

chi-en

Contents

Articles

Chi-squared test	1
Chi-squared distribution	2
Pearson's chi-squared test	10
Statistics	14

References

Article Sources and Contributors	25
Image Sources, Licenses and Contributors	26

Article Licenses

License	27
---------	----

Chi-squared test

A **chi-squared test**, also referred to as **chi-square test** or χ^2 test, is any statistical hypothesis test in which the sampling distribution of the test statistic is a chi-squared distribution when the null hypothesis is true. Also considered a chi-squared test is a test in which this is *asymptotically* true, meaning that the sampling distribution (if the null hypothesis is true) can be made to approximate a chi-squared distribution as closely as desired by making the sample size large enough.

Some examples of chi-squared tests where the chi-squared distribution is only approximately valid:

- Pearson's chi-squared test, also known as the chi-squared goodness-of-fit test or chi-squared test for independence. When the chi-squared test is mentioned without any modifiers or without other precluding context, this test is usually meant (for an exact test used in place of χ^2 , see Fisher's exact test).
- Yates's correction for continuity, also known as Yates' chi-squared test.
- Cochran–Mantel–Haenszel chi-squared test.
- McNemar's test, used in certain 2×2 tables with pairing
- Tukey's test of additivity
- The portmanteau test in time-series analysis, testing for the presence of autocorrelation
- Likelihood-ratio tests in general statistical modelling, for testing whether there is evidence of the need to move from a simple model to a more complicated one (where the simple model is nested within the complicated one).

One case where the distribution of the test statistic is an exact chi-squared distribution is the test that the variance of a normally distributed population has a given value based on a sample variance. Such a test is uncommon in practice because values of variances to test against are seldom known exactly.

Chi-squared test for variance in a normal population

If a sample of size n is taken from a population having a normal distribution, then there is a result (see distribution of the sample variance) which allows a test to be made of whether the variance of the population has a pre-determined value. For example, a manufacturing process might have been in stable condition for a long period, allowing a value for the variance to be determined essentially without error. Suppose that a variant of the process is being tested, giving rise to a small sample of n product items whose variation is to be tested. The test statistic T in this instance could be set to be the sum of squares about the sample mean, divided by the nominal value for the variance (i.e. the value to be tested as holding). Then T has a chi-squared distribution with $n - 1$ degrees of freedom. For example if the sample size is 21, the acceptance region for T for a significance level of 5% is the interval 9.59 to 34.17.

References

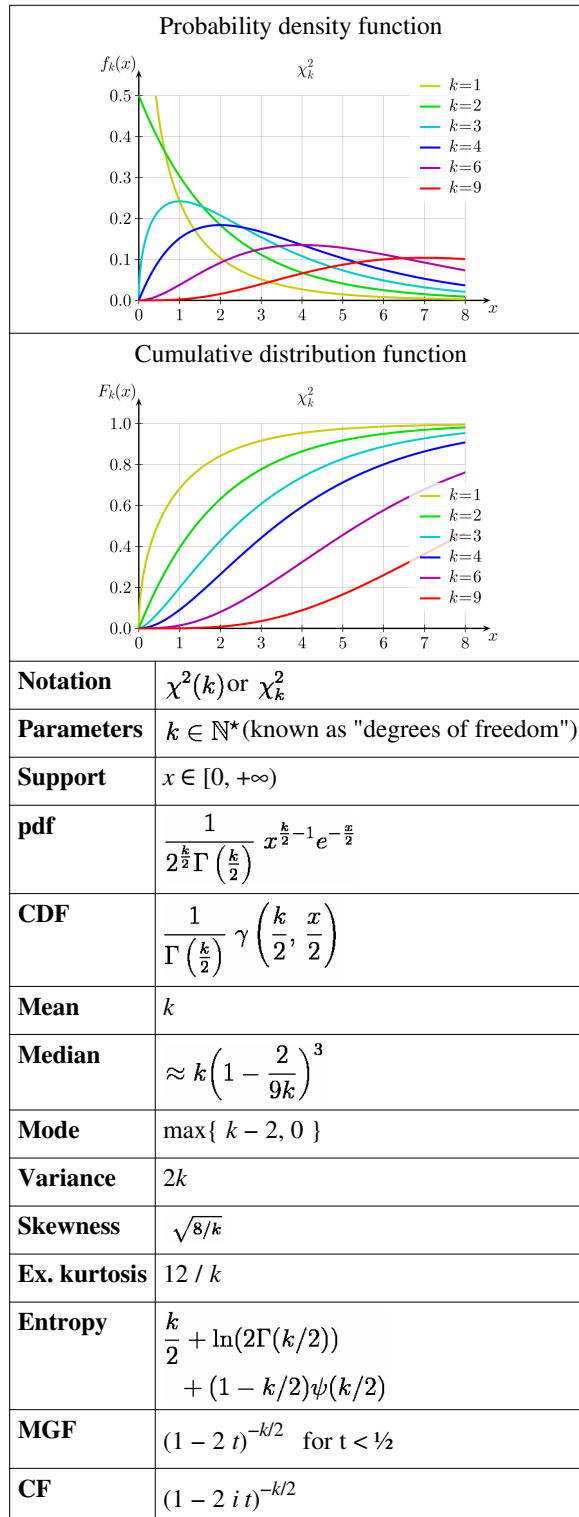
- Weisstein, Eric W., "Chi-Squared Test"^[1], *MathWorld*.
- Corder, G.W., Foreman, D.I. (2009). *Nonparametric Statistics for Non-Statisticians: A Step-by-Step Approach* Wiley, ISBN 978-0-470-45461-9
- Greenwood, P.E., Nikulin, M.S. (1996) *A guide to chi-squared testing*. Wiley, New York. ISBN 0-471-55779-X
- Nikulin, M.S. (1973). "Chi-squared test for normality". In: *Proceedings of the International Vilnius Conference on Probability Theory and Mathematical Statistics*, v.2, pp. 119–122.
- Bagdonavicius, V., Nikulin, M.S. (2011) "Chi-square goodness-of-fit test for right censored data". *The International Journal of Applied Mathematics and Statistics*, p. 30-50. Wikipedia:Citing sources#What information to include

References

[1] <http://mathworld.wolfram.com/Chi-SquaredTest.html>

Chi-squared distribution

This article is about the mathematics of the chi-squared distribution. For its uses in statistics, see chi-squared test. For the music group, see Chi2 (band).



In probability theory and statistics, the **chi-squared distribution** (also **chi-square** or **χ^2 -distribution**) with k degrees of freedom is the distribution of a sum of the squares of k independent standard normal random variables. A special case of the gamma distribution, it is one of the most widely used probability distributions in inferential statistics, e.g., in hypothesis testing or in construction of confidence intervals.^[1] When there is a need to contrast it with the noncentral chi-squared distribution, this distribution is sometimes called the **central chi-squared distribution**.

The chi-squared distribution is used in the common chi-squared tests for goodness of fit of an observed distribution to a theoretical one, the independence of two criteria of classification of qualitative data, and in confidence interval estimation for a population standard deviation of a normal distribution from a sample standard deviation. Many other statistical tests also use this distribution, like Friedman's analysis of variance by ranks.

History and name

This distribution was first described by the German statistician Friedrich Robert Helmert in papers of 1875/1876,^{[2][3]} where he computed the sampling distribution of the sample variance of a normal population. Thus in German this was traditionally known as the *Helmertsche* ("Helmertian") or "Helmert distribution".

The distribution was independently rediscovered by the English mathematician Karl Pearson in the context of goodness of fit, for which he developed his Pearson's chi-squared test, published in (Pearson 1900), with computed table of values published in (Elderton 1902), collected in (Pearson 1914, pp. xxxi–xxxiii, 26–28, Table XII). The name "chi-squared" ultimately derives from Pearson's shorthand for the exponent in a multivariate normal distribution with the Greek letter Chi, writing $-1/2\chi^2$ for what would appear in modern notation as $-1/2\mathbf{x}^T\boldsymbol{\Sigma}^{-1}\mathbf{x}$ ($\boldsymbol{\Sigma}$ being the covariance matrix).^[4] The idea of a family of "chi-squared distributions", however, is not due to Pearson but arose as a further development due to Fisher in the 1920s.^[2]

Definition

If Z_1, \dots, Z_k are independent, standard normal random variables, then the sum of their squares,

$$Q = \sum_{i=1}^k Z_i^2,$$

is distributed according to the **chi-squared distribution** with k degrees of freedom. This is usually denoted as

$$Q \sim \chi^2(k) \text{ or } Q \sim \chi_k^2.$$

The chi-squared distribution has one parameter: k — a positive integer that specifies the number of degrees of freedom (i.e. the number of Z_i 's)

Characteristics

Further properties of the chi-squared distribution can be found in the box at the upper right corner of this article.

Probability density function

The probability density function (pdf) of the chi-squared distribution is

$$f(x; k) = \begin{cases} \frac{x^{(k/2-1)} e^{-x/2}}{2^{k/2} \Gamma(\frac{k}{2})}, & x \geq 0; \\ 0, & \text{otherwise.} \end{cases}$$

where $\Gamma(k/2)$ denotes the Gamma function, which has closed-form values for integer k .

