# Mathematics Learning Support Centre

# Statistics: 1.3 The Chi-squared test for two-way tables

Rosie Shier. 2004.

## 1 Introduction

If we have two categorical variables we can look at the relationship between these variables by putting the data in a two-way table.

Example

A sample of 200 components is selected from the output of a factory that uses three different machines to manufacture these components. Each component in the sample is inspected to determine whether or not it is defective. The machine that produced the component is also recorded. The results are as follows:

|              | Machine |   |   |       |
| ------------ | ------- | - | - | ----- |
| Outcome      | A       | B | C | Total |
| Defective    | 8 (12.9%) | 6 (8.8%) | 12 (17.1%) | 26 (13.0%) |
| Non-defective | 54     | 62 | 58 | 174 |
| Total        | 62      | 68 | 70 | 200 |

The manager wishes to determine whether or not there is a relationship between the proportion of defectives and the machine used.

In general, the null hypothesis in a two-way table is that there is <u>no association</u> between the row variable and the column variable and the alternative hypothesis is that there is an association. In this case, this is equivalent to saying:

$H_0$: There are no differences between machines in the percentage of defectives produced.
$H_1$: There are differences . . .

To test the null hypothesis we compare the observed cell counts with **expected** cell counts calculated under the assumption that the null hypothesis is true. If the null hypothesis were true the row (or column) percentages would all be the same. Therefore:

$$
\begin{aligned}
\text{Expected cell count} \quad &= \quad \frac{\text{Row percentage}}{100} \times \text{Column total} \\
&= \quad \frac{\text{Row total}}{\text{Overall total}} \times \text{Column Total} \\
&= \quad \frac{\text{Row total} \times \text{Column Total}}{n}
\end{aligned}
$$

where $n$ is the overall total.

Looking at our example – expected cell counts:

| Outcome | Machine | | |
|---|---|---|---|
| | A | B | C |
| Defective | $\frac{26\times62}{200} = 8.06$ | $\frac{26\times68}{200} = 8.84$ | $\frac{26\times70}{200} = 9.10$ |
| Non-defective | $\frac{174\times62}{200} = 53.94$ | $\frac{174\times68}{200} = 59.16$ | $\frac{174\times70}{200} = 60.90$ |

Notice that (as specified by the null hypothesis) the expected row percentages are all the same (i.e. $\frac{8.06}{62} = \frac{8.84}{68} = \frac{9.10}{70} = 13\%$ = overall percentage of defectives).

# 2 Carrying out the chi-squared test

To test the null hypothesis we now compute a statistic that compares the entire set of observed counts with the set of expected counts. This statistic is called the **chi-squared statistic** and is given by:

$$\chi^2 = \sum \frac{(O-E)^2}{E} \qquad \text{where O = observed cell count \& E = expected cell count}$$

and the sum is over all $r \times c$ cells in the table, where $r$ = number of rows, $c$ = number of columns in the table.

In our example:

$$\chi^2 = \frac{(8-8.06)^2}{8.06} + \frac{(6-8.84)^2}{8.84} + \ldots + \frac{(62-59.16)^2}{59.16} + \frac{(58-60.90)^2}{60.90} = 2.11 \quad (3 \text{ d.p.})$$

The number of degrees of freedom is simply $(r-1) \times (c-1)$.
In our case this is $(3-1) \times (2-1) = 2$. The p-value for the chi-squared test in this case is $P(\chi^2_2 > 2.11)$.
Looking this up in tables of the chi-squared distribution gives $0.25 < p < 0.5$.
In this case we have no real evidence that the percentage of defectives varies from machine to machine.

# 3 Validity of chi-squared tests for two-way tables

Chi-squared tests are only valid when you have a reasonable sample size. The following guidelines can be used:

1. For 2 x 2 tables:

   - If the total sample size is greater than 40, $\chi^2$ can be used.
   - If the total sample size si between 20 and 40 and the smallest expected frequency is at least 5, $\chi^2$ can be used.
   - Otherwise Fisher's exact test must be used.

2. For other tables:

   - $\chi^2$ can be used if no more than 20% of the expected frequencies are less than 5 <u>and</u> none is less than 1.

# 4   Carrying out a chi-squared test in SPSS

Your data could be in one of two formats:

1. Individual data

| Machine | Outcome |
|:---:|:---:|
| A | 0 |
| A | 1 |
| B | 1 |
| C | 0 |
| A | 1 |
| B | 0 |
| ⋮ | ⋮ |

Where 0 stands for defective and 1 for non-defective. In this case, you would have 200 lines of data, one for each component.

2. Grouped data

| Machine | Outcome | Frequency |
|:---:|:---:|:---:|
| A | 0 | 8 |
| A | 1 | 54 |
| B | 0 | 6 |
| B | 1 | 62 |
| C | 0 | 12 |
| C | 1 | 58 |

In this case, you only have $r \times c = 6$ lines of data, one for each cell in the table.

If you have <u>grouped data</u>, you need to do the following before carrying out the chi-squared test:
— **Data**
— **Weight Cases**
— Select **Weight cases by** and choose your frequency variable as the **Frequency Variable**.

You are now ready to carry out the chi-squared test (the procedure is now the same whether you have individual or grouped data):

— **Analyze**
— **Descriptive Statistics**
— **Crosstabs**
— Choose your outcome variable as the **Row Variable** and your explanatory variable as the **Column Variable** (it actually doesn't matter which way round they are, but the commands below are based on them being this way around).
— Click on the **Cells** button and select **Column** under **Percentages**. (You can also ask SPSS to display expected cell counts by clicking on the appropriate button here.) Now click on **Continue**. — Click on the **Statistics** button and select **Chi-square** in the top left-hand corner. Now click on **Continue**.
— Click on **OK**

(See over for output)

Your output will look like this:

**Case Processing Summary**

| | Cases | | | | | |
|---|---|---|---|---|---|---|
| | Valid | | Missing | | Total | |
| | N | Percent | N | Percent | N | Percent |
| OUTCOME * MACHINE | 200 | 100.0% | 0 | .0% | 200 | 100.0% |

**OUTCOME * MACHINE Crosstabulation**

| | | MACHINE | | | |
|---|---|---|---|---|---|
| OUTCOME | | A | B | C | Total |
| Defective | Count | 8 | 6 | 12 | 26 |
| | % within MACHINE | 12.9% | 8.8% | 17.1% | 13.0% |
| Non-defective | Count | 54 | 62 | 58 | 174 |
| | % within MACHINE | 87.1% | 91.2% | 82.9% | 87.0% |
| Total | Count | 62 | 68 | 70 | 200 |
| | % within MACHINE | 100.0% | 100.0% | 100.0% | 100.0% |

**Chi-Square Tests**

| | Value | df | Asymp. Sig. (2-sided) |
|---|---|---|---|
| Pearson Chi-Square | 2.112 | 2 | .348 |
| Likelihood Ratio | 2.144 | 2 | .342 |
| Linear-by-Linear Association | .585 | 1 | .444 |
| N of Valid Cases | 200 | | |

a. 0 cells (.0%) have expected count less than 5.
The minimum expected count is 8.06.

We are interested in the top row of the last table — i.e. the **Pearson Chi-Square**. The last column in the table gives the p-value — in this case we have $p = 0.348$. We would normally round this to $p = 0.3$.