Testing Independence

Dipankar Bandyopadhyay

Department of Biostatistics, Virginia Commonwealth University

BIOS 625: Categorical Data & GLM

イロト イポト イヨト イヨト 三日

Testing Independence

- Previously, we looked at RR = OR = 1 to determine independence.
- Now, lets revisit the Pearson and Likelihood Ratio Chi-Squared tests.
- Pearson's Chi-Square

$$X^{2} = \sum_{i=1}^{2} \sum_{j=1}^{2} \frac{(O_{ij} - E_{ij})^{2}}{E_{ij}}$$

Likelihood Ratio Test

$$G^2 = \sum_{i=1}^2 \sum_{j=1}^2 O_{ij} \log(\frac{O_{ij}}{E_{ij}})$$

Since both X² and G² are distributed as approximately χ², in order to draw inference about the significance of both, we need the degrees of freedom.

- A way to think about degrees of freedom is to relate it to the number of "pieces" of information you need to complete a table.
- More specifically, Degrees of Freedom (df) equals

df = Number of cells - Number of Constraints - Number of Parameters Estimated

- First, lets consider Pearson's Chi-Square
- We will derive *df* for the Cross Sectional Design using this definition.

- For the general $I \times J$ contingency table, there are a total of IJ cells.
- Under the Multinomial sampling design, the only constraint is that $\sum p_{ij} = 1$ so there is only one constraint.
- Under the hypothesis on interest, we are interested in estimating the marginal probabilities.
 - Since the sample size is fixed, we only need to estimate *I* 1 marginal row probabilities.
 - Namely $p_1, p_2, ..., p_{(I-1)}$.
 - Likewise, we only need to estimate J-1 column marginals.
- Thus,

$$df = IJ - \text{Number of Constraints} - \text{Number of Parameters Estimated} df = IJ - 1 - ((I - 1) + (J - 1)) = IJ - I - J + 1 = (I - 1)(J - 1)$$

- Again, there are IJ cells in our $I \times J$ contingency table
- For the Prospective design, we have constraints that each rows probability sums to 1, so there are *I* constraints.
- Although we did not state it directly before, the hypothesis of interest is the "Homogeneity" hypothesis. That is, that H₀ = p_{ij} = p_{.j} for j = 1, 2, ... J. Therefore, there are J 1 estimated marginal probabilities.
- Then the DF equals,

$$df = IJ - I - (J - 1) = IJ - I - J + 1 = (I - 1)(J - 1)$$

- For the remaining study design (Case-Control), the degrees of freedom can be shown to be (I 1)(J 1).
- Therefore, regardless of the sample design, the df for any $I \times J$ contingency table using Pearson's Chi-Square is (I-1)(J-1).
- For the 2×2 tables we have been studying,

$$df = (2-1) \times (2-1) = 1$$

イロト イポト イヨト イヨト 三日

Likelihood Ratio Test

- If you recall, we described the *df* for the likelihood ratio test as the difference in the number of parameters estimated under the alternative minus the number estimated under the null.
- Under the multinomial sampling design, the alternative model is that p_{ij} ≠ p_i.p_{.j} and as such, ∑_i ∑_j p_{ij} = 1. Thus, there is only one constraint and we estimate IJ − 1 cell probabilities.
- Under the null, we have $p_{ij} = p_{i.}p_{.j}$ which is determined by (I-1) and (J-1) marginals. Thus, we only estimate [(I-1) + (J-1)] marginal probabilities.
- Thus, the DF of G^2 is

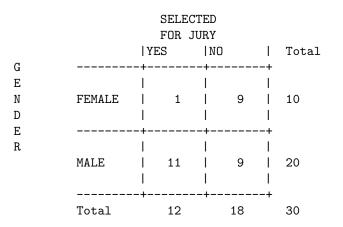
$$df = IJ - 1 - [(I - 1) + (J - 1)] = (I - 1)(J - 1)$$

- Pearson and the LRT have same limiting distribution. (both converge in distribution to χ^2 with df = (I 1)(J 1) as $n \to \infty$)
- Pearson's is score based
- LRT combines the information of the null and alternative hypotheses
- So which one is best?

