# Chapter 15

## Chi-Squared Tests

# A Common Theme…

| What to do? | Data Type? | Number of Categories? | Statistical Technique: |
|---|---|---|---|
| Describe a population | Nominal | Two or more | $\chi^2$ goodness of fit test |
| Compare two populations | Nominal | Two or more | $\chi^2$ test of a contingency table |
| Compare two or more populations | Nominal | -- | $\chi^2$ test of a contingency table |
| Analyze relationship between two variables | Nominal | -- | $\chi^2$ test of a contingency table |

One data type…

…Two techniques

# Two Techniques...

The first is a ***goodness-of-fit test*** applied to data produced by a ***multinomial experiment***, a generalization of a binomial experiment and is used to describe one population of data.

The second uses data arranged in a ***contingency table*** to determine whether two classifications of a population of nominal data are ***statistically independent***; this test can also be interpreted as a comparison of two or more populations.

In both cases, we use the chi-squared ($\chi^2$) distribution.

# Example 15.1

Two companies, A and B, have recently conducted aggressive advertising campaigns to maintain and possibly increase their respective shares of the market for cars. These two companies enjoy a dominant position in the market. Before the advertising campaigns began, the market share of company A was 45%, whereas company B had 40% of the market. Other competitors accounted for the remaining 15%.



Benz b180



TOYOTA Camry

# Example 15.1

- To determine whether these market shares changed after the advertising campaigns, a marketing analyst solicited the preferences of a random sample of 200 customers of cars. Of the 200 customers, 102 indicated a preference for company A's product, 82 preferred company B's car, and the remaining 16 preferred the products of one of the competitors. Can the analyst infer at the 5% significance level that customer preferences have changed from their levels before the advertising campaigns were launched?

我們不知道 *90* 人到 *102* 人是不是採樣誤差 還是真的支持度有上升

# Example 15.1

- To determine whether these market shares changed after the advertising campaigns, a marketing analyst solicited the preferences of a random sample of 200 customers of cars.

- If the market shares remain the same, each person in the sample will be

– A person who prefer A with a probability 45%

– A person who prefer B with a probability 40%

– A person who prefer other with a probability 15%

# The Multinomial Experiment...

Unlike a binomial experiment which only has two possible outcomes (e.g. heads or tails), a ***multinomial experiment***:

- Consists of a fixed number, **n**, of trials.
- Each trial can have one of **k** outcomes, called cells.
- Each probability $p_i$ remains constant.
- Our usual notion of probabilities holds, namely:

$$p_1 + p_2 + \ldots + p_k = 1, \text{ and}$$

- Each trial is ***independent*** of the other trials.

在這題裡面 *n = 200, k = 3, 45% like company A 40 % like company B etc*

# Chi-squared Goodness-of-Fit Test...

We test whether there is sufficient evidence to reject a *specified set* of values for $p_i$.

To illustrate, our null hypothesis is:

$$H_0: p_1 = a_1, p_2 = a_2, \ldots, p_k = a_k$$

where $a_1, a_2, \ldots, a_k$ are the values we want to test.

Our research hypothesis is:

$$H_1: \text{At least one } p_i \text{ is not equal to its specified value}$$

# Example 15.1

Two companies, A and B, have recently conducted aggressive advertising campaigns to maintain and possibly increase their respective shares of the market for cars. These two companies enjoy a dominant position in the market. Before the advertising campaigns began, the market share of company A was 45%, whereas company B had 40% of the market. Other competitors accounted for the remaining 15%.



Benz b180



TOYOTA Camry

# Example 15.1

- To determine whether these market shares changed after the advertising campaigns, a marketing analyst solicited the preferences of a random sample of 200 customers of cars. Of the 200 customers, 102 indicated a preference for company A's product, 82 preferred company B's car, and the remaining 16 preferred the products of one of the competitors. Can the analyst infer at the 5% significance level that customer preferences have changed from their levels before the advertising campaigns were launched?

# Example 15.1...

We compare market share *before* and *after* an advertising campaign to see if there is a *difference* (i.e. if the advertising was effective in improving market share). We hypothesize values for the parameters equal to the before-market share. That is,

$$H_0: p_1 = .45, p_2 = .40, p_3 = .15$$

The alternative hypothesis is a denial of the null. That is,

$$H_1: \text{At least one } p_i \text{ is not equal to its specified value}$$

# Example 15.1…

**Test Statistic**

If the null hypothesis is true, we would expect the number of customers selecting brand A, brand B, and other to be 200 times the proportions specified under the null hypothesis. That is,

