

UCLA STAT 13
Introduction to Statistical Methods for the
Life and Health Sciences

Instructor: Ivo Dinov,
 Asst. Prof. of Statistics and Neurology

Teaching Assistants:
 Brandi Shanata & Tiffany Head

University of California, Los Angeles, Fall 2007
http://www.stat.ucla.edu/~dinov/courses_students.html

Slide 1

Stat 13, UCLA, Ivo Dinov

Chapter 10
 Chi-Square Test
 Relative Risk/Odds Ratios

Slide 2

Stat 13, UCLA, Ivo Dinov

The χ^2 Goodness of Fit Test

- Let's start by considering analysis of a **single sample of categorical data**
- This is a hypothesis test, so we will be going over the four major HT parts:
 - #1 The general for of the hypotheses:
 - H_0 : probabilities are equal to some specified values
 - H_a : probabilities are not equal to some specified values
 - #2 The Chi-Square test statistic (p.393)
 - O – Observed frequency
 - E – Expected frequency (according to H_0)
 - For the goodness of fit test
 - $$\chi^2 = \sum \frac{(O - E)^2}{E}$$
 - df = # of categories – 1

Slide 3

Stat 13, UCLA, Ivo Dinov

The χ^2 Goodness of Fit Test

- Like other test statistics a smaller value for indicates that the data agree with H_0
- If there is disagreement from H_0 , the test stat will be large because the difference between the observed and expected values is large
- #3 P-value:
 - Table 9, p.686
 - http://socr.stat.ucla.edu/htmls/SOCR_Distributions.html
 - Uses df (similar idea to the t table)
 - After first n-1 categories have been specified, the last can be determined because the proportions must add to 1
 - One tailed distribution, not symmetric (different from t table)
- #4 Conclusion similar to other conclusions (TBD)

Slide 4

Stat 13, UCLA, Ivo Dinov

The χ^2 Goodness of Fit Test

Example: Mendel's pea experiment. Suppose a tall offspring is the event of interest and that the true proportion of tall peas (based on a 3:1 phenotypic ratio) is 3/4 or $p = 0.75$. He would like to show that his data follow this 3:1 phenotypic ratio.

The hypotheses (#1):

H_0 : $P(\text{tall}) = 0.75$ (No effect, follows a 3:1 phenotypic ratio)
 $P(\text{dwarf}) = 0.25$
 H_a : $P(\text{tall}) \neq 0.75$
 $P(\text{dwarf}) \neq 0.25$

Slide 5

Stat 13, UCLA, Ivo Dinov

The χ^2 Goodness of Fit Test

Suppose the data were:

$N = 1064$ (Total)
 Tall = 787 These are the O's (observed values)
 Dwarf = 277

	Observed	
	Observed	Expected
Tall	787	798
Dwarf	277	266

To calculate the E's (expected values), we will take the hypothesized proportions under H_0 and multiply them by the total sample size

Tall = $(0.75)(1064) = 798$ These are the E's (expected values),
 Dwarf = $(0.25)(1064) = 266$

Quick check to see if total = 1064

Slide 6

Stat 13, UCLA, Ivo Dinov

The χ^2 Goodness of Fit Test

Next calculate the test statistic (#2)

$$\chi^2_s = \frac{(787 - 798)^2}{798} + \frac{(277 - 266)^2}{266} = 0.152 + 0.455 = 0.607$$

The p-value (#3):

$$df = 2 - 1 = 1$$

$P > 0.20$, fail to reject H_0

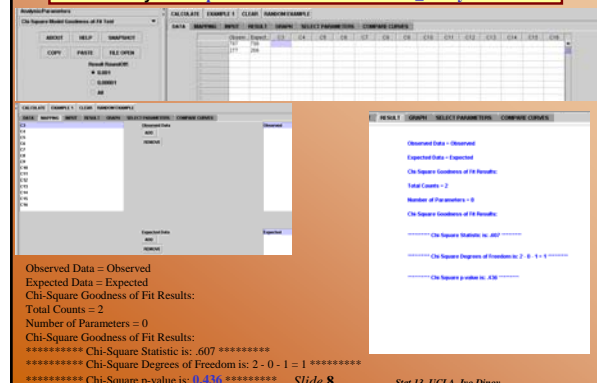
CONCLUSION: These data provide evidence that the true proportions of tall and dwarf offspring are not statistically significantly different from their hypothesized values of 0.75 and 0.25, respectively. In other words, these data are reasonably consistent with the Mendelian 3:1 phenotypic ratio.

