k - E(a_1 u_1 + \dots + a_k u_k)]^2 \\ &= E[a_1 (u_1 - \bar{u}_1) + \dots + a_k (u_k - \bar{u}_k)]^2 \quad \text{where } E(u_i) = \bar{u}_i \end{aligned}$$

By squaring the quantity in [] and considering the expected values of the terms in the series, the following result is obtained.

Theorem 3.5. The variance of u , a linear combination of random variables, is given by the following equation

$$\text{Var}(u) = \sum_i^k a_i^2 \sigma_i^2 + \sum_{i \neq j} a_i a_j \sigma_{ij}$$

where σ_i^2 is the variance of u_i and σ_{ij} is the covariance of u_i and u_j .

Theorems 3.3 and 3.5 are very useful because many estimates from probability samples are linear combinations of random variables.

Illustration 3.10. Suppose for a srs (simple random sample) that data have been obtained for two characteristics X and Y , the sample values being x_1, \dots, x_n and y_1, \dots, y_n . What is the variance of $\bar{x} - \bar{y}$? From the theory and results that have been presented one can proceed immediately to write the answer. From Theorem 3.5 we know that $\text{Var}(\bar{x} - \bar{y}) = \text{Var}(\bar{x}) + \text{Var}(\bar{y}) - 2\text{Cov}(\bar{x}, \bar{y})$. From the sampling specifications we know the variances of \bar{x} and \bar{y} and the covariance. See Equations (3.9) and (3.10). Thus, the following result is easily obtained:

$$\text{Var}(\bar{x} - \bar{y}) = \left(\frac{N-n}{N-1}\right) \left(\frac{1}{n}\right) (\sigma_X^2 + \sigma_Y^2 - 2\sigma_{XY}) \quad (3.13)$$

Some readers might be curious about the relationship between covariance and correlation. By definition the correlation between X and Y is

$$r_{XY} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

Therefore, one could substitute $r_{XY} \sigma_X \sigma_Y$ for σ_{XY} in Equation (3.13).

Exercise 3.13. In a statistical publication suppose you find 87 bushels per acre as the yield of corn in State A and 83 is the estimated yield for State B. The estimated standard errors are given as 1.5 and

2.0 bushels. You become interested in the standard error of the difference in yield between the two States and want to know how large the estimated difference is in relation to its standard error. Find the standard error of the difference. You may assume that the two yield estimates are independent because the sample selection in one State was completely independent of the other. Answer: 2.5.

Illustration 3.11. No doubt students who are familiar with sampling have already recognized the application of Theorems 3.3 and 3.5 to several sampling plans and methods of estimation. For example, for stratified random sampling, an estimator of the population total is

$$x' = N_1 \bar{x}_1 + \dots + N_k \bar{x}_k = \sum_{i=1}^k N_i \bar{x}_i$$

where N_i is the population number of sampling units in the i^{th} stratum and \bar{x}_i is the average per sampling unit of characteristic, X , from a sample of n_i sampling units from the i^{th} stratum. According to Theorem 3.3

$$E(x') = \sum_{i=1}^k N_i E(\bar{x}_i)$$

If the sampling is such that $E(\bar{x}_i) = \bar{X}_i$ for all strata, x' is an unbiased estimate of the population total. According to Theorem 3.5

$$\text{Var}(x') = N_1^2 \text{Var}(\bar{x}_1) + \dots + N_k^2 \text{Var}(\bar{x}_k) \quad (3.14)$$

There are no covariance terms in Equation (3.14) because the sample selection in one stratum is independent of another stratum. Assuming a srs from each stratum, Equation (3.14) becomes

$$\text{Var}(x') = N_1^2 \left(\frac{N_1 - n_1}{N_1 - 1} \right) \frac{\sigma_1^2}{n_1} + \dots + N_k^2 \left(\frac{N_k - n_k}{N_k - 1} \right) \frac{\sigma_k^2}{n_k}$$

where σ_i^2 is the variance of X among sampling units within the i^{th} stratum.

Illustration 3.12. Suppose x'_1, \dots, x'_k are independent estimates of the same quantity, T . That is, $E(x'_i) = T$. Let σ_i^2 be the variance of x'_i .

Consider a weighted average of the estimates, namely

$$x' = w_1 x'_1 + \dots + w_k x'_k \quad (3.15)$$

where $\sum w_i = 1$. Then

$$E(x') = w_1 E(x'_1) + \dots + w_k E(x'_k) = T \quad (3.16)$$

That is, for any set of weights where $\sum w_i = 1$ the expected value of x' is T . How should the weights be chosen?

The variance of x' is

$$\text{Var}(x') = w_1^2 \sigma_1^2 + \dots + w_k^2 \sigma_k^2$$

If we weight the estimates equally, $w_i = \frac{1}{k}$ and the variance of x' is

$$\text{Var}(x') = \frac{1}{k} \left[\frac{\sum \sigma_i^2}{k} \right] \quad (3.17)$$

which is the average variance divided by k . However, it is reasonable to give more weight to estimates having low variance. Using differential calculus we can find the weights which will minimize the variance of x' .

The optimum weights are inversely proportional to the variances of the estimates. That is, $w_i \propto \frac{1}{\sigma_i^2}$

As an example, suppose one has two independent unbiased estimates of the same quantity which originate from two different samples. The optimum weighting of the two estimates would be

$$\frac{\frac{1}{\sigma_1^2} x'_1 + \frac{1}{\sigma_2^2} x'_2}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}}$$

As another example, suppose x'_1, \dots, x'_k are the values of X in a sample of k sampling units selected with equal probability and with replacement. In this case each x'_i is an unbiased estimate of \bar{X} . If we let $w_i = \frac{1}{k}$, x' is \bar{x} , the simple average of the sample values. Notice, as one would expect, Equation (3.16) reduces to $E(\bar{x}) = \bar{X}$. Also, since each estimate, x'_i , is the same random variable that could be equal to any value in the set X_1, \dots, X_N , it is clear that all of the σ_i^2 's must be equal to $\sigma^2 = \frac{\sum (X_i - \bar{X})^2}{N}$. Hence, Equation (3.17) reduces to $\frac{\sigma^2}{n}$ which agrees with the first part of Section 3.5.1.

Exercise 3.14. If you equate x'_i in Equation (3.15) with $\frac{x_i}{p_i}$ in Section 3.5.2 and let $w_i = \frac{1}{n}$ and $k = n$, then x' in Equation (3.15) is the

same as $x' = \frac{\sum \frac{x_i}{p_i}}{n}$ in Section 3.5.2. Show that in this case Equation (3.17) becomes the same as Equation (3.12).

3.7 ESTIMATION OF VARIANCE

All of the variance formulas presented in previous sections have involved calculations from a population set of values. In practice, we have data for only a sample. Hence, we must consider means of estimating variances from sample data.

3.7.1 SIMPLE RANDOM SAMPLING

In Section 3.5.1, we found that the variance of the mean of a srs is

$$\text{Var}(\bar{x}) = \frac{N-n}{N-1} \frac{\sigma_X^2}{n} \quad (3.18)$$

where

$$\sigma_X^2 = \frac{\sum (X_i - \bar{X})^2}{N}$$

BEST COPY AVAILABLE

As an estimator of σ_X^2 , $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ seems like a natural first choice for consideration. However, when sampling finite populations, it is customary to define variance among units of the population as follows:

$$S^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2$$

and to use $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$ as an estimator of S^2 . A reason for this will become apparent when we find the expected value of s^2 as follows:

The formula for s^2 can be written in a form that is more convenient for finding $E(s^2)$. Thus,

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{\sum_{i=1}^n x_i^2 - n\bar{x}^2}{n-1}$$

and
$$E(s^2) = \frac{1}{n-1} [\sum_{i=1}^n E(x_i^2) - nE(\bar{x}^2)]$$

We have shown previously that x_i is a random variable that has an equal probability of being any value in the set X_1, \dots, X_N . Therefore

$$E(x_i^2) = \frac{1}{N} \sum_{i=1}^N X_i^2 \quad \text{and} \quad \sum_{i=1}^n E(x_i^2) = \frac{n \sum_{i=1}^N X_i^2}{N}$$

Hence,
$$E(s^2) = \frac{n}{n-1} \left[\frac{\sum_{i=1}^N X_i^2}{N} - E(\bar{x}^2) \right] \quad (3.19)$$

We know, by definition, that $\sigma_x^2 = E(\bar{x} - \bar{X})^2$ and it is easy to show that

$$E(\bar{x} - \bar{X})^2 = E(\bar{x}^2) - \bar{X}^2$$

Therefore,
$$E(\bar{x}^2) = \sigma_x^2 + \bar{X}^2$$

BEST COPY AVAILABLE

By substitution in Equation (3.19) we obtain

$$E(s^2) = \frac{n}{n-1} \left[\frac{\sum X_i^2}{N} - \bar{X}^2 - \sigma_x^2 \right]$$

By definition $\sigma_x^2 = \frac{\sum (X_i - \bar{X})^2}{N} = \frac{\sum X_i^2}{N} - \bar{X}^2$ and since the specified method of sampling was srs, $\sigma_x^2 = \frac{N-n}{N-1} \frac{\sigma_X^2}{n}$, we have $E(s^2) = \frac{n}{n-1} \left[\sigma_x^2 - \frac{N-n}{N-1} \frac{\sigma_X^2}{n} \right]$

which after simplification is

$$E(s^2) = \frac{N}{N-1} \sigma_X^2$$

Note from the above definitions of σ_X^2 and S^2 that

$$S^2 = \frac{N}{N-1} \sigma_X^2$$

Therefore $E(s^2) = S^2$

Since s^2 is an unbiased estimate of S^2 , we will now substitute $\frac{N-1}{N} s^2$ for σ_x^2 in Equation (3.18) which gives

$$\text{Var}(\bar{x}) = \frac{N-n}{N} \frac{S^2}{n} \quad (3.20)$$

Both Equations, (3.18) and (3.20), for the $\text{Var}(\bar{x})$ give identical results and both agree with $E(\bar{x} - \bar{X})^2$ as a definition of variance. We have shown that s^2 is an unbiased estimate of S^2 . Substituting s^2 for S^2 in Equation (3.20) we have

$$\text{var}(\bar{x}) = \frac{N-n}{N} \frac{s^2}{n} \quad (3.21)$$

as an estimate of the variance of \bar{x} . With regard to Equation (3.18), $\frac{N-1}{N} s^2$ is an unbiased estimate of σ_x^2 . When $\frac{N-1}{N} s^2$ is substituted for σ_x^2 , Equation (3.21) is obtained.

Since in Equation (3.20), $\frac{N-n}{N}$ is exactly 1 minus the sampling fraction and s^2 is an unbiased estimate of S^2 , there is some advantage to using

Equation (3.20) and $S^2 = \frac{\sum (X_i - \bar{X})^2}{N-1}$ as a definition of variance among sampling units in the population.

Exercise 3.15. For a small population of 4 elements suppose the values of X are $X_1 = 2$, $X_2 = 5$, $X_3 = 3$; and $X_4 = 6$. Consider simple random samples of size 2. There are six possible samples.

- (a) For each of the six samples calculate \bar{x} and s^2 . That is, find the sampling distribution of \bar{x} and the sampling distribution of s^2 .
- (b) Calculate S^2 , then find $\text{Var}(\bar{x})$ using Equation (3.20).
- (c) Calculate the variance among the six values of \bar{x} and compare the result with $\text{Var}(\bar{x})$ obtained in (b). The results should be the same.
- (d) From the sampling distribution of s^2 calculate $E(s^2)$ and verify that $E(s^2) = S^2$.

3.7.2 UNEQUAL PROBABILITY OF SELECTION

In Section 3.5.2, we derived a formula for the variance of the estimator x' where

$$x' = \frac{\sum \frac{x_i}{p_i}}{n} \quad (3.22)$$

The sampling was with unequal selection probabilities and with replacement.

