

32. STATISTICS

Revised September 2007 by G. Cowan (RHUL).

This chapter gives an overview of statistical methods used in high-energy physics. In statistics, we are interested in using a given sample of data to make inferences about a probabilistic model, *e.g.*, to assess the model's validity or to determine the values of its parameters. There are two main approaches to statistical inference, which we may call frequentist and Bayesian. In frequentist statistics, probability is interpreted as the frequency of the outcome of a repeatable experiment. The most important tools in this framework are parameter estimation, covered in Section 32.1, and statistical tests, discussed in Section 32.2. Frequentist confidence intervals, which are constructed so as to cover the true value of a parameter with a specified probability, are treated in Section 32.3.2. Note that in frequentist statistics one does not define a probability for a hypothesis or for a parameter.

Frequentist statistics provides the usual tools for reporting the outcome of an experiment objectively, without needing to incorporate prior beliefs concerning the parameter being measured or the theory being tested. As such, they are used for reporting most measurements and their statistical uncertainties in high-energy physics.

In Bayesian statistics, the interpretation of probability is more general and includes *degree of belief* (called subjective probability). One can then speak of a probability density function (p.d.f.) for a parameter, which expresses one's state of knowledge about where its true value lies. Bayesian methods allow for a natural way to input additional information, such as physical boundaries and subjective information; in fact they *require* the *prior* p.d.f. as input for the parameters, *i.e.*, the degree of belief about the parameters' values before carrying out the measurement. Using Bayes' theorem Eq. (31.4), the prior degree of belief is updated by the data from the experiment. Bayesian methods for interval estimation are discussed in Sections 32.3.1 and 32.3.2.6

Bayesian techniques are often used to treat systematic uncertainties, where the author's beliefs about, say, the accuracy of the measuring device may enter. Bayesian statistics also provides a useful framework for discussing the validity of different theoretical interpretations of the data. This aspect of a measurement, however, will usually be treated separately from the reporting of the result.

For many inference problems, the frequentist and Bayesian approaches give similar numerical answers, even though they are based on fundamentally different interpretations of probability. For small data samples, however, and for measurements of a parameter near a physical boundary, the different approaches may yield different results, so we are forced to make a choice. For a discussion of Bayesian vs. non-Bayesian methods, see References written by a statistician[1], by a physicist[2], or the more detailed comparison in Ref. [3].

Following common usage in physics, the word "error" is often used in this chapter to mean "uncertainty." More specifically it can indicate the size of an interval as in "the standard error" or "error propagation," where the term refers to the standard deviation of an estimator.

2 32. Statistics

32.1. Parameter estimation

Here we review the frequentist approach to *point estimation* of parameters. An *estimator* $\hat{\theta}$ (written with a hat) is a function of the data whose value, the *estimate*, is intended as a meaningful guess for the value of the parameter θ .

There is no fundamental rule dictating how an estimator must be constructed. One tries, therefore, to choose that estimator which has the best properties. The most important of these are (a) *consistency*, (b) *bias*, (c) *efficiency*, and (d) *robustness*.

(a) An estimator is said to be *consistent* if the estimate $\hat{\theta}$ converges to the true value θ as the amount of data increases. This property is so important that it is possessed by all commonly used estimators.

(b) The *bias*, $b = E[\hat{\theta}] - \theta$, is the difference between the expectation value of the estimator and the true value of the parameter. The expectation value is taken over a hypothetical set of similar experiments in which $\hat{\theta}$ is constructed in the same way. When $b = 0$, the estimator is said to be unbiased. The bias depends on the chosen metric, *i.e.*, if $\hat{\theta}$ is an unbiased estimator of θ , then $\hat{\theta}^2$ is not in general an unbiased estimator for θ^2 . If we have an estimate \hat{b} for the bias, we can subtract it from $\hat{\theta}$ to obtain a new $\hat{\theta}' = \hat{\theta} - \hat{b}$. The estimate \hat{b} may, however, be subject to statistical or systematic uncertainties that are larger than the bias itself, so that the new $\hat{\theta}'$ may not be better than the original.

