

yes to (3) two-sample problem? no to (4) underlying distribution normal or can central-limit theorem be assumed to hold? and yes to (5) underlying distribution binomial?

We now refer to the flowchart at the end of this chapter (p. 409). We answer yes to (1) are samples independent? (2) are all expected values  $\geq 5$ ? and (3)  $2 \times 2$  contingency table? This leads us to the box labeled “Use the two-sample test for binomial proportions or  $2 \times 2$  contingency-table methods if no confounding is present, or Mantel-Haenszel test if confounding is present.” In brief, a confounder is another variable that is potentially related to both the row and column classification variables, and it must be controlled for. We discuss methods for controlling for confounding in Chapter 13. In this chapter, we assume no confounding is present. Thus we use either the two-sample test for binomial proportions (Equation 10.3) or the equivalent chi-square test for  $2 \times 2$  contingency tables (Equation 10.5).

### REVIEW QUESTIONS 10A

- 1 What is a contingency table?
- 2 Suppose we have 50 ovarian-cancer cases and 100 controls, all of whom are age 50–54. Ten of the ovarian-cancer cases and 12 of the controls reached menarche (age when periods begin) at  $<11$  years.
  - (a) What test can be used to assess whether there is a significant association between early age at menarche and ovarian cancer?
  - (b) Perform the test in Review Question 10A.2a, and report a two-tailed  $p$ -value.

## 10.3 Fisher's Exact Test

In Section 10.2, we discussed methods for comparing two binomial proportions using either normal-theory or contingency-table methods. Both methods yield identical  $p$ -values. However, they require that the normal approximation to the binomial distribution be valid, which is not always the case, especially for small samples.

**Example 10.16** **Cardiovascular Disease, Nutrition** Suppose we want to investigate the relationship between high salt intake and death from cardiovascular disease (CVD). Groups of high- and low-salt users could be identified and followed over a long time to compare relative frequency of death from CVD in the two groups. In contrast, a much less expensive study would involve looking at death records, separating CVD deaths from non-CVD deaths, asking a close relative (such as a spouse) about the dietary habits of the deceased, and then comparing salt intake between people who died of CVD vs. people who died of other causes.

The latter type of study, a retrospective study, may be impossible to perform for a number of reasons. But if it is possible, it is almost always less expensive than the former type, a prospective study.

**Example 10.17** **Cardiovascular Disease, Nutrition** Suppose a retrospective study is done among men ages 50–54 in a specific county who died over a 1-month period. The investigators try to include approximately an equal number of men who died from CVD (the cases) and men who died from other causes (the controls). Of 35 people who died from CVD, 5 were on a high-salt diet before they died, whereas of 25 people who died from other causes 2 were on such a diet. These

data, presented in Table 10.9, are in the form of a  $2 \times 2$  contingency table, so the methods of Section 10.2 may be applicable.

**Table 10.9** Data concerning the possible association between cause of death and high salt intake

Cause of death	Type of diet		Total
	High salt	Low salt	
Non-CVD	2	23	25
CVD	5	30	35
Total	7	53	60

However, the expected values of this table are too small for such methods to be valid. Indeed,

$$E_{11} = 7(25)/60 = 2.92$$
$$E_{12} = 7(35)/60 = 4.08$$

thus two of the four cells have expected values less than 5. How should the possible association between cause of death and type of diet be assessed?

In this case, **Fisher’s exact test** can be used. This procedure gives exact levels of significance for any  $2 \times 2$  table but is only necessary for tables with small expected values, tables in which the standard chi-square test as given in Equation 10.5 is not applicable. For tables in which use of the chi-square test is appropriate, the two tests give very similar results. Suppose the probability that a man was on a high-salt diet given that his cause of death was noncardiovascular (non-CVD) =  $p_1$  and the probability that a man was on a high-salt diet given that his cause of death was cardiovascular (CVD) =  $p_2$ . We wish to test the hypothesis  $H_0: p_1 = p_2 = p$  vs.  $H_1: p_1 \neq p_2$ . Table 10.10 gives the general layout of the data.