- X^2 converges in distribution faster than G^2 .
- When n/IJ < 5 (less than 5 per cell), G^2 usually is not a good estimate.
- When I or J is large, Pearson's usually is valid when some $E_{ij} < 5$ but most are greater than 5.
- Therefore, for the general $I \times J$ table, you can usually just use Pearson's Chi Square.
- We will now develop a test for small samples.

Small Samples

Question: Is there Gender Bias in Jury Selection?



The sampling distribution for the this study design is the hypergeometric. However, we will adapt the study design into a small sample exact test.

- In this study, we COULD consider the column totals fixed by design (since the jury has to have 12 members), and the row totals random.
- Then, the columns are independent binomials.
- Using SAS

```
data one;
input sex $ jury $ count;
cards;
1FEMALE 1YES 1
1FEMALE 2NO
             9
2MALE 1YES 11
2MALE 2NO
             9
٠
,
proc freq;
 table sex*jury/expected chisq;
weight count;
run;
```

TABLE OF SEX BY JURY SEX JURY								
Frequency	1							
Expected	1							
Percent	1							
Row Pct	1							
Col Pct	1YES	2N0	1	Total				
	+	+	+					
1FEMALE	1	1	9	10				
	4	1	6					
	3.33		•	33.33				
	10.00	90.0	00					
	8.33	50.0	00					
	+	+	+					
2MALE	11	1	9	20				
	8	1 1	12					
	36.67	30.0	00	66.67				
	55.00	45.0	00					
	91.67	50.0	00					
	+	+	+					
Total	12	1	18	30				
	40.00	60.0	00	100.00				

STATISTICS FOR TABLE OF SEX BY JURY

Statistic	DF	Value	Prob
Chi-Square	1	5.6250	0.0177
Likelihood Ratio Chi-Square	1	6.3535	0.0117
Continuity Adj. Chi-Square	1	3.9063	0.0481
Mantel-Haenszel Chi-Square	1	5.4375	0.0197
Phi Coefficient		-0.4330	
Contingency Coefficient		0.3974	
Cramer's V		-0.4330	
WARNING: 25% of the cells hav	ve exp	ected counts	less
than 5. Chi-Square m	nay no	t be a valid	test.

 A rule of thumb in SAS is that the Large Sample approximations for the likelihood ratio and Pearson's Chi-Square are not very good if the sample size is small

WARNING: 25% of the cells have expected counts less than 5. Chi-Square may not be a valid test.

- Suppose for a cross sectional, prospective, or case-control design: some of the cell counts are small (so that *E_{ij}* < 5), and you want to make inferences about the OR.
- A popular technique with small samples is to fix both margins of the (2 × 2) table, and use 'Exact Tests' and confidence intervals.

Suppose, then, for:

- A prospective study (rows margins fixed) we further condition on the column margins
- A case-control study (column margins fixed) we further condition on the rows margins
- A cross sectional (total fixed) we condition on both row and column margins.
- In all cases, we have a conditional distribution with row and column margins fixed.

- What is the conditional distribution of Y₁₁ given both row and column margins are fixed.
- First note, unlike the other distributions discussed, since the margins are fixed and known, we will show that this conditional distribution is a function of only one unknown parameter
- This follows from what we have seen:
- If the total sample size is fixed (cross sectional), we have 3 unknown parameters, (p_{11}, p_{12}, p_{21})
- If one of the margins is fixed (prospective, or case-control study), we have two unknown parameters, (p₁, p₂) or (π₁, π₂)
- Intuitively, given we know both margins, if we know one cell count (say Y_{11}), then we can figure out the other 3 cell counts by subtraction. This implies that we can characterize the conditional distribution by 1 parameter.
- Thus, given the margins are fixed, we only need to consider one cell count as random, and, by convention Y₁₁ is usually chosen. (you could have chosen any of the 4 cell counts, though).

Can you complete all of the observed cell counts given the information available? Yes.

