

[Click Here](#)



How to calculate p value from test statistic

Function of the observed sample results Not to be confused with the P-factor. In null-hypothesis significance testing, the p-value[note 1] is the probability of obtaining test results at least as extreme as the result actually observed, under the assumption that the null hypothesis is correct.[2][3] A very small p-value means that such an extreme observed outcome would be very unlikely under the null hypothesis. Even though reporting p-values of statistical tests is common practice in academic publications of many quantitative fields, misinterpretation and misuse of p-values is widespread and has been a major topic in mathematics and metascience.[4][5] In 2016, the American Statistical Association (ASA) made a formal statement that "p-values do not measure the probability that the studied hypothesis is true, or that a p-value, or statistical significance, does not measure the size of an effect or the importance of a result" or "evidence regarding a model or hypothesis".[6] That said, 2019 task force by ASA has issued a statement on statistical significance and replicability, concluding with, "p-values and significance tests, when properly applied and interpreted, increase the rigor of the conclusions drawn from data".[7] In statistics, every conjecture concerning the unknown probability distribution of a collection of random variables representing the observed data

X

{\displaystyle X}

 in some study is called a statistical hypothesis. If we state one hypothesis only and the aim of the statistical test is to see whether this hypothesis is tenable, but not to investigate other specific hypotheses, then such a test is called a null hypothesis test. As our statistical hypothesis will, by definition, state some property of the distribution, the null hypothesis is the default hypothesis under which that property does not exist. The null hypothesis is typically that some parameter (such as a correlation or a difference between means) in the populations of interest is zero. Our hypothesis might specify the probability distribution of

X

{\displaystyle X}

 precisely, or it might only specify that it belongs to some class of distributions. Often, we reduce the data to a single numerical statistic, e.g.

T

{\displaystyle T}

, whose marginal probability distribution is closely connected to a main question of interest in the study. The p-value is used in the context of null hypothesis testing in order to quantify the statistical significance of a result, the result being the observed value of the chosen statistic

T

{\displaystyle T}

. [note 1] The lower the p-value is, the lower the probability of getting that result if the null hypothesis were true. A result is said to be statistically significant if it allows us to reject the null hypothesis. All other things being equal, smaller p-values are taken as stronger evidence against the null hypothesis. Loosely speaking, rejection of the null hypothesis implies that there is sufficient evidence against it. As a particular example, if a null hypothesis states that a certain summary statistic

T

{\displaystyle T}

 follows the standard normal distribution

N
(
0
,
1
)
,

{\displaystyle N(0,1)}

 then the rejection of this null hypothesis could mean that (i) the mean of

T

{\displaystyle T}

 is not 0, or (ii) the variance of

T

{\displaystyle T}

 is not 1, or (iii)

T

{\displaystyle T}

 is not normally distributed. Different tests of the same null hypothesis would be more or less sensitive to different alternatives. However, even if we do manage to reject the null hypothesis for all 3 alternatives, and even if we know that the distribution is normal and variance = 1, the null hypothesis test does not tell us which non-zero values of the mean are now most plausible. The more independent observations from the same probability distribution one has, the more accurate the test will be, and the higher the precision with which one will be able to determine the mean value and show that it is not equal to zero; but this will also increase the importance of evaluating the real-world or scientific relevance of this deviation. The p-value is the probability under the null hypothesis of obtaining a real-valued test statistic at least as extreme as the one obtained. Consider an observed test-statistic

t

{\displaystyle t}

 from unknown distribution

T

{\displaystyle T}

. Then the p-value

p

{\displaystyle p}

 is what the prior probability would be of observing a test-statistic value at least as "extreme" as

t

{\displaystyle t}

 if null hypothesis

H
0

{\displaystyle H_{0}}

 were true. That is:

p
=
Pr
(
T
≥
t
|

H

0

)

{\displaystyle p=\Pr(T\geq t\mid H_{0})}

 for a one-sided right-tail test-statistic distribution.

p
=
Pr
(
T
≤
t
|

H

0

)

{\displaystyle p=\Pr(T\leq t\mid H_{0})}

 for a one-sided left-tail test-statistic distribution.

p
=
2
min
(
Pr
(
T
≥
t
|

H

0

)
,
Pr
(
T
≤
t
|

H

0

)
)

{\displaystyle p=2\min(\Pr(T\geq t\mid H_{0}),\Pr(T\leq t\mid H_{0}))}

 for a two-sided test-statistic distribution. If the distribution of

T

{\displaystyle T}

 is symmetric about zero, then

p
=
Pr
(
|
T
|
≥
|
t
|
|

H

0

)
.