**Table 10.10** General layout of data for Fisher’s exact test example

Cause of death	Type of diet		Total
	High salt	Low salt	
Non-CVD	a	b	a + b
CVD	c	d	c + d
Total	a + c	b + d	n

For mathematical convenience, we assume the margins of this table are *fixed*; that is, the numbers of non-CVD deaths and CVD deaths are fixed at  $a + b$  and  $c + d$ , respectively, and the numbers of people on high- and low-salt diets are fixed at  $a + c$  and  $b + d$ , respectively. Indeed, it is difficult to compute exact probabilities unless one assumes fixed margins. The *exact* probability of observing the table with cells  $a$ ,  $b$ ,  $c$ ,  $d$  is as follows.

**Equation 10.7****Exact Probability of Observing a Table with Cells  $a, b, c, d$** 

$$Pr(a, b, c, d) = \frac{(a+b)!(c+d)!(a+c)!(b+d)!}{n!a!b!c!d!}$$

The formula in Equation 10.7 is easy to remember because the numerator is the product of the factorials of each of the row and column margins, and the denominator is the product of the factorial of the grand total and the factorials of the individual cells.

**Example 10.18**

Suppose we have the  $2 \times 2$  table shown in Table 10.11. Compute the exact probability of obtaining this table assuming the margins are fixed.

**Solution**

$$Pr(2, 5, 3, 1) = \frac{7!4!5!6!}{11!2!5!3!1!} = \frac{5040(24)(120)(720)}{39,916,800(2)(120)(6)} = \frac{1.0450944 \times 10^{10}}{5.7480192 \times 10^{10}} = .182$$

**Table 10.11****Hypothetical  $2 \times 2$  contingency table in Example 10.18**

2	5	7
3	1	4
5	6	11

## The Hypergeometric Distribution

Suppose we consider all possible tables with fixed row margins denoted by  $N_1$  and  $N_2$  and fixed column margins denoted by  $M_1$  and  $M_2$ . We assume the rows and columns have been rearranged so that  $M_1 \leq M_2$  and  $N_1 \leq N_2$ . We refer to each table by its (1, 1) cell because all other cells are then determined from the fixed row and column margins. Let the random variable  $X$  denote the cell count in the (1, 1) cell. The probability distribution of  $X$  is given by

**Equation 10.8**

$$Pr(X = a) = \frac{N_1!N_2!M_1!M_2!}{N!a!(N_1 - a)!(M_1 - a)!(M_2 - N_1 + a)!}, a = 0, \dots, \min(M_1, N_1)$$

and  $N = N_1 + N_2 = M_1 + M_2$ . This probability distribution is called the **hypergeometric distribution**.

It will be useful for our subsequent work on combining evidence from more than one  $2 \times 2$  table in Chapter 13 to refer to the expected value and variance of the hypergeometric distribution. These are as follows.

**Equation 10.9****Expected Value and Variance of the Hypergeometric Distribution**

Suppose we consider all possible tables with fixed row margins  $N_1, N_2$  and fixed column margins  $M_1, M_2$ , where  $N_1 \leq N_2, M_1 \leq M_2$ , and  $N = N_1 + N_2 = M_1 + M_2$ . Let the random variable  $X$  denote the cell count in the (1, 1) cell. The expected value and variance of  $X$  are

$$E(X) = \frac{M_1 N_1}{N}$$

$$Var(X) = \frac{M_1 M_2 N_1 N_2}{N^2 (N - 1)}$$

Thus the exact probability of obtaining a table with cells  $a, b, c, d$  in Equation 10.7 is a special case of the hypergeometric distribution, where  $N_1 = a + b$ ,  $N_2 = c + d$ ,  $M_1 = a + c$ ,  $M_2 = b + d$ , and  $N = a + b + c + d$ . We can evaluate this probability by calculator using Equation 10.7, or we can use the HYPGEOMDIST function of Excel. In the latter case, to evaluate  $Pr(a, b, c, d)$ , we specify HYPGEOMDIST( $a, a + b, a + c, N$ ). In words, the hypergeometric distribution evaluates the probability of obtaining  $a$  successes out of a sample of  $a + b$  observations, given that the total population (in this case, the two samples combined), is of size  $N$ , of which  $a + c$  observations are successes. Thus, to evaluate the exact probability in Table 10.11, we specify HYPGEOMDIST( $2, 7, 5, 11$ ) = .182, which is the probability of obtaining two successes in a sample of 7 observations given that the total population consists of 11 observations, of which 5 are successes. The hypergeometric distribution differs from the binomial distribution, because in the latter case, we simply evaluate the probability of obtaining  $a$  successes out of  $a + b$  observations, assuming that each outcome is independent. For the hypergeometric distribution, the outcomes are not independent because once a success occurs it is less likely that another observation will be a success, as the total number of successes is fixed (at  $a + c$ ). If  $N$  is large, the two distributions are very similar because there is only a slight deviation from independence for the hypergeometric.

