

## Topic 8

### Chi Square Tests

<b>Topics</b>	1. Introduction to Contingency Tables .....	2
	2. Introduction to the Contingency Table Hypothesis Test of No Association .....	4
	3. The Chi Square Test of No Association in an R x C Table .....	10
	4. ( <i>For Epidemiologists</i> ) Special Case: More on the 2x2 Table .....	20
	5. Hypotheses of Independence or No Association .....	21
	Appendix Relationship Between the Normal(0,1) and the Chi Square Distribution ...	22

## 1. Introduction to Contingency Tables

### We wish to Explore the Association Between Two Variables, Both Discrete

- **Example - Is smoking associated with low birth weight?**
- A common goal of many research studies is investigation of the association of 2 factors, both discrete; eg – smoking (yes/no) and low birth weight (yes/no).
- In Topic 8, our focus is in the setting of two categorical variables, such as smoking and low birth weight, and the use of **chi-square tests of association and homogeneity**.

### Introduction to Contingency Tables

- **Example -** Suppose we do a study to investigate the relationship between smoking and impairment of lung function, measured by forced vital capacity (FVC).
- Suppose  $n = 100$  people are selected for the study.
- For each person we note their smoking behavior (smoke or don't smoke) and their forced vital capacity, FVC (normal or abnormal).

	FVC		
	normal	abnormal	
smoke	a	b	a + b
don't smoke	c	d	c + d
	a + c	b + d	n = a + b + c + d

these are counts

Fixed by sample size

- One scenario is the following set of counts

	fvc		
	abn	normal	
smoke	50	0	50
don't smoke	0	50	50
	50	50	100

What can be said about the relationship between fvc and smoking?

- All 50 smokers have an abnormal FVC
- And all 50 non-smokers have normal FVC
- This is an illustration of a **perfect association** in that “once smoking status is known, FVC status is known also”

- Another scenario is the following set of counts

	fvc		
	abn	normal	
smoke	25	25	50
don't smoke	25	25	50
	50	50	100

- In this scenario, half (25) of the smokers have an abnormal FVC
- But we also observe that half (25) of the 50 *non*-smokers have an abnormal FVC, also.
- This similarity in the data suggest that there is **no association** between smoking status and FVC
- Put another way, the data suggest that lung function, as measured by FVC, is **independent** of smoking status.

## 2. Introduction to the Contingency Table Hypothesis Test of No Association

In Topic 7 (*Hypothesis Testing*), we used the idea of “proof by contradiction” to develop hypothesis tests. We do this in the contingency table setting, too.

**Example, continued** - Suppose we do a study to investigate the relationship between smoking and impairment of lung function, measured by forced vital capacity (FVC).

What are our **null** and **alternative** hypotheses?

- Consider the notation that says

$\pi_1$  = the proportion of smokers with abnormal fvc

$\pi_2$  = the proportion of non-smokers with abnormal fvc

- The **null hypothesis** is that of **independence, NO association**, and says the proportion with abnormal fvc is the same, regardless of smoking status.

- $\pi_1 = \pi_2$

- The **alternative hypothesis** is that of **association/dependence** and says the proportion with abnormal fvc will be a different number, depending on smoking status.

- $\pi_1 \neq \pi_2$

- Thus,

$H_0$ : There is **no association** between the two variables,  $\pi_1 = \pi_2$

*“we’ll argue proof-by-contradiction from here ...”*

$H_a$ : The two variables are **associated**  $\pi_1 \neq \pi_2$

*“so as to advance this hypothesis”*

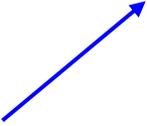
Recall also from Topic 7 that a **test statistic** (also called *pivotal quantity*) is a comparison of what the data are to what we **expected** *under the assumption that the null hypothesis is correct* –

### Introduction to **observed** versus **expected** counts.

- **Observed** counts are represented using the notation “**O**” or “**n**”.
- **Expected** counts are represented using the notation “**E**”