		Colı	umn	
		1	2	
Row	1	<i>Y</i> ₁₁		$Y_{1.}$
	2			Y ₂ .
		$Y_{\cdot 1}$	Y.2	$N = n_{}$

• **Question:** Then, what is the conditional distribution of Y₁₁ given both row and column margins are fixed.

$$P[Y_{11} = y_{11}|y_{1\cdot}, y_{\cdot 1}, y_{\cdot \cdot}, OR]$$

 After some tedious algebra, you can show it is non-central hypergeometric, i.e.,

$$P[Y_{11} = y_{11} | y_{1.}, y_{.1}, y_{..}, OR] = \frac{\begin{pmatrix} y_{.1} \\ y_{11} \end{pmatrix} \begin{pmatrix} y_{..} - y_{.1} \\ y_{1.} - y_{11} \end{pmatrix} (OR)^{y_{11}}}{\sum_{\ell=0}^{y_{.1}} \begin{pmatrix} y_{.1} \\ \ell \end{pmatrix} \begin{pmatrix} y_{..} - y_{.1} \\ y_{1.} - \ell \end{pmatrix} (OR)^{\ell}}$$

where, for all designs,

$$OR = \frac{O_{11}O_{22}}{O_{21}O_{12}},$$

• We denote the distribution of Y_{11} by

$$(Y_{11}|y_{1.}, y_{.1}, y_{..}) \sim HG(y_{..}, y_{.1}, y_{1.}, OR)$$

Notes about non-central hypergeometric

- Again, unlike the other distributions discussed, since the margins are fixed and known, the non-central hypergeometric is a function of only one unknown parameter, the OR.
- Thus, the conditional distribution given both margins is called non-central hypergeometric.
- Given both margins are fixed, if you know one of the 4 cells of the table, then you know all 4 cells (only one of the 4 counts in the table is non-redundant).
- Under the null H₀ OR=1, the non-central hypergeometric is called the central hypergeometric or just the hypergeometric.
- We will use the hypergeometric distribution (i.e., the non-central hypergeometric under H₀ OR=1) to obtain an 'Exact' Test for H₀ OR=1. This test is appropriate in small samples.

Let's consider the following table with both Row and Column totals fixed.

		Colı	umn	
		1	2	
Row	1	<i>Y</i> ₁₁	<i>Y</i> ₁₂	<i>Y</i> ₁ .
	2	Y_{21}	Y_{22}	<i>Y</i> ₂ .
		<i>Y</i> .1	Y.2	$N = Y_{}$

Many define the $\{1,1\}$ cell as the "Pivot Cell".

Before we consider the sampling distribution, lets consider the constraints on the Pivot Cell.

The Values L_1 and L_2

- We know that Y_{11} must not exceed the marginal totals, $Y_{.1}$ or $Y_{1.}$
- That is,

$$Y_{11} \leq Y_{\cdot 1}$$
 and $Y_{11} \leq Y_{1 \cdot 1}$

• Therefore, the largest value Y₁₁ can assume can be denoted as L₂ in which

$$L_2 = \min(Y_{\cdot 1}, Y_{1 \cdot})$$

- Similarly, the minimum value of Y_{11} is also constrained.
- It is harder to visualize, but the minimum value Y₁₁ can assume, denoted as L₁, is

$$L_1 = \max(0, Y_{1.} + Y_{.1} - Y_{..})$$

21 / 50

イロト イポト イヨト イヨト 三日

Suppose you observe the following marginal distribution.

		Colu	ımn	
		1	2	
Row	1	<i>y</i> ₁₁		6
	2			3
		5	4	9

- We want to determine L_1 and L_2
- So that we can determine the values the Pivot Cell can assume.
- The values in which the Pivot Cell can assume are used in the significance testing.

Based on the previous slide's table, L_1 and L_2 are

$$L_1 = \max(0, Y_{1.} + Y_{.1} - Y_{..})$$

= max(0, 6 + 5 - 9)
= max(0, 2)
= 2

and

$$L_2 = \min(Y_{1.}, Y_{.1}) \\ = \min(6, 5) \\ = 5$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへの

23 / 50

Therefore, the values that Y_{11} can assume are $\{2, 3, 4, 5\}$.