The basic strategy in testing the hypothesis  $H_0: p_1 = p_2$  vs.  $H_1: p_1 \neq p_2$  will be to enumerate all possible tables with the same margins as the observed table and to compute the exact probability for each such table based on the hypergeometric distribution. A method for accomplishing this is as follows.

#### Equation 10.10

#### Enumeration of All Possible Tables with the Same Margins as the Observed Table

- (1) Rearrange the rows and columns of the observed table so the smaller row total is in the first row and the smaller column total is in the first column.

Suppose that after the rearrangement, the cells in the observed table are  $a, b, c, d$ , as shown in Table 10.10.

- (2) Start with the table with 0 in the (1, 1) cell. The other cells in this table are then determined from the row and column margins. Indeed, to maintain the same row and column margins as the observed table, the (1, 2) element must be  $a + b$ , the (2, 1) cell must be  $a + c$ , and the (2, 2) element must be  $(c + d) - (a + c) = d - a$ .
- (3) Construct the next table by increasing the (1, 1) cell by 1 (i.e., from 0 to 1), decreasing the (1, 2) and (2, 1) cells by 1, and increasing the (2, 2) cell by 1.
- (4) Continue increasing and decreasing the cells by 1, as in step 3, until one of the cells is 0, at which point all possible tables with the given row and column margins have been enumerated. Each table in the sequence of tables is referred to by its (1, 1) element. Thus, the first table is the “0” table, the next table is the “1” table, and so on.

#### Example 10.19

**Cardiovascular Disease, Nutrition** Enumerate all possible tables with the same row and column margins as the observed data in Table 10.9.

**Solution**

The observed table has  $a = 2$ ,  $b = 23$ ,  $c = 5$ ,  $d = 30$ . The rows or columns do not need to be rearranged because the first row total is smaller than the second row total, and the first column total is smaller than the second column total. Start with the 0 table, which has 0 in the (1, 1) cell, 25 in the (1, 2) cell, 7 in the (2, 1) cell, and  $30 - 2$ , or 28, in the (2, 2) cell. The 1 table then has 1 in the (1, 1) cell,  $25 - 1 = 24$  in the (1, 2) cell,  $7 - 1 = 6$  in the (2, 1) cell, and  $28 + 1 = 29$  in the (2, 2) cell. Continue in this fashion until the 7 table is reached, which has 0 in the (2, 1) cell, at which point all possible tables with the given row and column margins have been enumerated. The set of hypergeometric probabilities in Table 10.12 can be easily evaluated using the recursive properties of Excel by (1) setting up a column with consecutive values from 0 to 7 (say from B1 to B8), (2) using the function HYPGEOMDIST to compute  $Pr(0) = \text{HYPGEOMDIST}(B1, 25, 7, 60)$  and placing it in C1, and then (3) dragging the cursor down column C to compute the remaining hypergeometric probabilities. See the Companion Website for more details on the use of the HYPGEOMDIST function. The collection of tables and their associated probabilities based on the hypergeometric distribution in Equation 10.8 are given in Table 10.12.

**Table 10.12** Enumeration of all possible tables with fixed margins and their associated probabilities, based on the hypergeometric distribution for Example 10.19

0	25	1	24	2	23	3	22
7	28	6	29	5	30	4	31
.017		.105		.252		.312	
4	21	5	20	6	19	7	18
3	32	2	33	1	34	0	35
.214		.082		.016		.001	

The question now is: What should be done with these probabilities to evaluate the significance of the results? The answer depends on whether a one-sided or a two-sided alternative is being used. In general, the following method can be used.