$$e_1 = 200(.45) = 90$$

$$e_2 = 200(.40) = 80$$

$$e_3 = 200(.15) = 30$$

In general, the ***expected frequency*** for each cell is given by

$$e_i = np_i$$

This expression is derived from the formula for the expected value of a binomial random variable, introduced in Section 7.4.
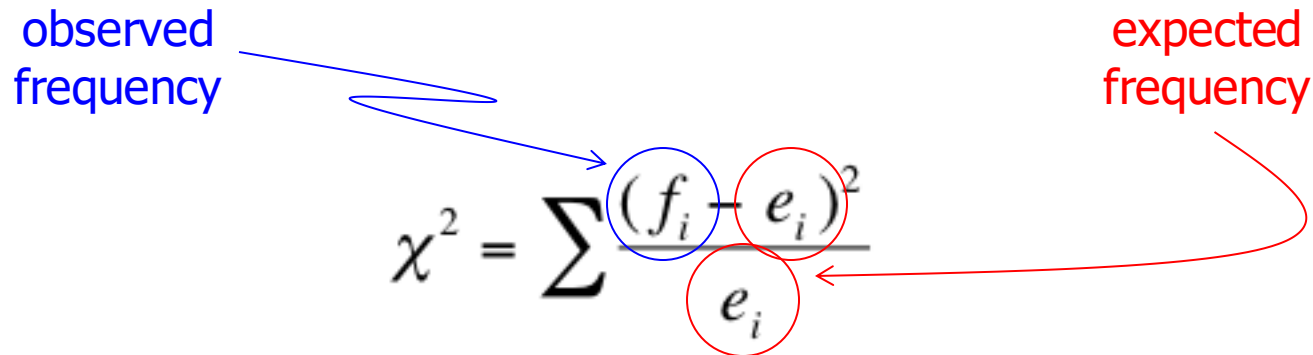
# Example 15.1...

If the expected frequencies and the observed frequencies are quite different, we would conclude that the null hypothesis is false, and we would reject it.

However, if the expected and observed frequencies are similar, we would not reject the null hypothesis.

The test statistic measures the similarity of the expected and observed frequencies.

# Chi-squared Goodness-of-Fit Test...

Our Chi-squared goodness of fit test statistic is given by:

observed frequency

expected frequency

$$\chi^2 = \sum \frac{(f_i - e_i)^2}{e_i}$$

<u>Note</u>: this statistic is *approximately* Chi-squared with k–1 degrees of freedom provided the sample size is large. The rejection region is: $\chi^2 > \chi^2_{\alpha, k-1}$

# Example 15.1...

In order to calculate our test statistic, we lay-out the data in a tabular fashion for easier calculation by hand:

| Company | Observed Frequency | Expected Frequency | Delta | Summation Component |
|---|---|---|---|---|
| | $f_i$ | $e_i$ | $(f_i - e_i)$ | $(f_i - e_i)^2/e_i$ |
| A | 102 | 90 | 12 | 1.60 |
| B | 82 | 80 | 2 | 0.05 |
| Others | 16 | 30 | -14 | 6.53 |
| Total | 200 | 200 | | **8.18** |

Check that these are equal

$$\chi^2 = \sum \frac{(f_i - e_i)^2}{e_i}$$

# Example 15.1...

Our rejection region is:

$$\chi^2 > \chi^2_{\alpha, k-1} = \chi^2_{.05, 3-1} = 5.99147$$

Since our test statistic is 8.18 which is greater than our critical value for Chi-squared, we reject $H_0$ in favor of $H_1$, that is,

*"There is sufficient evidence to infer that the proportions have changed since the advertising campaigns were implemented"*

# Example 15.1…

| ◇ | A | B | C |
|---|---|---|---|
| 1 | | Observed | Expected |
| 2 | | Frequency | Frequency |
| 3 | Company A | **102** | **90** |
| 4 | Company B | **82** | **80** |
| 5 | Others | **16** | **30** |
| 6 | | | |
| 7 | | p-value: | 0.01671136 |

=CHITEST(B3:B5,C3:C5)

**CHITEST**

| Actual_range | B3:B5 | | = {102;82;16} |
| Expected_range | C3:C5 | | = {90;80;30} |

= 0.016711358

Returns the test for independence: the value from the chi-squared distribution for the statistic and the appropriate degrees of freedom.

# =CHITEST(observed frequency, expected frequency)

# Required Conditions...

In order to use this technique, the sample size must be *large enough* so that the expected value for each cell is 5 or more. (i.e. $n \times p_i \geq 5$)

If the ***expected frequency*** is less than five, combine it with other cells to satisfy the condition.

| Company | Observed Frequency | Expected Frequency | Delta | Summation Component |
|---------|--------------------|--------------------|-------|---------------------|
|         | $f_i$ | $e_i$ | $(f_i - ei)$ | $(f_i - e_i)^2/e_i$ |
| A | 102 | 90 | 12 | 1.60 |
| B | 82 | 80 | 2 | 0.05 |
| Others | 16 | **3.5** | 12.5 | 6.53 |
| Total | 200 | 200 | | **8.18** |

# Example

- Acme Toy Company prints baseball cards. The company claims that 30% of the cards are rookies, 60% veterans but not All-Stars, and 10% are veteran All-Stars.

- Suppose a random sample of 100 cards has 50 rookies, 45 veterans, and 5 All-Stars. Is this consistent with Acme's claim? Use a 0.05 level of significance.