All Possible Contingency Tables

• Since each table is uniquely defined by the pivot cell, the following tables are all of the possible configurations.

$OR_{E} = 0.078$		Co	lumn		(
		1	2		
Row	1	2	4	6	F
	2	3	0	3	
		5	4	9	

OR = 0.5		Column		
		1	2	
Row	1	3	3	6
	2	2	1	3
		5	4	9

OR = 4		Со		
		1	2	
Row	1	4	2	6
	2	1	2	3
		5	4	9

$OR_{E} = 25.7$	**	Co	lumn	
		1	2	
Row	1	5	1	6
	2	0	3	3
		5	4	9

- Suppose the table observed is flagged with "**".
- How do we know if the Rows and Columns are independent?
- Note, as Y_{11} increases, so does the OR.

Test Statistics

• The probability of observing any given table is

$$P[Y_{11} = y_{11} | Y_{1.}, Y_{2.}, Y_{.1}, Y_{.2}] = \frac{\begin{pmatrix} y_{.1} \\ y_{11} \end{pmatrix} \begin{pmatrix} y_{.2} \\ y_{12} \end{pmatrix}}{\begin{pmatrix} y_{..} \\ y_{1.} \end{pmatrix}}$$

• The probability of observing our table is

$$P[Y_{11} = 5|6, 3, 5, 4] = \frac{\begin{pmatrix} 5 \\ 5 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \end{pmatrix}}{\begin{pmatrix} 9 \\ 6 \end{pmatrix}} = \frac{4}{84} = 0.0476$$

 We now need to develop tests to determine whether or not this arrangement supports or rejects independence.

Suppose we want to test

$$H_O:OR = 1$$
 or $E(Y_{11}) = y_{1.}y_{.1}/y_{..}$

versus

$$H_A: OR > 1$$
 or $E(Y_{11}) > y_{1.}y_{.1}/y_{..}$

- Let $y_{11,obs}$ be the observed value of Y_{11} ; we will reject the null in favor of the alternative if $y_{11,obs}$ is large (recall from the example, as Y_{11} increases, so does the OR).
- Then, the exact p-value (one-sided) is the sum of the table probabilities in which the pivot cell is greater than or equal to the Y_{11,obs}.

 Or more specifically, The exact p-value looks at the upper tail:

$$p - value = P[Y_{11} \ge y_{11,obs} | H_O: OR = 1]$$

$$= \sum_{\substack{\ell=y_{11,obs}}}^{L_2=\min(y_{\cdot 1},y_{1\cdot})} \frac{\begin{pmatrix} y_{\cdot 1} \\ \ell \end{pmatrix} \begin{pmatrix} y_{\cdot 2} \\ y_{1\cdot} - \ell \end{pmatrix}}{\begin{pmatrix} y_{\cdot 1} \\ y_{1\cdot} \end{pmatrix}}$$

- Note that ℓ increments the values of Y₁₁ to produce the tables as extreme (ℓ = Y_{11,obs} and more extreme (approaching L₂)
- Note $y_{1.} = y_{11} + y_{12}$ so $y_{12} = y_{1.} y_{11}$.

• Suppose we want to test

versus

$$H_O:OR = 1$$
 or $E(Y_{11}) = y_1 \cdot y_{\cdot 1} / y_{\cdot 1}$

$$H_A: OR < 1$$
 or $E(Y_{11}) < y_1.y_{.1}/y_{..}$

We will reject the null in favor of the alternative if y_{11,obs} is small.
Then, the exact p-value looks at the lower tail:

$$p - value = P[Y_{11} \le y_{11,obs}|H_0:OR = 1]$$

$$= \sum_{\ell=L_1=\max(0,y_1.+y_{\cdot 1}-y_{\cdot \cdot})}^{y_{\cdot 1}} \frac{\begin{pmatrix} y_{\cdot 1} \\ \ell \end{pmatrix} \begin{pmatrix} y_{\cdot 2} \\ y_{1.} - \ell \end{pmatrix}}{\begin{pmatrix} y_{\cdot} \\ y_{1.} \end{pmatrix}}$$