**Equation 10.11****Fisher's Exact Test: General Procedure and Computation of  $p$ -Value**

To test the hypothesis  $H_0: p_1 = p_2$  vs.  $H_1: p_1 \neq p_2$ , where the expected value of at least one cell is  $< 5$  when the data are analyzed in the form of a  $2 \times 2$  contingency table, use the following procedure:

- (1) Enumerate all possible tables with the same row and column margins as the observed table, as shown in Equation 10.10.
- (2) Compute the exact probability of each table enumerated in step 1, using either the computer or the formula in Equation 10.7.
- (3) Suppose the observed table is the  $a$  table and the last table enumerated is the  $k$  table.
  - (a) To test the hypothesis  $H_0: p_1 = p_2$  vs.  $H_1: p_1 \neq p_2$ , the  $p$ -value =  $2 \times \min[Pr(0) + Pr(1) + \dots + Pr(a), Pr(a) + Pr(a+1) + \dots + Pr(k), .5]$ .
  - (b) To test the hypothesis  $H_0: p_1 = p_2$  vs.  $H_1: p_1 < p_2$ , the  $p$ -value =  $Pr(0) + Pr(1) + \dots + Pr(a)$ .

(c) To test the hypothesis  $H_0: p_1 = p_2$  vs.  $H_1: p_1 > p_2$ , the  $p$ -value =  $Pr(a) + Pr(a+1) + \dots + Pr(k)$ .

For each of these three alternative hypotheses, the  $p$ -value can be interpreted as the probability of obtaining a table as extreme as or more extreme than the observed table.

### Example 10.20

**Cardiovascular Disease, Nutrition** Evaluate the statistical significance of the data in Example 10.17 using a two-sided alternative.

### Solution

We want to test the hypothesis  $H_0: p_1 = p_2$  vs.  $H_1: p_1 \neq p_2$ . Our table is the 2 table whose probability is .252 in Table 10.12. Thus, to compute the  $p$ -value, the smaller of the tail probabilities corresponding to the 2 table is computed and doubled. This strategy corresponds to the procedures for the various normal-theory tests studied in Chapters 7 and 8. First compute the left-hand tail area,

$$Pr(0) + Pr(1) + Pr(2) = .017 + .105 + .252 = .375$$

and the right-hand tail area,

$$Pr(2) + Pr(3) + \dots + Pr(7) = .252 + .312 + .214 + .082 + .016 + .001 = .878$$

Then  $p = 2 \times \min(.375, .878, .5) = 2(.375) = .749$

If a one-sided alternative of the form  $H_0: p_1 = p_2$  vs.  $H_1: p_1 < p_2$  is used, then the  $p$ -value equals

$$Pr(0) + Pr(1) + Pr(2) = .017 + .105 + .252 = .375$$

Thus the two proportions in this example are *not* significantly different with either a one-sided or two-sided test, and we *cannot* say, on the basis of this limited amount of data, that there is a significant association between salt intake and cause of death.

In most instances, computer programs are used to implement Fisher's exact test using statistical packages such as SAS. There are other possible approaches to significance testing in the two-sided case. For example, the approach used by SAS is to compute

$$p\text{-value (two-tailed)} = \sum_{\{i: Pr(i) \leq Pr(a)\}} Pr(i)$$

In other words, the two-tailed  $p$ -value using SAS is the sum of the probabilities of all tables whose probabilities are  $\leq$  the probability of the observed table. Using this approach, the two-tailed  $p$ -value would be

$$\begin{aligned} p\text{-value (two-tailed)} &= Pr(0) + Pr(1) + Pr(2) + Pr(4) + Pr(5) + Pr(6) + Pr(7) \\ &= .017 + .105 + .252 + .214 + .082 + .016 + .001 = .688 \end{aligned}$$

In this section, we learned about Fisher's exact test, which is used for comparing binomial proportions from two independent samples in  $2 \times 2$  tables with small expected counts ( $<5$ ). This is the two-sample analog to the exact one-sample binomial test given in Equation 7.44. If we refer to the flowchart at the end of this chapter (Figure 10.16, p. 409), we answer yes to (1) are samples independent? and no to (2) are all expected values  $\geq 5$ ? This leads us to the box labeled "Use Fisher's exact test."

disease–exposure relationships in a hypothesis-testing framework using the Mantel-Haenszel test. Finally, standardization can be based on stratification by factors other than age. For example, standardization by both age and sex is common. Similar methods can be used to obtain age–sex standardized risks and standardized *RRs* as given in Definition 13.15.