# Example

- $X^2 = 19.58$
- P-value=0.000056

# Example

256 visual artists were surveyed to find out their zodiac sign. The results were: Aries (29), Taurus (24), Gemini (22), Cancer (19), Leo (21), Virgo (18), Libra (19), Scorpio (20), Sagittarius (23), Capricorn (18), Aquarius (20), Pisces (23).
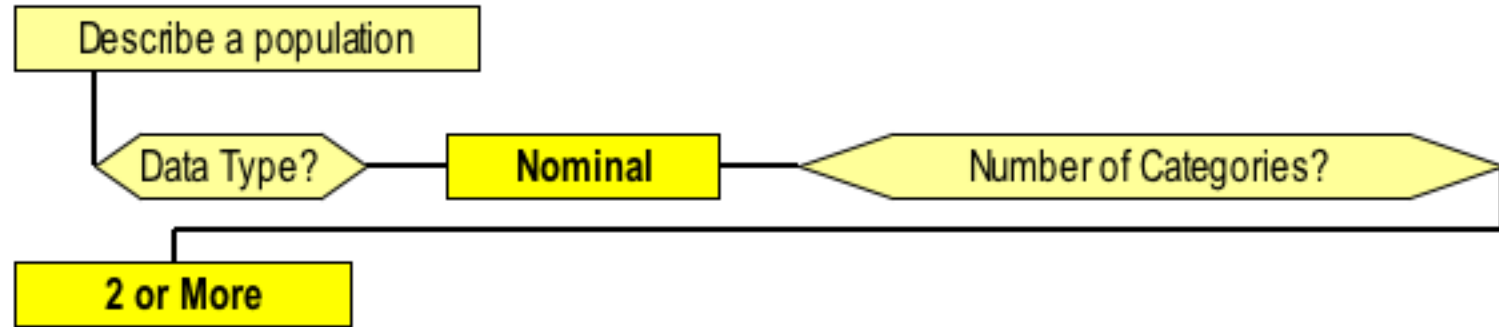
Test the null hypothesis that zodiac signs are evenly distributed across visual artists.

# Example

- $X^2=5.094$

- P-value=0.9265

# Identifying Factors...

Factors that Identify the Chi-Squared Goodness-of-Fit Test:



Test Statistic:
$$\chi^2 = \sum \frac{(f_i - e_i)^2}{e_i}$$

$e_i = (n)(p_i)$

Parameters of interest:

$$p_1, \ p_2, \ ..., \ p_k$$

Required Condition: $e_i \geq 5$

# A Common Theme...

| What to do? | Data Type? | Number of Categories? | Statistical Technique: |
|---|---|---|---|
| Describe a population | Nominal | Two or more | $\chi^2$ goodness of fit test |
| Compare two populations | Nominal | Two or more | $\chi^2$ test of a contingency table |
| Compare two or more populations | Nominal | -- | $\chi^2$ test of a contingency table |
| Analyze relationship between two variables | Nominal | -- | $\chi^2$ test of a contingency table |

One data type...

...Two techniques

# Example 15.2

The MBA program was experiencing problems scheduling their courses. The demand for the program's optional courses and majors was quite variable from one year to the next.

In desperation the dean of the business school turned to a statistics professor for assistance.

The statistics professor believed that the problem may be the variability in the academic background of the students and that the undergraduate degree affects the choice of major.

# Example 15.2

As a start he took a random sample of last year's MBA students and recorded the undergraduate degree and the major selected in the graduate program.

The undergraduate degrees were BA, BEng, BBA, and several others.

There are three possible majors for the MBA students, accounting, finance, and marketing. Can the statistician conclude that the undergraduate degree affects the choice of major?

# Two Techniques...

The first is a ***goodness-of-fit test*** applied to data produced by a ***multinomial experiment***, a generalization of a binomial experiment and is used to describe one population of data.

The second uses data arranged in a ***contingency table*** to determine whether two classifications of a population of nominal data are ***statistically independent***; this test can also be interpreted as a comparison of two or more populations.

In both cases, we use the chi-squared ($\chi^2$) distribution.

# Chi-squared Test of a Contingency Table

The ***Chi-squared test of a contingency table*** is used to:

   • determine whether there is enough evidence to infer that ***two nominal variables are related***, and

   • to infer that ***differences exist*** among two or more populations of nominal variables.

In order to use use these techniques, we need to classify the data according to two different criteria.

# Example 15.2

The MBA program was experiencing problems scheduling their courses. The demand for the program's optional courses and majors was quite variable from one year to the next.

In desperation the dean of the business school turned to a statistics professor for assistance.

The statistics professor believed that the problem may be the variability in the academic background of the students and that the undergraduate degree affects the choice of major.

# Example 15.2

As a start he took a random sample of last year's MBA students and recorded the undergraduate degree and the major selected in the graduate program.

The undergraduate degrees were BA, BEng, BBA, and several others.

There are three possible majors for the MBA students, accounting, finance, and marketing. Can the statistician conclude that the undergraduate degree affects the choice of major?

# Example 15.2

The data are stored in two columns. The first column consist of integers 1, 2, 3, and 4 representing the undergraduate degree where

1 = BA
2 = BEng
3 = BBA
4 = other

The second column lists the MBA major where

1= Accounting
2 = Finance
3 = Marketing

# Example 15.2

The problem objective is to determine whether two variables (undergraduate degree and MBA major) are related. Both variables are nominal. Thus, the technique to use is the chi-squared test of a contingency table. The alternative hypotheses specifies what we test. That is,

$H_1$: The two variables are **dependent**

The null hypothesis is a denial of the alternative hypothesis.

$H_0$: The two variables are **independent**.

# Test Statistic

The test statistic is the same as the one used to test proportions in the goodness-of-fit-test. That is, the test statistic is

$$\chi^2 = \sum \frac{(f_i - e_i)^2}{e_i}$$

Note however, that there is a major difference between the two applications. In this one the null does not specify the proportions $p_i$, from which we compute the expected values $e_i$, which we need to calculate the $\chi^2$ test statistic. That is, we cannot use

$$e = np_i$$

because we don't know the $p_i$ (they are not specified by the null hypothesis). It is necessary to estimate the $p_i$ from the data.

# Example 15.2

The first step is to count the number of students in each of the 12 combinations. The result is called a cross-classification table.