Fisher's Exact (2-sided) Test

• Suppose we want to test

$$H_O:OR = 1$$
 or $E(Y_{11}) = y_{1.}y_{.1}/y_{..}$

versus

$$H_A: OR \neq 1$$
 or $E(Y_{11}) \neq y_{1.}y_{.1}/y_{..}$

• The exact *p*-value here is the exact 2-sided *p*-value is

$$P\left[\begin{array}{c|c} \text{seeing a result as likely or} \\ \text{less likely than the observed} \\ \text{result in either direction} \end{array} \middle| \begin{array}{c} \mathsf{H}_0: \mathit{OR} = 1 \\ \text{H}_0: \mathit{OR} = 1 \\ \end{array} \right].$$

29 / 50

◆□> ◆□> ◆豆> ◆豆> ・豆 ・ のへの

In general, to calculate the 2-sided p-value,

Calculate the probability of the observed result under the null

$$\pi = P[Y_{11} = y_{11,obs} | H_O: OR = 1]$$

$$= \frac{\left(\begin{array}{c} y_{\cdot 1} \\ y_{11,obs} \end{array}\right) \left(\begin{array}{c} y_{\cdot \cdot} - y_{\cdot 1} \\ y_{1\cdot} - y_{11,obs} \end{array}\right)}{\left(\begin{array}{c} y_{\cdot \cdot} \\ y_{1\cdot} \end{array}\right)}$$

◆□ → ◆□ → ◆臣 → ◆臣 → □ 臣.

1 Recall, Y_{11} can take on the values

$$\max(0, y_{1.} + y_{.1} - y_{..}) \le Y_{11} \le \min(y_{1.}, y_{.1}),$$

Calculate the probabilities of all of these values,

$$\pi_{\ell} = P[Y_{11} = \ell | \mathsf{H}_O: OR = 1]$$

Sum the probabilities π_{ℓ} in (2.) that are less than or equal to the observed probability π in (1.)

.

$$p - value = \sum_{\ell=\max(0,y_1.+y_{\cdot 1}-y_{\cdot \cdot})}^{\min(y_1.,y_{\cdot 1})} \pi_{\ell} I(\pi_{\ell} \le \pi)$$

where

$$I(\pi_{\ell} \le \pi) = \begin{cases} 1 \text{ if } \pi_{\ell} \le \pi \\ 0 \text{ if } \pi_{\ell} > \pi \end{cases}$$

イロト イポト イヨト イヨト 三日

Using our example "By Hand"

Recall, $P(Y_{11,obs} = 5) = 0.0476$. Below are the calculations of the other three tables.

$$P[Y_{11} = 2|6, 3, 5, 4] = \frac{\begin{pmatrix} 5 \\ 2 \end{pmatrix} \begin{pmatrix} 4 \\ 4 \end{pmatrix}}{\begin{pmatrix} 9 \\ 6 \end{pmatrix}}$$

= $\frac{10}{84}$
= 0.1190
$$P[Y_{11} = 3|6, 3, 5, 4] = \frac{\begin{pmatrix} 5 \\ 3 \end{pmatrix} \begin{pmatrix} 4 \\ 3 \end{pmatrix}}{\begin{pmatrix} 9 \\ 6 \end{pmatrix}}$$

= $\frac{40}{84}$
= 0.4762

<ロト < 部ト < 国ト < 国ト = のへの 32/50

$$P[Y_{11} = 4|6, 3, 5, 4] = \frac{\begin{pmatrix} 5\\4 \end{pmatrix} \begin{pmatrix} 4\\2 \end{pmatrix}}{\begin{pmatrix} 9\\6 \end{pmatrix}}$$
$$= \frac{30}{84}$$
$$= 0.3571$$