In this section, we have introduced the concept of a confounding variable (*C*), a variable related to both the disease (*D*) and exposure (*E*) variables. Furthermore, we classified confounding variables as positive confounders if the associations between *C* and *D* and *C* and *E*, respectively, are in the same direction and as negative confounders if the associations between *C* and *D* and *C* and *E* are in opposite directions. We also discussed when it is or is not appropriate to control for a confounder, according to whether *C* is or is not in the causal pathway between *E* and *D*. Finally, because age is often an important confounding variable, it is reasonable to consider descriptive measures of proportions and relative risk that control for age. Age-standardized proportions and *RRs* are such measures.

REVIEW QUESTIONS 13B

- 1 Suppose we are interested in the association between smoking and bone density in women.
- (a) If body-mass index (BMI) is inversely associated with smoking and positively associated with bone density, then is BMI a positive confounder, a negative confounder, or neither?

(b) If alcohol intake is positively associated with smoking and is unrelated to bone density, then is alcohol intake a positive confounder, a negative confounder, or neither?
- 2 Suppose the age-specific risks of hypertension in adults with left ventricular hypertrophy (LVH) and controls are as shown in Table 13.10.

Table 13.10 Age-specific hypertension risks among patients with LVH and controls

	LVH		Control	
	Risk	<i>N</i>	Risk	<i>N</i>
40–49	.16	20	.14	50
50–59	.20	40	.18	40
60–69	.28	30	.20	36
70–79	.36	35	.25	29

- (a) Calculate the age-standardized risk of hypertension in each group using the total population in the two groups as the standard.
- (b) Calculate the age-standardized *RR* of hypertension for LVH patients vs. controls.

13.6 Methods of Inference for Stratified Categorical Data—The Mantel-Haenszel Test

Example 13.23

Cancer A 1985 study identified a group of 518 cancer cases ages 15–59 and a group of 518 age- and sex-matched controls by mail questionnaire [4]. The main purpose of the study was to look at the effect of passive smoking on cancer risk. The study

defined passive smoking as exposure to the cigarette smoke of a spouse who smoked at least one cigarette per day for at least 6 months. One potential confounding variable was smoking by the participants themselves (i.e., personal smoking) because personal smoking is related to both cancer risk and spouse smoking. Therefore, it was important to control for personal smoking before looking at the relationship between passive smoking and cancer risk.

To display the data, a  $2 \times 2$  table relating case-control status to passive smoking can be constructed for both nonsmokers and smokers. The data are given in Table 13.11 for nonsmokers and Table 13.12 for smokers.

**Table 13.11 Relationship of passive smoking to cancer risk among nonsmokers**

Case-control status	Passive smoker		Total
	Yes	No	
Case	120	111	231
Control	80	155	235
Total	200	266	466

Source: From Sandler et al., "Passive Smoking in Adulthood and Cancer Risk," *American Journal of Epidemiology*, 1985 121: 37–48. Reprinted by permission of Oxford University Press.

**Table 13.12 Relationship of passive smoking to cancer risk among smokers**

Case-control status	Passive smoker		Total
	Yes	No	
Case	161	117	278
Control	130	124	254
Total	291	241	532

Source: From Sandler et al., "Passive Smoking in Adulthood and Cancer Risk," *American Journal of Epidemiology*, 1985 121: 37–48. Reprinted by permission of Oxford University Press.

The passive-smoking effect can be assessed separately for nonsmokers and smokers. Indeed, we notice from Tables 13.11 and 13.12 that the *OR* in favor of a case being exposed to cigarette smoke from a spouse who smokes vs. a control is  $(120 \times 155)/(80 \times 111) = 2.1$  for nonsmokers, whereas the corresponding *OR* for smokers is  $(161 \times 124)/(130 \times 117) = 1.3$ . Thus for both subgroups the trend is in the direction of more passive smoking among cases than among controls. The key question is how to combine the results from the two tables to obtain an overall estimated *OR* and test of significance for the passive-smoking effect.

In general, the data are stratified into  $k$  subgroups according to one or more confounding variables to make the units within a stratum as homogeneous as possible. The data for each stratum consist of a  $2 \times 2$  contingency table relating exposure to disease, as shown in Table 13.13 for the  $i$ th stratum.