		<u>FVC</u>		
		Abnormal	Normal	
Smoke	$O_{11}$	$O_{12}$		$O_{1.}$
Don't smoke	$O_{21}$	$O_{22}$		$O_{2.}$
	$O_{.1}$	$O_{.2}$	$O_{..}$	

### How to read the subscripts -

- The **first** subscript tells you the “**row**” ( e.g.  $O_{21}$  is a cell count in row “2”)  

- The **second** subscript tells you “**column**” ( e.g.  $O_{21}$  is a cell count in col “1”)  

- Thus,  $O_{21}$  is the count for the cell in row “2” and column “1”
- A subscript replaced with a “dot” is tells you what has been totaled over.
  - Thus,  $O_{2.}$  is the total for row “2” taken over all the columns
  - Similarly  $O_{.1}$  is the total for column “1”, taken over all the rows
  - And  $O_{..}$  is the grand total taken over all rows and over all columns

- Here are the **observed** counts in another scenario

		<u>FVC</u>		
		Abnormal	Normal	
Smoke	O <sub>11</sub> =40	O <sub>12</sub> =10	O <sub>1.</sub> =50	
Don't smoke	O <sub>21</sub> =5	O <sub>22</sub> =45	O <sub>2.</sub> =50	
	O <sub>.1</sub> =45	O <sub>.2</sub> =55	O <sub>..</sub> =100	

- O<sub>21</sub> = 5 is # in row 2 column 1
- O<sub>12</sub> = 10 is # in row 1 column 2
- O<sub>1.</sub> = 50 is the row 1 total
- O<sub>.1</sub> = 45 is the column 1 total

What are the **expected** counts “E” under the assumption that the null hypothesis is true?

**Hint – You already have an intuition for this. If a fair coin is tossed 20 times, the expected number of heads is 10. 10 represents (20 tosses) x (.50 probability of heads on each toss) = 10.**

#### Solution:

- Recall that we are utilizing the notation that says

$\pi_1$  = the proportion of smokers with abnormal fvc

$\pi_2$  = the proportion of non-smokers with abnormal fvc

- Under the assumption that the null of **NO association/independence**, then
  - $\pi_1 = \pi_2 = \pi$  a common (*null hypothesis*) value
- The common  $\pi$  is estimated as the observed overall proportion of abnormal fvc.
 

**Hint – You already know this intuitively as = 45/100 = 0.45**  
**= (total # with abnormal fvc)/(grand total).**

$$\hat{\pi} = \frac{45}{100} = \frac{\text{column 1 total}}{\text{grand total}}, \text{ or a bit more formally ...}$$

$$\hat{\pi} = \frac{O_{11} + O_{21}}{O_{11} + O_{12} + O_{21} + O_{22}} = \frac{O_{.1}}{O_{..}} = \frac{40 + 5}{100} = 0.45$$

- Thus, under the assumption that  $H_0$  is true (meaning *no association, independence*), the proportion with abnormal fvc among smokers as well as among non-smokers should be the same as in the overall population, that is,

$$\pi_{1;\text{null}} = \pi_{2;\text{null}} = \hat{\pi} = 0.45$$

- So we expect 45% of the 50 smokers, or 22.5 persons, to have abnormal fvc, and we also expect 45% of the non-smokers, or 22.5 persons, to have abnormal fvc. *Yes, you are right. These expected counts are NOT whole integers. That's okay. Do NOT round.*

$$\text{Expected \# smokers w abnormal FVC} = (\#\text{Smokers})(\hat{\pi}) = (50)(.45) = 22.5 = E_{11}$$

$$\text{Expected \# NONsmokers w abnormal FVC} = (\#\text{NONSmokers})(\hat{\pi}) = (50)(.45) = 22.5 = E_{21}$$

- We also need to obtain the expected counts of **normal** fvc.
  - We expect 55% of the 50 smokers, or 27.5, to have **normal** fvc, and we also expect 55% of the non-smokers, or 27.5, to have **normal** fvc.
    - Expected # smokers w **normal** FVC =  $(\#\text{Smokers})(1-\hat{\pi}) = (50)(.55) = 27.5 = E_{12}$
    - Expected # NONsmokers w **normal** FVC =  $(\#\text{NONSmokers})(1-\hat{\pi}) = (50)(.55) = 27.5 = E_{22}$
- Thus the following **expected** counts emerge.