• Then, for
$$H_A: OR < 1$$
,
 $p-value = 0.1190 + 0.4762 + 0.3571 + 0.0476 = 1.0$

• for
$$H_A: OR > 1$$
,
 $p-value = 0.0476$

• For $H_A: OR \neq 1$, p-value = 0.0476 (we observed the most extreme arrangement)

Using SAS

```
data test;
 input row $ col$ count;
cards;
1row 1col 5
1row 2col 1
2row 1col 0
2row 2col 3
;
run;
proc freq;
tables row*col/exact;
weight count;
run;
```

Frequency							
Percent	L						
Row Pct	I.						
Col Pct	11	lcol	1:	2col	I	Total	
 1row		5	1	1		6	
IIOW	4	-		-		-	
			I	11.11	I	66.67	
		83.33	L	16.67			
	Т	100.00	I	25.00	I		
					-		
2row		0	L	3	Ι	3	
	Т	0.00	L	33.33	Ι	33.33	
	Т	0.00	T	100.00	Ι		
	Т	0.00	I	75.00	I		
	- -				-		
Total		5		4		9	
		55.56		44.44		100.00	

Statistics for Table of row Statistic	by col DF	Value	Prob
Chi-Square Likelihood Ratio Chi-Square Continuity Adj. Chi-Square Mantel-Haenszel Chi-Square	1 1 1 1	5.6250 6.9586 2.7563 5.0000	0.0177 0.0083 0.0969 0.0253
Manvol manufactor on pulle Phi Coefficient Contingency Coefficient Cramer's V WARNING: 100% of the cells than 5. Chi-Square	have exp	0.7906 0.6202 0.7906 mected count	s less

・ロト・(雪)・(当)・(雪)・(ロ)

Cell (1,1) Frequency (F)	5
Left-sided Pr <= F	1.0000
Right-sided Pr >= F	0.0476
Table Probability (P)	0.0476
Two-sided Pr <= P	0.0476
Sample Size = 9	

General Notes about Fisher's Exact Test

- Fisher's Exact p-values is one of the most frequently used p-values you will find in the medical literature (for "good studies")
- However, Cruess (1989) reviewed 201 scientific articles published during 1988 in *The American Journal of Tropical Medicine and Hygiene* and found 148 articles with at least one statistical error. The most common error was found to be the use of a large sample $\chi^2 p$ -value when the sample was too small for the approximation.
- Since the values of Y₁₁ is discrete (highly discrete given a small sample size such as in our example), the actual number of possible p-values is limited.
- For example, Given our example margins, {0.0476, 0.1666, 0.5237, 1.0} are our only potential values.

		1	Variable (7)	2
	1	Y ₁₁	Y ₁₂	Y ₁ .
Variable (X)	2	Y ₂₁	Y ₂₂	Y ₂ .

Variable (V)

Y.2

• The hypergeometric (when OR = 1) is symmetrically defined in the rows and columns.

 $Y_{.1}$

Y..

In particular, under $H_0: OR = 1$

$$P[Y_{11} = y_{11} | OR = 1] = \frac{\begin{pmatrix} y_{\cdot 1} \\ y_{11} \end{pmatrix} \begin{pmatrix} y_{\cdot 2} \\ y_{21} \end{pmatrix}}{\begin{pmatrix} y_{\cdot 1} \\ y_{11} \end{pmatrix}}$$
$$= \frac{\begin{pmatrix} y_{1 \cdot} \\ y_{11} \end{pmatrix} \begin{pmatrix} y_{2 \cdot} \\ y_{21} \end{pmatrix}}{\begin{pmatrix} y_{\cdot 1} \\ y_{11} \end{pmatrix}}$$

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Expected Value of Y_{11} under the null

- Recall, for the hypergeometric distribution, the margins $Y_{i.}$, $Y_{.j}$ and $Y_{..}$ are assumed known and fixed.
- From the theory of the hypergeometric distribution, under the null of no association, the mean is

$$E(Y_{ij}|OR=1) = \frac{y_{i}.y_{.j}}{y_{..}}$$

• For other distributions, we could not write the expected value in terms of the possibly random $Y_{i.}$ and/or $Y_{.j.}$. Since $(Y_{i.}, Y_{.j}, Y_{..})$ are known for the hypergeometric, we can write the expected value in terms of them.