**Table 13.13 Relationship of disease to exposure in the  $i$ th stratum**

Disease		Exposure		Total
		Yes	No	
	Yes	$a_i$	$b_i$	$a_i + b_i$
Disease	No	$c_i$	$d_i$	$c_i + d_i$
		$a_i + c_i$	$b_i + d_i$	$n_i$

Based on our work on Fisher's exact test, the distribution of  $a_i$  follows a **hypergeometric distribution**. The test procedure is based on a comparison of the observed number of units in the (1, 1) cell of each stratum (denoted by  $O_i = a_i$ ) with the



expected number of units in that cell (denoted by  $E_i$ ). The test procedure is the same regardless of order of the rows and columns; that is, which row (or column) is designated as the first row (or column) is arbitrary. Based on the hypergeometric distribution (Equation 10.9), the expected number of units in the (1, 1) cell of the  $i$ th stratum is given by

Equation 13.14

$$E_i = \frac{(a_i + b_i)(a_i + c_i)}{n_i}$$

The observed and expected numbers of units in the (1, 1) cell are then summed over all strata, yielding  $O = \sum_{i=1}^k O_i$ ,  $E = \sum_{i=1}^k E_i$ , and the test is based on  $O - E$ . Based on the hypergeometric distribution (Equation 10.9), the variance of  $O_i$  is given by

Equation 13.15

$$V_i = \frac{(a_i + b_i)(c_i + d_i)(a_i + c_i)(b_i + d_i)}{n_i^2(n_i - 1)}$$

Furthermore, the variance of  $O$  is denoted by  $V = \sum_{i=1}^k V_i$ . The test statistic is given by  $X_{MH}^2 = (|O - E| - .5)^2 / V$ , which should follow a chi-square distribution with 1 degree of freedom ( $df$ ) under the null hypothesis of no association between disease and exposure.  $H_0$  is rejected if  $X_{MH}^2$  is large. The abbreviation  $MH$  refers to Mantel-Haenszel; this procedure is known as the Mantel-Haenszel test and is summarized as follows.

Equation 13.16

#### Mantel-Haenszel Test

To assess the association between a dichotomous disease and a dichotomous exposure variable after controlling for one or more confounding variables, use the following procedure:

- (1) Form  $k$  strata, based on the level of the confounding variable(s), and construct a  $2 \times 2$  table relating disease and exposure within each stratum, as shown in Table 13.13.
- (2) Compute the total observed number of units ( $O$ ) in the (1, 1) cell over all strata, where

$$O = \sum_{i=1}^k O_i = \sum_{i=1}^k a_i$$

- (3) Compute the total expected number of units ( $E$ ) in the (1, 1) cell over all strata, where

$$E = \sum_{i=1}^k E_i = \sum_{i=1}^k \frac{(a_i + b_i)(a_i + c_i)}{n_i}$$

- (4) Compute the variance ( $V$ ) of  $O$  under  $H_0$ , where

$$V = \sum_{i=1}^k V_i = \sum_{i=1}^k \frac{(a_i + b_i)(c_i + d_i)(a_i + c_i)(b_i + d_i)}{n_i^2(n_i - 1)}$$

- (5) The test statistic is then given by

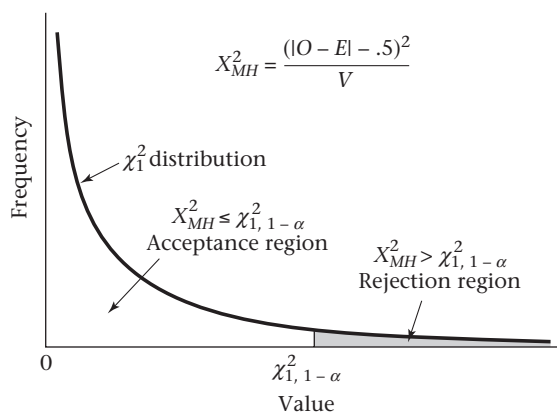
$$X_{MH}^2 = \frac{(|O - E| - .5)^2}{V}$$

which under  $H_0$  follows a chi-square distribution with 1  $df$ .