		<u>FVC</u>		
		Abnormal	Normal	
Smoke	$E_{11}=22.5$	$E_{12}=27.5$	$E_{1.}=50$	
Don't smoke	$E_{21}=22.5$	$E_{22}=27.5$	$E_{2.}=50$	
		$E_{.1}=45$	$E_{.2}=55$	$E_{..}=100$

- $E_{21}=22.5$   $E_{12}=27.5$
- $E_{1.}=50$   $E_{.1}=45$

- **IMPORTANT NOTE about row totals and column totals -**
  - **The expected row totals match the observed row totals.**
  - **The expected column totals match the observed column totals.**
  - **These totals have a special name - “**marginals**”.**

**The “marginals” are treated as fixed constants (“givens”).**

## What is the **test statistic** (pivotal quantity)?

It is a **chi square** statistic. Examination reveals that it is defined to be a function of the comparisons of observed and expected counts. It can also be appreciated as a kind of “signal-to-noise” construction.

$$\begin{aligned} \text{Chi Square}_{df} = \chi_{df}^2 &= \sum_{\text{all cells "i,j"}} \left[ \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \right] \\ &= \sum_{\text{all cells "i,j"}} \left[ \frac{(\text{Observed}_{ij} - \text{Expected}_{ij})^2}{\text{Expected}_{ij}} \right] \end{aligned}$$

- When the **null** hypothesis of **no association** is true, the observed and expected counts will be similar, their difference will be close to **zero**, resulting in a **SMALL** chi square statistic value.
  - **Example – A null hypothesis says that a coin is “fair”. Since a fair coin tossed 20 times is expected to yield 10 “heads”, the result of 20 tosses is likely to be a number of observed heads that is close to 10.**
- When the **alternative** hypothesis of an **association** is true the observed counts will be unlike the expected counts, their difference will be non zero and their squared difference will be positive, resulting in a **LARGE POSITIVE** chi square statistic value.
  - **Example – The alternative hypothesis says that a coin is “NOT fair”. An UNfair coin tossed 20 times is expected to yield a different number of heads that is NOT close to the null expectation value of 10.**
- Thus, evidence for the rejection of the null hypothesis is reflected in **LARGE POSITIVE** values of the chi square statistic.

### 3. The Chi Square Test of No Association in an R x C Table

For reasons not detailed here (see Appendix), the comparison of observed and expected counts defined on page 9 is, often, distributed **chi square** when the null is true.

- For one cell, when the null is true,

$$\frac{\left[ \begin{array}{cc} \text{Observed} & \text{Expected} \\ \text{Count} & - \text{Count} \end{array} \right]^2}{\text{Expected Count}} \text{ is distributed } \mathbf{Chi\ Square\ (df = 1)} \text{ approximately.}$$

- Summed over all cells in an R x C table, when the null is true,  
In a table that has “R” rows and “C” columns, the same calculation is repeated RC times and then summed

$$\sum_{i=1}^R \sum_{j=1}^C \frac{\left[ \begin{array}{cc} \text{Observed} & \text{Expected} \\ \text{Count (i, j)} & - \text{Count (i, j)} \end{array} \right]^2}{\text{Expected Count (i, j)}} \text{ is distributed } \mathbf{Chi\ Square\ (df = [R-1][C-1])}$$

approximately.

