# STA 570     Spring 2011

Lecture 25     *Tuesday, April 26*

8.2    Chi-Squared Test of Independence

# Chi-Squared Test of Independence

- ## Assumptions
  - – Two categorical variables
  - – Random sampling (perhaps stratified with respect to the categories of one variable)
  - – Expected cell count at least 5 in all cells

- ## Hypotheses
  - – Null hypothesis: Statistical independence of the two variables
  - – Alternative hypothesis: Statistical dependence

# Comparing Nominal Samples Chi-Squared Test of Independence

- Example: Family Structure and Sexual Activity

- Sociologists think that family structure may have an influence on sexual activity of teenagers

- 380 randomly selected females between 15 and 19 years of ages are asked to disclose
  - Family structure at age 14
  - Whether or not she has had sexual intercourse

- Response variable is binary (nominal)

# Observed and Expected Frequencies

| Sexual activity | Both parents | Single Parent | Parent and Stepparent | Nonparental Guardian | Total |
|---|---|---|---|---|---|
| Yes | 64 | 59 | 44 | 32 | 199 |
| No | 86 | 41 | 36 | 18 | 181 |
| Total | 150 | 100 | 80 | 50 | 380 |

- The expected frequency $f_e$ in a cell equals the product of row and column totals for that cell, divided by the total sample size

| Sexual activity | Both parents | Single Parent | Parent and Stepparent | Nonparental Guardian | Total |
|---|---|---|---|---|---|
| Yes | 78.6 | 52.4 | 41.9 | 26.2 | 199 |
| No | 71.4 | 47.6 | 38.1 | 23.8 | 181 |
| Total | 150 | 100 | 80 | 50 | 380 |

# Chi-Squared Test Statistic

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

- When the null hypothesis of independence is true, then the observed frequencies are close to the expected frequencies, so the chi-squared statistic takes a relatively small value

- A large value of the chi-squared statistic is evidence *against* the null hypothesis

- In order to quantify the evidence and calculate a P-value, we need the sampling distribution of the statistic

- Chi-Squared Distribution

# Chi-Squared Test of Independence

- Test Statistic

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

$$\text{where } f_e = \frac{(\text{Row total}) \cdot (\text{Column total})}{\text{Total sample size}}$$

- In our example:

# Chi-Squared Test of Independence

- ## P-Value

  - P = right-hand tail probability above the observed chi-squared value for chi-squared distribution with *df=(r-1)(c-1)*

  - For the chi-squared test, always use the right-hand tail probability!

- ## Report P-value, reject null hypothesis at alpha-level if P is less than alpha

- ## In our example, P-value =

# Degrees of Freedom

- *(r-1)x(c-1)*
- Given the row marginals and the column marginals, this is the number of frequencies that we need to determine all the other cell frequencies
- In our example, *(r-1)x(c-1)=(2-1)x(4-1)=3*
- This is the number of cell frequencies that are free to vary, because once we know row and column totals, they determine the remaining ones

  (the remaining ones are not free to vary)

| Sexual activity | Both parents | Single Parent | Parent and Stepparent | Nonparental Guardian | Total |
|---|---|---|---|---|---|
| Yes | 64 | 59 | 44 | | 199 |
| No | | | | | 181 |
| Total | 150 | 100 | 80 | 50 | 380 |

# Chi-Squared Test, Properties

- The chi-squared test treats the classifications as nominal

- Any reordering of rows or columns of the table leaves the value of the chi-squared test unchanged

- If either of the classifications is in fact ordinal, this information is not used

- If the response variable is in fact ordinal, one should use the Kruskal-Wallis test instead

# Chi-Squared Test, Misuse

- The chi-squared test should not be used when any of the expected frequencies is less than five
- For smaller sample sizes, there is a procedure that can be used
  - generalized version of Fisher's exact test
  - SAS: PROC FREQ, option EXACT
- Also, the test statistic must be calculated using the observed/expected frequencies, and not using percentages!
- This test can not be used when the samples are dependent.
- For example, when each row or each column has observations on the same subjects, the samples are dependent (McNemar's test can be used then)

# Special Case: Chi-Squared Test, 2x2 Table

- For the 2x2 table with large enough sample sizes, we can use
  - Either the test for a difference of proportions (using normal scores)
  - Or the chi-squared test for association
  - Fortunately, the two tests are equivalent

# Chi-Squared Test, 2x2 Table: Example

- 340 commercial motor vehicle drivers who had accidents in Kentucky from 1998 to 2002

- Two variables:
  - wearing a seat belt (y/n)
  - accident fatal (y/n)

| | | Accident Fatal | | |
|---|---|---|---|---|
| | | Yes | No | |
| Seat Belt | Yes | 30 | 212 | 242 |
| | No | 33 | 52 | 85 |
| | | 63 | 264 | 327 |

# Chi-Squared Test, 2x2 Table: Example

- Testing whether the two variables "Seat Belt" and "Fatal" are associated or independent is equivalent to

- testing whether the fatality rate is the same for the two groups "Seat Belt: Yes" and "Seat Belt: No"

- The row variable is explanatory, the column variable is response

- Calculate the p-value for both tests

# Chi-Squared Test, 2x2 Table

- For the 2x2-table, the chi-squared statistic is exactly the square of the z statistic

- Also, squaring z-scores for certain tail probabilities yields chi-squared scores with *df=1* for the same tail probabilities

- Squared normal = chi-squared with *df=1*

- *In short: For the special case of a 2x2 table, the* **chi-squared test for independence is equivalent to the test for equal proportions** *of two independent samples*

# Summary: Investigating Association Between Two Variables

| | | Response Variable | | |
|---|---|---|---|---|
| | | Unordered Categorical (Nominal) | Ordinal | Quantitative |
| Explanatory Variable | Nominal with 2 Levels | Comparing Proportions of 2 samples | Nonparametric Wilcoxon-Mann-Whitney Test | Comparing Means of 2 samples, t-test for independent samples |
| | Unordered Categorical (Nominal) More than 2 Levels | Analyzing Association, Chi Squared Tests | Nonparametric Kruskal-Wallis Test | ANOVA |
| | Ordinal | | *Spearman Rank Correlation* | |
| | Quantitative | *Logistic Regression* | | *Regression* |

# Quiz