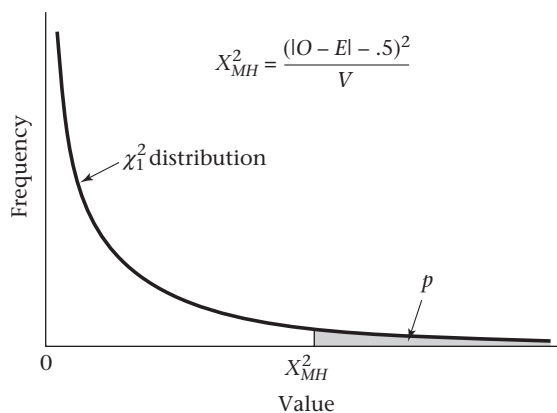
- (6) For a two-sided test with significance level  $\alpha$ ,  
 if  $X_{MH}^2 > \chi_{1,1-\alpha}^2$  then reject  $H_0$ .  
 if  $X_{MH}^2 \leq \chi_{1,1-\alpha}^2$  then accept  $H_0$ .
- (7) The exact  $p$ -value for this test is given by
- $$p = \Pr(\chi_1^2 > X_{MH}^2)$$
- (8) Use this test only if the variance  $V \geq 5$ .
- (9) Which row or column is designated as first is arbitrary. The test statistic  $X_{MH}^2$  and the assessment of significance are the same regardless of the order of the rows and columns.

The acceptance and rejection regions for the Mantel-Haenszel test are shown in Figure 13.1. The computation of the  $p$ -value for the Mantel-Haenszel test is illustrated in Figure 13.2.

**Figure 13.1** Acceptance and rejection regions for the Mantel-Haenszel test



**Figure 13.2** Computation of the  $p$ -value for the Mantel-Haenszel test



**Example 13.24** **Cancer** Assess the relationship between passive smoking and cancer risk using the data stratified by personal smoking status in Tables 13.11 and 13.12.

**Solution**

Denote the nonsmokers as stratum 1 and the smokers as stratum 2.

$O_1$  = observed number of nonsmoking cases who are passive smokers = 120

$O_2$  = observed number of smoking cases who are passive smokers = 161

Furthermore,

$$E_1 = \frac{231 \times 200}{466} = 99.1$$

$$E_2 = \frac{278 \times 291}{532} = 152.1$$

Thus the total observed and expected numbers of cases who are passive smokers are, respectively,

$$O = O_1 + O_2 = 120 + 161 = 281$$

$$E = E_1 + E_2 = 99.1 + 152.1 = 251.2$$

Therefore, more cases are passive smokers than would be expected based on their personal smoking habits. Now compute the variance to assess whether this difference is statistically significant.

$$V_1 = \frac{231 \times 235 \times 200 \times 266}{466^2 \times 465} = 28.60$$

$$V_2 = \frac{278 \times 254 \times 291 \times 241}{532^2 \times 531} = 32.95$$

$$\text{Therefore } V = V_1 + V_2 = 28.60 + 32.95 = 61.55$$

Thus the test statistic  $X_{MH}^2$  is given by

$$X_{MH}^2 = \frac{(|281 - 251.2| - .5)^2}{61.55} = \frac{858.17}{61.55} = 13.94 \sim \chi_1^2 \text{ under } H_0$$

Because  $\chi_{1,.999}^2 = 10.83 < 13.94 = X_{MH}^2$ , it follows that  $p < .001$ . Thus there is a highly significant positive association between case-control status and passive-smoking exposure, even after controlling for personal cigarette-smoking habit.

## Estimation of the Odds Ratio for Stratified Data

The Mantel-Haenszel method tests significance of the relationship between disease and exposure. However, it does not measure the strength of the association. Ideally, we would like a measure similar to the *OR* presented for a single  $2 \times 2$  contingency table in Definition 13.6. Assuming that the underlying *OR* is the same for each stratum, an estimate of the common underlying *OR* is provided by the Mantel-Haenszel estimator as follows.

**Equation 13.17**

**Mantel-Haenszel Estimator of the Common Odds Ratio for Stratified Data**

In a collection of  $k$   $2 \times 2$  contingency tables, where the table corresponding to the  $i$ th stratum is denoted as in Table 13.13, the Mantel-Haenszel estimator of the common *OR* is given by