### More on Degrees of Freedom (*What is the correct degrees of freedom?*)

Recall that “degrees of freedom” can be appreciated as the **number of independent pieces of information**. Here is how the idea works in contingency tables:

- In a 2x2 table

$x$	$n_1 - x$	$n_1$	$\therefore$ we have "freedom" to fill in only one of the cells
$n_3 - x$	$n_2 - (n_3 - x)$	$n_2$	
$n_3$	$n_4$	$n$	$\Rightarrow$ 1 degree of freedom

- Now see what happens in larger tables

$x$	$x$		$= 2 \text{ d.f.}$

$x$	$x$	$x$		$= 3 \text{ d.f.}$

$x$	$x$		$= 4 \text{ d.f.}$
$x$	$x$		

$x$	$x$	$x$	$x$		$= 4 \text{ d.f.}$

- In each scenario, the last column is not free and the last row is not free.

- More generally,

**In an R x C table**

**Degrees of freedom = (#rows - 1) x (#columns - 1)**  
**= (R - 1)(C - 1)**

We have the tools for computing the chi square test of association in a contingency table.

**Example**

Suppose we wish to investigate whether or not there is an association between income level and how regularly a person visits his or her doctor. Consider the following count data.

Last Consulted Physician				
Income	≤ 6 months	7-12 months	>12 months	Total
< \$6000	186	38	35	259
\$6000-\$9999	227	54	45	326
\$10,000-\$13,999	219	78	78	375
\$14,000-\$19,999	355	112	140	607
≥ \$20,000	653	285	259	1197
<b>Total</b>	<b>1640</b>	<b>567</b>	<b>557</b>	<b>2764</b>

The general layout uses either the “n<sub>ij</sub>” or “O<sub>ij</sub>” notation to represent the observed counts:

		<b>Columns, “j”</b>			
		<i>j</i> = 1	...	<i>j</i> = <i>C</i>	
<b>Rows, “I”</b>	<i>i</i> = 1	<b>O<sub>11</sub>=n<sub>11</sub></b>	...	<b>O<sub>1C</sub>=n<sub>1C</sub></b>	<b>N<sub>1.</sub> = O<sub>1.</sub></b>
	...	...			
	<i>i</i> = <i>R</i>	<b>O<sub>R1</sub>=n<sub>R1</sub></b>	...	<b>O<sub>RC</sub>=n<sub>RC</sub></b>	
	<b>N<sub>.1</sub> = O<sub>.1</sub></b>	...	<b>N<sub>.C</sub> = O<sub>.C</sub></b>		

## What are the $\pi_{ij}$ now?

$\pi_{ij}$  = the probability of being the combination that is income at level “i” and time since last consult at level “j”

Example:  $\pi_{11}$  = probability [ income is <\$6000 AND time since last visit is  $\leq$  6 mos]

$\pi_{i.}$  = the *overall (marginal)* probability that income is at level “i”

Example:  $\pi_{1.}$  = probability [ income is <\$6000 ]

$\pi_{.j}$  = the *overall (marginal)* probability that time since last visit is at level “j”

Example:  $\pi_{.1}$  = probability [ time since last visit is  $\leq$  6 months ]

## What is independence now?

Again, you already have an intuition for this.

Recall the example of tossing a fair coin two times.

Because the outcomes of the two tosses are independent,

Probability of “heads” on toss 1 and “heads” on toss 2 =  $(.50)(.50) = .25$

Now attach some notation to this intuition.

$\pi_{1.}$  = Probability of “heads” on toss 1

$\pi_{.2}$  = Probability of “heads” on toss 2

$\pi_{12}$  = Probability of “heads” on toss 1 and “heads” on toss 2

Independence  $\rightarrow$

$$\begin{aligned}\pi_{12} &= [ \text{probability heads on toss 1} ] \times [ \text{probability heads on toss 2} ] \\ &= [ \pi_{1.} ] [ \pi_{.2} ]\end{aligned}$$

Thus, under independence

$$\begin{array}{ccc} \pi_{ij} & = & [ \pi_{i.} ] [ \pi_{.j} ] \\ \downarrow & & \downarrow \quad \downarrow \\ \text{Pr[ “i” x “j” combination ]} & = & [\text{Marginal “i” prob}] \times [\text{Marginal “j”}] \end{array}$$