For derivations of the pdf in the cases of one, two and k degrees of freedom, see Proofs related to chi-squared distribution.

Differential equation

$$\left\{ 2x f'(x) + f(x)(-k + x + 2) = 0, f(1) = \frac{2^{-k/2}}{\sqrt{e}\Gamma\left(\frac{k}{2}\right)} \right\}$$

Cumulative distribution function

Its cumulative distribution function is:

$$F(x; k) = \frac{\gamma\left(\frac{k}{2}, \frac{x}{2}\right)}{\Gamma\left(\frac{k}{2}\right)} = P\left(\frac{k}{2}, \frac{x}{2}\right),$$

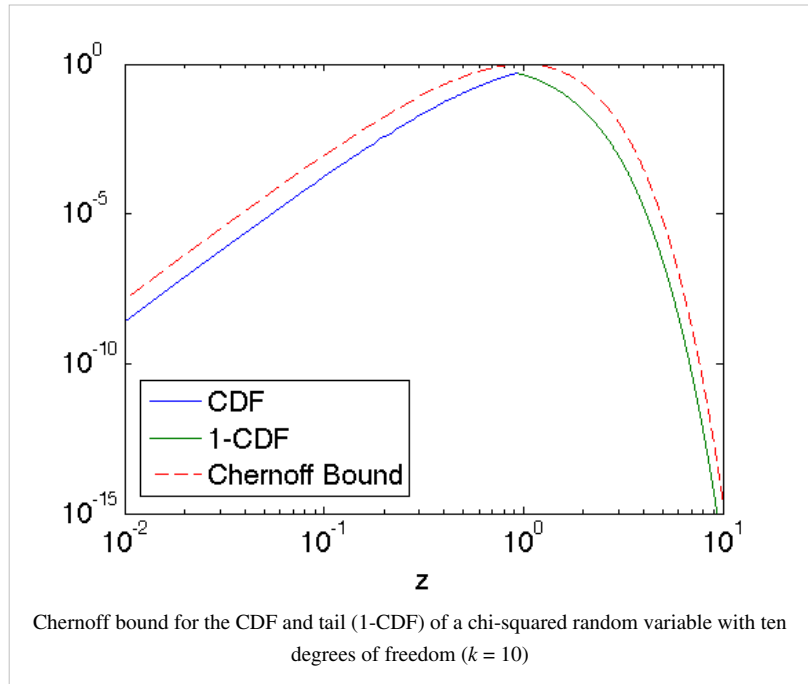
where $\gamma(s,t)$ is the lower incomplete Gamma function and $P(s,t)$ is the regularized Gamma function.

In a special case of $k = 2$ this function has a simple form:

$$F(x; 2) = 1 - e^{-\frac{x}{2}}$$

and the form is not much more complicated for other small even k .

Tables of the chi-squared cumulative distribution function are widely available and the function is included in many spreadsheets and all statistical packages.



Letting $z \equiv x/k$, Chernoff bounds on the lower and upper tails of the CDF may be obtained. For the cases when $0 < z < 1$ (which include all of the cases when this CDF is less than half):

$$F(zk; k) \leq (ze^{1-z})^{k/2}.$$

The tail bound for the cases when $z > 1$, similarly, is

$$1 - F(zk; k) \leq (ze^{1-z})^{k/2}.$$

For another approximation for the CDF modeled after the cube of a Gaussian, see under Noncentral chi-squared distribution.

Additivity

It follows from the definition of the chi-squared distribution that the sum of independent chi-squared variables is also chi-squared distributed. Specifically, if $\{X_i\}_{i=1}^n$ are independent chi-squared variables with $\{k_i\}_{i=1}^n$ degrees of freedom, respectively, then $Y = X_1 + \dots + X_n$ is chi-squared distributed with $k_1 + \dots + k_n$ degrees of freedom.

Sample mean

The sample mean of n i.i.d. chi-squared variables of degree k is distributed according to a gamma distribution with shape α and scale θ parameters:

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim \text{Gamma}(\alpha = n \cdot k/2, \theta = 2/n) \quad \text{where} \quad X_i \sim \chi^2(k)$$

Asymptotically, given that for a scale parameter α going to infinity, a Gamma distribution converges towards a Normal distribution with expectation $\mu = k\theta$ and variance $\sigma^2 = k\theta^2$, the sample mean converges towards:

$$\bar{X} \xrightarrow{n \rightarrow \infty} N(k, 2 \cdot k/n)$$

Note that we would have obtained the same result invoking instead the central limit theorem, noting that the expectation of the χ^2 is k , and its variance $2k$ (and hence the variance of the sample mean being $2k/n$).

Entropy

The differential entropy is given by

$$h = \int_{-\infty}^{\infty} f(x; k) \ln f(x; k) dx = \frac{k}{2} + \ln \left[2\Gamma\left(\frac{k}{2}\right) \right] + \left(1 - \frac{k}{2}\right) \psi\left[\frac{k}{2}\right],$$

where $\psi(x)$ is the Digamma function.

The chi-squared distribution is the maximum entropy probability distribution for a random variate X for which $E(X) = k$ and $E(\ln(X)) = \psi(k/2) + \log(2)$ are fixed. Since the chi-squared is in the family of gamma distributions, this can be derived by substituting appropriate values in the Expectation of the Log moment of Gamma. For derivation from more basic principles, see the derivation in moment generating function of the sufficient statistic.

Noncentral moments

The moments about zero of a chi-squared distribution with k degrees of freedom are given by^{[5][6]}

$$E(X^m) = k(k+2)(k+4)\cdots(k+2m-2) = 2^m \frac{\Gamma(m + \frac{k}{2})}{\Gamma(\frac{k}{2})}.$$

Cumulants

The cumulants are readily obtained by a (formal) power series expansion of the logarithm of the characteristic function:

$$\kappa_n = 2^{n-1}(n-1)! k$$

Asymptotic properties

By the central limit theorem, because the chi-squared distribution is the sum of k independent random variables with finite mean and variance, it converges to a normal distribution for large k . For many practical purposes, for $k > 50$ the distribution is sufficiently close to a normal distribution for the difference to be ignored. Specifically, if $X \sim \chi^2(k)$, then as k tends to infinity, the distribution of $(X - k)/\sqrt{2k}$ tends to a standard normal distribution. However, convergence is slow as the skewness is $\sqrt{8/k}$ and the excess kurtosis is $12/k$.

- The sampling distribution of $\ln(\chi^2)$ converges to normality much faster than the sampling distribution of χ^2 , as the logarithm removes much of the asymmetry. Other functions of the chi-squared distribution converge more rapidly to a normal distribution. Some examples are:
- If $X \sim \chi^2(k)$ then $\sqrt{2X}$ is approximately normally distributed with mean $\sqrt{2k-1}$ and unit variance (result credited to R. A. Fisher).
- If $X \sim \chi^2(k)$ then $\sqrt[3]{X/k}$ is approximately normally distributed with mean $1-2/(9k)$ and variance $2/(9k)$. This is known as the Wilson–Hilferty transformation.

Relation to other distributions

- As $k \rightarrow \infty$, $(\chi_k^2 - k)/\sqrt{2k} \xrightarrow{d} N(0, 1)$ (normal distribution)
- $\chi_k^2 \sim \chi_k^2(0)$ (Noncentral chi-squared distribution with non-centrality parameter $\lambda = 0$)
- If $X \sim F(\nu_1, \nu_2)$ then $Y = \lim_{\nu_2 \rightarrow \infty} \nu_1 X$ has the chi-squared distribution $\chi_{\nu_1}^2$
- As a special case, if $X \sim F(1, \nu_2)$ then $Y = \lim_{\nu_2 \rightarrow \infty} X$ has the chi-squared distribution χ_1^2
- $\|\mathbf{N}_{i=1, \dots, k}(0, 1)\|^2 \sim \chi_k^2$ (The squared norm of \mathbf{k} standard normally distributed variables is a chi-squared distribution with \mathbf{k} degrees of freedom)
- If $X \sim \chi^2(\nu)$ and $c > 0$, then $cX \sim \Gamma(k = \nu/2, \theta = 2c)$. (gamma distribution)
- If $X \sim \chi_k^2$ then $\sqrt{X} \sim \chi_k$ (chi distribution)
- If $X \sim \chi^2(2)$, then $X \sim \text{Exp}(1/2)$ is an exponential distribution. (See Gamma distribution for more.)
- If $X \sim \text{Rayleigh}(1)$ (Rayleigh distribution) then $X^2 \sim \chi^2(2)$
- If $X \sim \text{Maxwell}(1)$ (Maxwell distribution) then $X^2 \sim \chi^2(3)$
- If $X \sim \chi^2(\nu)$ then $\frac{1}{X} \sim \text{Inv-}\chi^2(\nu)$ (Inverse-chi-squared distribution)
- The chi-squared distribution is a special case of type 3 Pearson distribution
- If $X \sim \chi^2(\nu_1)$ and $Y \sim \chi^2(\nu_2)$ are independent then $\frac{X}{X+Y} \sim \text{Beta}(\frac{\nu_1}{2}, \frac{\nu_2}{2})$ (beta distribution)
- If $X \sim U(0, 1)$ (uniform distribution) then $-2 \log(X) \sim \chi^2(2)$
- $\chi^2(6)$ is a transformation of Laplace distribution
- If $X_i \sim \text{Laplace}(\mu, \beta)$ then $\sum_{i=1}^n \frac{2|X_i - \mu|}{\beta} \sim \chi^2(2n)$
- chi-squared distribution is a transformation of Pareto distribution
- Student's t-distribution is a transformation of chi-squared distribution
- Student's t-distribution can be obtained from chi-squared distribution and normal distribution
- Noncentral beta distribution can be obtained as a transformation of chi-squared distribution and Noncentral chi-squared distribution
- Noncentral t-distribution can be obtained from normal distribution and chi-squared distribution

A chi-squared variable with k degrees of freedom is defined as the sum of the squares of k independent standard normal random variables.