We found that the variance of x' was given by

$$\text{Var}(x') = \frac{\sum p_i \left(\frac{x_i}{p_i} - \bar{X} \right)^2}{n} \quad (3.23)$$

As a formula for estimating $\text{Var}(x')$ from a sample one might be inclined, as a first guess, to try a formula of the same form as Equation (3.23) but

that does not work. Equation (3.23) is a weighted average of the squares of deviations $(\frac{x_i}{p_i} - X)^2$ which reflects the unequal selection probabilities.

If one applied the same weighting system in a formula for estimating variance from a sample he would in effect be applying the weights twice; first, in the selection process itself and second, to the sample data.

The unequal probability of selection is already incorporated into the sample itself.

As in some of the previous discussion, look at the estimator as follows:

$$x' = \frac{\frac{x_1}{p_1} + \dots + \frac{x_n}{p_n}}{n} = \frac{x'_1 + \dots + x'_n}{n} \text{ where } x'_i = \frac{x_i}{p_i}$$

Each x'_i is an independent unbiased estimate of the population total. Since each value of x'_i receives an equal weight in determining x' it appears that the following formula for estimating $\text{Var}(x')$ might work:

$$\text{var}(x') = \frac{s^2}{n} \quad (3.24)$$

where

$$s^2 = \frac{\sum_{i=1}^n (x'_i - x')^2}{n-1}$$

By following an approach similar to that used in Section 3.7.1, one can prove that

$$E(s^2) = \sum_{i=1}^N P_i \left(\frac{x_i}{p_i} - X \right)^2$$

That is, Equation (3.24) does provide an unbiased estimate of $\text{Var}(x')$ in Equation (3.23). The proof is left as an exercise.

Exercise 3.16. Reference is made to Exercise 3.1, Illustration 3.7, and Exercise 3.11. In Illustration 3.7 the sampling distribution of x'

(See Equation (3.22)) is given for samples of 2 from the population of 4 elements that was given in Exercise 3.1.

- (a) Compute $\text{var}(\bar{x}) = \frac{s^2}{n}$ (Equation (3.24)) for each of the 10 possible samples.
- (b) Compute the expected value of $\text{var}(\bar{x})$ and compare it with the result obtained in Exercise 3.11. The results should be the same. Remember, when finding the expected value of $\text{var}(\bar{x})$, that the \bar{x} 's do not occur with equal frequency.

3.8 RATIO OF TWO RANDOM VARIABLES

In sampling theory and practice one frequently encounters estimates that are ratios of random variables. It was pointed out earlier that $E(\frac{u}{w}) \neq \frac{E(u)}{E(w)}$ where u and w are random variables. Formulas for the expected value of a ratio and for the variance of a ratio will now be presented without derivation. The formulas are approximations:

$$E\left(\frac{u}{w}\right) \approx \frac{\bar{u}}{\bar{w}} + \frac{\bar{u}}{\bar{w}} \left[\frac{\sigma_u^2}{\bar{u}^2} - \frac{\rho_{uw} \sigma_u \sigma_w}{\bar{u}\bar{w}} \right] \quad (3.25)$$

$$\text{Var}\left(\frac{u}{w}\right) \approx \left[\frac{\bar{u}}{\bar{w}}\right]^2 \left[\frac{\sigma_u^2}{\bar{u}^2} + \frac{\sigma_w^2}{\bar{w}^2} - \frac{2\rho_{uw} \sigma_u \sigma_w}{\bar{u}\bar{w}} \right] \quad (3.26)$$

where

$$\bar{u} = E(u)$$

$$\bar{w} = E(w)$$

$$\sigma_u^2 = E(u - \bar{u})^2$$

$$\sigma_w^2 = E(w - \bar{w})^2$$

and

$$\rho_{uw} = \frac{\sigma_{uw}}{\sigma_u \sigma_w} \quad \text{where } \sigma_{uw} = E(u - \bar{u})(w - \bar{w})$$

For a discussion of the conditions under which Equations (3.25) and (3.26) are good approximations, reference is made to Hansen, Hurwitz, and

Madow. 2/ The conditions are usually satisfied with regard to estimates from sample surveys. As a rule of thumb the variance formula is usually accepted as satisfactory if the coefficient of variation of the variable in the denominator is less than 0.1; that is, if $\frac{\sigma_w}{\bar{w}} < 0.1$. In other words, this condition states that the coefficient of variation of the estimate in the denominator should be less than 10 percent. A larger coefficient of variation might be tolerable before becoming concerned about Equation (3.26) as an approximation.

The condition $\frac{\sigma_w}{\bar{w}} < 0.1$ is more stringent than necessary for regarding the bias of a ratio as negligible. With few exceptions in practice the bias of a ratio is ignored. Some of the logic for this will appear in the illustration below. To summarize, the conditions when Equations (3.25) and (3.26) are not good approximations are such that the ratio is likely to be of questionable value owing to large variance.

If u and w are linear combinations of random variables, the theory presented in previous sections applies to u and to w . Assuming u and w are estimates from a sample, to estimate $\text{Var}(\frac{u}{w})$ take into account the sample design and substitute in Equation (3.26) estimates of \bar{u} , \bar{w} , σ_u^2 , σ_w^2 , and ρ_{uw} . Ignore Equation (3.25) unless there is reason to believe the bias of the ratio might be important relative to its standard error.

It is of interest to note the similarity between $\text{Var}(u-w)$ and $\text{Var}(\frac{u}{w})$. According to Theorem 3.5,

$$\text{Var}(u-w) = \sigma_u^2 + \sigma_w^2 - 2\rho_{uw} \sigma_u \sigma_w$$

2/ Hansen, Hurwitz, and Madow, Sample Surveys Methods and Theory, Volume I, Chapter 4, John Wiley and Sons, 1953.

By definition the relative variance of an estimate is the variance of the estimate divided by the square of its expected value. Thus, in terms of the relative variance of a ratio, Equation (3.26) can be written

$$\text{Rel Var}\left(\frac{u}{w}\right) = \frac{\sigma_u^2}{\bar{u}^2} + \frac{\sigma_w^2}{\bar{w}^2} - 2\rho_{uw} \frac{\sigma_u \sigma_w}{\bar{u} \bar{w}}$$

The similarity is an aid to remembering the formula for $\text{Var}\left(\frac{u}{w}\right)$.

Illustration 3.13. Suppose one has a simple random sample of n elements from a population of N . Let \bar{x} and \bar{y} be the sample means for characteristics X and Y . Then, $u = \bar{x}$, $w = \bar{y}$,

$$\sigma_u^2 = \frac{N-n}{N} \frac{S_X^2}{n} \quad \text{and} \quad \sigma_w^2 = \frac{N-n}{N} \frac{S_Y^2}{n}$$

Notice that the condition discussed above, $\frac{\sigma_w}{\bar{w}} < 0.1$, is satisfied if the sample is large enough so

$$\frac{N-n}{N} \frac{S_Y^2}{n \bar{y}^2} < 0.1^2$$

Substituting in Equation (3.26) we obtain the following as the variance of the ratio:

$$\text{Var}\left(\frac{\bar{x}}{\bar{y}}\right) = \left(\frac{N-n}{N}\right) \left(\frac{1}{n}\right) \frac{\bar{x}^2}{\bar{y}^2} \left[\frac{S_X^2}{\bar{x}^2} + \frac{S_Y^2}{\bar{y}^2} - \frac{2\rho_{XY} S_X S_Y}{\bar{X} \bar{Y}} \right]$$

The bias of $\frac{\bar{x}}{\bar{y}}$ as an estimate of $\frac{\bar{X}}{\bar{Y}}$ is given by the second term of Equation (3.25). For this illustration it becomes

$$\left(\frac{N-n}{N}\right) \left(\frac{1}{n}\right) \frac{\bar{x}}{\bar{y}} \left[\frac{S_Y^2}{\bar{y}^2} - \frac{\rho_{XY} \sigma_X \sigma_Y}{\bar{X} \bar{Y}} \right]$$

As the size of the sample increases, the bias decreases as $\frac{1}{n}$ whereas the standard error of the ratio decreases at a slower rate, namely $\frac{1}{\sqrt{n}}$.

Thus, we need not be concerned about a possibility of the bias becoming important relative to sampling error as the size of the sample increases. A possible exception occurs when several ratios are combined. An example is stratified random sampling when many strata are involved and separate ratio estimates are made for the strata. This is discussed in the books on sampling.

9.9 CONDITIONAL EXPECTATION

The theory for conditional expectation and conditional variance of a random variable is a very important part of sampling theory, especially in the theory for multistage sampling. The theory will be discussed with reference to two-stage sampling.

The notation that will be used in this and the next section is as follows:

M is the number of psu's (primary sampling units) in the population.

m is the number of psu's in the sample.

N_i is the total number of elements in the i^{th} psu.

$N = \sum_{i=1}^M N_i$ is the total number of elements in the population.

n_i is the sample number of elements from the i^{th} psu.

$n = \sum_{i=1}^m n_i$ is the total number of elements in the sample.

$$\bar{n} = \frac{n}{m}$$

X_{ij} is the value of X for the j^{th} element in the i^{th} psu. It

refers to an element in the population, that is, $j = 1, \dots, N_i$,

and $i = 1, \dots, M$.

x_{ij} is the value of X for the j^{th} element in the sample from the i^{th} psu in the sample, that is, the indexes i and j refer to the set of psu's and elements in the sample.

$X_{i.} = \sum_j x_{ij}$ is the population total for the i^{th} psu.

$\bar{X}_{i.} = \frac{X_{i.}}{N_i}$ is the average of X for all elements in the i^{th} psu.

$\bar{X}_{..} = \frac{\sum_i \sum_j x_{ij}}{N} = \frac{\sum_i X_{i.}}{N}$ is the average of all N elements.

$\bar{X}_{.} = \frac{\sum_i X_{i.}}{M}$ is the average of the psu totals. Be sure to note the difference between $\bar{X}_{..}$ and $\bar{X}_{.}$.

$x_{i.} = \sum_j x_{ij}$ is the sample total for the i^{th} psu in the sample.

$\bar{x}_{i.} = \frac{x_{i.}}{n_i}$ is the average for the n_i elements in the sample from the i^{th} psu.

$\bar{x}_{..} = \frac{\sum_i \sum_j x_{ij}}{n}$ is the average for all elements in the sample.

Assume simple random sampling, equal probability of selection without replacement, at both stages. Consider the sample of n_i elements from the i^{th} psu. We know from Section 3.3 that $\bar{x}_{i.}$ is an unbiased estimate of the psu mean $\bar{X}_{i.}$; that is, $E(\bar{x}_{i.}) = \bar{X}_{i.}$ and for a fixed i (a specified psu) $EN_i \bar{x}_{i.} = N_i E(\bar{x}_{i.}) = N_i \bar{X}_{i.} = X_{i.}$. But, owing to the first stage of sampling,

$\sum_i \bar{x}_i$ must be treated as a random variable. Hence, it is necessary to become involved with the expected value of an expected value.

First, consider X as a random variable, in the context of single-stage sampling, which could equal any one of the values X_{ij} in the population set of $N = \sum_i^M N_i$. Let $P(ij)$ be the probability of selecting the j^{th} element in the i^{th} psu; that is, $P(ij)$ is the probability of X being equal to X_{ij} . By definition

$$E(X) = \sum_i^M \sum_j^1 P(ij) X_{ij} \quad (3.27)$$

Now consider the selection of an element as a two-step procedure:

(1) selected a psu with probability $P(i)$, and (2) selected an element within the selected psu with probability $P(j|i)$. In words, $P(j|i)$ is the probability of selecting the j^{th} element in the i^{th} psu given that the i^{th} psu has already been selected. Thus, $P(ij) = P(i)P(j|i)$. By substitution, Equation (3.27) becomes

$$E(X) = \sum_i^M \sum_j^1 P(i)P(j|i) X_{ij}$$

or

$$E(X) = \sum_i^M P(i) \sum_j^1 P(j|i) X_{ij} \quad (3.28)$$

By definition, $\sum_j^1 P(j|i) X_{ij}$ is the expected value of X for a fixed value of i . It is called "conditional expectation."