**Under independence  $\pi_{ij} = [\pi_i] [\pi_j]$**

*Example, continued-*

$\pi_i$  = Probability that income is level "i"

$\pi_j$  = Probability that time since last visit is at level "j"

$\pi_{ij}$  = Probability income is level "i" AND time since last visit is at level "j"

**Under Independence,**

$$\pi_{ij} = [\pi_i] [\pi_j]$$

### *Assumptions of Chi Square Test of NO Association*

1. The contingency table of count data is a random sample from some population
2. The cross-classification of each individual is independent of the cross-classification of all other individuals.

### *Null and Alternative Hypotheses*

$$H_0 : \pi_{ij} = \pi_i \pi_j$$

$$H_A : \pi_{ij} \neq \pi_i \pi_j$$

### *Null Hypothesis Estimates of the $\pi_{ij}$*

$$\hat{\pi}_{ij} = \hat{\pi}_i \hat{\pi}_j \quad \text{by independence and where}$$

$$\hat{\pi}_i = \frac{n_{i.}}{n} = \frac{\text{row "i" total}}{\text{grand total}}$$

$$\hat{\pi}_j = \frac{n_{.j}}{n} = \frac{\text{column "j" total}}{\text{grand total}}$$

**Null Hypothesis Expected Counts  $E_{ij}$** 

$$E_{ij} = (\# \text{ trials})[\hat{\pi}_{ij} \text{ under null}] = (n)\hat{\pi}_{i.}\hat{\pi}_{.j} = \frac{[\text{row "i" total}][\text{column "j" total}]}{n}$$

**Test Statistic (Pivotal Quantity)**

Just as we did before...

For each cell, departure of the observed counts from the null hypothesis expected counts is obtained using:

$$\frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

The **chi square test statistic of association** is the sum of these over all the cells in the table:

$$\text{Test Statistic} = \sum_{i=1}^R \sum_{j=1}^C \left[ \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \right]$$

**Behavior of the Test Statistic under the assumption of the null hypothesis**

Under the null hypothesis

$$\text{Test Statistic} = \sum_{i=1}^R \sum_{j=1}^C \left[ \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \right] \text{ is distributed } \chi_{df=(R-1)(C-1)}^2$$

### **Decision Rule**

The null hypothesis is rejected for large values of the test statistic. Thus, evidence for rejection of the null hypothesis is reflected in the following (all will occur)

- **LARGE** value of test statistic
- **SMALL** value of achieved significance (p-value)
- **Test statistic value that EXCEEDS CRITICAL VALUE** threshold

### **Computations**

(1) For each cell, compute

$$E_{ij} = \frac{[\text{row "i" total}][\text{column "j" total}]}{n}$$

(2) For each cell, compute

$$\frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

**Example, continued -**Observed Counts (*this is just the table on page 12 again with the "O" notation provided*)**Last Consulted Physician**

<b>Income</b>	<b>≤ 6 months</b>	<b>7-12 months</b>	<b>&gt;12 months</b>	<b>Total</b>
< \$6000	$O_{11} = 186$	$O_{12} = 38$	$O_{13} = 35$	$O_{1.} = 259$
\$6000-\$9999	$O_{21} = 227$	$O_{22} = 54$	$O_{23} = 45$	$O_{2.} = 326$
\$10,000-\$13,999	$O_{31} = 219$	$O_{32} = 78$	$O_{33} = 78$	$O_{3.} = 375$
\$14,000-\$19,999	$O_{41} = 355$	$O_{42} = 112$	$O_{43} = 140$	$O_{4.} = 607$
≥ \$20,000	$O_{51} = 653$	$O_{52} = 285$	$O_{53} = 259$	$O_{5.} = 1197$
<b>Total</b>	$O_{.1} = 1640$	$O_{.2} = 567$	$O_{.3} = 557$	$O_{..} = 2764$

Expected Counts – *note that each entry is (row total)(column total)/(grand total)***Last Consulted Physician**