If Y is a k -dimensional Gaussian random vector with mean vector μ and rank k covariance matrix C , then $X = (Y-\mu)^T C^{-1} (Y-\mu)$ is chi-squared distributed with k degrees of freedom.

The sum of squares of statistically independent unit-variance Gaussian variables which do *not* have mean zero yields a generalization of the chi-squared distribution called the noncentral chi-squared distribution.

If Y is a vector of k i.i.d. standard normal random variables and A is a $k \times k$ idempotent matrix with rank $k-n$ then the quadratic form $Y^T A Y$ is chi-squared distributed with $k-n$ degrees of freedom.

The chi-squared distribution is also naturally related to other distributions arising from the Gaussian. In particular,

- Y is F-distributed, $Y \sim F(k_1, k_2)$ if $Y = \frac{X_1/k_1}{X_2/k_2}$ where $X_1 \sim \chi^2(k_1)$ and $X_2 \sim \chi^2(k_2)$ are statistically independent.
- If X is chi-squared distributed, then \sqrt{X} is chi distributed.
- If $X_1 \sim \chi^2_{k_1}$ and $X_2 \sim \chi^2_{k_2}$ are statistically independent, then $X_1 + X_2 \sim \chi^2_{k_1+k_2}$. If X_1 and X_2 are not independent, then $X_1 + X_2$ is not chi-squared distributed.

Generalizations

The chi-squared distribution is obtained as the sum of the squares of k independent, zero-mean, unit-variance Gaussian random variables. Generalizations of this distribution can be obtained by summing the squares of other types of Gaussian random variables. Several such distributions are described below.

Linear combination

If X_1, \dots, X_n are iid chi square random variables and $a_1, \dots, a_n \in \mathbb{R}_{>0}$, then a closed expression for the distribution of $X = \sum_{i=1}^n a_i X_i$ is not known. However, there exists a computationally efficient algorithm to calculate the pdf f_X to arbitrary precision.

Chi-squared distributions

Noncentral chi-squared distribution

Main article: Noncentral chi-squared distribution

The noncentral chi-squared distribution is obtained from the sum of the squares of independent Gaussian random variables having unit variance and *nonzero* means.

Generalized chi-squared distribution

Main article: Generalized chi-squared distribution

The generalized chi-squared distribution is obtained from the quadratic form $z'Az$ where z is a zero-mean Gaussian vector having an arbitrary covariance matrix, and A is an arbitrary matrix.

Gamma, exponential, and related distributions

The chi-squared distribution $X \sim \chi^2(k)$ is a special case of the gamma distribution, in that $X \sim \Gamma(k/2, 1/2)$ using the rate parameterization of the gamma distribution (or $X \sim \Gamma(k/2, 2)$ using the scale parameterization of the gamma distribution) where k is an integer.

Because the exponential distribution is also a special case of the Gamma distribution, we also have that if $X \sim \chi^2(2)$, then $X \sim \text{Exp}(1/2)$ is an exponential distribution.

The Erlang distribution is also a special case of the Gamma distribution and thus we also have that if $X \sim \chi^2(k)$ with even k , then X is Erlang distributed with shape parameter $k/2$ and scale parameter $1/2$.

Applications

The chi-squared distribution has numerous applications in inferential statistics, for instance in chi-squared tests and in estimating variances. It enters the problem of estimating the mean of a normally distributed population and the problem of estimating the slope of a regression line via its role in Student's t-distribution. It enters all analysis of variance problems via its role in the F-distribution, which is the distribution of the ratio of two independent chi-squared random variables, each divided by their respective degrees of freedom.

Following are some of the most common situations in which the chi-squared distribution arises from a Gaussian-distributed sample.

- if X_1, \dots, X_n are i.i.d. $N(\mu, \sigma^2)$ random variables, then $\sum_{i=1}^n (X_i - \bar{X})^2 \sim \sigma^2 \chi_n^2$ where $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$.
- The box below shows some statistics based on $X_i \sim \text{Normal}(\mu_i, \sigma_i^2)$, $i = 1, \dots, k$, independent random variables that have probability distributions related to the chi-squared distribution:

Name	Statistic
chi-squared distribution	$\sum_{i=1}^k \left(\frac{X_i - \mu_i}{\sigma_i} \right)^2$
noncentral chi-squared distribution	$\sum_{i=1}^k \left(\frac{X_i}{\sigma_i} \right)^2$
chi distribution	$\sqrt{\sum_{i=1}^k \left(\frac{X_i - \mu_i}{\sigma_i} \right)^2}$
noncentral chi distribution	$\sqrt{\sum_{i=1}^k \left(\frac{X_i}{\sigma_i} \right)^2}$

Table of χ^2 value vs p-value

The p-value is the probability of observing a test statistic *at least* as extreme in a chi-squared distribution. Accordingly, since the cumulative distribution function (CDF) for the appropriate degrees of freedom (*df*) gives the probability of having obtained a value *less extreme* than this point, subtracting the CDF value from 1 gives the p-value. The table below gives a number of p-values matching to χ^2 for the first 10 degrees of freedom.

A low p-value indicates greater statistical significance, i.e. greater confidence that the observed deviation from the null hypothesis is significant. A p-value of 0.05 is often used as a bright-line cutoff between significant and not-significant results.

Degrees of freedom (df)	χ^2 value ^[7]										
1	0.004	0.02	0.06	0.15	0.46	1.07	1.64	2.71	3.84	6.64	10.83
2	0.10	0.21	0.45	0.71	1.39	2.41	3.22	4.60	5.99	9.21	13.82
3	0.35	0.58	1.01	1.42	2.37	3.66	4.64	6.25	7.82	11.34	16.27
4	0.71	1.06	1.65	2.20	3.36	4.88	5.99	7.78	9.49	13.28	18.47
5	1.14	1.61	2.34	3.00	4.35	6.06	7.29	9.24	11.07	15.09	20.52
6	1.63	2.20	3.07	3.83	5.35	7.23	8.56	10.64	12.59	16.81	22.46
7	2.17	2.83	3.82	4.67	6.35	8.38	9.80	12.02	14.07	18.48	24.32
8	2.73	3.49	4.59	5.53	7.34	9.52	11.03	13.36	15.51	20.09	26.12
9	3.32	4.17	5.38	6.39	8.34	10.66	12.24	14.68	16.92	21.67	27.88
10	3.94	4.87	6.18	7.27	9.34	11.78	13.44	15.99	18.31	23.21	29.59
P value (Probability)	0.95	0.90	0.80	0.70	0.50	0.30	0.20	0.10	0.05	0.01	0.001

References

- [1] NIST (2006). Engineering Statistics Handbook - Chi-Squared Distribution (<http://www.itl.nist.gov/div898/handbook/eda/section3/eda3666.htm>)
- [2] Hald 1998, pp. 633–692, 27. Sampling Distributions under Normality.
- [3] F. R. Helmert, " Ueber die Wahrscheinlichkeit der Potenzsummen der Beobachtungsfehler und über einige damit im Zusammenhange stehende Fragen (http://gdz.sub.uni-goettingen.de/dms/load/img/?PPN=PPN599415665_0021&DMDID=DMDLOG_0018)", *Zeitschrift für Mathematik und Physik* 21 (http://gdz.sub.uni-goettingen.de/dms/load/toc/?PPN=PPN599415665_0021), 1876, S. 102–219
- [4] R. L. Plackett, *Karl Pearson and the Chi-Squared Test*, International Statistical Review, 1983, 61f. (<http://www.jstor.org/stable/1402731?seq=3>) See also Jeff Miller, Earliest Known Uses of Some of the Words of Mathematics (<http://jeff560.tripod.com/c.html>).
- [5] Chi-squared distribution (<http://mathworld.wolfram.com/Chi-SquaredDistribution.html>), from MathWorld, retrieved Feb. 11, 2009
- [6] M. K. Simon, *Probability Distributions Involving Gaussian Random Variables*, New York: Springer, 2002, eq. (2.35), ISBN 978-0-387-34657-1
- [7] Chi-Squared Test (<http://www2.lv.psu.edu/jxm57/irp/chisquar.html>) Table B.2. Dr. Jacqueline S. McLaughlin at The Pennsylvania State University. In turn citing: R.A. Fisher and F. Yates, *Statistical Tables for Biological Agricultural and Medical Research*, 6th ed., Table IV

Further reading

- Hald, Anders (1998). *A history of mathematical statistics from 1750 to 1930*. New York: Wiley. ISBN 0-471-17912-4.
- Elderton, William Palin (1902). "Tables for Testing the Goodness of Fit of Theory to Observation". *Biometrika* **1** (2): 155–163. doi: 10.1093/biomet/1.2.155 (<http://dx.doi.org/10.1093/biomet/1.2.155>).