• Thus, the null $H_0: OR = 1$ can be rewritten as

$$\mathsf{H}_0: \mathsf{E}(Y_{ij}|OR=1) = \frac{y_{i.}y_{.j}}{y_{..}},$$

Recall, for all other distribution discussed, under no association,

$$E_{ij} = \frac{[i^{th} \text{ row total } (y_{i\cdot})] \cdot [j^{th} \text{ column total } (y_{\cdot j})]}{[\text{total sample size } (y_{\cdot \cdot})]},$$

is the estimate of $E(Y_{ij})$ under the null of no association

 However, under independence, *E_{ij}* is the exact conditional mean (not an estimate) since *y_i*. and *y_i* are both fixed.

Miscellaneous notes regarding X^2 Test

• Suppose we have the following

$$p_1 = .4$$

and

$$p_2 = .6$$

- where p_1 and p_2 are the true success rates for a prospective study.
- Thus, the true odds ratio is

$$OR = \frac{.40 \cdot .80}{.20 \cdot .60} = 2\frac{2}{3} = 2.666$$

43 / 50

(日) (圖) (E) (E) (E)

• Suppose we randomized 50 subjects (25 in each group) and observe the following table

	Success	Failure	Total
Group 1	10	15	25
Group 2	5	20	25
Total	15	35	50

• And use SAS to test $p_1 = p_2$

```
options nocenter;
data one;
 input row col count;
cards;
 1 1 10
 1 2 15
2 1 5
2 2 20
 ;
run;
proc freq data=one;
 tables row*col/chisq measures;
 weight count;
run;
```

The FREQ Procedure Fisher's Exact Test			
Cell (1,1) Frequency (F)	10		
Left-sided Pr <= F	0.9689		
Right-sided Pr >= F	0.1083		
Table Probability (P)	0.0772		
Two-sided Pr <= P	0.2165		
Estimates of the	Relative Risk	(Row1/Row2)	
Type of Study	Value	95% Confidence	Limits
Case-Control (Odds Ratio)	2.6667	0.7525	9.4497
Cohort (Coll Risk)	2.0000	0.7976	5.0151
Cohort (Col2 Risk) Sample Size = 50	0.7500	0.5153	1.0916

- For this trial, we would fail to reject the null hypothesis (p=0.2165).
- However, our estimated odds ratio is 2.6666 and relative risk is 2.0
- What would happen if our sample size was larger?

```
data two;
 input row col count;
 cards;
1 1 40
1 2 60
2 1 20
2 2 80
;
run;
proc freq data=two;
  tables row*col/chisq measures;
  weight count;
run;
```

Fisher's Exact lest			
Cell (1,1) Frequency (F)	40		
Left-sided Pr <= F	0.9995		
Right-sided Pr >= F	0.0016		
Table Probability (P)	0.0010		
Two-sided Pr <= P	0.0032		
Estimates of the	Relative Risk	(Row1/Row2)	
Type of Study	Value	95% Confidence	Limits
Case-Control (Odds Ratio)	2.6667	1.4166	5.0199
Cohort (Col1 Risk)	2.0000	1.2630	3.1670
Cohort (Col2 Risk)	0.7500	0.6217	0.9048
Sample Size = 200			

Fichor's Exact Tost

Moral of the Story?

- Both examples have the exact same underlying probability distribution
- Both examples have the exact same estimates for OR and RR
- The statistical significance differed
- A Chi-square (or as presented Fisher's exact)'s *p*-value does not indicate how strong an association is in the data (i.e., a smaller p-value, say < 0.001, does not mean there is a "strong" treatment effect)
- It simply indicates that you have evidence for the alternative (i.e., $p_1 \neq p_2$).
- You must use a measure of association to quantify this difference