Slide 7

Stat 13, UCLA, Jon Dineen

The χ^2 Goodness of Fit Test

SOCR Analysis: http://socr.ucla.edu/htmls/SOCR_Analyses.html



Slide 8

Stat 13, UCLA, Jon Dineen

The χ^2 Goodness of Fit Test

● Tips for calculating χ^2 (p.393):

- Use the SOCR Resource (www.socr.ucla.edu)
- The table of observed frequencies must include ALL categories, so that the sum of the Observed's is equal to the total number of observations
- The O's must be absolute, rather than relative frequencies (i.e., counts not percentages)
- Can round each part to a minimum of 2 decimal places, if you aren't using your calculator's memory

Slide 9

Stat 13, UCLA, Jon Dineen

Compound Hypotheses

- The hypotheses for the t-test contained one assertion: that the means were equal or not.
- The goodness of fit test can contain more than one assertion (e.g., $a=a_0$, $b=b_0$, ..., $c=c_0$)
 - this is called a **compound hypothesis**
 - The alternative hypothesis is non-directional, it measures deviations in all directions (*at least one* probability differs from its hypothesized value)

Slide 10

Stat 13, UCLA, Jon Dineen

Directionality

- **RECALL:** dichotomous – having two categories
- If the categorical variable is dichotomous, H_0 is not compound, so we can specify a directional alternative
 - when one category goes up the other must go down
 - **RULE OF THUMB:** when $df = 1$, the alternative can be specified as directional

Slide 11

Stat 13, UCLA, Jon Dineen

Directionality

Example: A hotspot is defined as a 10 km² area that is species rich (heavily populated by the species of interest). Suppose in a study of butterfly hotspots in a particular region, the number of butterfly hotspots in a sample of 2,588, 10 km² areas is 165. In theory, 5% of the areas should be butterfly hotspots. Do the data provide evidence to suggest that the number of butterfly hotspots is increasing from the theoretical standards? Test using $\alpha = 0.01$.

Slide 12

Stat 13, UCLA, Jon Dineen

Directionality

H_0 : $p(\text{hotspot}) = 0.05$

$p(\text{other spot}) = 0.95$

H_a : $p(\text{hotspot}) > 0.05$

$p(\text{other spot}) < 0.95$

	Hotspot	Other spot	Total
Observed	165	2423	2588
Expected	$(0.05)(2588)$ = 129.4	$(0.95)(2588)$ = 2458.6	2588

$$\chi^2_s = \frac{(165 - 129.4)^2}{129.4} + \frac{(2423 - 2458.6)^2}{2458.6}$$

$$= 9.79 + 0.52 = 10.31$$

Slide 13

Stat 13, UCLA, Jon Dineen

Directionality

$df = 2 - 1 = 1$

$0.001 < p < 0.01$, however because of directional alternative the p-value needs to be divided by 2 (* see note at top of table 9)

Therefore, $0.0005 < p < 0.005$; Reject H_0 .

CONCLUSION: These data provide evidence that in this region the number of butterfly hotspots is increasing from theoretical standards (ie. greater than 5%).

Slide 14

Stat 13, UCLA, Jon Dineen

Goodness of Fit Test, in general

- The expected cell counts can be determined by:

- Pre-specified proportions set-up in the experiment

- For example: 5% hot spots, 95% other spots

- Implied

- For example: Of 250 births at a local hospital is there evidence that there is a gender difference in the proportion of males and females? Without further information this implies that we are looking for $P(\text{males}) = 0.50$ and $P(\text{females}) = 0.50$.