External links

- Hazewinkel, Michiel, ed. (2001), "Chi-squared distribution" (<http://www.encyclopediaofmath.org/index.php?title=p/c022100>), *Encyclopedia of Mathematics*, Springer, ISBN 978-1-55608-010-4
- Calculator for the pdf, cdf and quantiles of the chi-squared distribution (<http://calculus-calculator.com/statistics/chi-squared-distribution-calculator.html>)
- Earliest Uses of Some of the Words of Mathematics: entry on Chi squared has a brief history (<http://jeff560.tripod.com/c.html>)
- Course notes on Chi-Squared Goodness of Fit Testing (<http://www.stat.yale.edu/Courses/1997-98/101/chigf.htm>) from Yale University Stats 101 class.
- *Mathematica* demonstration showing the chi-squared sampling distribution of various statistics, e.g. $\sum x^2$, for a normal population (<http://demonstrations.wolfram.com/StatisticsAssociatedWithNormalSamples/>)
- Simple algorithm for approximating cdf and inverse cdf for the chi-squared distribution with a pocket calculator (<http://www.jstor.org/stable/2348373>)

Pearson's chi-squared test

Pearson's chi-squared test (χ^2) is a statistical test applied to sets of categorical data to evaluate how likely it is that any observed difference between the sets arose by chance. It is suitable for unpaired data from large samples. It is the most widely used of many chi-squared tests (Yates, likelihood ratio, portmanteau test in time series, etc.) – statistical procedures whose results are evaluated by reference to the chi-squared distribution. Its properties were first investigated by Karl Pearson in 1900. In contexts where it is important to improve a distinction between the test statistic and its distribution, names similar to *Pearson χ -squared* test or statistic are used.

It tests a null hypothesis stating that the frequency distribution of certain events observed in a sample is consistent with a particular theoretical distribution. The events considered must be mutually exclusive and have total probability 1. A common case for this is where the events each cover an outcome of a categorical variable. A simple example is the hypothesis that an ordinary six-sided die is "fair", i. e., all six outcomes are equally likely to occur.

Definition

Pearson's chi-squared test is used to assess two types of comparison: tests of goodness of fit and tests of independence.

- A test of **goodness of fit** establishes whether or not an observed frequency distribution differs from a theoretical distribution.
- A **test of independence** assesses whether paired observations on two variables, expressed in a contingency table, are independent of each other (e.g. polling responses from people of different nationalities to see if one's nationality is related to the response).

The procedure of the test includes the following steps:

1. Calculate the chi-squared test statistic, χ^2 , which resembles a normalized sum of squared deviations between observed and theoretical frequencies (see below).
2. Determine the degrees of freedom, df , of that statistic, which is essentially the number of frequencies reduced by the number of parameters of the fitted distribution.
3. Compare χ^2 to the critical value from the chi-squared distribution with df degrees of freedom, which in many cases gives a good approximation of the distribution of χ^2 .

Test for fit of a distribution

Discrete uniform distribution

In this case N observations are divided among n cells. A simple application is to test the hypothesis that, in the general population, values would occur in each cell with equal frequency. The "theoretical frequency" for any cell (under the null hypothesis of a discrete uniform distribution) is thus calculated as

$$E_i = \frac{N}{n},$$

and the reduction in the degrees of freedom is $p = 1$, notionally because the observed frequencies O_i are constrained to sum to N .

Other distributions

When testing whether observations are random variables whose distribution belongs to a given family of distributions, the "theoretical frequencies" are calculated using a distribution from that family fitted in some standard way. The reduction in the degrees of freedom is calculated as $p = s + 1$, where s is the number of co-variates used in fitting the distribution. For instance, when checking a three-co-variate Weibull distribution, $p = 4$, and when checking a normal distribution (where the parameters are mean and standard deviation), $p = 3$. In other words, there will be $n - p$ degrees of freedom, where n is the number of categories.

It should be noted that the degrees of freedom are not based on the number of observations as with a Student's t or F-distribution. For example, if testing for a fair, six-sided die, there would be five degrees of freedom because there are six categories/parameters (each number). The number of times the die is rolled will have absolutely no effect on the number of degrees of freedom.

Calculating the test-statistic

The value of the test-statistic is

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

where

χ^2 = Pearson's cumulative test statistic, which asymptotically approaches a χ^2 distribution.

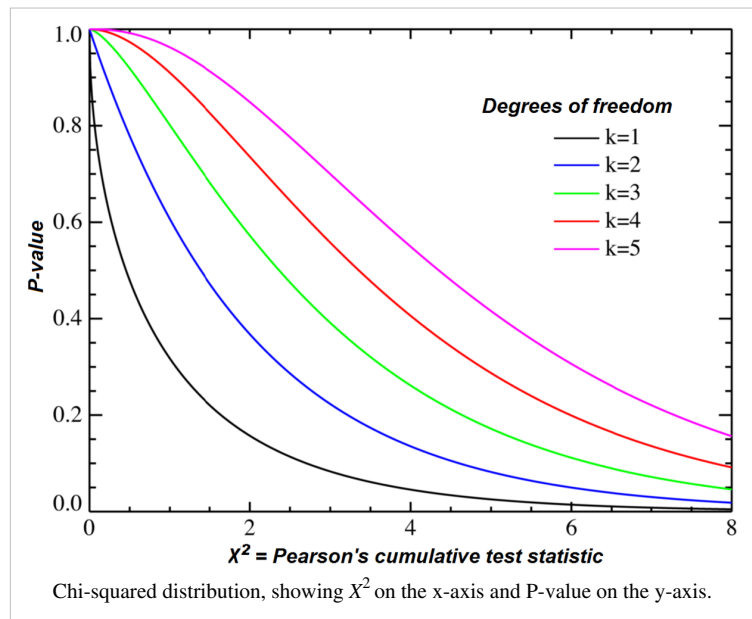
O_i = an observed frequency;

E_i = an expected (theoretical) frequency, asserted by the null hypothesis;

n = the number of cells in the table.

The chi-squared statistic can then be used to calculate a p-value by comparing the value of the statistic to a chi-squared distribution. The number of degrees of freedom is equal to the number of cells n , minus the reduction in degrees of freedom, p .

The result about the numbers of degrees of freedom is valid when the original data are multinomial and hence the estimated parameters are efficient for minimizing the chi-squared statistic. More generally however, when maximum likelihood estimation does not coincide with minimum chi-squared estimation, the distribution will lie somewhere between a chi-squared distribution with $n - 1 - p$ and $n - 1$ degrees of freedom (See for instance Chernoff and Lehmann, 1954).



Bayesian method

For more details on this topic, see Categorical distribution § With a conjugate prior.

In Bayesian statistics, one would instead use a Dirichlet distribution as conjugate prior. If one took a uniform prior, then the maximum likelihood estimate for the population probability is the observed probability, and one may compute a credible region around this or another estimate.

Test of independence

In this case, an "observation" consists of the values of two outcomes and the null hypothesis is that the occurrence of these outcomes is statistically independent. Each observation is allocated to one cell of a two-dimensional array of cells (called a contingency table) according to the values of the two outcomes. If there are r rows and c columns in the table, the "theoretical frequency" for a cell, given the hypothesis of independence, is

$$E_{i,j} = \frac{\left(\sum_{n_c=1}^c O_{i,n_c}\right) \cdot \left(\sum_{n_r=1}^r O_{n_r,j}\right)}{N},$$

where N is the total sample size (the sum of all cells in the table). With the term "frequencies" this page does not refer to already normalised values.

The value of the test-statistic is

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{i,j} - E_{i,j})^2}{E_{i,j}}.$$

Fitting the model of "independence" reduces the number of degrees of freedom by $p = r + c - 1$. The number of degrees of freedom is equal to the number of cells rc , minus the reduction in degrees of freedom, p , which reduces to $(r - 1)(c - 1)$.

For the test of independence, also known as the test of homogeneity, a chi-squared probability of less than or equal to 0.05 (or the chi-squared statistic being at or larger than the 0.05 critical point) is commonly interpreted by applied workers as justification for rejecting the null hypothesis that the row variable is independent of the column variable. The alternative hypothesis corresponds to the variables having an association or relationship where the structure of this relationship is not specified.

Assumptions

The chi-squared test, when used with the standard approximation that a chi-squared distribution is applicable, has the following assumptions: Wikipedia:Citation needed

- Simple random sample – The sample data is a random sampling from a fixed distribution or population where every collection of members of the population of the given sample size has an equal probability of selection. Variants of the test have been developed for complex samples, such as where the data is weighted. Other forms can be used such as purposive sampling^[1]
- Sample size (whole table) – A sample with a sufficiently large size is assumed. If a chi squared test is conducted on a sample with a smaller size, then the chi squared test will yield an inaccurate inference. The researcher, by using chi squared test on small samples, might end up committing a Type II error.
- Expected cell count – Adequate expected cell counts. Some require 5 or more, and others require 10 or more. A common rule is 5 or more in all cells of a 2-by-2 table, and 5 or more in 80% of cells in larger tables, but no cells with zero expected count. When this assumption is not met, Yates's Correction is applied.
- Independence – The observations are always assumed to be independent of each other. This means chi-squared cannot be used to test correlated data (like matched pairs or panel data). In those cases you might want to turn to McNemar's test.