{\displaystyle p=\Pr(|T|\geq |t|\mid H_{0}).}

 The error that a practising statistician would consider the more important to avoid (which is a subjective judgment) is called the error of the first kind. The first demand of the mathematical theory is to deduce such test criteria as would ensure that the probability of committing an error of the first kind would equal (or approximately equal, or not exceed) a preassigned number α, such as α = 0.05 or 0.01, etc. This number is called the level of significance.—Jerzy Neyman, "The Emergence of Mathematical Statistics"[8] In a significance test, the null hypothesis

H
0

{\displaystyle H_{0}}

 is rejected if the p-value is less than equal to a predefined threshold value

α

{\displaystyle \alpha }

, which is referred to as the alpha level or significance level. α is not derived from the data, but rather is set by the researcher before examining the data. α (

{\displaystyle \alpha }

) is commonly set to 0.05, though lower alpha levels are sometimes used. The 0.05 value (equivalent to 1/20 chances) was originally proposed by R. Fisher in 1925 in his famous book entitled "Statistical Methods for Research Workers".[9] Different p-values based on independent sets of data can be combined, for instance using Fisher's combined probability test. The p-value is a function of the chosen test statistic

T

{\displaystyle T}

 and is therefore a random variable. If the null hypothesis fixes the probability distribution of

T

{\displaystyle T}

 precisely (e.g.

H
0
:
θ
=
θ
0

{\displaystyle H_{0}:\theta =\theta _{0}}

, where

θ

{\displaystyle \theta }

 is the only parameter), and if that distribution is continuous, then when the null-hypothesis is true, the p-value is uniformly distributed between 0 and 1. Regardless of the truth of the

H
0

{\displaystyle H_{0}}

, the p-value is not fixed; if the same test is repeated indefinitely with fresh data, it will typically obtain a different p-value in each iteration. Usually only a single p-value relating to a hypothesis is observed, so the p-value is interpreted by a significance test, and no effort is made to estimate the distribution it was drawn from. When a collection of p-values are available (e.g. when considering a group of studies on the same subject), the distribution of p-values is sometimes called a p-curve.[10] A p-curve can be used to assess the reliability of scientific literature, such as by detecting publication bias or p-hacking.[10][11] In parametric hypothesis testing problems, a simple or point hypothesis refers to a hypothesis where the parameter's value is assumed to be a single number. In contrast, in a composite hypothesis the parameter's value is given by a set of numbers. When the null hypothesis is composite (or the distribution of the statistic is discrete), then when the null-hypothesis is true the probability of obtaining a p-value less than or equal to any number between 0 and 1 is still less than or equal to that number. In other words, it remains the case that very small values are relatively unlikely if the null-hypothesis is true, and that a significance test at level α (

{\displaystyle \alpha }

) is obtained by rejecting the null-hypothesis if the p-value is less than or equal to α (

{\displaystyle \alpha }

). [12][13] For example, when testing the null hypothesis that a distribution is normal with a mean less than or equal to zero against the alternative that the mean is greater than zero (

H
0
:
μ
≤
0

{\displaystyle H_{0}:\mu \leq 0}

, variance known), the null hypothesis does not specify the exact probability distribution of the appropriate test statistic. In this example that would be the Z-statistic belonging to the one-sided one-sample Z-test. For each possible value of the theoretical mean, the Z-test statistic has a different probability distribution. In these circumstances the p-value is defined by taking the least favorable null-hypothesis case, which is typically on the border between null and alternative. This definition ensures the complementarity of p-values and alpha-levels: α = 0.05 (

{\displaystyle \alpha =0.05}

 means one only rejects the null hypothesis if the p-value is less than or equal to 0.05 (

{\displaystyle 0.05}

), and the hypothesis test will indeed have a maximum type-I error rate of 0.05 (