<b>Income</b>	<b>≤ 6 months</b>	<b>7-12 months</b>	<b>&gt;12 months</b>	<b>Total</b>
< \$6000	$E_{11} = \frac{(259)(1640)}{2764} = 153.68$	$E_{12} = 53.13$	$E_{13} = 52.19$	$E_{1.} = 259$
\$6000-\$9999	$E_{21} = 193.43$	$E_{22} = 66.87$	$E_{23} = 65.70$	$E_{2.} = 326$
\$10,000-\$13,999	$E_{31} = 222.50$	$E_{32} = 76.93$	$E_{33} = 75.57$	$E_{3.} = 375$
\$14,000-\$19,999	$E_{41} = 360.16$	$E_{42} = 124.52$	$E_{43} = 122.32$	$E_{4.} = 607$
≥ \$20,000	$E_{51} = 710.23$	$E_{52} = 245.55$	$E_{53} = \frac{(1197)(557)}{2764} = 241.22$	$E_{5.} = 1197$
<b>Total</b>	$E_{.1} = 1640$	$E_{.2} = 567$	$E_{.3} = 557$	$E_{..} = 2764$

$$\chi^2_{(R-1)(C-1)} = \sum_{\text{all cells}} \left[ \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \right] = \frac{(186 - 153.68)^2}{153.68} + \dots + \frac{(259 - 241.22)^2}{241.22} = 47.90$$

with degrees of freedom =  $(R-1)(C-1) = (5-1)(3-1) = 8$

- **Statistical Decision Using Significance Level (p-value) Approach**  
<http://www.anu.edu.au/nceph/surfstat/surfstat-home/tables/chi.php>

**p-value = Probability [ Chi square with df=8  $\geq$  47.90 ]  $\ll$  .0001**  
**Such an extremely small p-value is statistically significant  $\rightarrow$**   
**Reject the null hypothesis.**

*To use: enter values in the boxes for "df" and chi square value". Click on RIGHT arrow"*

**SurfStat chi-squared calculator**

The calculator interface shows two chi-square distribution curves. The left curve is shaded red, and the right curve is white with a red tail. Below the curves are input fields for df (8),  $\chi^2$  value (47.90), and probability (empty). There are also right and left arrow buttons between the input fields.

*You should then see "0" in the box labeled "probability"; thus, probability  $\ll$  .0001*

**SurfStat chi-squared calculator**

The calculator interface shows two chi-square distribution curves. The left curve is shaded red, and the right curve is white with a red tail. Below the curves are input fields for df (8),  $\chi^2$  value (47.90), and probability (0). There are also right and left arrow buttons between the input fields.

### Statistical Decision Using Critical Region Approach with type I error=.05

<http://www.anu.edu.au/nceph/surfstat/surfstat-home/tables/chi.php>

The critical region approach with type I error = .05 tells us to reject the null hypothesis for values of the test statistic that are greater than the 95<sup>th</sup> percentile. Equivalently, it tells us to reject the null hypothesis for values of the test statistic that are in the upper 5% of the distribution.

$\chi^2_{.95;df=8} = 15.51$  is our critical value.

Observed statistic = 47.90 >>  $\chi^2_{.95;df=8} = 15.51 \rightarrow$

Reject the null hypothesis.

*To use: Click the radio dial button for “area to the left”. Then enter values in the boxes for “df” and probability”. Then click on LEFT arrow”*

The screenshot shows a web-based calculator for chi-square distributions. It features two radio buttons: the left one is selected, indicating 'area to the left', and the right one is unselected, indicating 'area to the right'. Below the radio buttons are three input fields: 'd.f.' with the value '8', 'χ² value' which is currently empty, and 'probability' with the value '0.95'. There are also left and right arrow buttons between the input fields.

*You should then see the value 15.51 in to box labeled chi square value*

This screenshot shows the same calculator interface as above, but the 'χ² value' input field now contains the number '15.51', which is the critical value for a chi-square distribution with 8 degrees of freedom and a probability of 0.95.