A test that relies on different assumptions is Fisher's exact test; if its assumption of fixed marginal distributions is met it is substantially more accurate in obtaining a significance level, especially with few observations. In the vast majority of applications this assumption will not be met, and Fisher's exact test will be over conservative and not have correct coverage. [Wikipedia:Citation needed](#)

Examples

Goodness of fit

In this context, the frequencies of both theoretical and empirical distributions are unnormalised counts, and for a chi-squared test the total sample sizes N of both these distributions (sums of all cells of the corresponding contingency tables) have to be the same.

For example, to test the hypothesis that a random sample of 100 people has been drawn from a population in which men and women are equal in frequency, the observed number of men and women would be compared to the theoretical frequencies of 50 men and 50 women. If there were 44 men in the sample and 56 women, then

$$\chi^2 = \frac{(44 - 50)^2}{50} + \frac{(56 - 50)^2}{50} = 1.44.$$

If the null hypothesis is true (i.e., men and women are chosen with equal probability), the test statistic will be drawn from a chi-squared distribution with one degree of freedom (because if the male frequency is known, then the female frequency is determined).

Consultation of the chi-squared distribution for 1 degree of freedom shows that the probability of observing this difference (or a more extreme difference than this) if men and women are equally numerous in the population is approximately 0.23. This probability is higher than conventional criteria for statistical significance (0.001 or 0.05), so normally we would not reject the null hypothesis that the number of men in the population is the same as the number of women (i.e., we would consider our sample within the range of what we'd expect for a 50/50 male/female ratio.)

Problems

The approximation to the chi-squared distribution breaks down if expected frequencies are too low. It will normally be acceptable so long as no more than 20% of the events have expected frequencies below 5. Where there is only 1 degree of freedom, the approximation is not reliable if expected frequencies are below 10. In this case, a better approximation can be obtained by reducing the absolute value of each difference between observed and expected frequencies by 0.5 before squaring; this is called Yates's correction for continuity.

In cases where the expected value, E , is found to be small (indicating a small underlying population probability, and/or a small number of observations), the normal approximation of the multinomial distribution can fail, and in such cases it is found to be more appropriate to use the G-test, a likelihood ratio-based test statistic. When the total sample size is small, it is necessary to use an appropriate exact test, typically either the binomial test or (for contingency tables) Fisher's exact test. This test uses the conditional distribution of the test statistic given the marginal totals; however, it does not assume that the data were generated from an experiment in which the marginal totals are fixed and is valid whether or not that is the case.

Notes

[1] . See 'Discovering Statistics Using SPSS' by Andy Field for assumptions on Chi Square. -

References

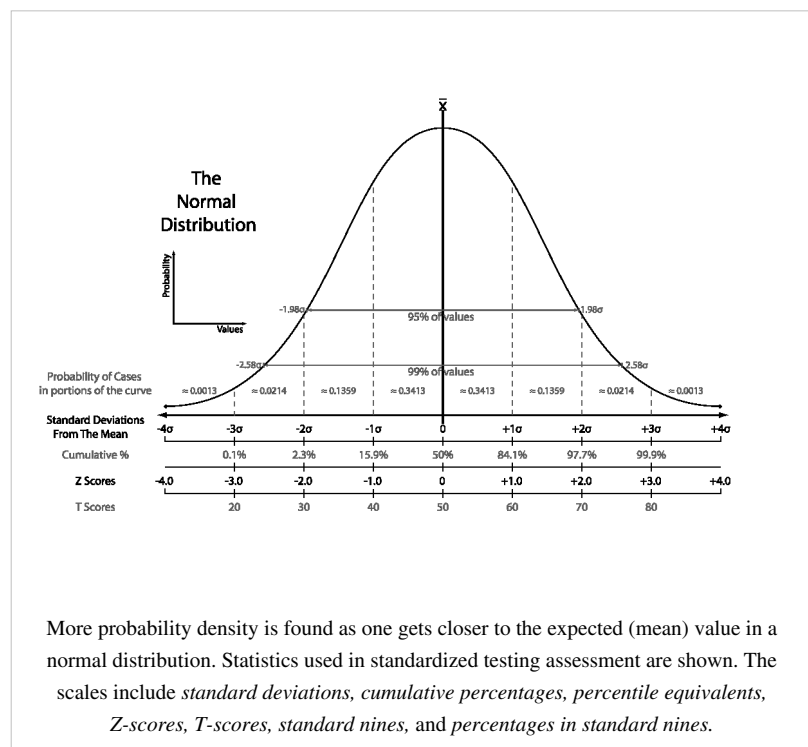
- Chernoff, H.; Lehmann, E. L. (1954). "The Use of Maximum Likelihood Estimates in χ^2 Tests for Goodness of Fit". *The Annals of Mathematical Statistics* **25** (3): 579–586. doi: 10.1214/aoms/1177728726 (<http://dx.doi.org/10.1214/aoms/1177728726>).
- Plackett, R. L. (1983). "Karl Pearson and the Chi-Squared Test". *International Statistical Review* (International Statistical Institute (ISI)) **51** (1): 59–72. doi: 10.2307/1402731 (<http://dx.doi.org/10.2307/1402731>). JSTOR 1402731 (<http://www.jstor.org/stable/1402731>).
- Greenwood, P.E.; Nikulin, M.S. (1996). *A guide to chi-squared testing*. New York: Wiley. ISBN 0-471-55779-X.

Statistics

Statistics is the study of the collection, organization, analysis, interpretation and presentation of data.^[1] It deals with all aspects of data including the planning of data collection in terms of the design of surveys and experiments. When analyzing data, it is possible to use one of two statistics methodologies: descriptive statistics or inferential statistics.

Scope

Statistics is a mathematical body of science that pertains to the collection, analysis, interpretation or explanation, and presentation of data,^[2] or as a branch of mathematics.^[3] Some consider statistics to be a distinct mathematical science rather than a branch of mathematics. Wikipedia:Vagueness



Mathematical statistics

Mathematical statistics is the application of mathematics to statistics, which was originally conceived as the science of the state — the collection and analysis of facts about a country: its economy, land, military, population, and so forth. Mathematical techniques which are used for this include mathematical analysis, linear algebra, stochastic analysis, differential equations, and measure-theoretic probability theory.

Overview

In applying statistics to e.g. a scientific, industrial, or societal problem, it is necessary to begin with a population or process to be studied. Populations can be diverse topics such as "all persons living in a country" or "every atom composing a crystal".

Ideally, statisticians compile data about the entire population (an operation called census). This may be organized by governmental statistical institutes. **Descriptive statistics** can be used to summarize the population data. Numerical descriptors include mean and standard deviation for continuous data types (like income), while frequency and percentage are more useful in terms of describing categorical data (like race).

When a census is not feasible, a chosen subset of the population called a sample is studied. Once a sample that is representative of the population is determined, data is collected for the sample members in an observational or experimental setting. Again, descriptive statistics can be used to summarize the sample data. However, the drawing of the sample has been subject to an element of randomness, hence the established numerical descriptors from the sample are also due to uncertainty. In order to still draw meaningful conclusions about the entire population, **inferential statistics** is needed. It uses patterns in the sample data to draw inferences about the population represented, accounting for randomness. These inferences may take the form of: answering yes/no questions about the data (hypothesis testing), estimating numerical characteristics of the data (estimation), describing associations within the data (correlation) and modeling relationships within the data (for example, using regression analysis). Inference can extend to forecasting, prediction and estimation of unobserved values either in or associated with the population being studied; it can include extrapolation and interpolation of time series or spatial data, and can also include data mining.

Data collection

Sampling

In case census data cannot be collected, statisticians collect data by developing specific experiment designs and survey samples. Statistics itself also provides tools for prediction and forecasting the use of data through statistical models. To use a sample as a guide to an entire population, it is important that it truly represent the overall population. Representative sampling assures that inferences and conclusions can safely extend from the sample to the population as a whole. A major problem lies in determining the extent that the sample chosen is actually representative. Statistics offers methods to estimate and correct for any random trending within the sample and data collection procedures. There are also methods of experimental design for experiments that can lessen these issues at the outset of a study, strengthening its capability to discern truths about the population.

Sampling theory is part of the mathematical discipline of probability theory. Probability is used in "mathematical statistics" (alternatively, "statistical theory") to study the sampling distributions of sample statistics and, more generally, the properties of statistical procedures. The use of any statistical method is valid when the system or population under consideration satisfies the assumptions of the method. The difference in point of view between classic probability theory and sampling theory is, roughly, that probability theory starts from the given parameters of a total population to deduce probabilities that pertain to samples. Statistical inference, however, moves in the opposite direction—inductively inferring from samples to the parameters of a larger or total population.