Slide 15

Stat 13, UCLA, Jon Dineen

Goodness of Fit Test, in general

- Goodness of fit tests can be compound (i.e., Have more than 2 categories):

- For example: Of 250 randomly selected CP college students is there evidence to show that there is a difference in area of home residence, defined as: Northern California (North of SLO); Southern California (In SLO or South of SLO); or Out of State? Without further information this implies that we are looking for $P(\text{N.Cal}) = 0.33$, $P(\text{S.Cal}) = 0.33$, and $P(\text{Out of State}) = 0.33$.

- <http://socr.stat.ucla.edu/Applets.dir/SOCRCurveFitter.html>

Slide 16

Stat 13, UCLA, Jon Dineen

Example – Polynomial model fitting

<http://www.socr.ucla.edu/Applets.dir/SOCRCurveFitter.html>

Slide 17

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

- We will now consider analysis of two samples of categorical data
- This type of analysis utilizes tables, called contingency tables
 - Contingency tables focus on the dependency or association between column and row variables

Slide 18

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

Example: Suppose 200 randomly selected cancer patients were asked if their primary diagnosis was Brain cancer and if they owned a cell phone before their diagnosis. The results are presented in the table below:

		Brain cancer		
		Yes	No	Total
Cell Phone	Yes	18	80	98
	No	7	95	102
	Total	25	175	200

Slide 19

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

- Does it seem like there is an association between brain cancer and cell phone use?

■ How could we tell quickly?

Of the brain cancer patients $18/25 = 0.72$, owned a cell phone before their diagnosis.

$\hat{P}(CP|BC) = 0.72$, estimated probability of owning a cell phone given that the patient has brain cancer.

Of the other cancer patients, $80/175 = 0.46$, owned a cell phone before their diagnosis.

$\hat{P}(CP|NBC) = 0.46$, estimated probability of owning a cell phone given that the patient has another cancer.

		Brain cancer		
		Yes	No	Total
Cell Phone	Yes	18	80	98
	No	7	95	102
	Total	25	175	200

Slide 20

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

- The goal: We want to analyze the association, if any, between brain cancer and cell phone use
- This is a 2 X 2 table because there are two possible outcomes for each variable (each variable is dichotomous)
- Consider the following population parameters:
 $P(CP|BC)$ = true probability of owning a cell phone (CP) given that the patient had brain cancer (BC) is estimated by
 $\hat{P} = (CP|BC) = 0.72$

 $P(CP|NBC)$ = true probability of owning a cell phone given that the patient had another cancer, is estimated by
 $\hat{P} = (CP|NBC) = 0.46$

Slide 21

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

- The general form of a hypothesis test for a contingency table:
 ■ #1 The hypotheses:
 H_0 : there is no association between variable 1 and variable 2 (independence)
 H_a : there is an association between variable 1 and variable 2 (dependence)
 NOTE: Using symbols can be tricky, be careful and read section 10.3

Slide 22

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

- #2 The test statistic:

Expected cell counts can be calculated by

$$E = \frac{(\text{row total})(\text{column total})}{\text{grand total}}$$

$$\chi^2_s = \sum \frac{(O - E)^2}{E}$$

with $df = (\# \text{ rows} - 1)(\# \text{ col} - 1)$

- #3 p-value and #4 conclusion are similar to the goodness of fit test.

Slide 23

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

Example: Brain cancer (cont')

Test to see if there is an association between brain cancer and cell phone use using $\alpha = 0.05$

H_0 : there is no association between brain cancer and cell phone (using notation $P(CP|BC) = P(CP|NBC)$)

H_a : there is an association between brain cancer and cell phone (using notation $P(CP|BC) \neq P(CP|NBC)$)

$(98)(25)/200$

		Brain cancer		
		Yes	No	Total
Cell Phone	Yes	18 (12.25)	80 (85.75)	98
	No	7 (12.75)	95 (89.25)	102
	Total	25	175	200

Slide 24

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

$$\chi^2_s = \frac{(18-12.25)^2}{12.25} + \frac{(7-12.75)^2}{12.75} + \frac{(80-85.75)^2}{85.75} + \frac{(95-89.25)^2}{89.25}$$

$$= 2.699 + 2.539 + 0.386 + 0.370 = 6.048$$

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

$0.01 < p < 0.02$, reject H_0 .