{\displaystyle 0.05}

). The p-value is widely used in statistical hypothesis testing, specifically in null hypothesis significance testing. In this method, before conducting the study, one first chooses a model (the null hypothesis) and the alpha level α (most commonly 0.05). After analyzing the data, if the p-value is less than α, that is taken to mean that the observed data is sufficiently inconsistent with the null hypothesis for the null hypothesis to be rejected. However, that does not prove that the null hypothesis is false. The p-value does not, in itself, establish probabilities of hypotheses. Rather, it is a tool for deciding whether to reject the null hypothesis.[14] Main article: Misuse of p-values According to the ASA, there is widespread agreement that p-values are often misused and misinterpreted.[3] One practice that has been particularly criticized is accepting the alternative hypothesis for any p-value nominally less than 0.05 without other supporting evidence. Although p-values are helpful in assessing how incompatible the data are with a specified statistical model, contextual factors must also be considered, such as "the design of a study, the quality of the measurements, the external evidence for the phenomenon under study, and the validity of assumptions that underlie the data analysis".[3] Another concern is that the p-value is often misunderstood as being the probability that the null hypothesis is true.[3][15] Some statisticians have proposed abandoning p-values and focusing more on other inferential statistics,[3] such as confidence intervals,[16][17] likelihood ratios,[18][9] or Bayes factors.[20][21][22] but there is heated debate on the feasibility of these alternatives.[23][24] Others have suggested to remove fixed significance thresholds and to interpret p-values as continuous indices of the strength of evidence against the null hypothesis.[25][26] Yet others suggested to report alongside p-values the prior probability of a real effect that would be required to obtain a false positive risk (i.e. the probability that there is no real effect) below a pre-specified threshold (e.g. 5%).[27] That said, in 2019 a task force by ASA had convened to consider the use of statistical methods in scientific studies, specifically hypothesis tests and p-values, and their connection to replicability.[7] It states that "Different measures of uncertainty can complement one another; no single measure serves all purposes", citing p-value as one of these measures. They also stress that p-values can provide valuable information when considering the specific value as well as when compared to some threshold. In general, it stresses that "p-values and significance tests, when properly applied and interpreted, increase the rigor of the conclusions drawn from data". Usually,

T

{\displaystyle T}

 is a test statistic. A test statistic is the output of a scalar function of all the observations. This statistic provides a single number, such as a t-statistic or an F-statistic. As such, the test statistic follows a distribution determined by the function used to define that test statistic and the distribution of the input observational data. For the important case in which the data are hypothesized to be a random sample from a normal distribution, depending on the nature of the test statistic and the hypotheses of interest about its distribution, different null hypothesis tests have been developed. Some such tests are the z-test for hypotheses concerning the mean of a normal distribution with known variance, the t-test based on Student's t-distribution of a suitable statistic for hypotheses concerning the mean of a normal distribution when the variance is unknown, the F-test based on the F-distribution of yet another statistic for hypotheses concerning the variance. For data of other nature, for instance, categorical (discrete) data, test statistics might be constructed whose null hypothesis distribution is based on normal approximations to appropriate statistics obtained by invoking the central limit theorem for large samples, as in the case of Pearson's chi-squared test. Thus computing a p-value requires a null hypothesis, a test statistic (together with deciding whether the researcher is performing a one-tailed test or a two-tailed test), and data. Even though computing the test statistic on given data may be easy, computing the sampling distribution under the null hypothesis, and then computing its cumulative distribution function (CDF) is often a difficult problem. Today, this computation is done using statistical software, often via numeric methods (rather than exact formulae), but, in the early and mid 20th century, this was instead done via tables of values, and one interpolated or extrapolated p-values from these discrete values[citation needed]. Rather than using a table of p-values, Fisher instead inverted the CDF, publishing a list of values of the test statistic for given fixed p-values; this corresponds to computing the quantile function (inverse CDF). Main article: Checking whether a coin is fair As an example of a statistical test, an experiment is performed to determine whether a coin flip is fair (equal chance of landing heads or tails) or unfairly biased (one outcome being more likely than the other). Suppose that the experimental results show the coin turning up heads 14 times out of 20 total flips. The full data

X

{\displaystyle X}

 would be a sequence of twenty times the symbol "H" or "T". The statistic on which one might focus could be the total number

T

{\displaystyle T}

 of heads. The null hypothesis is that the coin is fair, and coin tosses are independent of one another. If a right-tailed test is considered, which would be the case if one is actually interested in the possibility that the coin is biased towards falling heads, then the p-value of this result is the chance of a fair coin landing on heads at least 14 times out of 20 flips. That probability can be computed from binomial coefficients as

Pr
(
14
heads
)
+
Pr
(
15
heads
)
+
⋯
+
Pr
(
20
heads
)
=
1
2

20

(
20
14
)
+
(
20
15
)
+
⋯
+
(
20
20
)
]
=
60
460
1
048
576
=
0.058.