# Example 15.2

| Undergrad Degree | MBA Major | | | |
|---|---|---|---|---|
| | Accounting | Finance | Marketing | Total |
| BA | 31 | 13 | 16 | 60 |
| BEng | 8 | 16 | 7 | 31 |
| BBA | 12 | 10 | 17 | 39 |
| Other | 10 | 5 | 7 | 22 |
| Total | 61 | 44 | 47 | **152** |

# Example 15.2

If the null hypothesis is true (Remember we always start with this assumption.) and the two nominal variables are independent, then, for example

P(BA and Accounting) = [P(BA)] [P(Accounting)]

Since we don't know the values of P(BA) or P(Accounting)

We need to use the data to estimate the probabilities.

# Test Statistic

There are 152 students of which 61 who have chosen accounting as their MBA major. Thus, we estimate the probability of accounting as

$$P(\text{Accounting}) \approx \frac{61}{152} = .401$$

Similarly

$$P(\text{BA}) \approx \frac{60}{152} = .395$$

# Example 15.2…

If the null hypothesis is true

$$P(BA \text{ and Accounting}) = (60/152)(61/152)$$

Now that we have the probability we can calculate the expected value. That is,

$$E(BA \text{ and Accounting}) = 152(60/152)(61/152)$$
$$= (60)(61)/152 = 24.08$$

We can do the same for the other 11 cells.

# Example 15.2

We can now compare *observed* with *expected* frequencies…

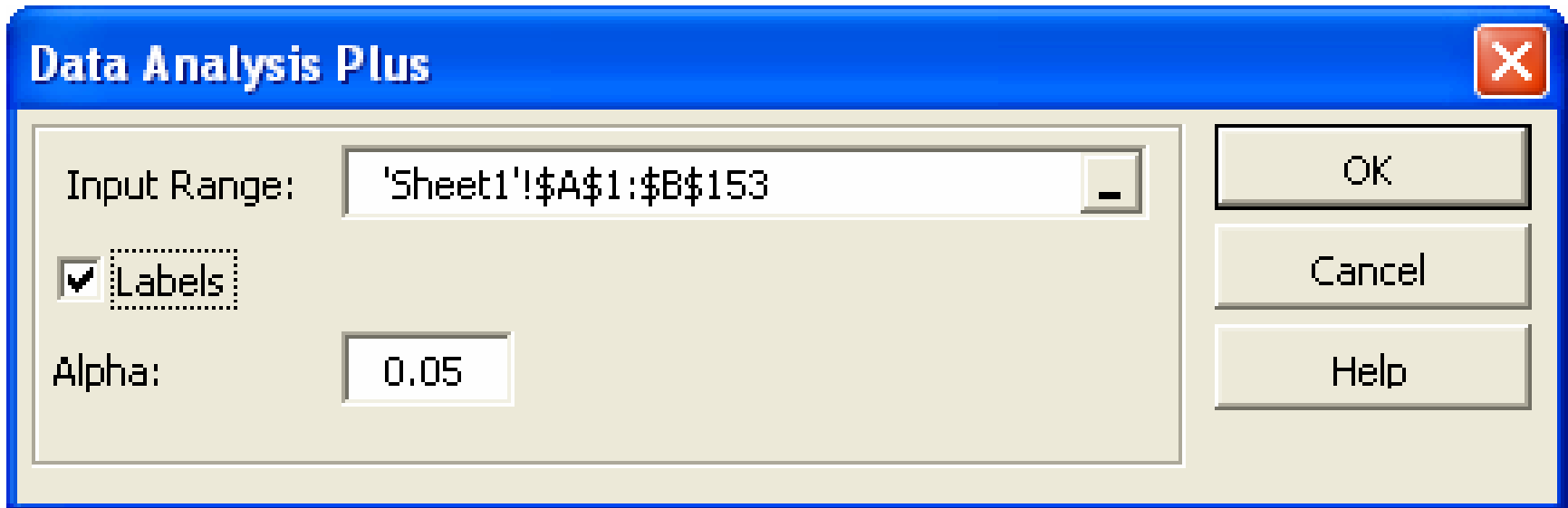| Undergrad Degree | MBA Major | | | | | |
|---|---|---|---|---|---|---|
| | Accounting | | Finance | | Marketing | |
| BA | 31 | 24.08 | 13 | 17.37 | 16 | 18.55 |
| BEng | 8 | 12.44 | 16 | 8.97 | 7 | 9.59 |
| BBA | 12 | 15.65 | 10 | 11.29 | 17 | 12.06 |
| Other | 10 | 8.83 | 5 | 6.37 | 7 | 6.80 |

*df = (4-1)(3-1)=6*

and calculate our test statistic:     d.f = (r-1)*(c-1)

$$\chi^2 = \frac{(31-24.08)^2}{24.08} + \frac{(13-17.37)^2}{17.37} + \ldots + \frac{(7-6.80)^2}{6.80} = 14.70$$

# Example 15.2...

Using Excel : Click Add-Ins, Data Analysis Plus, Contingency Table [if the table has already been prepared] or Contingency Table (Raw Data) [if the table has not been completed]