Let $E_2(X|i) = \sum_j^1 P(j|i) X_{ij}$ where $E_2(X|i)$ is the form of notation we

will be using to designate conditional expectation. To repeat, $E_2(X|i)$ means the expected value of X for a fixed i . The subscript 2 indicates

that the conditional expectation applies to the second stage of sampling. E_1 and E_2 will refer to expectation at the first and second stages, respectively.

Substituting $E_2(X|i)$ in Equation (3.28) we obtain

$$E(X) = \sum_{i=1}^M P(i) E_2(X|i) \quad (3.29)$$

There is one value of $E_2(X|i)$ for each of the M psu's. In fact $E_2(X|i)$ is a random variable where the probability of $E_2(X|i)$ is $P(i)$. Thus the right-hand side of Equation (3.29) is, by definition, the expected value of $E_2(X|i)$. This leads to the following theorem:

Theorem 3.6. $E(X) = E_1 E_2(X|i)$

Suppose $P(j|i) = \frac{1}{N_i}$ and $P(i) = \frac{1}{M}$. Then,

$$E_2(X|i) = \sum_{j=1}^{N_i} \left(\frac{1}{N_i}\right) X_{ij} = \bar{X}_{i.}$$

and
$$E(X) = E_1(\bar{X}_{i.}) = \sum_{i=1}^M \left(\frac{1}{M}\right) (\bar{X}_{i.}) = \frac{\sum \bar{X}_{i.}}{M}.$$

In this case $E(X)$ is an unweighted average of the psu averages. It is important to note that, if $P(i)$ and $P(j|i)$ are chosen in such a way that $P(ij)$ is constant, every element has the same chance of selection. This point will be discussed later.

Theorem 3.3 dealt with the expected value of a linear combination of random variables. There is a corresponding theorem for conditional expectation. Assume the linear combination is

$$U = a_1 u_1 + \dots + a_k u_k = \sum_{t=1}^k a_t u_t$$

where a_1, \dots, a_k are constants and u_1, \dots, u_k are random variables. Let $E(U|c_i)$ be the expected value of U under a specified condition, c_i , where c_i is one of the conditions out of a set of M conditions that could occur. The theorem on conditional expectation can then be stated symbolically as follows:

Theorem 3.7. $E(U|c_i) = a_1 E(u_1|c_i) + \dots + a_k E(u_k|c_i)$ ✓

$$\text{or } E(U|c_i) = \sum_t^k a_t E(u_t|c_i)$$

Compare Theorems 3.7 and 3.3 and note that Theorem 3.7 is like Theorem 3.3 except that conditional expectation is applied. Assume c is a random event and that the probability of the event c_i occurring is $P(i)$. Then $E(U|c_i)$ is a random variable and by definition the expected value of $E(U|c_i)$ is $\sum_i^M P(i) E(U|c_i)$ which is $E(U)$. Thus, we have the following theorem:

Theorem 3.8. The expected value of U is the expected value of the conditional expected value of U , which in symbols is written as follows:

$$E(U) = E[E(U|c_i)] \quad (3.30)$$

Substituting the value of $E(U|c_i)$ from Theorem 3.7 in Equation (3.30) we have

$$E(U) = E[a_1 E(u_1|c_i) + \dots + a_k E(u_k|c_i)] = E[\sum_t^k a_t E(u_t|c_i)] \quad (3.31)$$

Illustration 3.14. Assume two-stage sampling with simple random sampling at both stages. Let x' , defined as follows, be the estimator of the population total:

$$x' = \frac{M}{m} \sum_i^m \frac{N_i}{n_i} \sum_j^{n_i} x_{ij} \quad (3.32)$$

Exercise 3.17. Examine the estimator, x' , Equation (3.32). Express it in other forms that might help show its logical structure. For example, for a fixed i what is $\frac{N_i}{n_i} \sum_j x_{ij}$? Does it seem like a reasonable way of estimating the population total?

To display x' as a linear combination of random variables it is convenient to express it in the following form:

$$x' = \left[\frac{M}{m} \frac{N_1}{n_1} x_{11} + \dots + \frac{M}{m} \frac{N_1}{n_1} x_{1n_1} \right] + \dots + \left[\frac{M}{m} \frac{N_m}{n_m} x_{m1} + \dots + \frac{M}{m} \frac{N_m}{n_m} x_{mn_m} \right] \quad (3.33)$$

Suppose we want to find the expected value of x' to determine whether it is equal to the population total. According to Theorem 3.8,

$$E(x') = E_1 E_2 (x' | i) \quad (3.34)$$

$$E(x') = E_1 E_2 \left\{ \left[\frac{M}{m} \sum_i \frac{N_i}{n_i} \sum_j x_{ij} \right] | i \right\} \quad (3.35)$$

Equations (3.34) and (3.35) are obtained simply by substituting x' as the random variable in (3.30). The c_i now refers to any one of the m psu's in the sample. First we must solve the conditional expectation, $E_2(x' | i)$. Since $\frac{M}{m}$ and $\frac{N_i}{n_i}$ are constant with respect to the conditional expectation, and making use of Theorem 3.7, we can write

$$E_2(x' | i) = \frac{M}{m} \sum_i \frac{N_i}{n_i} \sum_j E_2(x_{ij} | i) \quad (3.36)$$

We know for any given psu in the sample that x_{ij} is an element in a simple random sample from the psu and according to Section 3.3 its expected value is the psu mean, \bar{X}_i . That is,

$$E_2(x_{ij} | i) = \bar{X}_i.$$

and
$$\sum_{j=1}^{n_1} E_2(x_{1j}|1) = n_1 \bar{X}_1. \quad (3.37)$$

Substituting the result from Equation (3.37) in Equation (3.36) gives

$$E_2(x'|1) = \frac{M}{m} \sum_{i=1}^m N_i \bar{X}_1. \quad (3.38)$$

Next we need to find the expected value of $E_2(x'|1)$. In Equation (3.38), N_i is a random variable, as well as \bar{X}_1 , associated with the first stage of sampling. Accordingly, we will take $X_{i.} = N_i \bar{X}_1$ as the random variable which gives in lieu of Equation (3.38).

$$E_2(x'|1) = \frac{M}{m} \sum_{i=1}^m X_{i.}$$

Therefore,

$$E(x') = E_1\left[\frac{M}{m} \sum_{i=1}^m X_{i.}\right]$$

From Theorem 3.3

$$E_1\left[\frac{M}{m} \sum_{i=1}^m X_{i.}\right] = \frac{M}{m} \sum_{i=1}^m E_1(X_{i.})$$

Since

$$\sum_{i=1}^m E_1(X_{i.}) = \sum_{i=1}^m \frac{\sum X_{i.}}{M}$$

$$E_1\left[\frac{M}{m} \sum_{i=1}^m X_{i.}\right] = \frac{M}{m} \sum_{i=1}^m X_{i.}$$

Therefore, $E(x') = \sum_{i=1}^M X_{i.} = X_{..}$ This shows that x' is an unbiased estimator of the population total.

3.10 CONDITIONAL VARIANCE

Conditional variance refers to the variance of a variable under a specified condition or limitation. It is related to conditional probability and to conditional expectation.

To find the variance of \bar{x} (See Equation (3.32) or (3.33)) the following important theorem will be used:

Theorem 3.9. The variance of \bar{x} is given by

$$V(\bar{x}) = V_1 E_2(\bar{x}^2 | 1) + E_1 V_2(\bar{x}^2 | 1).$$

where V_1 is the variance for the first stage of sampling and V_2 is the "conditional" variance for the second stage.

We have discussed $E_2(\bar{x}^2 | 1)$ and noted there is one value of $E_2(\bar{x}^2 | 1)$ for each psu in the population. Hence $V_1 E_2(\bar{x}^2 | 1)$ is simply the variance of the M values of $E_2(\bar{x}^2 | 1)$.

In Theorem 3.9 the conditional variance, $V_2(\bar{x}^2 | 1)$, by definition is

$$V_2(\bar{x}^2 | 1) = E_2\{[\bar{x} - E_2(\bar{x} | 1)]^2 | 1\}$$

To understand $V_2(\bar{x}^2 | 1)$ think of \bar{x} as a linear combination of random variables (see Equation (3.33)). Consider the variance of \bar{x} when i is held constant. All terms (random variables) in the linear combination are now constant except those originating from sampling within the i^{th} psu. Therefore, $V_2(\bar{x}^2 | 1)$ is associated with variation among elements in the i^{th} psu. $V_2(\bar{x}^2 | 1)$ is a random variable with M values in the set, one for each psu. Therefore, $E_1 V_2(\bar{x}^2 | 1)$ by definition is

$$E_1 V_2(\bar{x}^2 | 1) = \sum_i^M P(i) V_2(\bar{x}^2 | i)$$

That is, $E_1 V_2(\bar{x}^2 | 1)$ is an average of M values of $V_2(\bar{x}^2 | i)$ weighted by $P(i)$, the probability that the i^{th} psu had of being in the sample.

Three illustrations of the application of Theorem 3.9 will be given.

In each case there will be five steps in finding the variance of \bar{x} :

Step 1, find $E_2(\bar{x}^2 | 1)$

Step 2, find $V_1 E_2(\bar{x}^2 | 1)$

Step 3, find $V_2(x'|i)$

Step 4, find $E_1 V_2(x'|i)$

Step 5, combine results from Steps 2 and 4.

Illustration 3.15. This is a simple illustration, selected because we know what the answer is from previous discussion and a linear combination of random variables is not involved. Suppose x' in Theorem 3.9 is simply the random variable X where X has an equal probability of being any one of the X_{ij} values in the set of $N = \sum_{i=1}^M N_i$. We know that the variance of X can be expressed as follows:

$$V(x') = \frac{1}{N} \sum_{ij}^M (X_{ij} - \bar{X}_{..})^2 \quad (3.39)$$

In the case of two-stage sampling, an equivalent method of selecting a value of X is to select a psu first and then select an element within the psu, the condition being that $P(ij) = P(i)P(j|i) = \frac{1}{N}$. This condition is satisfied by letting $P(i) = \frac{N_i}{N}$ and $P(j|i) = \frac{1}{N_i}$. We now want to find $V(X)$ by using Theorem 3.9 and check the result with Equation (3.39).

Step 1. From the random selection specifications we know that $E_2(x'|i) = \bar{X}_{i.}$. Therefore,

$$\text{Step 2. } V_1 E_2(x'|i) = V_1(\bar{X}_{i.})$$

We know that $\bar{X}_{i.}$ is a random variable that has a probability of $\frac{N_i}{N}$ of being equal to the i^{th} value in the set $\bar{X}_1, \dots, \bar{X}_M$. Therefore, by definition of the variance of a random variable,

$$V_1 E(x'|i) = \sum_{i=1}^M \frac{N_i}{N} (\bar{X}_{i.} - \bar{X}_{..})^2 \quad (3.40)$$

where

$$\bar{X}_{..} = \sum_{i=1}^M \frac{N_i}{N} \bar{X}_{i.} = \frac{\sum X_{ij}}{N}$$

Step 3. By definition

$$V_2(x'|i) = \sum_j \frac{N_1}{N} (x_{1j} - \bar{x}_{1.})^2$$

Step 4. Since each value of $V_2(x'|i)$ has a probability $\frac{N_1}{N}$

$$E_1 V_2(x'|i) = \sum_i \frac{N_1}{N} \sum_j \frac{N_1}{N} (x_{1j} - \bar{x}_{1.})^2 \quad (3.41)$$

Step 5. From Equations (3.40) and (3.41) we obtain

$$V(x') = \frac{1}{N} \left[\sum_i N_1 (\bar{x}_{1.} - \bar{x}_{..})^2 + \sum_i \sum_j \frac{N_1}{N} (x_{1j} - \bar{x}_{1.})^2 \right] \quad (3.42)$$

The fact that Equations (3.42) and (3.39) are the same is verified by Equation (1.10) in Chapter I.