Experimental and observational studies

A common goal for a statistical research project is to investigate causality, and in particular to draw a conclusion on the effect of changes in the values of predictors or independent variables on dependent variables or response. There are two major types of causal statistical studies: experimental studies and observational studies. In both types of studies, the effect of differences of an independent variable (or variables) on the behavior of the dependent variable are observed. The difference between the two types lies in how the study is actually conducted. Each can be very effective. An experimental study involves taking measurements of the system under study, manipulating the system, and then taking additional measurements using the same procedure to determine if the manipulation has modified the values of the measurements. In contrast, an observational study does not involve experimental manipulation. Instead, data are gathered and correlations between predictors and response are investigated. While the tools of data analysis work best on data from randomized studies, they are also applied to other kinds of data – like natural experiments and observational studies^[4] – for which a statistician would use a modified, more structured estimation method (e.g., Difference in differences estimation and instrumental variables, among many others) that will produce consistent estimators.

Experiments

The basic steps of a statistical experiment are:

1. Planning the research, including finding the number of replicates of the study, using the following information: preliminary estimates regarding the size of treatment effects, alternative hypotheses, and the estimated experimental variability. Consideration of the selection of experimental subjects and the ethics of research is necessary. Statisticians recommend that experiments compare (at least) one new treatment with a standard treatment or control, to allow an unbiased estimate of the difference in treatment effects.
2. Design of experiments, using blocking to reduce the influence of confounding variables, and randomized assignment of treatments to subjects to allow unbiased estimates of treatment effects and experimental error. At this stage, the experimenters and statisticians write the *experimental protocol* that shall guide the performance of the experiment and that specifies the *primary analysis* of the experimental data.
3. Performing the experiment following the experimental protocol and analyzing the data following the experimental protocol.
4. Further examining the data set in secondary analyses, to suggest new hypotheses for future study.
5. Documenting and presenting the results of the study.

Experiments on human behavior have special concerns. The famous Hawthorne study examined changes to the working environment at the Hawthorne plant of the Western Electric Company. The researchers were interested in determining whether increased illumination would increase the productivity of the assembly line workers. The researchers first measured the productivity in the plant, then modified the illumination in an area of the plant and checked if the changes in illumination affected productivity. It turned out that productivity indeed improved (under the experimental conditions). However, the study is heavily criticized today for errors in experimental procedures, specifically for the lack of a control group and blindness. The Hawthorne effect refers to finding that an outcome (in this case, worker productivity) changed due to observation itself. Those in the Hawthorne study became more productive not because the lighting was changed but because they were being observed.

Observational study

An example of an observational study is one that explores the correlation between smoking and lung cancer. This type of study typically uses a survey to collect observations about the area of interest and then performs statistical analysis. In this case, the researchers would collect observations of both smokers and non-smokers, perhaps through a case-control study, and then look for the number of cases of lung cancer in each group.

Type of data

Main articles: Statistical data type and Levels of measurement

Various attempts have been made to produce a taxonomy of levels of measurement. The psychophysicist Stanley Smith Stevens defined nominal, ordinal, interval, and ratio scales. Nominal measurements do not have meaningful rank order among values, and permit any one-to-one transformation. Ordinal measurements have imprecise differences between consecutive values, but have a meaningful order to those values, and permit any order-preserving transformation. Interval measurements have meaningful distances between measurements defined, but the zero value is arbitrary (as in the case with longitude and temperature measurements in Celsius or Fahrenheit), and permit any linear transformation. Ratio measurements have both a meaningful zero value and the distances between different measurements defined, and permit any rescaling transformation.

Because variables conforming only to nominal or ordinal measurements cannot be reasonably measured numerically, sometimes they are grouped together as categorical variables, whereas ratio and interval measurements are grouped together as quantitative variables, which can be either discrete or continuous, due to their numerical nature. Such distinctions can often be loosely correlated with data type in computer science, in that dichotomous categorical variables may be represented with the Boolean data type, polytomous categorical variables with arbitrarily assigned integers in the integral data type, and continuous variables with the real data type involving floating point computation. But the mapping of computer science data types to statistical data types depends on which categorization of the latter is being implemented.

Other categorizations have been proposed. For example, Mosteller and Tukey (1977)^[5] distinguished grades, ranks, counted fractions, counts, amounts, and balances. Nelder (1990)^[6] described continuous counts, continuous ratios, count ratios, and categorical modes of data. See also Chrisman (1998),^[7] van den Berg (1991).^[8]

The issue of whether or not it is appropriate to apply different kinds of statistical methods to data obtained from different kinds of measurement procedures is complicated by issues concerning the transformation of variables and the precise interpretation of research questions. "The relationship between the data and what they describe merely reflects the fact that certain kinds of statistical statements may have truth values which are not invariant under some transformations. Whether or not a transformation is sensible to contemplate depends on the question one is trying to answer" (Hand, 2004, p. 82).^[9]

Terminology and theory of inferential statistics

Statistics, estimators and pivotal quantities

Consider an independent identically distributed (iid) random variables with a given probability distribution: standard statistical inference and estimation theory defines a random sample as the random vector given by the column vector of these iid variables.^[10] The population being examined is described by a probability distribution which may have unknown parameters.

A statistic is random variable which is a function of the random sample, but *not a function of unknown parameters*. The probability distribution of the statistics, though, may have unknown parameters.

Consider now a function of the unknown parameter: an estimator is a statistic used to estimate such function. Commonly used estimators include sample mean, unbiased sample variance and sample covariance.

A random variable which is a function of the random sample and of the unknown parameter, but whose probability distribution *does not depend on the unknown parameter* is called a pivotal quantity or pivot. Widely used pivots include the z-score, the chi square statistic and Student's t-value.

Between two estimators of a given parameter, the one with lower mean squared error is said to be more efficient. Furthermore an estimator is said to be unbiased if its expected value is equal to the true value of the unknown parameter which is being estimated and asymptotically unbiased if its expected value converges at the limit to the true value of such parameter.

Other desirable properties for estimators include: UMVUE estimators which have the lowest variance for all possible values of the parameter to be estimated (this is usually an easier property to verify than efficiency) and consistent estimators which converges in probability to the true value of such parameter.

This still leaves the question of how to obtain estimators in a given situation and carry the computation, several methods have been proposed: the method of moments, the maximum likelihood method, the least squares method and the more recent method of estimating equations.

Null hypothesis and alternative hypothesis

Interpretation of statistical information can often involve the development of a null hypothesis in that the assumption is that whatever is proposed as a cause has no effect on the variable being measured.

The best illustration for a novice is the predicament encountered by a jury trial. The null hypothesis, H_0 , asserts that the defendant is innocent, whereas the alternative hypothesis, H_1 , asserts that the defendant is guilty. The indictment comes because of suspicion of the guilt. The H_0 (status quo) stands in opposition to H_1 and is maintained unless H_1 is supported by evidence "beyond a reasonable doubt". However, "failure to reject H_0 " in this case does not imply innocence, but merely that the evidence was insufficient to convict. So the jury does not necessarily *accept* H_0 but *fails to reject* H_0 . While one can not "prove" a null hypothesis, one can test how close it is to being true with a power test, which tests for type II errors.

What statisticians call a alternative hypothesis is simply an hypothesis which contradicts the null hypothesis.

Error

Working from a null hypothesis two basic forms of error are recognized:

- Type I errors where the null hypothesis is falsely rejected giving a "false positive".
- Type II errors where the null hypothesis fails to be rejected and an actual difference between populations is missed giving a "false negative".

Standard deviation refers to the extent to which individual observations in a sample differ from a central value, such as the sample or population mean, while Standard error refers to an estimate of difference between sample mean and population mean.

A statistical error is the amount by which an observation differs from its expected value, a residual is the amount an observation differs from the value the estimator of the expected value assumes on a given sample (also called prediction).

Mean squared error is used for obtaining efficient estimators, a widely used class of estimators. Root mean square error is simply the square root of mean squared error.

Many statistical methods seek to minimize the residual sum of squares, and these are called "methods of least squares" in contrast to Least absolute deviations. The later gives equal weight to small and big errors, while the former gives more weight to large errors. Residual sum of squares is also differentiable, which provides a handy property for doing regression. Least squares applied to linear regression is called ordinary least squares method and least squares applied to nonlinear regression is called non-linear least squares. Also in a linear regression model the non deterministic part of the model is called error term, disturbance or more simply noise.

Measurement processes that generate statistical data are also subject to error. Many of these errors are classified as random (noise) or systematic (bias), but other important types of errors (e.g., blunder, such as when an analyst reports incorrect units) can also be important. The presence of missing data and/or censoring may result in biased estimates and specific techniques have been developed to address these problems.^[11]

Interval estimation

Main article: Interval estimation

Most studies only sample part of a population, so results don't fully represent the whole population. Any estimates obtained from the sample only approximate the population value. Confidence intervals allow statisticians to express how closely the sample estimate matches the true value in the whole population. Often they are expressed as 95% confidence intervals. Formally, a 95% confidence interval for a value is a range where, if the sampling and analysis were repeated under the same conditions (yielding a different dataset), the interval would include the true (population) value in 95% of all possible cases. This does *not* imply that the probability that the true value is in the confidence interval is 95%. From the frequentist perspective, such a claim does not even make sense, as the true value is not a random variable. Either the true value is or is not within the given interval. However, it is true that, before any data are sampled and given a plan for how to construct the confidence interval, the probability is 95% that the yet-to-be-calculated interval will cover the true value: at this point, the limits of the interval are yet-to-be-observed random variables. One approach that does yield an interval that can be interpreted as having a given probability of containing the true value is to use a credible interval from Bayesian statistics: this approach depends on a different way of interpreting what is meant by "probability", that is as a Bayesian probability.