**Data Analysis Plus**

Input Range:    'Sheet1'!$A$1:$B$153

☑ Labels

Alpha:    0.05

OK

Cancel

Help

# Example 15.2...

The printout below was produced from file Xm15-02 using the Contingency Table (Raw Data) command

|   | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | **Contingency Table** | | | | | |
| 2 | | | | | | |
| 3 | | *Degree* | | | | |
| 4 | *MBA Major* | | 1 | 2 | 3 | TOTAL |
| 5 | | 1 | 31 | 13 | 16 | 60 |
| 6 | | 2 | 8 | 16 | 7 | 31 |
| 7 | | 3 | 12 | 10 | 17 | 39 |
| 8 | | 4 | 10 | 5 | 7 | 22 |
| 9 | | TOTAL | 61 | 44 | 47 | 152 |
| 10 | | | | | | |
| 11 | | | | | | |
| 12 | | chi-squared Stat | | | 14.7019 | |
| 13 | | df | | | 6 | |
| 14 | | p-value | | | 0.0227 | |
| 15 | | chi-squared Critical | | | 12.5916 | |

# Example 15.2...

The p-value is .0227. There is enough evidence to infer that the MBA major and the undergraduate degree are related.

We can also interpret the results of this test in two other ways.

1. There is enough evidence to infer that there are differences in MBA major between the four undergraduate categories.

2. There is enough evidence to infer that there are differences in undergraduate degree between the majors.
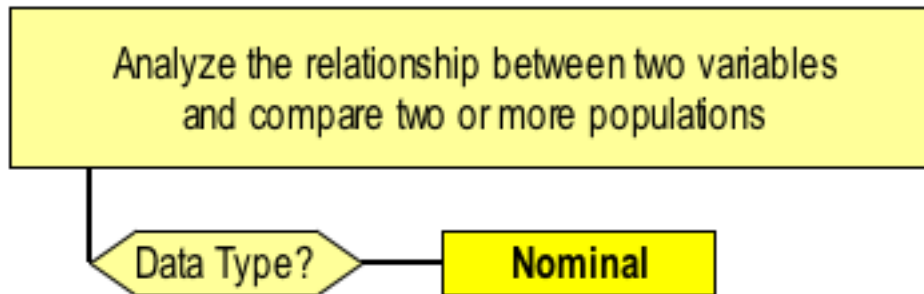
# Required Condition – Rule of Five...

In a contingency table where one or more cells have **expected values** of <span style="color:red">**less than 5**</span>, we need to combine rows or columns to satisfy the rule of five.

**Note:** by doing this, the degrees of freedom must be changed as well.

# Identifying Factors...

Factors that identify the Chi-squared test of a contingency table:



Analyze the relationship between two variables and compare two or more populations

Data Type? — **Nominal**

Test Statistic:

$$\chi^2 = \sum \frac{(f_i - e_i)^2}{e_i}$$

Parameters of interest:

$$p_1, \ p_2, \ ..., \ p_k$$

Required Condition: $e_i \geq 5$

$$e_{ij} = \frac{Row \ i \ total \times Column \ j \ total}{Sample \ size}$$

## Table 15.1 Statistical Techniques for Nominal Data

| Problem Objective | Categories | Statistical Technique |
|---|---|---|
| Describe a population | 2 | z-test of p or the chi-squared goodness-of-fit test |
| Describe a population | More than 2 | Chi-squared goodness-of-fit test |
| Compare two populations | 2 | z-test $p_1 - p_2$ or chi-squared test of a contingency table |
| Compare two populations | More than 2 | Chi-squared test of a contingency table |
| Compare more than two populations | 2 or more | Chi-squared test of a contingency table |
| Analyze the relationship between two variables | 2 or more | Chi-squared test of a contingency table |

# Chi-Squared test for Normality

- The goodness of fit Chi-squared test can be used to determine if data were drawn from any distribution.

# Chi-Squared test for Normality

# Chi-Squared test for Normality

- The goodness of fit Chi-squared test can be used to determined if data were drawn from any distribution.

- The general procedure:

  – Hypothesize on the parameter values of the distribution we test (i.e. $\mu = \mu_0, \sigma = \sigma_0$ for the normal distribution).

  – For the variable tested X specify disjoint ranges that cover all its possible values.

  – Build a Chi squared statistic that (aggregately) compares the expected frequency under $H_0$ and the actual frequency of observations that fall in each range.

  – Run a goodness of fit test based on the multinomial experiment.

# Chi-Squared test for Normality

- **Testing for normality in Example 12.1**

  For a sample size of n=50 (see Xm12-01) ,the sample mean was 460.38 with standard error of 38.83.

- Can we infer from the data provided that this sample was drawn from a **normal distribution** with $\mu = 460.38$ and $\sigma = 38.83$? Use 5% significance level.

# $\chi^2$ test for normality

**Solution**

First let us select z values that define each cell (expected frequency > 5 for each cell.)