CONCLUSION: These data show that there is a statistically significant association between brain cancer and cell phone use in patients that have been previously diagnosed with cancer.

Slide 25

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

Results of Chi-Square Test for Independent or Homogeneity

Number of Rows = 2

Number of Columns = 2

Column 1 Column 2 Row Total

Row 1 18.0 (12.25000) 80.0 (85.75000) 98

Row 2 7.0 (12.75000) 95.0 (89.25000) 102

Col Total 25 175 200

Degrees of Freedom = 1

Pearson Chi-Square Statistics = 6.04813

Slide 26

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

● Output:

Chi-Square Test: C1, C2

Expected counts are printed below observed counts
Chi-Square contributions are printed below expected counts

	C1	C2	Total
1	18	80	98
	12.25	85.75	
	2.699	0.386	

	C1	C2	Total
2	7	95	102
	12.75	89.25	
	2.593	0.370	

	C1	C2	Total
Total	25	175	200

Chi-Sq = 6.048, DF = 1, P-Value = 0.014

Slide 27

Stat 13, UCLA, Jon Dineen

The χ^2 Test for the 2 X 2 Contingency Table

● NOTE: $df = 1$, we could have carried this out as a one-tailed test

■ The probability that a patient with brain cancer owned a cell phone is greater than the probability that another cancer patient owned a cell phone

□ $H_a: P(CP|BC) > P(CP|NBC)$

■ Why didn't we carry this out as a one tailed test?

● CAUTION: Association does not imply Causality!

Slide 28

Stat 13, UCLA, Jon Dineen

Computational Notes

1. Contingency table is useful for calculations, but not nice for presentation in reports.

2. When calculating observed values should be absolute frequencies, not relative frequencies. Also sum of observed values should equal grand total.

● To eyeball a contingency table for differences, check for proportionality of columns:

■ If the columns are nearly proportional then the data seem to agree with H_0

■ If the columns are not proportional then the data seem to disagree with H_0

Slide 29

Stat 13, UCLA, Jon Dineen

Independence and Association in the 2x2 Contingency Table

● There are two main contexts for contingency tables:

■ Two independent samples with a dichotomous observed variable

■ One sample with two dichotomous observed variables

NOTE: The χ^2 test procedure is the same for both situations

Example: Vitamin E. Subjects treated with either vitamin E or placebo for two years, then evaluated for a reduction in plaque from their baseline (Yes or No).

■ Any study involving a dichotomous observed variable and completely randomized allocation to two treatments can be viewed this way

Example: Brain cancer and cell phone use. One sample, cancer patients, two observed variables: brain cancer (yes or no) and cell phone use (yes or no)

Slide 30

Stat 13, UCLA, Jon Dineen

Independence and Association in the 2x2 Contingency Table

- When a dataset is viewed as a single sample with two observed variables, the relationship between the variables is thought of as independence or association.
 - H_0 : independence (no association) between the variables
 - H_a : dependence (association) between the variables

● χ^2 is often called a test of independence or a test of association.

NOTE: If columns and rows are interchanged test statistic will be the same

Slide 31

Stat 13, UCLA, Jon Dineen

The r X k Contingency Table

- We now consider tables that are larger than a 2x2 (more than 2 groups or more than 2 categories), called r x k contingency tables

● Testing procedure is the same as the 2x2 contingency table, just more work and no possibility for a directional alternative

- The goal of an r x k contingency table is to investigate the relationship between the row and column variables

● NOTE: H_0 is a compound hypothesis because it contains more than one independent assertion

- This will be true for all r x k tables larger than 2x2
- In other words, the alternative hypothesis for r x k tables larger than 2x2, will always be non-directional.