{\displaystyle {\begin{aligned}& \Pr(14{\text{ heads}})+\Pr(15{\text{ heads}})+\dots +\Pr(20{\text{ heads}})\\& ={\frac {1}{2}}^{20}[({\binom {20}{14}}+{\binom {20}{15}})+\dots +{\binom {20}{20}}]{\text{right}}={\frac {60}{460}}\approx 0.058.\end{aligned}}}

 This probability is the p-value, considering only extreme results that favor heads. This is called a one-tailed test. However, one might be interested in deviations in either direction, favoring either heads or tails. The two-tailed p-value, which considers deviations favoring either heads or tails, may instead be calculated. As the binomial distribution is symmetrical for a fair coin, the two-sided p-value is simply twice the above calculated single-sided p-value: the two-sided p-value is 0.115. In the above example: Null hypothesis (H0): The coin is fair, with Pr(heads) = 0.5. Test statistic: Number of heads. Alpha level (designated threshold of significance): 0.05. Observation O: 14 heads out of 20 flips. Two-tailed p-value of observation O given H0 = 2 × min(Pr.no. of heads ≥ 14 heads), Pr(no. of heads ≤ 14 heads) = 2 × min(0.058, 0.978) = 2 × 0.058 = 0.115. The Pr.no. of heads ≤ 14 heads = 1 − Pr(no. of heads ≥ 14 heads) + Pr(no. of head = 14) = 1 − 0.058 + 0.036 = 0.978; however, the symmetry of this binomial distribution makes it an unnecessary computation to find the smaller of the two probabilities. Here, the calculated p-value exceeds 0.05, meaning that the data falls within the range of what would happen 95% of the time, if the coin were fair. Hence, the null hypothesis is not rejected at the 0.05 level. However, had one more head been obtained, the resulting p-value (two-tailed) would have been 0.0414 (4.14%), in which case the null hypothesis would be rejected at the 0.05 level. The difference between the two meanings of "extreme" appear when we consider a sequential hypothesis testing, or optional stopping, for the fairness of the coin. In general, optional stopping changes how p-value is calculated.[28][29] Suppose we design the experiment as follows: Flip the coin twice. If both comes up heads or tails, end the experiment. Else, flip the coin 4 more times. This experiment has 7 types of outcomes: 2 heads, 2 tails, 3 heads 1 tail, ..., 1 head 5 tails. We now calculate the p-value of the "3 heads 3 tails" outcome. If we use the test statistic heads / tails (

{\displaystyle {\text{heads}}/{\text{tails}}}

), then under the null hypothesis is exactly 1 for two-sided p-value, and exactly 19 / 32 (

{\displaystyle 19/32}

) for one-sided left-tail p-value, and same for one-sided right-tail p-value. If we consider every outcome that has equal or lower probability than "3 heads 3 tails" as "at least as extreme", then the p-value is exactly 1 / 2. (

{\displaystyle 1/2.}

) However, suppose we have planned to simply flip the coin 6 times no matter what happens, then the second definition of p-value would mean that the p-value of "3 heads 3 tails" is exactly 1. Thus, the "at least as extreme" definition of p-value is deeply contextual and depends on what the experimenter planned to do even in situations that did not occur. John Arbuthnot Pierre-Simon Laplace Carl Pearson Ronald Fisher P-value computations date back to the 1700s, where they were computed for the human sex ratio at birth, and used to compute statistical significance compared to the null hypothesis of equal probability of male and female births.[30] John Arbuthnot studied this question in 1710,[31][32][33][34] and examined birth records in London for each of the 82 years from 1629 to 1710. In every year, the number of males born in London exceeded the number of females. Considering more male or more female births as equally likely, the probability of the observed outcome is 1/282, or about 1 in 4,836,000,000,000,000,000,000; in modern terms, the p-value. This is vanishingly small, leading Arbuthnot to conclude that this was not due to chance, but to divine providence: "From whence it follows, that it is Art, not Chance, that governs.". In modern terms, he rejected the null hypothesis of equally likely male and female births at the