Illustration 3.16. Find the variance of the estimator x' given by Equation (3.32) assuming simple random sampling at both stages of sampling.

Step 1. Theorem 3.7 is applicable. That is,

$$E_2(x'|i) = \sum_{ij} \frac{mn_1}{N} E_2 \left[\frac{M}{m} \frac{N_1}{n_1} x_{1j} | i \right]$$

which means "sum the conditional expected values of each of the n terms in Equation (3.33)."

With regard to any one of the terms in Equation (3.33), the conditional expectation is

$$E_2 \left[\frac{M}{m} \frac{N_1}{n_1} x_{1j} | i \right] = \frac{M}{m} \frac{N_1}{n_1} E_2(x_{1j} | i) = \frac{M}{m} \frac{N_1}{n_1} \bar{x}_{1.} = \frac{M}{m} \frac{X_{1.}}{n_1}$$

Therefore

$$E_2(x'|i) = \sum_{ij} \frac{mn_1}{N} \frac{M}{m} \frac{X_{1.}}{n_1} \quad (3.43)$$

With reference to Equation (3.43) and summing with respect to j , we have

$$\sum_{j=1}^n \frac{1}{m} \frac{X_{j.}}{n_j} = \frac{1}{m} \sum_{i=1}^M X_{i.}$$

Hence Equation (3.43) becomes

$$E_2(x'|1) = \frac{1}{m} \sum_{i=1}^m X_{i.} \quad (3.44)$$

Step 2. Find $V_1 E_2(x'|1)$. This is simple because $\frac{1}{m} \sum_{i=1}^m X_{i.}$ in Equation (3.44) is the mean of a random sample of m from the set of psu totals X_1, \dots, X_M . Therefore,

$$V_1 E_2(x'|1) = M^2 \left(\frac{M-m}{M-1} \right) \frac{\sigma_{b1}^2}{m} \quad (3.45)$$

where

$$\sigma_{b1}^2 = \frac{\sum_{i=1}^M (X_{i.} - \bar{X}_.)^2}{M} \quad \text{and} \quad \bar{X}_. = \frac{\sum_{i=1}^M X_{i.}}{M}$$

In the subscript to σ^2 , the "b" indicates between psu variance and "1" distinguishes this variance from between psu variances in later illustrations.

Step 3. Finding $V_2(x'|1)$, is more involved because the conditional variance of a linear combination of random variables must be derived.

However, this is analogous to using Theorem 3.5 for finding the variance of a linear combination of random variables. Theorem 3.5 applies except that $V(u|1)$ replaces $V(u)$ and conditional variance and conditional covariance replace the variances and covariances in the formula for $V(u)$. As the solution proceeds, notice that the strategy is to shape the problem so previous results can be used.

Look at the estimator x' , Equation (3.33), and determine whether any covariances exist. An element selected from one psu is independent of an

element selected from another; but within a psu the situation is the same as the one we had when finding the variance of the mean of a simple random sample. This suggests writing \bar{x} in terms of \bar{x}_i , because the \bar{x}_i 's are independent. Accordingly, we will start with

$$\bar{x} = \frac{M}{m} \sum_{i=1}^m N_i \bar{x}_i.$$

Hence

$$V_2(\bar{x}|1) = V_2\left(\left[\frac{M}{m} \sum_{i=1}^m N_i \bar{x}_i\right]|1\right)$$

Since the \bar{x}_i 's are independent

$$V_2(\bar{x}|1) = \frac{M^2}{m^2} \sum_{i=1}^m V_2(N_i \bar{x}_i|1)$$

and since N_i is constant with regard to the conditional variance

$$V_2(\bar{x}|1) = \frac{M^2}{m^2} \sum_{i=1}^m N_i^2 V_2(\bar{x}_i|1) \quad (3.46)$$

Since the sampling within each psu is simple random sampling

$$V_2(\bar{x}_i|1) = \left(\frac{N_i - n_i}{N_i - 1}\right) \frac{\sigma_i^2}{n_i} \quad (3.47)$$

where

$$\sigma_i^2 = \sum_{j=1}^{N_i} \frac{1}{N_i} (x_{ij} - \bar{x}_i)^2$$

Step 4. After substituting the value of $V_2(\bar{x}_i|1)$ in Equation (3.46), and then applying Theorem 3.3, we have

$$E_1 V_2(\bar{x}|1) = \frac{M^2}{m^2} \sum_{i=1}^m E_1 \left[N_i^2 \frac{N_i - n_i}{N_i - 1} \frac{\sigma_i^2}{n_i} \right]$$

Since the first stage of sampling was simple random sampling and each psu had an equal chance of being in the sample,

$$E_1 \left[N_1^2 \frac{N_1 - n_1}{N_1 - 1} \frac{\sigma_1^2}{n_1} \right] = \frac{1}{M} \sum_1^M N_1^2 \frac{N_1 - n_1}{N_1 - 1} \frac{\sigma_1^2}{n_1}$$

Hence

$$E_1 V_2(x'|1) = \frac{M}{m} \sum_1^M N_1^2 \frac{N_1 - n_1}{N_1 - 1} \frac{\sigma_1^2}{n_1} \quad (3.48)$$

Step 5. Combining Equation (3.48) and Equation (3.45) the answer is

$$V(x') = M^2 \frac{M-m}{M-1} \frac{\sigma_{b1}^2}{m} + \frac{M}{m} \sum_1^M N_1^2 \frac{N_1 - n_1}{N_1 - 1} \frac{\sigma_1^2}{n_1} \quad (3.49)$$

Illustration 3.17. The sampling specifications are: (1) at the first stage select m psu's with replacement and probability $P(1) = \frac{N_1}{N}$, and (2) at the second stage a simple random sample of \bar{n} elements is to be selected from each of the m psu's selected at the first stage. This will give a sample of $n = m\bar{n}$ elements. Find the variance of the sample estimate of the population total.

The estimator needs to be changed because the psu's are not selected with equal probability. Sample values need to be weighted by the reciprocals of their probabilities of selection if the estimator is to be unbiased. Let

$P'(ij)$ be the probability of element ij being in the sample,

$P'(i)$ be the relative frequency of the i^{th} psu being in a sample of m , and let

$P'(j|i)$ equal the conditional probability of element ij being in the sample given that the i^{th} psu is already in the sample.

Then

$$P'(ij) = P'(i)P'(j|i)$$

According to the sampling specifications $P'(i) = m \frac{N_1}{N}$. This probability was described as relative frequency because "probability of being

in a sample of m psu's" is subject to misinterpretation. The i^{th} psu can appear in a sample more than once and it is counted every time it appears. That is, if the i^{th} psu is selected more than once, a sample of \bar{n} is selected within the i^{th} psu every time that it is selected. By substitution

$$P'(i_j) = \left[m \frac{N_i}{N} \right] \frac{\bar{n}}{N_i} = \frac{m\bar{n}}{N} = \frac{n}{N} \quad (3.50)$$

Equation (3.50) means that every element has an equal probability of being in the sample. Consequently, the estimator is very simple,

$$x' = \frac{N}{m\bar{n}} \sum_{i,j} \sum x_{ij} \quad (3.51)$$

Exercise 3.18. Show that x' , Equation (3.51), is an unbiased estimator of the population total.

In finding $V(x')$ our first step was to solve for $E_2(x'|1)$.

Step 1. By definition

$$E_2(x'|1) = E_2 \left\{ \left[\frac{N}{m\bar{n}} \sum_{i,j} \sum x_{ij} \right] | 1 \right\}$$

Since 1 is constant with regard to E_2 ,

$$E_2(x'|1) = \frac{N}{m\bar{n}} \sum_{i,j} \sum E_2(x_{ij}|1) \quad (3.52)$$

Proceeding from Equation (3.52) to the following result is left as an exercise:

$$E_2(x'|1) = \frac{N}{m} \sum_i \bar{x}_i \quad (3.53)$$

Step 2. From Equation (3.53) we have

$$V_1 E_2(x'|1) = V_1 \left(\frac{N}{m} \sum_i \bar{x}_i \right)$$

Since the $\bar{x}_{1.}$'s are independent

BEST COPY AVAILABLE

$$V_1 E_2(x'|1) = \frac{N^2}{m} \sum_{i=1}^m V_1(\bar{x}_{1.})$$

Because the first stage of sampling is sampling with probability proportional to N_i and with replacement,

$$V_1(\bar{x}_{1.}) = \sum_{i=1}^M \frac{N_i}{N} (\bar{x}_{1.} - \bar{x}_{..})^2 \quad (3.54)$$

Let

$$V_1(\bar{x}_{1.}) = \sigma_{b2}^2$$

Then

$$V_1 E_2(x'|1) = \frac{N^2}{m} (m \sigma_{b2}^2) = \frac{N^2}{m} \sigma_{b2}^2 \quad (3.55)$$

Exercise 3.19. Prove that $E(\bar{x}_{1.}) = \bar{x}_{..}$ which shows that it is appropriate to use $\bar{x}_{..}$ in Equation (3.54).

Step 3. To find $V_2(x'|1)$, first write the estimator as

$$x' = \frac{N}{m} \sum_{i=1}^m \bar{x}_{1.} \quad (3.56)$$

Then, since the $\bar{x}_{1.}$'s are independent

$$V_2(x'|1) = \frac{N^2}{m} \sum_{i=1}^m V_2(\bar{x}_{1.})$$

and

$$V_2(\bar{x}_{1.}) = \frac{N_i - \bar{n}}{N_i - 1} \frac{\sigma_i^2}{\bar{n}}$$

where

$$\sigma_i^2 = \sum_{j=1}^{N_i} \frac{1}{N_i} (x_{1j} - \bar{x}_{1.})^2$$

Therefore

$$V_2(x'|1) = \frac{N^2}{m^2} \sum_i^m \frac{N_1 - \bar{n}}{N_1 - 1} \frac{\sigma_1^2}{\bar{n}}$$

Step 4.

$$E_1 V_2(x'|1) = \frac{N^2}{m^2} \frac{1}{\bar{n}} \sum_i^m E_1 \left(\frac{N_1 - \bar{n}}{N_1 - 1} \sigma_1^2 \right)$$

Since the probability of $V_2(x'|1)$ is $\frac{N_1}{N}$

$$E_1 V_2(x'|1) = \frac{N^2}{m^2} \frac{1}{\bar{n}} \sum_i^m \left[\sum_i^M \frac{N_1}{N} \left(\frac{N_1 - \bar{n}}{N_1 - 1} \right) \sigma_1^2 \right]$$

which becomes

$$E_1 V_2(x'|1) = \frac{N^2}{mn} \sum_i^M \frac{N_1}{N} \left(\frac{N_1 - \bar{n}}{N_1 - 1} \right) \sigma_1^2 \quad (3.57)$$

Step 5. Combining Equation (3.55) and Equation (3.57) we have the

answer

$$V(x') = N^2 \left[\frac{\sigma_b^2}{m} + \frac{1}{mn} \sum_i^M \frac{N_1}{N} \left(\frac{N_1 - \bar{n}}{N_1 - 1} \right) \sigma_1^2 \right] \quad (3.58)$$

CHAPTER IV. THE DISTRIBUTION OF AN ESTIMATE

4.1 PROPERTIES OF SIMPLE RANDOM SAMPLES

The distribution of an estimate is a primary basis for judging the accuracy of an estimate from a sample survey. But an estimate is only one number. How can one number have a distribution? Actually, "distribution of an estimate" is a phrase that refers to the distribution of all possible estimates that might occur under repetition of a prescribed sampling plan and estimator (method of estimation). Thanks to theory and empirical testing of the theory, it is not necessary to generate physically the distribution of an estimate by selecting numerous samples and making an estimate from each. However, to have a tangible distribution of an estimate as a basis for discussion, an illustration has been prepared.