Slide 32

Stat 13, UCLA, Jon Dineen

The r X k Contingency Table

Example: Many factors are considered when purchasing earthquake insurance. One factor of interest may be location with respect to a major earthquake fault. Suppose a survey was mailed to California residents in four counties (data shown below). Is there a statistically significant association between county of residence and purchase of earthquake insurance? Test using $\alpha = 0.05$.

		County				
		Contra Costa CC	Santa Clara SC	Los Angeles LA	San Bernardino SB	Total
Earthquake Insurance	Yes	117	222	133	109	581
	No	404	334	204	263	1205
	Total	521	556	337	372	1786

Slide 33

Stat 13, UCLA, Jon Dineen

The r X k Contingency Table

H_0 : There is no association between Earthquake insurance and county of residence in California.

$$\begin{cases} P(Y|CC) = P(Y|SC) = P(Y|LA) = P(Y|SB) \\ P(N|CC) = P(N|SC) = P(N|LA) = P(N|SB) \end{cases}$$

H_a : There is an association between Earthquake insurance and county of residence in California.

The probability of having earthquake insurance is not the same in each county.

Slide 34

Stat 13, UCLA, Jon Dineen

The r X k Contingency Table

Chi-Square Test: C1, C2, C3, C4

<http://socr.stat.ucla.edu/Applets.dir/ChiSquareTable.html>

Expected counts are printed below observed counts

Chi-Square contributions are printed below expected counts

	C1	C2	C3	C4	Total
1	117 169.49 16.253	222 180.87 9.352	133 109.63 4.982	109 121.01 1.193	581
2	404 351.51 7.837	334 375.13 4.509	204 227.37 2.402	263 250.99 0.575	1205
Total	521	556	337	372	1786

Chi-Sq = 47.105, DF = 3, P-Value = 0.000

Slide 35

Stat 13, UCLA, Jon Dineen

The r X k Contingency Table

$p = 0.000 < 0.05$, reject H_0 .

CONCLUSION: These data show that there is a statistically significant association between purchase of earthquake insurance and county of residence in California.

Slide 36

Stat 13, UCLA, Jon Dineen

Applicability of Methods

- Conditions for validity of the χ^2 test:
 1. Design conditions
 - for a goodness of fit, it must be reasonable to regard the data as a random sample of categorical observations from a large population.
 - for a contingency table, it must be appropriate to view the data in one of the following ways:
 - as two or more independent random samples, observed with respect to a categorical variable
 - as one random sample, observed with respect to two categorical variables
- * for either type of test, the observations within a sample must be independent of one another.

Slide 37

Stat 13, UCLA, Jon Dineen

Applicability of Methods

- Conditions for validity of the χ^2 test (cont'):
 2. Sample conditions
 - critical values for table 9 only work if each expected value ≥ 5
 3. Form of H_0
 - for goodness of fit, H_0 specifies values
 - for contingency table, H_0 : row and column are not associated or use notation

Slide 38

Stat 13, UCLA, Jon Dineen

Verification of Conditions

- Data consisting of several samples need to be independent sample.
 - If the design contains blocking or pairing the samples are not independent
- Try to reduce bias
- Only simple random sampling
 - No pairing for the version we are learning, although there is a paired Chi-Square test (section 10.8)
- No hierarchical structure
- Check expected cell counts

Slide 39

Stat 13, UCLA, Jon Dineen

CI for the difference between probabilities

- Chi-Square tests for contingency tables tell us if there is an association or not between categories.
 - They tell us that there is a difference, but is it an important difference?
 - They do not give us any information as to the magnitude of any differences between probabilities
- For this we will calculate a confidence interval for the difference between probabilities