p
=
1/282

 significance level. This and other work by Arbuthnot is credited as "... the first use of significance tests ..."[35] the first example of reasoning about statistical significance.[36] and "... perhaps the first published report of a nonparametric test ..."[32] specifically the sign test; see details at Sign test § History. The same question was later addressed by Pierre-Simon Laplace, who instead used a parametric test, modeling the number of male births with a binomial distribution:[37] In the 1770s Laplace considered the statistics of almost half a million births. The statistics showed an excess of boys compared to girls. He concluded by calculation of a p-value that the excess was a real, but unexplained, effect. The p-value was first formally introduced by Karl Pearson, in his Pearson's chi-squared test.[38] using the chi-squared distribution and noted as capital P.[38] The p-values for the chi-squared distribution (for various values of

χ
2

 and degrees of freedom), now notated as

P
, were calculated in (Elderton 1902), collected in (Pearson 1914, pp. xxxi–xxxi, 26–28, Table XII). Ronald Fisher formalized and popularized the use of the p-value in statistics.[39][40] with it playing a central role in his approach to the subject.[41] In his highly influential book Statistical Methods for Research Workers (1925), Fisher proposed the level

p
=
0.05

, or a 1 in 20 chance of being exceeded by chance, as a limit for statistical significance, and applied this to a normal distribution (as a two-tailed test), thus yielding the rule of two standard deviations (on a normal distribution) for statistical significance (see 68–95–99.77 rule).[42][note 3][43] He then computed a table of values, similar to Elderton but, importantly, reversed the roles of

χ
2

 and

p

. That is, rather than computing

p

 for different values of

χ
2

 (and degrees of freedom

n

), he computed values of

χ
2

 that yield specified p-values, specifically 0.99, 0.98, 0.95, 0.90, 0.80, 0.70, 0.50, 0.30, 0.20, 0.10, 0.05, 0.02, and 0.01.[44] That allowed computed values of

χ
2

 to be compared against cutoffs and encouraged the use of p-values (especially 0.05, 0.02, and 0.01) as cutoffs, instead of computing and reporting p-values themselves. The same type of tables were then compiled in (Fisher & Yates 1938), which cemented the approach.[43] As an illustration of the application of p-values to the design and interpretation of experiments, in his following book The Design of Experiments (1935), Fisher presented the lady tasting tea experiment.[45] which is the archetypal example of the p-value. To evaluate a lady's claim that she (Muriel Bristol) could distinguish by taste how tea is prepared (first adding the milk to the cup, then the tea, or first tea, then milk), she was sequentially presented with 8 cups: 4 prepared one way, 4 prepared the other, and asked to determine the preparation of each cup (knowing that there were 4 of each). In that case, the null hypothesis was that she had no special ability; the test was Fisher's exact test, and the p-value

1

(

4

1
)

=
1

70

≈
0.014.

{\displaystyle 1/{\binom {3}{4}}=1/70\approx 0.014.}

 so Fisher was willing to reject the null hypothesis (consider the outcome highly unlikely to be due to chance) if all were classified correctly. (In the actual experiment, Bristol correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this level of significance.[46] Fisher also underlined the interpretation of

p

, as the long-run proportion of values at least as extreme as the data, assuming the null hypothesis is true. In later editions, Fisher explicitly contrasted the use of the p-value for statistical inference in science with the Neyman-Pearson method, which he terms "Acceptance Procedures".[47] Fisher emphasizes that while fixed levels such as 5%, 2%, and 1% are convenient, the exact p-value can be used, and the strength of evidence can and will be revised with further experimentation. In contrast, decision procedures require a clear-cut decision, yielding an irreversible action, and the procedure is based on costs of error, which, he argues, are inapplicable to scientific research. The E-value can refer to two concepts, both of which are related to the p-value and both of which play a role in multiple testing. First, it corresponds to a generic, more robust alternative to the p-value that can deal with optional continuation of experiments. Second, it is also used to abbreviate "expect value", which is the expected number of times that one expects to obtain a test statistic at least as extreme as the one that was actually observed if one assumes that the null hypothesis is true.[48] This expect-value is the product of the number of tests and the p-value. The q-value is the analog of the p-value with respect to the false discovery rate (FDR) and the multiple hypothesis testing to maintain statistical power while minimizing the false positive rate.[50] The Probability of Direction (pd) is the Bayesian numerical equivalent of the p-value.[51] It corresponds to the proportion of the posterior distribution that is the median's correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this level of significance.[46] Fisher also underlined the interpretation of