In principle confidence intervals can be symmetrical or asymmetrical. An interval can be asymmetrical because it works as lower or upper bound for a parameter (left-sided interval or right sided interval), but it can also be asymmetrical because the two sided interval is built violating symmetry around the estimate. Sometimes the bounds for a confidence interval are reached asymptotically and these are used to approximate the true bounds.

Significance

Main article: Statistical significance

Statistics rarely give a simple Yes/No type answer to the question asked of them. Interpretation often comes down to the level of statistical significance applied to the numbers and often refers to the probability of a value accurately rejecting the null hypothesis (sometimes referred to as the p-value).

Referring to statistical significance does not necessarily mean that the overall result is significant in real world terms. For example, in a large study of a drug it may be shown that the drug has a statistically significant but very small beneficial effect, such that the drug is unlikely to help the patient noticeably.

Criticisms arise because the hypothesis testing approach forces one hypothesis (the null hypothesis) to be "favored," and can also seem to exaggerate the importance of minor differences in large studies. A difference that is highly statistically significant can still be of no practical significance, but it is possible to properly formulate tests in account for this. (See also criticism of hypothesis testing.)

One response involves going beyond reporting only the significance level to include the *p*-value when reporting whether a hypothesis is rejected or accepted. The *p*-value, however, does not indicate the size of the effect. A better and increasingly common approach is to report confidence intervals. Although these are produced from the same calculations as those of hypothesis tests or *p*-values, they describe both the size of the effect and the uncertainty surrounding it.

Examples

Some well-known statistical tests and procedures are:

- Analysis of variance (ANOVA)
- Chi-squared test
- Correlation
- Factor analysis
- Mann–Whitney U
- Mean square weighted deviation (MSWD)
- Pearson product-moment correlation coefficient
- Regression analysis
- Spearman's rank correlation coefficient
- Student's t -test
- Time series analysis

Misuse of statistics

Main article: Misuse of statistics

Misuse of statistics can produce subtle, but serious errors in description and interpretation—subtle in the sense that even experienced professionals make such errors, and serious in the sense that they can lead to devastating decision errors. For instance, social policy, medical practice, and the reliability of structures like bridges all rely on the proper use of statistics.

Even when statistical techniques are correctly applied, the results can be difficult to interpret for those lacking expertise. The statistical significance of a trend in the data—which measures the extent to which a trend could be caused by random variation in the sample—may or may not agree with an intuitive sense of its significance. The set of basic statistical skills (and skepticism) that people need to deal with information in their everyday lives properly is referred to as statistical literacy.

There is a general perception that statistical knowledge is all-too-frequently intentionally misused by finding ways to interpret only the data that are favorable to the presenter.^[12] A mistrust and misunderstanding of statistics is associated with the quotation, "There are three kinds of lies: lies, damned lies, and statistics". Misuse of statistics can be both inadvertent and intentional, and the book *How to Lie with Statistics* outlines a range of considerations. In an attempt to shed light on the use and misuse of statistics, reviews of statistical techniques used in particular fields are conducted (e.g. Warne, Lazo, Ramos, and Ritter (2012)).^[13]

Ways to avoid misuse of statistics include using proper diagrams and avoiding bias. Misuse can occur when conclusions are overgeneralized and claimed to be representative of more than they really are, often by either deliberately or unconsciously overlooking sampling bias. Bar graphs are arguably the easiest diagrams to use and understand, and they can be made either by hand or with simple computer programs. Unfortunately, most people do not look for bias or errors, so they are not noticed. Thus, people may often believe that something is true even if it is not well represented. To make data gathered from statistics believable and accurate, the sample taken must be representative of the whole. According to Huff, "The dependability of a sample can be destroyed by [bias]... allow yourself some degree of skepticism."

To assist in the understanding of statistics Huff proposed a series of questions to be asked in each case:

- Who says so? (Does he/she have an axe to grind?)
 - How does he/she know? (Does he/she have the resources to know the facts?)
 - What's missing? (Does he/she give us a complete picture?)
 - Did someone change the subject? (Does he/she offer us the right answer to the wrong problem?)
 - Does it make sense? (Is his/her conclusion logical and consistent with what we already know?)
-

Misinterpretation: correlation

The concept of correlation is particularly noteworthy for the potential confusion it can cause. Statistical analysis of a data set often reveals that two variables (properties) of the population under consideration tend to vary together, as if they were connected. For example, a study of annual income that also looks at age of death might find that poor people tend to have shorter lives than affluent people. The two variables are said to be correlated; however, they may or may not be the cause of one another. The correlation phenomena could be caused by a third, previously unconsidered phenomenon, called a lurking variable or confounding variable. For this reason, there is no way to immediately infer the existence of a causal relationship between the two variables. (See Correlation does not imply causation.)

History of statistical science

Main articles: History of statistics and Founders of statistics

R A Fisher is the founder of Statistics.

Statistical methods date back at least to the 5th century BC.

Some scholars pinpoint the origin of statistics to 1663, with the publication of *Natural and Political Observations upon the Bills of Mortality* by John Graunt.^[14] Early applications of statistical thinking revolved around the needs of states to base policy on demographic and economic data, hence its *stat-* etymology. The scope of the discipline of statistics broadened in the early 19th century to include the collection and analysis of data in general. Today, statistics is widely employed in government, business, and natural and social sciences.

Its mathematical foundations were laid in the 17th century with the development of the probability theory by Blaise Pascal and Pierre de Fermat. Mathematical probability theory arose from the study of games of chance, although the concept of probability was already examined in medieval law and by philosophers such as Juan Caramuel.^[15] The method of least squares was first described by Adrien-Marie Legendre in 1805.

The modern field of statistics emerged in the late 19th and early 20th century in three stages. The first wave, at the turn of the century, was led by the work of Sir Francis Galton and Karl Pearson, who transformed statistics into a rigorous mathematical discipline used for analysis, not just in science, but in industry and politics as well. Galton's contributions to the field included introducing the concepts of standard deviation, correlation, regression and the application of these methods to the study of the variety of human characteristics – height, weight, eyelash length among others.^[16] Pearson developed the Correlation coefficient, defined as a product-moment, the method of moments for the fitting of distributions to samples and the



Karl Pearson, the founder of mathematical statistics.

Pearson's system of continuous curves, among many other things. Galton and Pearson founded *Biometrika* as the first journal of mathematical statistics and biometry, and the latter founded the world's first university statistics department at University College London.

The second wave of the 1910s and 20s was initiated by William Gosset, and reached its culmination in the insights of Sir Ronald Fisher, who wrote the textbooks that were to define the academic discipline in universities around the world. Fisher's most important publications were his 1916 seminal paper *The Correlation between Relatives on the Supposition of Mendelian Inheritance* and his classic 1925 work *Statistical Methods for Research Workers*. His paper was the first to use the statistical term, variance. He developed rigorous experimental models and also originated the concepts of sufficiency, ancillary statistics, Fisher's linear discriminator and Fisher information.

The final wave, which mainly saw the refinement and expansion of earlier developments, emerged from the collaborative work between Egon Pearson and Jerzy Neyman in the 1930s. They introduced the concepts of "Type II" error, power of a test and confidence intervals. Jerzy Neyman in 1934 showed that stratified random sampling was in general a better method of estimation than purposive (quota) sampling.^[17] Today, statistical methods are applied in all fields that involve decision making, for making accurate inferences from a collated body of data and for making decisions in the face of uncertainty based on statistical methodology. The use of modern computers has expedited large-scale statistical computations, and has also made possible new methods that are impractical to perform manually.

Trivia

Applied statistics, theoretical statistics and mathematical statistics

"Applied statistics" comprises descriptive statistics and the application of inferential statistics.^[18] Wikipedia:Verifiability *Theoretical statistics* concerns both the logical arguments underlying justification of approaches to statistical inference, as well encompassing *mathematical statistics*. Mathematical statistics includes not only the manipulation of probability distributions necessary for deriving results related to methods of estimation and inference, but also various aspects of computational statistics and the design of experiments.

Machine learning and data mining

Statistics has many ties to machine learning and data mining.