$z_1 = -1$; $P(z < -1) = p_1 = .1587$; $e_1 = np_1 = 50(.1587) = \boxed{7.94}$

$z_2 = 0$; $P(-1 < z < 0) = p_2 = .3413$; $e_2 = np_2 = 50(.3413) = \boxed{17.07}$
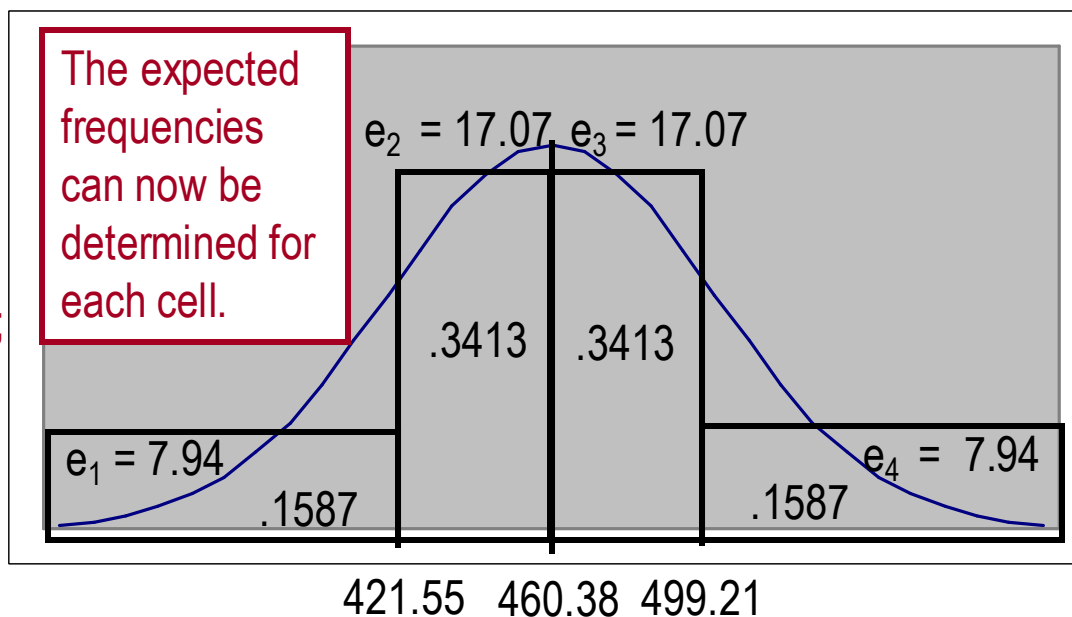
$z_3 = 1$; $P(0 < z < 1) = p_3 = .3413$; $e_3 = \boxed{17.07}$

$\quad\quad\quad P(z > 1) = p_4 = .1587$; $e_4 = \boxed{7.94}$

The cell boundaries are calculated from the corresponding z values **under $H_0$**.

$z_1 = (x_1 - 460.38)/38.83 = -1$;
$x_1 = 421.55$

The expected frequencies can now be determined for each cell.

$e_2 = 17.07$  $e_3 = 17.07$

.3413   .3413

$e_1 = 7.94$

$e_4 = 7.94$

.1587   .1587

421.55   460.38   499.21

# $\chi^2$ test for normality

– The test statistic

$$\chi^2 = \frac{(10 - 7.94)^2}{7.94} + \frac{(13 - 17.07)^2}{17.07} + \frac{(19 - 17.07)^2}{17.07} + \frac{(8 - 7.94)^2}{7.94} = 1.72$$

# $\chi^2$ test for normality

- The test statistic

$$\chi^2 = \frac{(10 - 7.94)^2}{7.94} + \frac{(13 - 17.07)^2}{17.07} + \frac{(19 - 17.07)^2}{17.07} + \frac{(8 - 7.94)^2}{7.94} = 1.72$$

- The rejection region

$\chi^2 > \chi^2_{\alpha, k-1-L}$ where L is the number of parameters estimated from the data.

$$\chi^2_{\alpha, k-3} = \chi^2_{.05, 4-3} = 3.84146$$

Conclusion: There is insufficient evidence to conclude at 5% significance level that the data are not normally distributed.

# Chi-Squared test for Normality

- **Sample size > 220**
  - $Z < -2$        0.0228
  - $-2 < Z \leqq -1$   0.1359
  - $-1 < Z \leqq 0$   0.3413
  - $0 < Z \leqq 1$    0.1359
  - $1 < Z \leqq 2$    0.0228

- **Sample size between 80 and 220**
  - $Z \leqq -1.5$        0.0668
  - $-1.5 < Z \leqq -0.5$   0.2417
  - $-0.5 < Z \leqq 0.5$   0.3829
  - $0.5 < Z \leqq 1.5$    0.2417
  - $Z > 1.5$        0.0668

# Chi-Squared test for Normality

- **Sample size < 80**
  - $Z \leqq -1$        0.1587
  - $-1 < Z \leqq 0$    0.3413
  - $0 < Z \leqq 1$     0.3413
  - $Z > 1$              0.1587

# Example

- The Anger and Heart Disease study

– A study followed a random sample of 8474 people with normal blood pressure for about four years. All the individuals were free of heart disease at the beginning of the study.

– Each person took the Spielberger Trait Anger Scale Test, which measures how prone a person is to sudden anger.

– Researchers also recorded whether each individual developed coronary heart disease (CHD). This includes people who had heart attacks and those who needed medical treatment for heart disease.

# Example

- The data

- – CHD: coronary heart disease

|  | Low Anger | Moderate Anger | High Anger | Total |
|---|---|---|---|---|
| CHD | 53 | 110 | 27 | 190 |
| NO CHD | 3057 | 4621 | 606 | 8284 |
| Total | 3110 | 4731 | 633 | 8474 |

– What is the relationship between anger and CHD status?

# Example

- The data

– First, we can eyeball the frequency distribution of the data:

| | Low Anger | Moderate Anger | High Anger | Total |
|---|---|---|---|---|
| CHD | 27.9% | 57.9% | 14.2% | 100.0% |
| NO CHD | 36.9% | 55.8% | 7.3% | 100.0% |
| Total | 36.7% | 55.8% | 7.5% | 100.0% |

– Does the relationship between anger and CHD status exist?