Slide 40

Stat 13, UCLA, Jon Dineen

CI for the difference between probabilities

- A 95% confidence interval for $p_1 - p_2$

$$(\tilde{p}_1 - \tilde{p}_2) \pm z_{0.025} (SE_{\tilde{p}_1 - \tilde{p}_2})$$

Where

$$\tilde{p}_1 = \frac{y_1 + 1}{n_1 + 2} \quad \tilde{p}_2 = \frac{y_2 + 1}{n_2 + 2}$$

$$SE_{\tilde{p}_1 - \tilde{p}_2} = \sqrt{\frac{\tilde{p}_1(1 - \tilde{p}_1)}{n_1 + 2} + \frac{\tilde{p}_2(1 - \tilde{p}_2)}{n_2 + 2}}$$

Sample	
1	2
y_1	y_2
$n_1 - y_1$	$n_2 - y_2$
n_1	n_2

Slide 41

Stat 13, UCLA, Jon Dineen

CI for the difference between probabilities

Example: Brain cancer continued

		Brain cancer		
		Yes	No	Total
Cell Phone	Yes	18	80	98
	No	7	95	102
Total		25	175	200

Calculate a 95% confidence interval for the difference in cell phone use between brain cancer and other cancer patients

$$\tilde{p}_1 = \frac{18 + 1}{25 + 2} = 0.704 \quad \tilde{p}_2 = \frac{80 + 1}{175 + 2} = .458$$

Slide 42

Stat 13, UCLA, Jon Dineen

CI for the difference between probabilities

95% CI continued...

$$SE_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{0.704(0.296)}{25+2} + \frac{0.457(0.543)}{175+2}} = \sqrt{0.009} = 0.095$$

$$(0.704 - 0.458) \pm 1.96(0.095) \\ = 0.246 \pm 0.186 = (0.06, 0.432)$$

We are 95% confident that the difference in the proportion of cell phone ownership between patients with brain cancer and those without brain cancer, is between 6% and 43%.

Slide 43

Stat 13, UCLA, Jon Dineen

CI for the difference between probabilities

What does this mean?

Does this seem like a significant difference?

Can we say that based on this data it appears that owning a cell phone increases the probability of brain cancer?

Slide 44

Stat 13, UCLA, Jon Dineen

Relative Risk

- The chi-square test is often referred to as a test of independence
- Another measure of dependence is relative risk
 - Allows researchers to compare probabilities in terms of their ratio (p_1 / p_2) rather than their difference ($p_1 - p_2$)
 - widely used in studies of public health
- In general a relative risk of 1 indicates that the probabilities of two events are the same.
 - A relative risk > 1 implies that there is increased risk
 - A relative risk < 1 implies that there is decreased risk

Slide 45

Stat 13, UCLA, Jon Dineen

Relative Risk

Example: Brain Cancer and cell phone use (continued)

Cell Phone	Brain cancer		
	Yes	No	Total
Yes	18	80	98
No	7	95	102
Total	25	175	200

Thinking in terms of conditional probability again, but switching the conditional probability around...

$$\hat{P} = (BC|CP) = 18/98 = 0.184$$

$$\hat{P} = (BC|NCP) = 7/102 = 0.069$$

So the relative risk is $0.184 / 0.069 = 2.67$

The risk of having brain cancer is more than 2.5 times greater for cell phone owners when compared to non-cell phone owners.

Slide 46

Stat 13, UCLA, Jon Dineen

Odds Ratio

- Another way to compare two probabilities is in terms of odds
- If an event takes place with probability p , then the odds in favor of the event are $p / (1 - p)$
 - If event A|B has $p = 1/2$, then the odds are $(1/2) / (1/2) = 1$ or 1 to 1 (the probability that event A|B occurs is equal to the probability that it does not occur)
 - If event A|C has $p = 3/4$, then the odds are $(3/4) / (1/4) = 3$ or 3 to 1 (the probability that event A|C occurs is three times as large as the probability that it does not occur)

Slide 47

Stat 13, UCLA, Jon Dineen

Odds Ratio: OR

- The **odds ratio** is the ratio of odds for two probabilities

$$\hat{\theta} = \frac{\hat{P}(A|B)}{1 - \hat{P}(A|B)} \div \frac{\hat{P}(A|C)}{1 - \hat{P}(A|C)}$$