Illustration 4.1. Consider simple random samples of 4 from an assumed population of 8 elements. There are $\frac{N!}{n!(N-n)!} = \frac{8!}{4!4!} = 70$ possible samples. In Table 4.1, the sample values for all of the 70 possible samples of four are shown. The 70 samples were first listed in an orderly manner to facilitate getting all of them accurately recorded. The mean, \bar{x} , for each sample was computed and the samples were then arrayed according to the value of \bar{x} for purposes of presentation in Table 4.1. The distribution of \bar{x} is the 70 values of \bar{x} shown in Table 4.1, including the fact that each of the 70 values of \bar{x} has an equal probability of being the estimate. These 70 values have been arranged as a frequency distribution in Table 4.2.

As discussed previously, one of the properties of simple random sampling is that the sample average is an unbiased estimate of the population average; that is, $E(\bar{x}) = \bar{X}$. This means that the distribution of

Table 4.1--Samples of four elements from a population of eight ^{1/}

Sample number	Values of x_i	\bar{x}	s^2	Sample number	Values of x_i	\bar{x}	s^2
1c	2,1,6,4	3.25	4.917	36s	1,6,8,9	6.00	12.667
2	2,1,4,7	3.50	7.000	37s	1,4,8,11	6.00	19.333
3	2,1,4,8	3.75	9.583	38s	2,6,8,9	6.25	9.583
4	2,1,6,7	4.00	8.667	39s	2,4,8,11	6.25	16.250
5	2,1,4,9	4.00	12.667	40s	1,6,7,11	6.25	16.917
6	2,1,6,8	4.25	10.917	41s	1,4,11,9	6.25	20.917
7	2,1,6,9	4.50	13.667	42	1,7,8,9	6.25	12.917
8	2,1,4,11	4.50	20.333	43cs	6,4,7,8	6.25	2.917
9cs	2,1,7,8	4.50	12.333	44s	2,6,7,11	6.50	13.667
10	1,6,4,7	4.50	7.000	45s	2,4,11,9	6.50	17.667
11s	2,1,7,9	4.75	14.917	46	2,7,8,9	6.50	9.667
12	2,6,4,7	4.75	4.917	47s	1,6,8,11	6.50	17.667
13	1,6,4,8	4.75	8.917	48s	6,4,7,9	6.50	4.333
14	2,1,6,11	5.00	20.667	49s	2,6,8,11	6.75	14.250
15s	2,1,8,9	5.00	16.667	50s	1,6,11,9	6.75	18.917
16	2,6,4,8	5.00	6.667	51	1,7,8,11	6.75	17.583
17	1,6,4,9	5.00	11.333	52s	6,4,8,9	6.75	4.917
18s	1,4,7,8	5.00	10.000	53s	2,6,11,9	7.00	15.333
19s	2,1,7,11	5.25	21.583	54	2,7,8,11	7.00	14.000
20	2,6,4,9	5.25	8.917	55	1,7,11,9	7.00	18.667
21s	2,4,7,8	5.25	7.583	56s	6,4,7,11	7.00	8.667
22s	1,4,7,9	5.25	12.250	57	4,7,8,9	7.00	4.667
23s	2,1,8,11	5.50	23.000	58	2,7,11,9	7.25	14.917
24s	2,4,7,9	5.50	9.667	59	1,8,11,9	7.25	18.917
25	1,6,4,11	5.50	17.667	60s	6,4,8,11	7.25	8.917
26s	1,6,7,8	5.50	9.667	61	2,8,11,9	7.50	15.000
27s	1,4,8,9	5.50	13.667	62cs	6,4,11,9	7.50	9.667
28cs	2,1,11,9	5.75	24.917	63	6,7,8,9	7.50	1.667
29	2,6,4,11	5.75	14.917	64	4,7,8,11	7.50	8.333
30s	2,6,7,8	5.75	6.917	65	4,7,11,9	7.75	8.917
31s	2,4,8,9	5.75	10.917	66	6,7,8,11	8.00	4.667
32s	1,6,7,9	5.75	11.583	67	4,8,11,9	8.00	8.667
33s	1,4,7,11	5.75	18.250	68	6,7,11,9	8.25	4.917
34s	2,6,7,9	6.00	8.667	69	6,8,11,9	8.50	4.333
35s	2,4,7,11	6.00	15.333	70c	7,8,11,9	8.75	2.917

^{1/} Values of X for the population of eight elements are $X_1 = 2$, $X_2 = 1$, $X_3 = 6$, $X_4 = 4$, $X_5 = 7$, $X_6 = 8$, $X_7 = 11$, $X_8 = 9$; $\bar{X} = 6.00$; and

$$s^2 = \frac{\sum (X_i - \bar{X})^2}{N-1} = 12.$$

Table 4.2--Sampling distribution of \bar{x}

\bar{x}	Relative frequency of \bar{x}		
	Simple random sampling Illustration 4.1	Cluster sampling Illustration 4.2	Stratified random sampling Illustration 4.2
3.25	1	1	
3.50	1		
3.75	1		
4.00	2		
4.25	1		
4.50	4	1	1
4.75	3		1
5.00	5		2
5.25	4		3
5.50	5		4
5.75	6	1	5
6.00	4		4
6.25	6	1	5
6.50	5		4
6.75	4		3
7.00	5		2
7.25	3		1
7.50	4	1	1
7.75	1		
8.00	2		
8.25	1		
8.50	1		
8.75	1	1	
Total	70	6	36
Expected value of \bar{x}	6.00	6.00	6.00
Variance of \bar{x}	1.50	3.29	0.49

\bar{x} is centered on \bar{X} . If the theory is correct, the average of \bar{x} for the 70 samples, which are equally likely to occur, should be equal to the population average, 6.00. The average of the 70 samples does equal 6.00.

From the theory of expected values, we also know that the variance of \bar{x} is given by

$$s_{\bar{x}}^2 = \frac{N-n}{N} \frac{S^2}{n} \quad (4.1)$$

where

$$S^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2$$

With reference to Illustration 4.1 and Table 4.1, $S^2 = 12.00$ and $s_{\bar{x}}^2 = \frac{8-4}{8} \frac{12}{4} = 1.5$. The formula (4.1) can be verified by computing the variance among the 70 values of \bar{x} as follows:

$$\frac{(3.25-6.00)^2 + (3.50-6.00)^2 + \dots + (8.75-6.00)^2}{70} = 1.5$$

Since S^2 is a population parameter, it is usually unknown. Fortunately, as discussed in Chapter 3, $E(s^2) = S^2$ where

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

In Table 4.1, the value of s^2 is shown for each of the 70 samples. The average of the 70 values of s^2 is equal to S^2 . The fact that $E(s^2) = S^2$ is another important property of simple random samples. In practice s^2 is used as an estimate of S^2 . That is,

$$s_{\bar{x}}^2 = \frac{N-n}{N} \frac{s^2}{n}$$

is an unbiased estimate of the variance of \bar{x} .

To recapitulate, we have just verified three important properties of simple random samples:

$$(1) E(\bar{x}) = \bar{X}$$

$$(2) S_{\bar{x}} = \sqrt{\frac{N-n}{N}} \frac{S}{\sqrt{n}}$$

$$(3) E(s^2) = S^2$$

The standard error of \bar{x} , namely $S_{\bar{x}}$, is a measure of how much \bar{x} varies under repeated sampling from \bar{X} . Incidentally, notice that Equation (4.1) shows how the variance of \bar{x} is related to the size of the sample. Now we need to consider the form or shape of the distribution of \bar{x} .

Definition 4.1. The distribution of an estimate is often called the sampling distribution. It refers to the distribution of all possible values of an estimate that could occur under a prescribed sampling plan.

4.2 SHAPE OF THE SAMPLING DISTRIBUTION

For random sampling there is a large volume of literature on the distribution of an estimate which we will not attempt to review. In practice, the distribution is generally accepted as being normal (See Figure 4.1) unless the sample size is "small." The theory and empirical tests show that the distribution of an estimate approaches the normal distribution rapidly as the size of the sample increases. The closeness of the distribution of an estimate to the normal distribution depends on:

- (1) the distribution of X (i.e., the shape of the frequency distribution of the values of X in the population being sampled),
- (2) the form of the estimator,
- (3) the sample design, and
- (4) the sample size.

It is not possible to give a few simple, exact guidelines for deciding when the degree of approximation is good enough. In practice, it is generally a matter of working as though the distribution of an estimate is normal but being mindful of the possibility that the distribution might differ

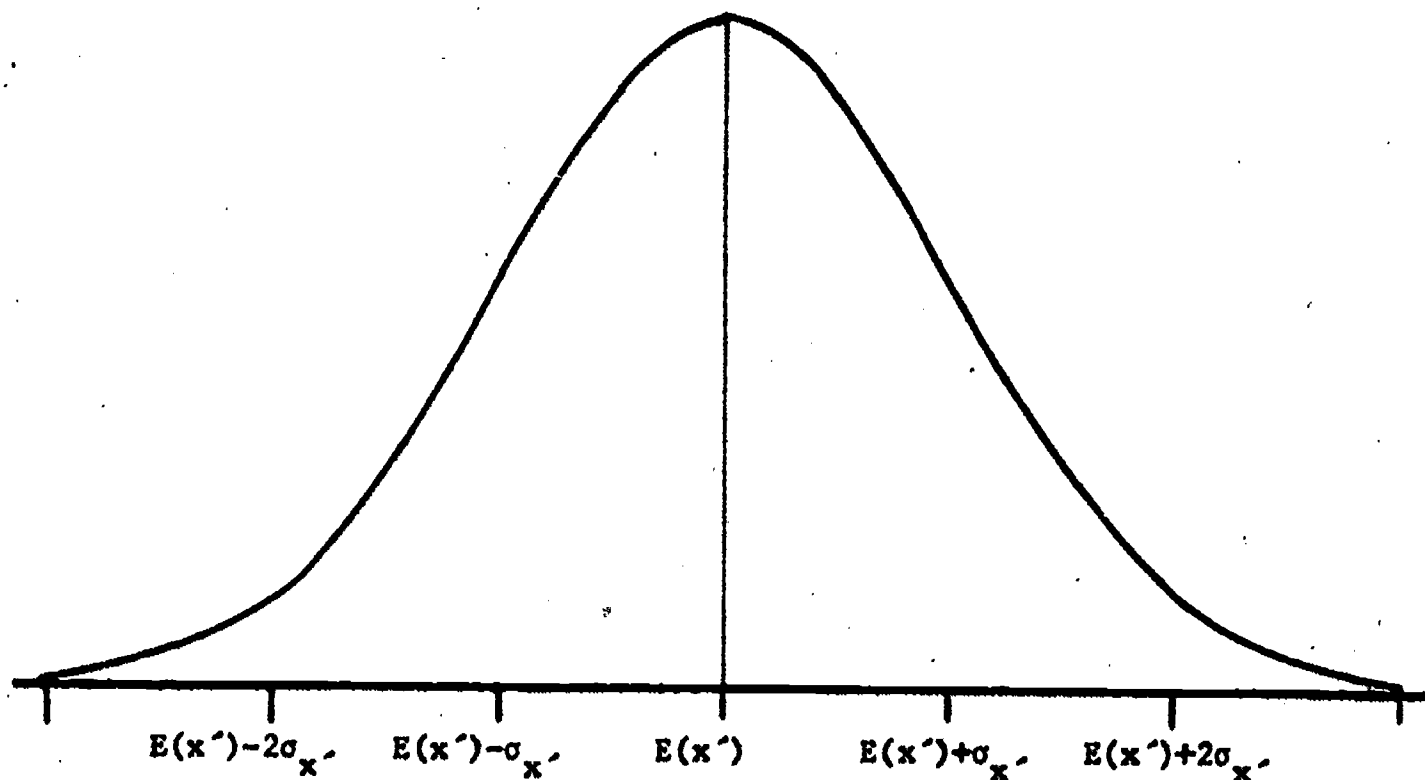


Figure 4.1--Distribution of an estimate (normal distribution)

considerably from normal when the sample is very small and the population distribution is highly skewed. ^{3/}

It is very fortunate that the sampling distribution is approximately normal as it gives a basis for probability statements about the precision of an estimate. As notation, x' will be the general expression for any estimate, and $\sigma_{x'}$ is the standard error of x' .