Statistics in society

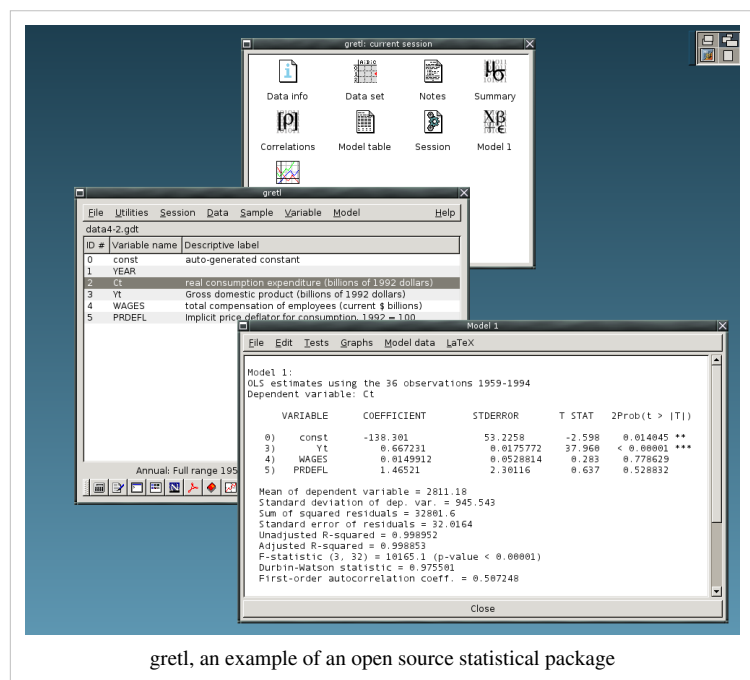
Statistics is applicable to a wide variety of academic disciplines, including natural and social sciences, government, and business. Statistical consultants can help organizations and companies that don't have in-house expertise relevant to their particular questions.

Statistical computing

Main article: Computational statistics

The rapid and sustained increases in computing power starting from the second half of the 20th century have had a substantial impact on the practice of statistical science. Early statistical models were almost always from the class of linear models, but powerful computers, coupled with suitable numerical algorithms, caused an increased interest in nonlinear models (such as neural networks) as well as the creation of new types, such as generalized linear models and multilevel models.

Increased computing power has also led to the growing popularity of computationally



gretl, an example of an open source statistical package

intensive methods based on resampling, such as permutation tests and the bootstrap, while techniques such as Gibbs sampling have made use of Bayesian models more feasible. The computer revolution has implications for the future of statistics with new emphasis on "experimental" and "empirical" statistics. A large number of both general and special purpose statistical software are now available.

Statistics applied to mathematics or the arts

Traditionally, statistics was concerned with drawing inferences using a semi-standardized methodology that was "required learning" in most sciences. This has changed with use of statistics in non-inferential contexts. What was once considered a dry subject, taken in many fields as a degree-requirement, is now viewed enthusiastically. Wikipedia:Avoid weasel words Initially derided by some mathematical purists, it is now considered essential methodology in certain areas.

- In number theory, scatter plots of data generated by a distribution function may be transformed with familiar tools used in statistics to reveal underlying patterns, which may then lead to hypotheses.
- Methods of statistics including predictive methods in forecasting are combined with chaos theory and fractal geometry to create video works that are considered to have great beauty.
- The process art of Jackson Pollock relied on artistic experiments whereby underlying distributions in nature were artistically revealed. Wikipedia:Citation needed With the advent of computers, statistical methods were applied to formalize such distribution-driven natural processes to make and analyze moving video art. Wikipedia:Citation needed
- Methods of statistics may be used predicatively in performance art, as in a card trick based on a Markov process that only works some of the time, the occasion of which can be predicted using statistical methodology.
- Statistics can be used to predicatively create art, as in the statistical or stochastic music invented by Iannis Xenakis, where the music is performance-specific. Though this type of artistry does not always come out as expected, it does behave in ways that are predictable and tunable using statistics.

Specialized disciplines

Main article: List of fields of application of statistics

Statistical techniques are used in a wide range of types of scientific and social research, including: biostatistics, computational biology, computational sociology, network biology, social science, sociology and social research. Some fields of inquiry use applied statistics so extensively that they have specialized terminology. These disciplines include:

- Actuarial science (assesses risk in the insurance and finance industries)
- Applied information economics
- Biostatistics
- Business statistics
- Chemometrics (for analysis of data from chemistry)
- Data mining (applying statistics and pattern recognition to discover knowledge from data)
- Demography
- Econometrics
- Energy statistics
- Engineering statistics
- Epidemiology
- Geography and Geographic Information Systems, specifically in Spatial analysis
- Image processing

Medical Statistics

- Psychological statistics
- Reliability engineering
- Social statistics

In addition, there are particular types of statistical analysis that have also developed their own specialised terminology and methodology:

- Bootstrap / Jackknife resampling
- Multivariate statistics
- Statistical classification
- Structured data analysis (statistics)
- Structural equation modelling
- Survey methodology
- Survival analysis
- Statistics in various sports, particularly baseball - known as 'Sabmetrics' - and cricket

Statistics form a key basis tool in business and manufacturing as well. It is used to understand measurement systems variability, control processes (as in statistical process control or SPC), for summarizing data, and to make data-driven decisions. In these roles, it is a key tool, and perhaps the only reliable tool.

References

- [1] Dodge, Y. (2006) *The Oxford Dictionary of Statistical Terms*, OUP. ISBN 0-19-920613-9
- [2] Moses, Lincoln E. (1986) *Think and Explain with Statistics*, Addison-Wesley, ISBN 978-0-201-15619-5 . pp. 1–3
- [3] Hays, William Lee, (1973) *Statistics for the Social Sciences*, Holt, Rinehart and Winston, p.xii, ISBN 978-0-03-077945-9
- [4] Freedman, D.A. (2005) *Statistical Models: Theory and Practice*, Cambridge University Press. ISBN 978-0-521-67105-7
- [5] Mosteller, F., & Tukey, J. W. (1977). *Data analysis and regression*. Boston: Addison-Wesley.
- [6] Nelder, J. A. (1990). The knowledge needed to computerise the analysis and interpretation of statistical information. In *Expert systems and artificial intelligence: the need for information about data*. Library Association Report, London, March, 23–27.
- [7] Chrisman, Nicholas R. (1998). Rethinking Levels of Measurement for Cartography. *Cartography and Geographic Information Science*, vol. 25 (4), pp. 231–242
- [8] van den Berg, G. (1991). *Choosing an analysis method*. Leiden: DSWO Press
- [9] Hand, D. J. (2004). *Measurement theory and practice: The world through quantification*. London, UK: Arnold.
- [10] P. Elio, *Probabilità e Statistica*, Esculapio 2007
- [11] Rubin, Donald B.; Little, Roderick J. A., *Statistical analysis with missing data*, New York: Wiley 2002
- [12] Huff, Darrell (1954) *How to Lie with Statistics*, WW Norton & Company, Inc. New York, NY. ISBN 0-393-31072-8
- [13] Warne, R. Lazo, M., Ramos, T. and Ritter, N. (2012). Statistical Methods Used in Gifted Education Journals, 2006–2010. *Gifted Child Quarterly*, 56(3) 134–149.
- [14] Willcox, Walter (1938) "The Founder of Statistics". *Review of the International Statistical Institute* 5(4):321–328.
- [15] J. Franklin, *The Science of Conjecture: Evidence and Probability before Pascal*, Johns Hopkins Univ Pr 2002
- [16] Galton F (1877) Typical laws of heredity. *Nature* 15: 492–553
- [17] Neyman, J (1934) On the two different aspects of the representative method: The method of stratified sampling and the method of purposive selection. *Journal of the Royal Statistical Society* 97 (4) 557–625
- [18] Anderson, D.R.; Sweeney, D.J.; Williams, T.A.. (1994) *Introduction to Statistics: Concepts and Applications*, pp. 5–9. West Group. ISBN 978-0-314-03309-3

Image Sources, Licenses and Contributors

File:Chi-square pdf.svg *Source:* http://en.wikipedia.org/w/index.php?title=File:Chi-square_pdf.svg *License:* Creative Commons Attribution 3.0 *Contributors:* Geek3

File:Chi-square cdf.svg *Source:* http://en.wikipedia.org/w/index.php?title=File:Chi-square_cdf.svg *License:* Creative Commons Attribution 3.0 *Contributors:* Geek3

File:Chernoff XS CDF.png *Source:* http://en.wikipedia.org/w/index.php?title=File:Chernoff_XS_CDF.png *License:* GNU Free Documentation License *Contributors:* Willem

File:Chi-square distributionCDF-English.png *Source:* http://en.wikipedia.org/w/index.php?title=File:Chi-square_distributionCDF-English.png *License:* Public Domain *Contributors:* Mikael Häggström

File:The Normal Distribution.svg *Source:* http://en.wikipedia.org/w/index.php?title=File:The_Normal_Distribution.svg *License:* Public Domain *Contributors:* Abdull, CarolSpears, Inductiveload, Trijnstel, 青子守歌, 17 anonymous edits

File:Mw160883.jpg *Source:* <http://en.wikipedia.org/w/index.php?title=File:Mw160883.jpg> *License:* Public Domain *Contributors:* Julian Felsenburgh

File:Gretl screenshot.png *Source:* http://en.wikipedia.org/w/index.php?title=File:Gretl_screenshot.png *License:* GNU General Public License *Contributors:* Cathy Richards, Den fjättrade ankan, Hannibal, Marcus Cyron, WikipediaMaster, 2 anonymous edits

License

Creative Commons Attribution-Share Alike 3.0
[//creativecommons.org/licenses/by-sa/3.0/](https://creativecommons.org/licenses/by-sa/3.0/)