#### Example, continued -

These data provide statistically significant evidence that time since last visit to the doctor is NOT independent of income, that there is an association between income and frequency of visit to the doctor.

**Important note!** What we've learned is that there **is** an association, but **not its nature**. This will be considered further in PubHlth 640, *Intermediate Biostatistics*.

#### 4. (For Epidemiologists) Special Case: More on the 2x2 Table

Many epidemiology texts use a different notation for representing counts in the same chi square test of no association in a 2x2 table. Counts are “a”, “b”, “c”, and “d” as follows.

		2 <sup>nd</sup> Classification Variable		
		1	2	
1 <sup>st</sup> Classification	1	a	b	a + b
	2	c	d	c + d
		a + c	b + d	n

The calculation for the chi square test that you’ve learned as being given by

$$\chi^2 = \sum_{\text{all cells}} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

is the same calculation as the following shortcut formula

$$\chi^2 = \frac{n(ad-bc)^2}{(a+c)(b+d)(c+d)(a+b)}$$

when the notation for the cell entries is the “a”, “b”, “c”, and “d”, above.

## 5. Hypotheses of Independence or No Association

“Independence”, “No Association”, “Homogeneity of Proportions” are alternative wordings for the same thing.

For example,

- (1) “Length of time since last visit to physician” is independent of “income” means that income has no bearing on the elapsed time between visits to a physician. The expected elapsed time is the same regardless of income level.
- (2) There is no association between coffee consumption and lung cancer means that an individual’s likelihood of lung cancer is not affected by his or her coffee consumption.
- (3) The equality of probability of success on treatment (experimental versus standard of care) in a randomized trial of two groups is a test of homogeneity of proportions.

The hypotheses of “independence”, “no association”, “homogeneity of proportions” are equivalent wordings of the same null hypothesis in an analysis of contingency table data.

**Appendix**  
**Relationship Between the Normal(0,1) and the Chi Square Distributions**  
*For the interested reader .....*

**This appendix explains how it is reasonable to use a continuous probability model distribution (the chi square) for the analysis of discrete (counts) data, in particular, investigations of association in a contingency table.**

- Previously (see Topic 6, *Estimation*), we obtained a chi square random variable when working with a function of the sample variance  $S^2$ .
- It is also possible to obtain a chi square random variable as the square of a Normal(0,1) variable. *Recall that this is what we have so far ...*

IF	THEN	Has a Chi Square Distribution with DF =
Z has a distribution that is Normal (0,1)	$Z^2$	1
X has a distribution that is Normal ( $\mu, \sigma^2$ ), so that  Z - score = $\frac{X - \mu}{\sigma}$	{ Z-score } <sup>2</sup>	1
$X_1, X_2, \dots, X_n$ are each distributed Normal ( $\mu, \sigma^2$ ) and are independent, so that  $\bar{X}$ is Normal ( $\mu, \sigma^2/n$ ) and  Z - score = $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$	{ Z-score } <sup>2</sup>	1
$X_1, X_2, \dots, X_n$ are each distributed Normal ( $\mu, \sigma^2$ ) and are independent and we calculate  $S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n - 1}$	$\frac{(n - 1)S^2}{\sigma^2}$	(n-1)

**Our new formulation of a chi square random variable comes from working with a Bernoulli, the sum of independent Bernoulli random variables, and the central limit theorem. What we get is a great result. The chi square distribution for a continuous random variable can be used as a good model for the analysis of discrete data, namely data in the form of counts.**