Figure 4.1 is a graphical representation of the sampling distribution of an estimate. It is the normal distribution. In the mathematical equation for the normal distribution of a variable there are two parameters: the average value of the variable, and the standard error of the variable.

^{3/} For a good discussion of the distribution of a sample estimate, see Vol. I, Chapter 1, Hansen, Hurwitz, and Madow. Sample Survey Methods and Theory, John Wiley and Sons, 1953.

Suppose \bar{x} is an estimate from a probability sample. The characteristics of the sampling distribution of \bar{x} are specified by three things: (1) the expected value of \bar{x} , $E(\bar{x})$, which is the mean of the distribution; (2) the standard error of \bar{x} , $\sigma_{\bar{x}}$, and (3) the assumption that the distribution is normal. If \bar{x} is normally distributed, two-thirds of the values that \bar{x} could equal are between $[E(\bar{x}) - \sigma_{\bar{x}}]$ and $[E(\bar{x}) + \sigma_{\bar{x}}]$, 95 percent of the possible values of \bar{x} are between $[E(\bar{x}) - 2\sigma_{\bar{x}}]$ and $[E(\bar{x}) + 2\sigma_{\bar{x}}]$, and 99.7 percent of the estimates are within $3\sigma_{\bar{x}}$ from $E(\bar{x})$.

Exercise 4.1. With reference to Illustration 4.1, find $E(\bar{x}) - \sigma_{\bar{x}}$ and $E(\bar{x}) + \sigma_{\bar{x}}$. Refer to Table 4.2 and find the proportion of the 70 values of \bar{x} that are between $E(\bar{x}) - \sigma_{\bar{x}}$ and $E(\bar{x}) + \sigma_{\bar{x}}$. How does this compare with the expected proportion assuming the sampling distribution of \bar{x} is normal? The normal approximation is not expected to be close, owing to the small size of the population and of the sample. Also compute $E(\bar{x}) - 2\sigma_{\bar{x}}$ and $E(\bar{x}) + 2\sigma_{\bar{x}}$ and find the proportion of the 70 values of \bar{x} that are between these two limits.

4.3 SAMPLE DESIGN

There are many methods of designing and selecting samples and of making estimates from samples. Each sampling method and estimator has a sampling distribution. Since the sampling distribution is assumed to be normal, alternative methods are compared in terms of $E(\bar{x})$ and $\sigma_{\bar{x}}$ (or $\sigma_{\bar{x}}^2$).

For simple random sampling, we have seen, for a sample of n , that every possible combination of n elements has an equal chance of being the sample selected. Some of these possible combinations (samples) are much better than others. It is possible to introduce restrictions in sampling so some of the combinations cannot occur or so some combinations have a

higher probability of occurrence than others. This can be done without introducing bias in the estimate \bar{x} and without losing a basis for estimating $\sigma_{\bar{x}}$. Discussion of particular sample designs is not a primary purpose of this chapter. However, a few simple illustrations will be used to introduce the subject of design and to help develop concepts of sampling variation.

Illustration 4.2. Suppose the population of 8 elements used in Table 4.1 is arranged so it consists of four sampling units as follows:

<u>Sampling Unit</u>	<u>Elements</u>	<u>Values of X</u>	<u>Sample Unit Total</u>
1	1,2	$X_1 = 2, X_2 = 1$	3
2	3,4	$X_3 = 6, X_4 = 4$	10
3	5,6	$X_5 = 7, X_6 = 8$	15
4	7,8	$X_7 = 11, X_8 = 9$	20

For sampling purposes the population now consists of four sampling units rather than eight elements. If we select a simple random sample of two sampling units from the population of four sampling units, it is clear that the sampling theory for simple random sampling applies. This illustration points out the importance of making a clear distinction between a sampling unit and an element that a measurement pertains to. A sampling unit corresponds to a random selection and it is the variation among sampling units (random selections) that determines the sampling error of an estimate. When the sampling units are composed of more than one element, the sampling is commonly referred to as cluster sampling because the elements in a sampling unit are usually close together geographically.

For a simple random sample of 2 sampling units, the variance of \bar{x}_c , where \bar{x}_c is the sample average per sampling unit, is

$$s_{\bar{x}_c}^2 = \frac{N-n}{N} \frac{S_c^2}{n} = 13.17$$

where

$$N = 4, n = 2, \text{ and } S_c^2 = \frac{(3-12)^2 + (10-12)^2 + (15-12)^2 + (20-12)^2}{3} = \frac{158}{3}$$

Instead of the average per sampling unit one will probably be interested in the average per element, which is $\bar{x} = \frac{\bar{x}_c}{2}$, since there are two elements in each sampling unit. The variance of \bar{x} is one-fourth of the variance of \bar{x}_c . Hence, the variance of \bar{x} is $\frac{13.17}{4} = 3.29$.

There are only six possible random samples as follows:

Sample	Sampling Units	Sample average per sampling unit, \bar{x}_c	s_c^2
1	1,2	6.5	24.5
2	1,3	9.0	72.0
3	1,4	11.5	144.5
4	2,3	12.5	12.5
5	2,4	15.0	50.0
6	3,4	17.5	12.5

where $s_c^2 = \frac{\sum (x_i - \bar{x}_c)^2}{n-1}$ and x_i is a sampling unit total. Be sure to notice that s_c^2 (which is the sample estimate of S_c^2) is the variance among sampling units in the sample, not the variance among individual elements in the sample. From the list of six samples, it is easy to verify that \bar{x}_c is an unbiased estimate of the population average per sampling unit and that s_c^2 is an unbiased estimate of $\frac{158}{3}$, the variance among the four sampling

units in the population. Also, the variance among the six values of \bar{x} is 13.17 which agrees with the formula.

The six possible cluster samples are among the 70 samples listed in Table 4.1. Their sample numbers in Table 4.1 are 1, 9, 28, 43, 62, and 70. A "c" follows these sample numbers. The sampling distribution for the six samples is shown in Table 4.2 for comparison with simple random sampling. It is clear from inspection that random selection from these six is less desirable than random selection from the 70. For example, one of the two extreme averages, 3.25 or 8.75, has a probability of $\frac{1}{3}$ of occurring for the cluster sampling and a probability of only $\frac{1}{35}$ when selecting a simple random sample of four elements. In this illustration, the sampling restriction (clustering of elements) increased the sampling variance from 1.5 to 3.29.

It is of importance to note that the average variance among elements within the four clusters is only 1.25. (Students should compute the within cluster variances and verify 1.25). This is much less than 12.00, the variance among the 8 elements of the population. In reality, the variance among elements within clusters is usually less than the variance among all elements in the population, because clusters (sampling units) are usually composed of elements that are close together and elements that are close together usually show a tendency to be alike.

Exercise 4.2. In Illustration 4.2, if the average variance among elements within clusters had been greater than 12.00, the sampling variance for cluster sampling would have been less than the sampling variance for a simple random sample of elements. Repeat what was done in Illustration 4.2

using as sampling units elements 1 and 6, 2 and 5, 3 and 8, and 4 and 7.

Study the results.

Illustration 4.3. Perhaps the most common method of sampling is to assign sampling units of a population to groups called strata. A simple random sample is then selected from each stratum. Suppose the population used in Illustration 4.1 is divided into two strata as follows:

Stratum 1 $X_1 = 2, X_2 = 1, X_3 = 6, X_4 = 4$

Stratum 2 $X_5 = 7, X_6 = 8, X_7 = 11, X_8 = 9$

The sampling plan is to select a simple random sample of two elements from each stratum. There are 36 possible samples of 4, two from each stratum. These 36 samples are identified in Table 4.1 by an s after the sample number so you may compare the 36 possible stratified random samples with the 70 simple random samples and with the six cluster samples. Also, see Table 4.2.

Consider the variance of \bar{x} . We can write

$$\bar{x} = \frac{\bar{x}_1 + \bar{x}_2}{2}$$

where \bar{x}_1 is the sample average for stratum 1 and \bar{x}_2 is the average for stratum 2. According to Theorem 3.5

$$S_{\bar{x}}^2 = \left(\frac{1}{4}\right)(S_{\bar{x}_1}^2 + S_{\bar{x}_2}^2 + 2S_{\bar{x}_1\bar{x}_2})$$

We know the covariance, $S_{\bar{x}_1\bar{x}_2}$, is zero because the sampling from one stratum is independent of the sampling from the other stratum. And, since the sample within each stratum is a simple random sample,

$$S_{\bar{x}_1}^2 = \frac{N_1 - n_1}{N_1} \frac{S_1^2}{n_1} \quad \text{where} \quad S_1^2 = \frac{\sum_{i=1}^{N_1} (X_{1i} - \bar{X}_1)^2}{N_1 - 1}$$

The subscript "1" refers to stratum 1. $S_{x_2}^2$ is of the same form as $S_{x_1}^2$.

Therefore,

$$S_{\bar{x}}^2 = \frac{1}{4} \left[\frac{N_1 - n_1}{N_1} \frac{S_1^2}{n_1} + \frac{N_2 - n_2}{N_2} \frac{S_2^2}{n_2} \right]$$

Since

$$\frac{N_1 - n_1}{N_1} = \frac{N_2 - n_2}{N_2} = \frac{1}{2}, \text{ and } n_1 = n_2 = 2,$$

$$S_{\bar{x}}^2 = \frac{1}{8} \left[\frac{S_1^2 + S_2^2}{2} \right] = \frac{1}{8} \left[\frac{4.92 + 2.92}{2} \right] = 0.49$$

The variance, 0.49, is comparable to 1.5 in Illustration 4.1 and to 3.29 in Illustration 4.2.

In Illustration 4.2, the sampling units were groups of two elements and the variance among these groups (sampling units) appeared in the formula for the variance of \bar{x} . In Illustration 4.3, each element was a sampling unit but the selection process (randomization) was restricted to taking one stratum (subset) at a time, so the sampling variance was determined by variability within strata. As you study sampling plans, form mental pictures of the variation which the sampling error depends on. With experience and accumulated knowledge of what the patterns of variation in various populations are like, one can become expert in judging the efficiency of alternative sampling plans in relation to specific objectives of a survey.

If the population and the samples in the above illustrations had been larger, the distributions in Table 4.2 would have been approximately normal. Thus, since the form of the distribution of an estimate from a probability sample survey is accepted as being normal, only two attributes of an estimate need to be evaluated, namely its expected value and its variance.

In the above illustrations ideal conditions were implicitly assumed. Such conditions do not exist in the real world so the theory must be extended to fit, more exactly, actual conditions. There are numerous sources of error or variation to be evaluated. The nature of the relationship between theory and practice is a major governing factor determining the rate of progress toward improvement of the accuracy of survey results.

We will now extend error concepts toward more practical settings.