# Example

- The data

– Or, we can eyeball the frequency distribution of the data in another way:

|  | Low Anger | Moderate Anger | High Anger | Total |
|---|---|---|---|---|
| CHD | 1.7% | 2.3% | 4.3% | 2.2% |
| NO CHD | 98.3% | 97.7% | 95.7% | 97.8% |
| Total | 100.0% | 100.0% | 100.0% | 100.0% |

– Does the relationship between anger and CHD status exist?

# Example

- Hypothesis

– H0: there is no association between anger level and heart disease in the population of people with normal blood pressure

– H1: there is an association between anger level and heart disease in the population of people with normal blood pressure

or

– H0: anger and heart disease are independent in the population of people with normal blood pressure

– H1: anger and heart disease are not independent in the population of people with normal blood pressure.

# Example

- Analysis
– Chi-square statistic: $\chi^2 = 16.077$
– P-value is less than 0.0005


- Conclusion
– Because the P-value is clearly less than $\alpha = 0.05$, we reject H0 and conclude that anger level and heart disease are associated in the population of people with normal blood pressure.

# Logistic Regression

Why logistic regression?

- There are many important research topics for which the dependent variable is "limited" (i.e., nominal).

- For example: voting, morbidity or mortality, and participation data is not continuous or distributed normally.

- Binary logistic regression is a type of regression analysis where the dependent variable is a dummy variable: coded 0 (did not vote) or 1(did vote)
  - A special case of multinomial experiment

# Logistic Regression

- If using a linear regression…

  $Y = \gamma + \varphi X + e$ ; where $Y = (0, 1)$

- e is not normally distributed because Y takes on only two values

- The predicted probabilities can be greater than 1 or less than 0

# Example

- Hurricane evacuations

Did you evacuate your home to go someplace safer before Hurricane Dennis (Floyd) hit?

- 1 YES

- 2 NO

- 3 DON'T KNOW

- 4 REFUSED

# Example: Data

| EVAC | PETS | MOBLHOME | TENURE | EDUC |
|------|------|----------|--------|------|
| 0 | 1 | 0 | 16 | 16 |
| 0 | 1 | 0 | 26 | 12 |
| 0 | 1 | 1 | 11 | 13 |
| 1 | 1 | 1 | 1 | 10 |
| 1 | 0 | 0 | 5 | 12 |
| 0 | 0 | 0 | 34 | 12 |
| 0 | 0 | 0 | 3 | 14 |
| 0 | 1 | 0 | 3 | 16 |
| 0 | 1 | 0 | 10 | 12 |
| 0 | 0 | 0 | 2 | 18 |
| 0 | 0 | 0 | 2 | 12 |
| 0 | 1 | 0 | 25 | 16 |
| 1 | 1 | 1 | 20 | 12 |

# Regression results

| Dependent Variable: EVAC | | |
|---|---|---|
| Variable | B | t-value |
| (Constant) | 0.190 | 2.121 |
| PETS | -0.137 | -5.296 |
| MOBLHOME | 0.337 | 8.963 |
| TENURE | -0.003 | -2.973 |
| EDUC | 0.003 | 0.424 |
| FLOYD | 0.198 | 8.147 |
| $R^2$ | 0.145 | |
| F-stat | 36.010 | |

# Regression results

Predicted value outside of the supposed [0,1]

## Descriptive Statistics

|  | N | Minimum | Maximum | Mean | Std. Deviation |
|---|---|---|---|---|---|
| Unstandardized Predicted Value | 1070 | -.0849 | .7602 | .24299 | .16325 |
| Valid N (listwise) | 1070 | | | | |

16.66

# Logistic Regression Model

The "logit" model solves these problems:

$$\ln[p/(1-p)] = \alpha + \beta X + e$$

- p is the probability that the event Y occurs, p(Y=1)
- p/(1-p) is the "odds ratio"
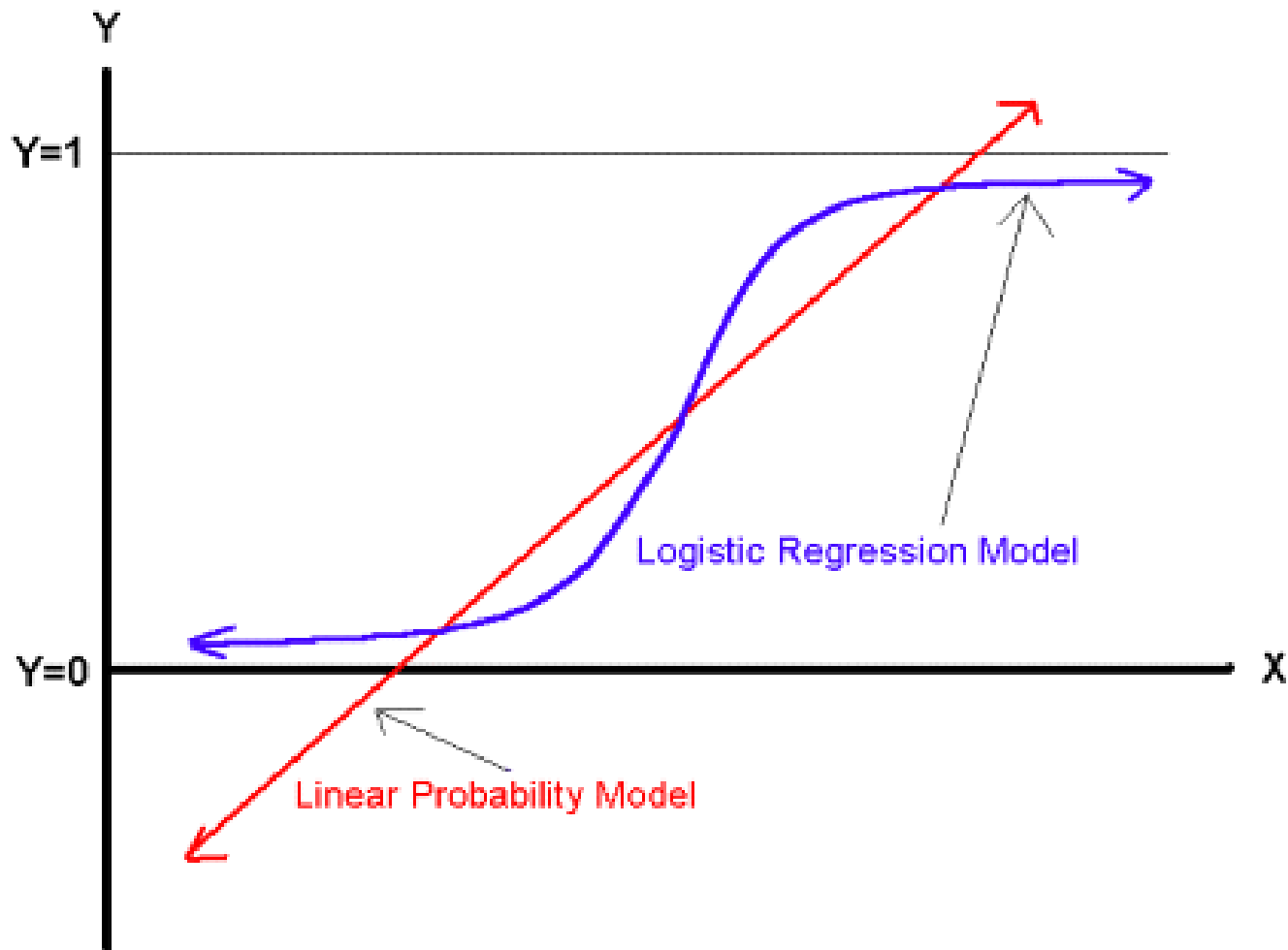- ln[p/(1-p)] is the log odds ratio, or "logit"

# Logistic Regression Model

- The logistic distribution constrains the estimated probabilities to lie between 0 and 1.

- The estimated probability is:

$$p = 1/[1 + \exp(-\alpha - \beta X)]$$

- if you let $\alpha + \beta X = 0$, then p = .50
- as $\alpha + \beta X$ gets really big, p approaches 1
- as $\alpha + \beta X$ gets really small, p approaches 0

Comparing the LP and Logit Models

# Maximum Likelihood Estimation (MLE)

- MLE is a statistical method for estimating the coefficients of a model.

- The likelihood function (L) measures the probability of observing the particular set of dependent variable values ($p_1$, $p_2$, ..., $p_n$) that occur in the sample:

  $$L = Prob (p_1 * p_2 * * * p_n)$$

- The higher the L, the higher the probability of observing the ps in the sample.

# Maximum Likelihood Estimation (MLE)

- MLE involves finding the coefficients ($\alpha$, $\beta$) that makes the log of the likelihood function (LL < 0) as large as possible

- Or, finds the coefficients that make -2 times the log of the likelihood function (-2LL) as small as possible

- The maximum likelihood estimates solve the following condition:

$$\{Y - p(Y=1)\}X_i = 0$$

summed over all observations, i = 1,…,n

# Interpreting Coefficients

$$\ln[p/(1-p)] = \alpha + \beta X + e$$

- The slope coefficient ($\beta$) is interpreted as the rate of change in the "log odds" as X changes … not very useful.

- An interpretation of the logit coefficient which is usually more intuitive is the "odds ratio"

# Odds ratio

- Since:

$$[p/(1-p)] = \exp(\alpha + \beta X)$$

$\exp(\beta)$ is the effect of the independent variable on the "odds ratio"

# Odds ratio (from R output)

| Variable | B | Exp(B) | 1/Exp(B) |
|----------|------|--------|----------|
| PETS | -0.6593 | 0.5172 | 1.933 |
| MOBLHOME | 1.5583 | 4.7508 | |
| TENURE | -0.0198 | 0.9804 | 1.020 |
| EDUC | 0.0501 | 1.0514 | |
| Constant | -0.916 | | |

"Households with pets are 1.933 times more likely to evacuate than those without pets."

# test statistics

The Wald statistic for the $\beta$ coefficient is:

$$\text{Wald} = [\beta / se_{\beta}]^2$$

which is distributed chi-square with 1 degree of freedom.

# Model output

| Variable | B | S.E. | Wald | R | Sig | t-value |
|---|---|---|---|---|---|---|
| PETS | -0.6593 | 0.2012 | 10.732 | -0.1127 | 0.0011 | -3.28 |
| MOBLHOME | 1.5583 | 0.2874 | 29.39 | 0.1996 | 0 | 5.42 |
| TENURE | -0.0198 | 0.008 | 6.1238 | -0.0775 | 0.0133 | -2.48 |
| EDUC | 0.0501 | 0.0468 | 1.1483 | 0.0000 | 0.2839 | 1.07 |
| Constant | -0.916 | 0.69 | 1.7624 | 1 | 0.1843 | -1.33 |