	<p><math>Z_1, Z_2, \dots, Z_n</math> are each Bernoulli with probability of event = <math>\pi</math>.</p> $E[Z_i] = \mu = \pi$ $\text{Var}[Z_i] = \sigma^2 = \pi(1 - \pi)$ <p style="text-align: center;">↓</p>	
	<p><b>1.</b> The net number of events <math>X = \sum_{i=1}^n Z_i</math> is Binomial (<math>N, \pi</math>)</p> <p><b>2.</b> We learned previously that the distribution of the <u>average</u> of the <math>Z_i</math> is well described as Normal(<math>\mu, \sigma^2/n</math>).</p> <p style="text-align: center;">Apply this notion here: By convention,</p> $\bar{Z} = \frac{\sum_{i=1}^n Z_i}{n} = \frac{X}{n} = \bar{X}$ <p style="text-align: center;">↓</p>	
	<p><b>3.</b> So perhaps the distribution of the <u>sum</u> is also well described as Normal. At least approximately</p> <p>If <math>\bar{X}</math> is described well as Normal (<math>\mu, \sigma^2/n</math>)</p> <p>Then <math>X = n\bar{X}</math> is described well as Normal (<math>n\mu, n\sigma^2</math>)</p> <p style="text-align: center;">↓</p>	
	<p style="text-align: center;">Exactly: <math>X</math> is distributed Binomial(<math>n, \pi</math>)</p> <p style="text-align: center;">Approximately: <math>X</math> is distributed Normal (<math>n\mu, n\sigma^2</math>)</p> <p style="text-align: center;">Where: <math>\mu = \pi</math> and <math>\sigma^2 = \pi(1 - \pi)</math></p>	

## Putting it all together ...

IF	THEN	Comment
<p><b>X has a distribution that is <u>Binomial</u> (<math>n, \pi</math>) <u>exactly</u></b></p>	<p><b>X has a distribution that is <u>Normal</u> (<math>n\mu, n\sigma^2</math>) <u>approximately</u>, where</b></p> $\mu = \pi$ $\sigma^2 = \pi(1-\pi)$ <p style="text-align: center;">↓</p>	
	$Z\text{-score} = \frac{X - E(X)}{SD(X)}$ $= \frac{X - n\mu}{\sqrt{n\sigma}}$ $= \frac{X - n\pi}{\sqrt{n\pi(1-\pi)}}$ <p><b>is approx. Normal(0,1)</b></p> <p style="text-align: center;">↓</p>	
	<p><b>{ Z-score }<sup>2</sup> has distribution that is well described as Chi Square with df = 1.</b></p>	<p><b>We arrive at a continuous distribution model (chi square) approximation for count data.</b></p>

Thus, the { Z-score }<sup>2</sup> that is distributed approximately Chi Square (df=1) is the (O-E)<sup>2</sup>/E introduced previously.

- Preliminaries

$$\begin{aligned} X &= \text{"Observed"} = O \\ n\pi &= \text{"Expected"} = E \end{aligned}$$

- As n gets larger and larger

$$n\pi(1-\pi) \rightarrow n\pi(1) = \text{"Expected"} = E$$

- Upon substitution,

$$\{\text{Z-Score}\}^2 = \left\{ \frac{X-n\pi}{\sqrt{n\pi(1-\pi)}} \right\}^2 \rightarrow \left\{ \frac{X-n\pi}{\sqrt{n\pi(1)}} \right\}^2 = \left\{ \frac{O-E}{\sqrt{E}} \right\}^2 = \frac{(O-E)^2}{E}$$

Thus,

- For **one cell**, when the *null hypothesis is true*, the **central limit theorem** gives us

$$\frac{\left[ \begin{array}{cc} \text{Observed} & \text{Expected} \\ \text{Count} & - \text{Count} \end{array} \right]^2}{\text{Expected Count}} \text{ is Chi Square (df = 1) approximately.}$$

- For **RC cells**, when the *null hypothesis is true*, the **central limit theorem** and the **definition of the chi square distribution** give us

$$\sum_{i=1}^R \sum_{j=1}^C \frac{\left[ \begin{array}{cc} \text{Observed} & \text{Expected} \\ \text{Count}_{ij} & - \text{Count}_{ij} \end{array} \right]^2}{\text{Expected Count}_{ij}} \text{ is Chi Square [df=(R-1)(C-1)] approx.}$$