p

, as the long-run proportion of values at least as extreme as the data, assuming the null hypothesis is true. In later editions, Fisher explicitly contrasted the use of the p-value for statistical inference in science with the Neyman-Pearson method, which he terms "Acceptance Procedures".[47] Fisher emphasizes that while fixed levels such as 5%, 2%, and 1% are convenient, the exact p-value can be used, and the strength of evidence can and will be revised with further experimentation. In contrast, decision procedures require a clear-cut decision, yielding an irreversible action, and the procedure is based on costs of error, which, he argues, are inapplicable to scientific research. The E-value can refer to two concepts, both of which are related to the p-value and both of which play a role in multiple testing. First, it corresponds to a generic, more robust alternative to the p-value that can deal with optional continuation of experiments. Second, it is also used to abbreviate "expect value", which is the expected number of times that one expects to obtain a test statistic at least as extreme as the one that was actually observed if one assumes that the null hypothesis is true.[48] This expect-value is the product of the number of tests and the p-value. The q-value is the analog of the p-value with respect to the false discovery rate (FDR) and the multiple hypothesis testing to maintain statistical power while minimizing the false positive rate.[50] The Probability of Direction (pd) is the Bayesian numerical equivalent of the p-value.[51] It corresponds to the proportion of the posterior distribution that is the median's correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this level of significance.[46] Fisher also underlined the interpretation of

p

, as the long-run proportion of values at least as extreme as the data, assuming the null hypothesis is true. In later editions, Fisher explicitly contrasted the use of the p-value for statistical inference in science with the Neyman-Pearson method, which he terms "Acceptance Procedures".[47] Fisher emphasizes that while fixed levels such as 5%, 2%, and 1% are convenient, the exact p-value can be used, and the strength of evidence can and will be revised with further experimentation. In contrast, decision procedures require a clear-cut decision, yielding an irreversible action, and the procedure is based on costs of error, which, he argues, are inapplicable to scientific research. The E-value can refer to two concepts, both of which are related to the p-value and both of which play a role in multiple testing. First, it corresponds to a generic, more robust alternative to the p-value that can deal with optional continuation of experiments. Second, it is also used to abbreviate "expect value", which is the expected number of times that one expects to obtain a test statistic at least as extreme as the one that was actually observed if one assumes that the null hypothesis is true.[48] This expect-value is the product of the number of tests and the p-value. The q-value is the analog of the p-value with respect to the false discovery rate (FDR) and the multiple hypothesis testing to maintain statistical power while minimizing the false positive rate.[50] The Probability of Direction (pd) is the Bayesian numerical equivalent of the p-value.[51] It corresponds to the proportion of the posterior distribution that is the median's correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this level of significance.[46] Fisher also underlined the interpretation of

p

, as the long-run proportion of values at least as extreme as the data, assuming the null hypothesis is true. In later editions, Fisher explicitly contrasted the use of the p-value for statistical inference in science with the Neyman-Pearson method, which he terms "Acceptance Procedures".[47] Fisher emphasizes that while fixed levels such as 5%, 2%, and 1% are convenient, the exact p-value can be used, and the strength of evidence can and will be revised with further experimentation. In contrast, decision procedures require a clear-cut decision, yielding an irreversible action, and the procedure is based on costs of error, which, he argues, are inapplicable to scientific research. The E-value can refer to two concepts, both of which are related to the p-value and both of which play a role in multiple testing. First, it corresponds to a generic, more robust alternative to the p-value that can deal with optional continuation of experiments. Second, it is also used to abbreviate "expect value", which is the expected number of times that one expects to obtain a test statistic at least as extreme as the one that was actually observed if one assumes that the null hypothesis is true.[48] This expect-value is the product of the number of tests and the p-value. The q-value is the analog of the p-value with respect to the false discovery rate (FDR) and the multiple hypothesis testing to maintain statistical power while minimizing the false positive rate.[50] The Probability of Direction (pd) is the Bayesian numerical equivalent of the p-value.[51] It corresponds to the proportion of the posterior distribution that is the median's correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this level of significance.[46] Fisher also underlined the interpretation of