4.4 RESPONSE ERROR

So far, we have discussed sampling under implicit assumptions that measurements are obtained from all n elements in a sample and that the measurement for each element is without error. Neither assumption fits, exactly, the real world. In addition, there are "coverage" errors of various kinds. For example, for a farm survey a farm is defined but application of the definition involves some degree of ambiguity about whether particular enterprises satisfy the definition. Also, two persons might have an interest in the same farm tract giving rise to the possibility that the tract might be counted twice (included as a part of two farms) or omitted entirely.

Partly to emphasize that error in an estimate is more than a matter of sampling, statisticians often classify the numerous sources of error into one of two general classes: (1) Sampling errors which are errors associated with the fact that one has measurements for a sample of elements rather than measurements for all elements in the population, and (2) non-sampling errors--errors that occur whether sampling is involved or not. Mathematical error models can be very complex when they include a term for

each of many sources of error and attempt to represent exactly the real world. However, complicated error models are not always necessary, depending upon the purposes.

For purposes of discussion, two oversimplified response-error models will be used. This will introduce the subject of response error and give some clues regarding the nature of the impact of response error on the distribution of an estimate. For simplicity, we will assume that a measurement is obtained for each element in a random sample and that no ambiguity exists regarding the identity or definition of an element. Thus, we will be considering sampling error and response error simultaneously.

Illustration 4.4. Let T_1, \dots, T_N be the "true values" of some variable for the N elements of a population. The mention of true values raises numerous questions about what is a true value. For example, what is your true weight? How would you define the true weight of an individual? We will refrain from discussing the problem of defining true values and simply assume that true values do exist according to some practical definition. When an attempt is made to ascertain T_i , some value other than T_i might be obtained. Call the actual value obtained X_i . The difference, $e_i = X_i - T_i$, is the response error for the i^{th} element. If the characteristic, for example, is a person's weight, the observed weight, X_i , for the i^{th} individual depends upon when and how the measurement is taken. However, for simplicity, assume that X_i is always the value obtained regardless of the conditions under which the measurement is taken. In other words, assume that the response error, e_i , is constant for the i^{th} element. In this hypothetical case, we are actually sampling a population set of values X_1, \dots, X_N instead of a set of true values T_1, \dots, T_N .

Under the conditions as stated, the sampling theory applies exactly to the set of population values X_1, \dots, X_N . If a simple random sample of elements is selected and measurements for all elements in the sample are

obtained, then $E(\bar{x}) = \bar{X}$. That is, if the purpose is to estimate $\bar{T} = \frac{1}{N} \sum_{i=1}^N T_i$, the estimate is biased unless \bar{T} happens to be equal to \bar{X} . The bias is $\bar{X} - \bar{T}$ which is appropriately called "response bias."

Rewrite $e_i = X_i - T_i$ as follows:

$$X_i = T_i + e_i \quad (4.2)$$

Then, the mean of a simple random sample may be expressed as

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{\sum_{i=1}^n (t_i + e_i)}{n},$$

or, as
$$\bar{x} = \bar{t} + \bar{e}.$$

From the theory of expected values, we have

$$E(\bar{x}) = E(\bar{t}) + E(\bar{e})$$

Since $E(\bar{x}) = \bar{X}$ and $E(\bar{t}) = \bar{T}$ it follows that

$$\bar{X} = \bar{T} + E(\bar{e})$$

Thus, \bar{x} is a biased estimate of \bar{T} unless $E(\bar{e}) = 0$, where $E(\bar{e}) = \frac{1}{N} \sum_{i=1}^N e_i$.

That is, $E(\bar{e})$ is the average of the response errors, e_i , for the whole population.

For simple random sampling the variance of \bar{x} is

$$s_{\bar{x}}^2 = \frac{N-n}{N} \frac{s_X^2}{n} \quad \text{where} \quad s_X^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2$$

How does the response error affect the variance of X and of \bar{x} ? We have already written the observed value for the i^{th} element as being equal to

its true value plus a response error, that is, $X_i = T_i + e_i$. Assuming random sampling, T_i and e_i are random variables. We can use Theorem 3.5 from Chapter III and write

$$S_X^2 = S_T^2 + S_e^2 + 2S_{T,e} \quad (4.3)$$

where S_X^2 is the variance of X , S_T^2 is the variance of T , S_e^2 is the response variance (that is, the variance of e), and $S_{T,e}$ is the covariance of T and e . The terms on the right-hand side of Equation (4.3) cannot be evaluated unless data on X_i and T_i are available; however, the equation does show how the response error influences the variance of X and hence of \bar{X} .

As a numerical example, assume a population of five elements and the following values for T and X :

	<u>T_i</u>	<u>X_i</u>	<u>e_i</u>
	23	26	3
	13	12	-1
	17	23	6
	25	25	0
	<u>7</u>	<u>9</u>	<u>2</u>
Average	17	19	2

Students may wish to verify the following results, especially the variance of e and the covariance of T and e :

$$S_X^2 = 62.5 \quad S_T^2 = 54.0 \quad S_e^2 = 7.5 \quad S_{T,e} = 0.5$$

As a verification of Equation (4.3) we have

$$62.5 = 54.0 + 7.5 + (2)(0.5)$$

From data in a simple random sample one would compute $s_x^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$

and use $\frac{N-n}{N} \frac{s_x^2}{n}$ as an estimate of the variance of \bar{x} . Is it clear that s_x^2 is an unbiased estimate of S_X^2 rather than of S_T^2 and that the impact of variation in e_i is included in s_x^2 ?

To summarize, response error caused a bias in \bar{x} as an estimate of \bar{T} that was equal to $\bar{x} - \bar{T}$. In addition, it was a source of variation included in the standard error of \bar{x} . To evaluate bias and variance attributable to response error, information on X_i and T_i must be available.

Illustration 4.5. In this case, we assume that the response error for a given element is not constant. That is, if an element were measured on several occasions, the observed values for the i^{th} element could vary even though the true value, T_i , remained unchanged. Let the error model be

$$X_{ij} = T_i + e'_{ij}$$

where X_{ij} is the observed value of X for the i^{th} element when the observation is taken on a particular occasion, j ,

T_i is the true value of X for the i^{th} element,

and e'_{ij} is the response error for the i^{th} element on a particular occasion, j .

Assume, for any given element, that the response error, e'_{ij} , is a random variable. We can let $e'_{ij} = \bar{e}_i + e_{ij}$, where \bar{e}_i is the average value of e_{ij} for a fixed i , that is, $\bar{e}_i = E(e'_{ij} | i)$. This divides the response error for the i^{th} element into two components: a constant component, \bar{e}_i , and a variable component, e_{ij} . By definition, the expected value of e_{ij} is zero for any given element. That is, $E(e_{ij} | i) = 0$.

Substituting $\bar{e}_i + e_{ij}$ for e'_{ij} , the model becomes

$$x_{ij} = T_i + \bar{e}_i + e_{ij} \quad (4.4)$$

The model, Equation (4.4), is now in a good form for comparison with the model in Illustration 4.4. In Equation (4.4), \bar{e}_i , like e_i in Equation (4.2) is constant for a given element. Thus, the two models are alike except for the added term, e_{ij} , in Equation (4.4) which allows for the possibility that the response error for the i^{th} element might not be constant.

Assume a simple random sample of n elements and one observation for each element. According to the model, Equation (4.4), we may now write the sample mean as follows:

$$\bar{x} = \frac{\sum t_i}{n} + \frac{\sum \bar{e}_i}{n} + \frac{\sum e_{ij}}{n}$$

Summation with respect to j is not needed as there is only one observation for each element in the sample. Under the conditions specified the expected value of \bar{x} may be expressed as follows:

$$E(\bar{x}) = \bar{T} + \bar{e}$$

where $\bar{T} = \frac{\sum T_i}{N}$ and $\bar{e} = \frac{\sum \bar{e}_i}{N}$

The variance of \bar{x} is complicated unless some further assumptions are made. Assume that all covariance terms are zero. Also, assume that the conditional variance of e_{ij} is constant for all values of i ; that is, let $V(e_{ij}|i) = S_e^2$. Then, the variance of \bar{x} is

$$S_{\bar{x}}^2 = \frac{N-n}{N} \frac{S_T^2}{n} + \frac{N-n}{N} \frac{S_e^2}{n} + \frac{S_e^2}{n}$$

where
$$s_T^2 = \frac{\sum (T_i - \bar{T})^2}{N-1}, \quad s_e^2 = \frac{\sum (\bar{e}_i - \bar{e})^2}{N-1},$$

and s_e^2 is the conditional variance of e_{ij} , that is, $V(e_{ij}|1)$. For this model the variance of \bar{x} does not diminish to zero as $n \rightarrow \infty$. However, assuming

N is large, the variance of \bar{x} , which becomes $\frac{s_e^2}{N}$, is probably negligible.

Definition 4.2. Mean-Square Error. In terms of the theory of expected values the mean-square error of an estimate, x' , is $E(x' - T)^2$ where T is the target value, that is, the value being estimated. From the theory it is easy to show that

$$E(x' - T)^2 = [E(x') - T]^2 + E[x' - E(x')]^2$$

Thus, the mean-square error, mse, can be expressed as follows:

$$\text{mse} = B^2 + \sigma_{x'}^2 \quad (4.5)$$

$$\text{where } B = E(x') - T \quad (4.6)$$

$$\text{and } \sigma_{x'}^2 = E[x' - E(x')]^2 \quad (4.7)$$

Definition 4.3. Bias. In Equation (4.5), B is the bias in x' as an estimate of T .

Definition 4.4. Precision. The precision of an estimate is the standard error of the estimate, namely, $\sigma_{x'}$ in Equation (4.7).

Precision is a measure of repeatability. Conceptually, it is a measure of the dispersion of estimates that would be generated by repetition of the same sampling and estimation procedures many times under the same conditions. With reference to the sampling distribution, it is a measure of the dispersion of the estimates from the center of the distribution and

does not include any indication of where the center of the distribution is in relation to a target.

In Illustrations 4.1, 4.2, and 4.3, the target value was implicitly assumed to be \bar{X} ; that is, T was equal to \bar{X} . Therefore, B was zero and the mean-square error of x' was the same as the variance of x' . In Illustrations 4.4 and 4.5 the picture was broadened somewhat by introducing response error and examining, theoretically, the impact of response error on $E(x')$ and $\sigma_{x'}$. In practice many factors have potential for influencing the sampling distribution of x' . That is, the data in a sample are subject to error that might be attributed to several sources.

From sample data an estimate, x' , is computed and an estimate of the variance of x' is also computed. How does one interpret the results? In Illustrations 4.4 and 4.5 we found that response error could be divided into bias and variance. The error from any source can, at least conceptually, be divided into bias and variance. An estimate from a sample is subject to the combined influence of bias and variance corresponding to each of the several sources of error. When an estimate of the variance of x' is computed from sample data, the estimate is a combination of variances that might be identified with various sources. Likewise the difference between $E(x')$ and T is a combination of biases that might be identified with various sources.

