

# Chi-square tests



## Overview

Chi-square ( $\chi^2$ ) tests are used to analyze categorical data. There are two main types:

Test	Purpose
Goodness of Fit	Does the data fit a specified distribution?
Test of Independence	Are two categorical variables independent?

## The chi-square test statistic

For both tests, the test statistic is:

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

where  $O$  = observed frequency and  $E$  = expected frequency.

**Interpretation:** Large values of  $\chi^2$  indicate a poor fit between observed and expected values, leading to rejection of  $H_0$ .

## Goodness of fit test

**Use when:** Testing whether observed data follows a hypothesized distribution.

**Hypotheses:**

- $H_0$ : The data follows the specified distribution
- $H_a$ : The data does not follow the specified distribution

**Degrees of freedom:**  $df = k - 1$  where  $k$  is the number of categories.

**Example:** A die is rolled 120 times with the following results. Is the die fair?

Outcome	1	2	3	4	5	6
Observed ( $O$ )	25	17	15	23	24	16
Expected ( $E$ )	20	20	20	20	20	20

If the die is fair, each outcome has probability  $\frac{1}{6}$ , so  $E = 120 \times \frac{1}{6} = 20$  for each.

$$\begin{aligned}\chi^2 &= \frac{(25 - 20)^2}{20} + \frac{(17 - 20)^2}{20} + \frac{(15 - 20)^2}{20} + \frac{(23 - 20)^2}{20} + \frac{(24 - 20)^2}{20} + \frac{(16 - 20)^2}{20} \\ &= \frac{25 + 9 + 25 + 9 + 16 + 16}{20} = \frac{100}{20} = 5.0\end{aligned}$$

Student Learning Support, Teaching and Learning Centre

[studentlearning@ontariotechu.ca](mailto:studentlearning@ontariotechu.ca)  
[ontariotechu.ca/studentlearning](http://ontariotechu.ca/studentlearning)



This document is licensed under Attribution-NonCommercial 4.0 International (CC BY-NC 4.0).

Degrees of freedom:  $df = 6 - 1 = 5$

Critical value at  $\alpha = 0.05$ :  $\chi_{0.05,5}^2 = 11.07$

Since  $\chi^2 = 5.0 < 11.07$ , we fail to reject  $H_0$ . There is not enough evidence to conclude the die is unfair.

## Test of independence

**Use when:** Testing whether two categorical variables are independent (no association).

**Hypotheses:**

- $H_0$ : The variables are independent
- $H_a$ : The variables are not independent (they are associated)

**Expected frequencies:** For each cell in the table:

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

**Degrees of freedom:**  $df = (r - 1)(c - 1)$  where  $r = \text{rows}$ ,  $c = \text{columns}$ .

**Example:** A survey asks 400 people about exercise habits and stress levels. Is there an association?

**Observed counts:**

	Low Stress	Medium Stress	High Stress	Row Total
Exercises regularly	60	50	30	140
Does not exercise	40	80	140	260
Column Total	100	130	170	400

**Calculate expected frequencies:**

$$\begin{aligned} E_{11} &= \frac{140 \times 100}{400} = 35 & E_{12} &= \frac{140 \times 130}{400} = 45.5 & E_{13} &= \frac{140 \times 170}{400} = 59.5 \\ E_{21} &= \frac{260 \times 100}{400} = 65 & E_{22} &= \frac{260 \times 130}{400} = 84.5 & E_{23} &= \frac{260 \times 170}{400} = 110.5 \end{aligned}$$

**Expected counts:**

	Low Stress	Medium Stress	High Stress
Exercises regularly	35	45.5	59.5
Does not exercise	65	84.5	110.5

**Calculate  $\chi^2$ :**

$$\begin{aligned} \chi^2 &= \frac{(60 - 35)^2}{35} + \frac{(50 - 45.5)^2}{45.5} + \frac{(30 - 59.5)^2}{59.5} + \frac{(40 - 65)^2}{65} + \frac{(80 - 84.5)^2}{84.5} + \frac{(140 - 110.5)^2}{110.5} \\ &= \frac{625}{35} + \frac{20.25}{45.5} + \frac{870.25}{59.5} + \frac{625}{65} + \frac{20.25}{84.5} + \frac{870.25}{110.5} \\ &= 17.86 + 0.45 + 14.63 + 9.62 + 0.24 + 7.88 = 50.68 \end{aligned}$$

Student Learning Support, Teaching and Learning Centre

[studentlearning@ontariotechu.ca](mailto:studentlearning@ontariotechu.ca)

[ontariotechu.ca/studentlearning](http://ontariotechu.ca/studentlearning)



This document is licensed under Attribution-NonCommercial 4.0 International (CC BY-NC 4.0).

Degrees of freedom:  $df = (2 - 1)(3 - 1) = 2$

Critical value at  $\alpha = 0.05$ :  $\chi_{0.05,2}^2 = 5.99$

Since  $\chi^2 = 50.68 > 5.99$ , we reject  $H_0$ . There is significant evidence of an association between exercise and stress levels.

## Conditions for chi-square tests

- Data must be counts (frequencies), not percentages
- Observations must be independent
- Expected frequency in each cell should be at least 5
- If expected frequencies are too small, combine categories

## Chi-square critical values (right-tail)

$df$	$\alpha = 0.10$	$\alpha = 0.05$	$\alpha = 0.025$	$\alpha = 0.01$
1	2.706	3.841	5.024	6.635
2	4.605	5.991	7.378	9.210
3	6.251	7.815	9.348	11.345
4	7.779	9.488	11.143	13.277
5	9.236	11.070	12.833	15.086
6	10.645	12.592	14.449	16.812
7	12.017	14.067	16.013	18.475
8	13.362	15.507	17.535	20.090
9	14.684	16.919	19.023	21.666
10	15.987	18.307	20.483	23.209

**Decision rule:** Reject  $H_0$  if  $\chi^2 > \chi_{\alpha,df}^2$  (critical value from table).

**Note:** Chi-square tests are always right-tailed because we're looking for large deviations from expected values.