p

, as the long-run proportion of values at least as extreme as the data, assuming the null hypothesis is true. In later editions, Fisher explicitly contrasted the use of the p-value for statistical inference in science with the Neyman-Pearson method, which he terms "Acceptance Procedures".[47] Fisher emphasizes that while fixed levels such as 5%, 2%, and 1% are convenient, the exact p-value can be used, and the strength of evidence can and will be revised with further experimentation. In contrast, decision procedures require a clear-cut decision, yielding an irreversible action, and the procedure is based on costs of error, which, he argues, are inapplicable to scientific research. The E-value can refer to two concepts, both of which are related to the p-value and both of which play a role in multiple testing. First, it corresponds to a generic, more robust alternative to the p-value that can deal with optional continuation of experiments. Second, it is also used to abbreviate "expect value", which is the expected number of times that one expects to obtain a test statistic at least as extreme as the one that was actually observed if one assumes that the null hypothesis is true.[48] This expect-value is the product of the number of tests and the p-value. The q-value is the analog of the p-value with respect to the false discovery rate (FDR) and the multiple hypothesis testing to maintain statistical power while minimizing the false positive rate.[50] The Probability of Direction (pd) is the Bayesian numerical equivalent of the p-value.[51] It corresponds to the proportion of the posterior distribution that is the median's correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this level of significance.[46] Fisher also underlined the interpretation of

p

, as the long-run proportion of values at least as extreme as the data, assuming the null hypothesis is true. In later editions, Fisher explicitly contrasted the use of the p-value for statistical inference in science with the Neyman-Pearson method, which he terms "Acceptance Procedures".[47] Fisher emphasizes that while fixed levels such as 5%, 2%, and 1% are convenient, the exact p-value can be used, and the strength of evidence can and will be revised with further experimentation. In contrast, decision procedures require a clear-cut decision, yielding an irreversible action, and the procedure is based on costs of error, which, he argues, are inapplicable to scientific research. The E-value can refer to two concepts, both of which are related to the p-value and both of which play a role in multiple testing. First, it corresponds to a generic, more robust alternative to the p-value that can deal with optional continuation of experiments. Second, it is also used to abbreviate "expect value", which is the expected number of times that one expects to obtain a test statistic at least as extreme as the one that was actually observed if one assumes that the null hypothesis is true.[48] This expect-value is the product of the number of tests and the p-value. The q-value is the analog of the p-value with respect to the false discovery rate (FDR) and the multiple hypothesis testing to maintain statistical power while minimizing the false positive rate.[50] The Probability of Direction (pd) is the Bayesian numerical equivalent of the p-value.[51] It corresponds to the proportion of the posterior distribution that is the median's correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this level of significance.[46] Fisher also underlined the interpretation of

p

, as the long-run proportion of values at least as extreme as the data, assuming the null hypothesis is true. In later editions, Fisher explicitly contrasted the use of the p-value for statistical inference in science with the Neyman-Pearson method, which he terms "Acceptance Procedures".[47] Fisher emphasizes that while fixed levels such as 5%, 2%, and 1% are convenient, the exact p-value can be used, and the strength of evidence can and will be revised with further experimentation. In contrast, decision procedures require a clear-cut decision, yielding an irreversible action, and the procedure is based on costs of error, which, he argues, are inapplicable to scientific research. The E-value can refer to two concepts, both of which are related to the p-value and both of which play a role in multiple testing. First, it corresponds to a generic, more robust alternative to the p-value that can deal with optional continuation of experiments. Second, it is also used to abbreviate "expect value", which is the expected number of times that one expects to obtain a test statistic at least as extreme as the one that was actually observed if one assumes that the null hypothesis is true.[48] This expect-value is the product of the number of tests and the p-value. The q-value is the analog of the p-value with respect to the false discovery rate (FDR) and the multiple hypothesis testing to maintain statistical power while minimizing the false positive rate.[50] The Probability of Direction (pd) is the Bayesian numerical equivalent of the p-value.[51] It corresponds to the proportion of the posterior distribution that is the median's correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this level of significance.[46] Fisher also underlined the interpretation of