Cell Phone	Brain cancer		
	Yes: A	No	Total
Yes	18	80	98 B
No	7	95	102 C
Total	25	175	200

- In general an OR and it's relationship to 1 is similar to relative risk

- An OR = 1 indicates that the probabilities of two events are the same
- An OR > 1 implies that there is increased risk
- An OR < 1 implies that there is decreased risk

Slide 48

Stat 13, UCLA, Jon Dineen

Odds Ratio

Example: Brain Cancer and cell phone use
(continued)

Cell Phone	Brain cancer		
	Yes	No	Total
Yes	18	80	98
No	7	95	102
Total	25	175	200

Calculate the odds of having brain cancer (A) for cell phone owners (B) compared to non-cell phone (C) owners.

$$\hat{p} = (BC|CP) = 18/98 = 0.184 \quad \hat{\theta} = \frac{0.184}{0.069} = \frac{0.225}{0.074} = 3.04$$

$$\hat{p} = (BC|NCP) = 7/102 = 0.069 \quad \frac{0.931}{0.543}$$

The odds of having brain cancer is about 3 times greater for cell phone owners when compared to non-cell phone owners.

Slide 49

Stat 13, UCLA, Jon Dineen

Odds Ratio

- We could have compared the odds of owning a cell phone given that a patient had brain cancer versus an other cancer (i.e., the *column-wise* probabilities)

$$\hat{p} (CP|BC) = 18/25 = 0.72 \text{ versus } \hat{p} (CP|NBC) = 80/175 = 0.457$$

However this does not seem as important scientifically

- But if we did calculate the OR of owning a cell phone given that a patient had brain cancer versus an other cancer we'd get:

$$\hat{\theta} = \frac{0.72}{0.457} = \frac{2.57}{0.842} = 3.05$$

- Note that this OR comes out to be approximately equal!

Slide 50

Stat 13, UCLA, Jon Dineen

Odds Ratio

- Shortcut formula for an odds ratio:

$$\hat{\theta} = \frac{n_{11}n_{22}}{n_{12}n_{21}}$$

Now it is easier to see why the OR would be the same for the *row-wise* and *column-wise* probabilities!

Where the table structure looks like:

n_{11}	n_{12}
n_{21}	n_{22}

Slide 51

Stat 13, UCLA, Jon Dineen

Relative Risk vs. Odds Ratio

- The formula and reasoning for the relative risk is a little bit easier to follow

- In most cases the two measures are roughly equal to each other

- Odds ratios have an advantage over relative risk because they can be calculated no matter the row or column comparison

- Relative risk runs into problems when the study design is a cohort study or a case-control design

- Odds ratios are an approximation of relative risk

$$OR = RR \cdot (1-P_2)/(1-P_1)$$

Slide 52

Stat 13, UCLA, Jon Dineen

Relative Risk vs. Odds Ratio

Example: Suppose a group of 200 people who have experienced a heart attack and 200 with no heart attack were asked if they were ever smokers.

Ever Smoker (SMK)?	Heart Attack (HA)?	
	Yes	No
Yes	33	18
No	167	182
Total	200	200

- We can reasonably calculate $\hat{p}(\text{SMK}|HA) = 33/200 = 0.165$ and $\hat{p}(\text{SMK}|NHA) = 18/200 = 0.09$

- However, the *row-wise* probabilities (incidence of heart attacks given that someone is a smoker or non-smoker) should not be estimated

- Because the number of subject with and without heart attacks were predetermined in the study design
- We have no information about the incidence of heart attacks

Slide 53

Stat 13, UCLA, Jon Dineen

Relative Risk vs. Odds Ratio

- If we calculate the *column-wise* probabilities and the odds of smoking for those with heart attacks compared to no heart attacks, using $\hat{p}(\text{SMK}|HA) = 0.165$ and $\hat{p}(\text{SMK}|NHA) = 0.09$

OR comes out to about 2.0

- Had we incorrectly calculated the *rowwise* probabilities and then the odds of heart attacks for people who smoke versus non-smokers, using $\hat{p}(HA|\text{SMK}) = 33/51 = 0.65$ and $\hat{p}(HA|\text{NSMK}) = 167/349 = 0.48$

