

Tutorials and worked examples for simulation, curve fitting, statistical analysis, and plotting. https://simfit.org.uk https://simfit.silverfrost.com

The chi-square test on observed and expected frequencies is based on forming a test statistic that, in the limit of a very large sample size, becomes asymptotic to a chi-square distribution, with the number of degrees of freedom dependent on the number of categories, and also on the number of parameters estimated from the sample.

To be precise, it is assumed that the user has counted the frequency of occurrence of k observations partitioned into n categories with O_i in category i, and also knows the frequencies E_i expected under the null hypothesis that the observations are consistent with the expected frequencies given by the assumed distribution. This allows the calculation of C defined as

$$C = \frac{(O_1 - E_1)^2}{E_1} + \frac{(O_2 - E_2)^2}{E_2} + \dots + \frac{(O_n - E_n)^2}{E_n}$$

which has an approximate chi-square distribution with ν degrees of freedom given by

$$v = n - 1 - m$$

where $v \ge 2$, and *m* is the number of parameters estimated from the sample. The reason for subtracting 1 + m from *n* to get the degrees of freedom will be clear by considering the identity

$$\sum_{i=1}^{n} E_i = \sum_{i=1}^{n} O_i = k$$

which reduces the effective numbers of terms in the calculation of *C* to n - 1. Similarly, every further equation of constraint can be considered to reduce by one the effective number of terms in the calculation of *C*.

It is usually recommended that the expected values are at least 5 and, if this cannot be realized, then categories could be combined until this condition is met. Alternatively, if the total number of observations is k as above, and the number of categories is not fixed by other considerations, then the number of bins n used to partition the data is sometimes suggested as

$$n \approx k^{0.4}$$
,

but obviously this all depends on the shape of the assumed distribution.

To illustrate this test consider the next table, which records the results from one hundred observations on the number of heads resulting from tossing five different coins. Clearly there are six categories, as the number of heads per toss of the five coins can only be 0, 1, 2, 3, 4, or 5, but note that $100^{0.4} \approx 6$ anyway in this case.

Number of Heads	0	1	2	3	4	5
Observed	3	16	36	32	11	2
Expected	4.0	17.9	32.6	29.6	13.5	2.4

There were 238 heads in all from the total of 500 tosses, so the expected frequencies were calculated using a binomial distribution with binomial N = 5 and estimated parameter $\hat{p} = 238/500 = 0.476$, and therefore 4 degrees of freedom.

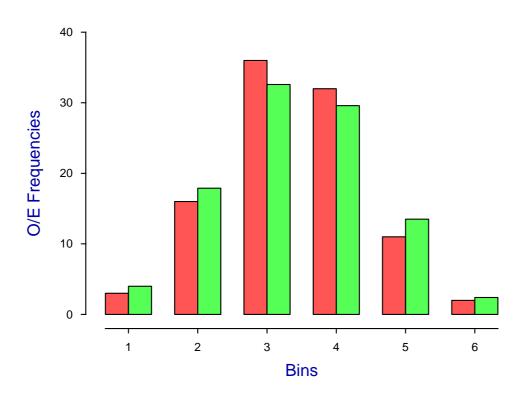
Choose [A/Z] from the main SIMFIT menu, open program **chisqd**, select chi-square test on observed and expected frequencies, then analyze the above data contained in the test files chisqd.tf2 and chisqd.tf3, with one parameter estimated from the sample to get these results.

Number of partitions (bins)	6	
Number of degrees of freedom	4	
Chi-square test statistic C	1.531	
$P(\chi^2 \ge C)$	0.8212	Consider accepting H0
Upper tail 5% critical point	9.488	
Upper tail 1% critical point	13.28	

SIMFIT first displays a warning that the expected frequencies for 0 and 5 heads are below 5, and so these two categories could be combined if it was thought necessary. However, in this case the p value of 0.8212 is much larger than 0.05, so the conclusion is that the null hypothesis of a binomial distribution with parameters N = 5, and p = 0.476 cannot be rejected.

Note that $SIMF_IT$ also lists the 1% and 5% upper tail critical points, as this is how the test results were analyzed in the past by looking up tables of critical points, before the availability of computers made this unnecessary.

The most widely used technique to display the agreement between the observed and expected frequencies is a bar chart, as in the next figure.



Observed and Expected Frequencies