p

, as the long-run proportion of values at least as extreme as the data, assuming the null hypothesis is true. In later editions, Fisher explicitly contrasted the use of the p-value for statistical inference in science with the Neyman-Pearson method, which he terms "Acceptance Procedures".[47] Fisher emphasizes that while fixed levels such as 5%, 2%, and 1% are convenient, the exact p-value can be used, and the strength of evidence can and will be revised with further experimentation. In contrast, decision procedures require a clear-cut decision, yielding an irreversible action, and the procedure is based on costs of error, which, he argues, are inapplicable to scientific research. The E-value can refer to two concepts, both of which are related to the p-value and both of which play a role in multiple testing. First, it corresponds to a generic, more robust alternative to the p-value that can deal with optional continuation of experiments. Second, it is also used to abbreviate "expect value", which is the expected number of times that one expects to obtain a test statistic at least as extreme as the one that was actually observed if one assumes that the null hypothesis is true.[48] This expect-value is the product of the number of tests and the p-value. The q-value is the analog of the p-value with respect to the false discovery rate (FDR) and the multiple hypothesis testing to maintain statistical power while minimizing the false positive rate.[50] The Probability of Direction (pd) is the Bayesian numerical equivalent of the p-value.[51] It corresponds to the proportion of the posterior distribution that is the median's correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this level of significance.[46] Fisher also underlined the interpretation of

p

, as the long-run proportion of values at least as extreme as the data, assuming the null hypothesis is true. In later editions, Fisher explicitly contrasted the use of the p-value for statistical inference in science with the Neyman-Pearson method, which he terms "Acceptance Procedures".[47] Fisher emphasizes that while fixed levels such as 5%, 2%, and 1% are convenient, the exact p-value can be used, and the strength of evidence can and will be revised with further experimentation. In contrast, decision procedures require a clear-cut decision, yielding an irreversible action, and the procedure is based on costs of error, which, he argues, are inapplicable to scientific research. The E-value can refer to two concepts, both of which are related to the p-value and both of which play a role in multiple testing. First, it corresponds to a generic, more robust alternative to the p-value that can deal with optional continuation of experiments. Second, it is also used to abbreviate "expect value", which is the expected number of times that one expects to obtain a test statistic at least as extreme as the one that was actually observed if one assumes that the null hypothesis is true.[48] This expect-value is the product of the number of tests and the p-value. The q-value is the analog of the p-value with respect to the false discovery rate (FDR) and the multiple hypothesis testing to maintain statistical power while minimizing the false positive rate.[50] The Probability of Direction (pd) is the Bayesian numerical equivalent of the p-value.[51] It corresponds to the proportion of the posterior distribution that is the median's correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this level of significance.[46] Fisher also underlined the interpretation of

p

, as the long-run proportion of values at least as extreme as the data, assuming the null hypothesis is true. In later editions, Fisher explicitly contrasted the use of the p-value for statistical inference in science with the Neyman-Pearson method, which he terms "Acceptance Procedures".[47] Fisher emphasizes that while fixed levels such as 5%, 2%, and 1% are convenient, the exact p-value can be used, and the strength of evidence can and will be revised with further experimentation. In contrast, decision procedures require a clear-cut decision, yielding an irreversible action, and the procedure is based on costs of error, which, he argues, are inapplicable to scientific research. The E-value can refer to two concepts, both of which are related to the p-value and both of which play a role in multiple testing. First, it corresponds to a generic, more robust alternative to the p-value that can deal with optional continuation of experiments. Second, it is also used to abbreviate "expect value", which is the expected number of times that one expects to obtain a test statistic at least as extreme as the one that was actually observed if one assumes that the null hypothesis is true.[48] This expect-value is the product of the number of tests and the p-value. The q-value is the analog of the p-value with respect to the false discovery rate (FDR) and the multiple hypothesis testing to maintain statistical power while minimizing the false positive rate.[50] The Probability of Direction (pd) is the Bayesian numerical equivalent of the p-value.[51] It corresponds to the proportion of the posterior distribution that is the median's correctly classified all 8 cups.) Fisher reiterated the

p
=
0.05

 threshold and explained its rationale, stating:[46] It is usual and convenient for experimenters to take 5 per cent as a standard level of significance, in the sense that they are prepared to ignore all results which fail to reach this standard, and, by this means, to eliminate from further discussion the greater part of the fluctuations which chance causes have introduced into their experimental results. He also applies this threshold to the design of experiments, noting that had only 6 cups been presented (3 of each), a perfect classification would have only yielded a p-value of

1

(

6

3
)

=
1

20
=
0.05.

{\displaystyle 1/{\binom {6}{3}}=1/20=0.05.}

 which would not have met this