However, the OR comes out to about 2.0

- Remember that the OR will come out to be approximately equal for row and column comparisons

Slide 54

Stat 13, UCLA, Jon Dineen

Relative Risk vs. Odds Ratio

- Because these estimates of the odds ratio are the same for column-wise and row-wise probabilities (see p. 449)
- And we know that the odds ratio is an approximation of relative risk
- We can say that we estimate the relative risk of a heart attack is about 2 twice as great for those who smoke versus who do not smoke
 - Without incorrectly calculating the *row-wise* probabilities

Slide 55

Stat 13, UCLA, Jon Dineen

Odds Ratio Confidence Interval

- Common to report odds ratios along with their CI
- One problem with our estimate of the odds ratio $\hat{\theta}$ is that its sampling distribution is not normal!!!
 - To solve this we take the log of $\hat{\theta}$ and so that the sampling distribution of $\ln(\hat{\theta})$ is normally distributed

SE of $\ln(\hat{\theta})$:
$$SE_{\ln(\hat{\theta})} = \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}}$$

Where the table structure looks like:
$$\begin{array}{c|c} n_{11} & n_{12} \\ \hline n_{21} & n_{22} \end{array}$$

Slide 56

Stat 13, UCLA, Jon Dineen

Odds Ratio Confidence Interval

- A $100(1 - \alpha)\%$ CI for $\ln(\hat{\theta})$ is

table structure to remember
$$\begin{array}{c|c} n_{11} & n_{12} \\ \hline n_{21} & n_{22} \end{array}$$

$$\ln(\hat{\theta}) \pm Z_{\alpha/2}(SE_{\ln(\hat{\theta})})$$

Where $\hat{\theta} = \frac{n_{11}n_{22}}{n_{12}n_{21}}$ and $SE_{\ln(\hat{\theta})} = \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}}$

HINT: You can use your t table to find certain values of $Z_{\alpha/2}$

Slide 57

Stat 13, UCLA, Jon Dineen

Odds Ratio Confidence Interval

Example: Smoking and heart attack (continued)

		Heart Attack?	
		Yes	No
Ever Smoker?	Yes	33	18
	No	167	182
	Total	200	200

Calculate a 90% confidence interval for odds of smoking for heart attack subjects and non-heart attack subjects.

$$\begin{aligned} \hat{\theta} &= \frac{n_{11}n_{22}}{n_{12}n_{21}} = \frac{33 \cdot 182}{18 \cdot 167} = 1.998 \\ SE_{\ln(\hat{\theta})} &= \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}} \\ &= \sqrt{\frac{1}{33} + \frac{1}{18} + \frac{1}{167} + \frac{1}{182}} \\ &= \sqrt{0.09734} = 0.3120 \end{aligned}$$

Slide 58

Stat 13, UCLA, Jon Dineen

Odds Ratio Confidence Interval

So the 90% CI for $\ln(\hat{\theta})$ is

$$\ln(\hat{\theta}) = \ln(1.998) = 0.6921$$

$$\ln(\hat{\theta}) \pm Z_{\alpha/2}(SE_{\ln(\hat{\theta})})$$

$$0.6921 \pm Z_{0.05}(0.3120) =$$

$$0.6921 \pm 1.645(0.3120) =$$

$$(0.1789, 1.2053)$$

But right now this is transformed data (natural log) so we need to untransform it by taking the exponent of the CI

$$(e^{0.1789}, e^{1.2053}) = (1.196, 3.338)$$

Slide 59

Stat 13, UCLA, Jon Dineen

Odds Ratio Confidence Interval

We are confident at the 0.10 level that the true odds of smoking for heart attack subjects and non-heart attack subjects are between 1.196 and 3.338

- So what does this actually mean?
- Does the zero rule work here? *One rule?*
- What if the CI came out to be (0.196, 1.338)?

Slide 60

Stat 13, UCLA, Jon Dineen