Figure 4.2 illustrates the sampling distribution of x' for four different cases: A, no bias and low standard error; B, no bias and large standard error; C, large bias and low standard error; and D, large bias and large standard error. The accuracy of an estimator is sometimes defined as the square root of the mean-square error of the estimator. According

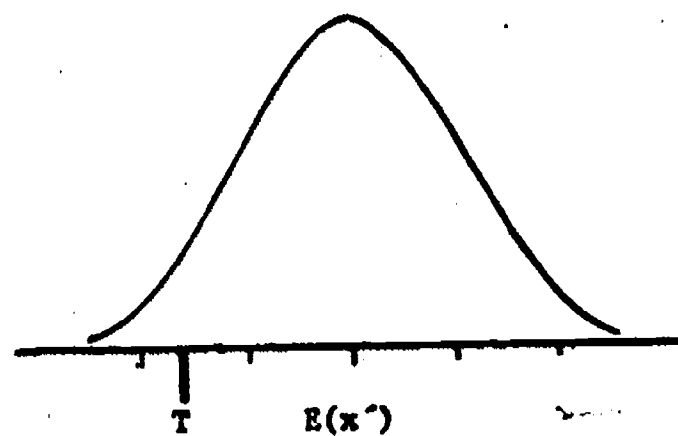
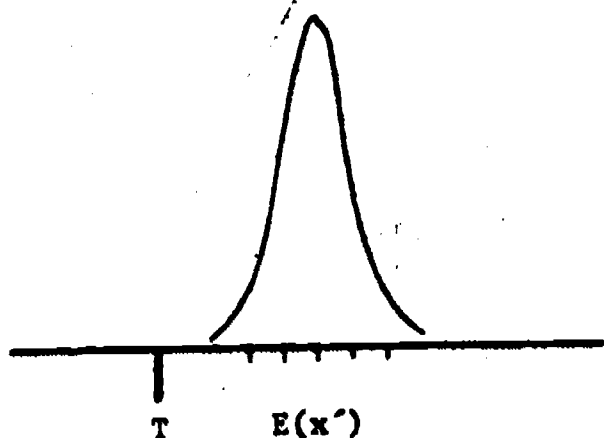
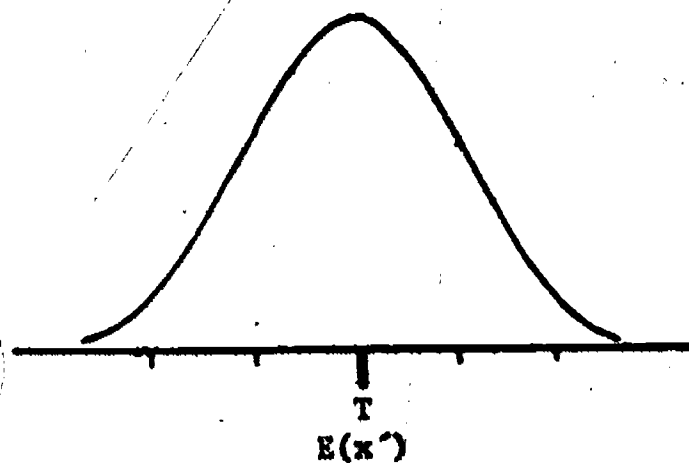
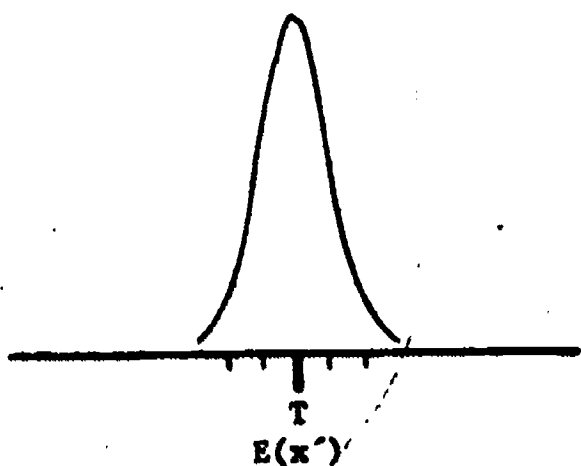


Figure 4.2--Examples of four sampling distributions

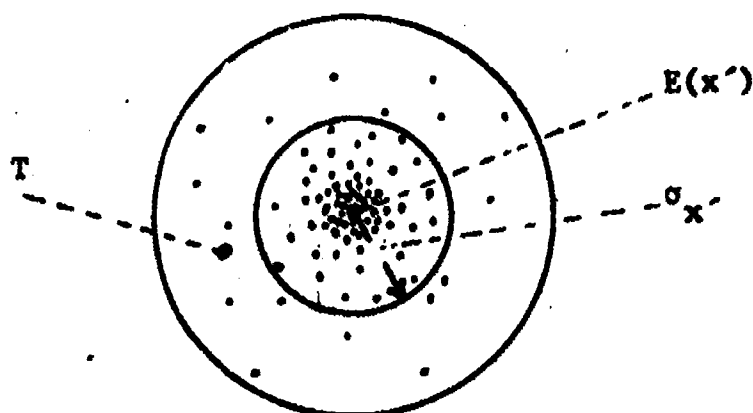


Figure 4.3--Sampling distribution--
Each small dot corresponds to an estimate

to that definition, we could describe estimators having the four sampling distributions in Figure 4.2 as follows: In case A the estimator is precise and accurate; in B the estimator lacks precision and is therefore inaccurate; in C the estimator is precise but inaccurate because of bias, and in D the estimator is inaccurate because of bias and low precision.

Unfortunately, it is generally not possible to determine, exactly, the magnitude of bias in an estimate, or of a particular component of bias. However, evidence of the magnitude of bias is often available from general experience, from knowledge of how well the survey processes were performed, and from special investigations. The author accepts a point of view that the mean-square error is an appropriate concept of accuracy to follow. In that context, the concern becomes a matter of the magnitude of the mse and the size of B relative to σ_x . That viewpoint is important because it is not possible to be certain that B is zero. Our goal should be to prepare survey specifications and to conduct survey operations so B is small in relation to σ_x . Or, one might say we want the mse to be minimum for a given cost of doing the survey. Ways of getting evidence on the magnitude of bias is a major subject and is outside the scope of this publication.

As indicated in the previous paragraph, it is important to know something about the magnitude of the bias, B , relative to the standard error, σ_x . The standard error is controlled primarily by the design of a sample and its size. For many survey populations, as the size of the sample increases, the standard error becomes small relative to the bias. In fact, the bias might be larger than the standard error even for samples of moderate size, for example a few hundred cases, depending upon the circumstances. The point is that if the mean-square error is to be small, both

B and σ_x must be small. The approaches for reducing B are very different from the approaches for reducing σ_x . The greater concern about non-sampling error is bias rather than impact on variance. In the design and selection of samples and in the processes of doing the survey an effort is made to prevent biases that are "sampling" in origin. However, in survey work one must be constantly aware of potential biases and on the alert to minimize biases as well as random error (that is, σ_x).

The above discussion puts a census in the same light as a sample. Results from both have a mean-square error. Both are surveys with reference to use of results. Uncertain inferences are involved in the use of results from a census as well as from a sample. The only difference is that in a census one attempts to get a measurement for all N elements, but making $n = N$ does not reduce the mse to zero. Indeed, as the sample size increases, there is no positive assurance that the mse will always decrease; because, as the variance component of the mse decreases, the bias component might increase. This can occur especially when the population is large and items on the questionnaire are such that simple, accurate answers are difficult to obtain. For a large sample or a census, compared to a small sample, it might be more difficult to control factors that cause bias. Thus, it is possible for a census to be less accurate (have a larger mse) than a sample wherein the sources of error are more adequately controlled. Much depends upon the kind of information being collected.

4.5 BIAS AND STANDARD ERROR

The words "bias," "biased," and "unbiased" have a wide variety of meaning among various individuals. As a result, much confusion exists,

especially since the terms are often used loosely. Technically, it seems logical to define the bias in an estimate as being equal to B in Equation (4.6), which is the difference between the expected value of an estimate and the target value. But, except for hypothetical cases, numerical values do not exist for either $E(x')$ or the target T . Hence, defining an unbiased estimate as one where $B = E(x') - T = 0$ is of little, if any, practical value unless one is willing to accept the target as being equal to $E(x')$. From a sampling point of view there are conditions that give a rational basis for accepting $E(x')$ as the target. However, regardless of how the target is defined, a good practical interpretation of $E(x')$ is needed.

It has become common practice among survey statisticians to call an estimate unbiased when it is based on methods of sampling and estimation that are "unbiased." For example, in Illustration 4.4, \bar{x} would be referred to as an unbiased estimate--unbiased because the method of sampling and estimation was unbiased. In other words, since \bar{x} was an unbiased estimate of \bar{X} , \bar{x} could be interpreted as an unbiased estimate of the result that would have been obtained if all elements in the population had been measured.

In Illustration 4.5 the expected value of \bar{x} is more difficult to describe. Nevertheless, with reference to the method of sampling and estimation, \bar{x} was "unbiased" and could be called an unbiased estimate even though $E(\bar{x})$ is not equal to \bar{T} .

The point is that a simple statement which says, "the estimate is unbiased" is incomplete and can be very misleading, especially if one is not familiar with the context and concepts of bias. Calling an estimate unbiased is equivalent to saying the estimate is an unbiased estimate of

its expected value. Regardless of how "bias" is defined or used, $E(x')$ is the mean of the sampling distribution of x ; and this concept of $E(x')$ is very important because $E(x')$ appears in the standard error, $\sigma_{x'}$, of x' as well as in B. See Equations (4.6) and (4.7).

As a simple concept or picture of the error of an estimate from a survey, the writer likes the analogy between an estimate and a shot at a target with a gun or an arrow. Think of a survey being replicated many times using the same sampling plan, but a different sample for each replication. Each replication would provide an estimate that corresponds to a shot at a target.

In Figure 4.3, each dot corresponds to an estimate from one of the replicated samples. The center of the cluster of dots is labeled $E(x')$ because it corresponds to the expected value of an estimate. Around the point $E(x')$ a circle is drawn which contains two-thirds of the points. The radius of this circle corresponds to $\sigma_{x'}$, the standard error of the estimate. The outer circle has a radius of two standard errors and contains 95 percent of the points. The target is labeled T. The distance between T and $E(x')$ is bias, which in the figure is greater than the standard error.

In practice, we usually have only one estimate, x' , and an estimate, $s_{x'}$, of the standard error of x' . With reference to Figure 4.3, this means one point and an estimate of the radius of the circle around $E(x')$ that would contain two-thirds of the estimates in repeated samplings. We do not know the value of $E(x')$; that is, we do not know where the center of the circles is. However, when we make a statement about the standard error of x' , we are expressing a degree of confidence about how close a

particular estimate prepared from a survey is to $E(x')$; that is, how close one of the points in Figure 4.3 probably is to the unknown point $E(x')$. A judgment as to how far $E(x')$ is from T is a matter of how T is defined and assessment of the magnitude of biases associated with various sources of error.

Unfortunately, it is not easy to make a short, rigorous, and complete interpretative statement about the standard error of x' . If the estimated standard error of x' is three percent, one could simply state that fact and not make an interpretation. It does not help much to say, for example, that the odds are about two out of three that the estimate is within three percent of its expected value, because a person familiar with the concepts already understands that and it probably does not help the person who is unfamiliar with the concepts. Suppose one states, "the standard error of x' means the odds are two out of three that the estimate is within three percent of the value that would have been obtained from a census taken under identically the same conditions." That is a good type of statement to make but, when one engages considerations of the finer points, interpretation of "a census taken under identically the same conditions" is needed--especially since it is not possible to take a census under identically the same conditions.

In summary, think of a survey as a fully defined system or process including all details that could affect an estimate, including: the method of sampling; the method of estimation; the wording of questions; the order of the questions on the questionnaire; interviewing procedures; selection, training, and supervision of interviewers; and editing and processing of

data. Conceptually, the sampling is then replicated many times, holding all specifications and conditions constant. This would generate a sampling distribution as illustrated in Figures 4.2 or 4.3. We need to recognize that a change in any of the survey specifications or conditions, regardless of how trivial the change might seem, has a potential for changing the sampling distribution, especially the expected value of \bar{x} . Changes in survey plans, even though the definition of the parameters being estimated remains unchanged, often result in discrepancies that are larger than the random error that can be attributed to sampling.

The points discussed in the latter part of this chapter were included to emphasize that much more than a well designed sample is required to assure accurate results. Good survey planning and management calls for evaluation of errors from all sources and for trying to balance the effort to control error from various sources so the mean-square error will be within acceptable limits as economically as possible.