(c) *Efficiency* is the inverse of the ratio of the variance $V[\hat{\theta}]$ to its minimum possible value. Under rather general conditions, the minimum variance is given by the Rao-Cramér-Frechet bound,

$$\sigma_{\min}^2 = \left(1 + \frac{\partial b}{\partial \theta}\right)^2 / I(\theta), \quad (32.1)$$

where

$$I(\theta) = E \left[\left(\frac{\partial}{\partial \theta} \sum_i \ln f(x_i; \theta) \right)^2 \right] \quad (32.2)$$

is the *Fisher information*. The sum is over all data, assumed independent, and distributed according to the p.d.f. $f(x; \theta)$, b is the bias, if any, and the allowed range of x must not depend on θ .

The *mean-squared error*,

$$\text{MSE} = E[(\hat{\theta} - \theta)^2] = V[\hat{\theta}] + b^2, \quad (32.3)$$

is a convenient quantity which combines the uncertainties in an estimate due to bias and variance.

(d) *Robustness* is the property of being insensitive to departures from assumptions in the p.d.f., *e.g.*, owing to uncertainties in the distribution's tails.

For some common estimators, the properties above are known exactly. More generally, it is possible to evaluate them by Monte Carlo simulation. Note that they will often depend on the unknown θ .

32.1.1. Estimators for mean, variance and median :

Suppose we have a set of N independent measurements, x_i , assumed to be unbiased measurements of the same unknown quantity μ with a common, but unknown, variance σ^2 . Then

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i \quad (32.4)$$

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu})^2 \quad (32.5)$$

are unbiased estimators of μ and σ^2 . The variance of $\hat{\mu}$ is σ^2/N and the variance of $\hat{\sigma}^2$ is

$$V[\hat{\sigma}^2] = \frac{1}{N} \left(m_4 - \frac{N-3}{N-1} \sigma^4 \right), \quad (32.6)$$

where m_4 is the 4th central moment of x . For Gaussian distributed x_i , this becomes $2\sigma^4/(N-1)$ for any $N \geq 2$, and for large N , the standard deviation of $\hat{\sigma}$ (the “error of the error”) is $\sigma/\sqrt{2N}$. Again, if the x_i are Gaussian, $\hat{\mu}$ is an efficient estimator for μ , and the estimators $\hat{\mu}$ and $\hat{\sigma}^2$ are uncorrelated. Otherwise the arithmetic mean (32.4) is not necessarily the most efficient estimator; this is discussed in more detail in Sec. 8.7 [4].

If σ^2 is known, it does not improve the estimate $\hat{\mu}$, as can be seen from Eq. (32.4); however, if μ is known, substitute it for $\hat{\mu}$ in Eq. (32.5) and replace $N-1$ by N to obtain an estimator of σ^2 still with zero bias but smaller variance. If the x_i have different, known variances σ_i^2 , then the weighted average

$$\hat{\mu} = \frac{1}{w} \sum_{i=1}^N w_i x_i \quad (32.7)$$

is an unbiased estimator for μ with a smaller variance than an unweighted average; here $w_i = 1/\sigma_i^2$ and $w = \sum_i w_i$. The standard deviation of $\hat{\mu}$ is $1/\sqrt{w}$.

As an estimator for the median x_{med} , one can use the value \hat{x}_{med} such that half the x_i are below and half above (the sample median). If the sample median lies between two observed values, it is set by convention halfway between them. If the p.d.f. of x has the form $f(x - \mu)$ and μ is both mean and median, then for large N the variance of the sample median approaches $1/[4Nf^2(0)]$, provided $f(0) > 0$. Although estimating the median can often be more difficult computationally than the mean, the resulting estimator is generally more robust, as it is insensitive to the exact shape of the tails of a distribution.

4 32. Statistics

32.1.2. The method of maximum likelihood :

Suppose we have a set of N measured quantities $\mathbf{x} = (x_1, \dots, x_N)$ described by a joint p.d.f. $f(\mathbf{x}; \boldsymbol{\theta})$, where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)$ is set of n parameters whose values are unknown. The *likelihood function* is given by the p.d.f. evaluated with the data \mathbf{x} , but viewed as a function of the parameters, *i.e.*, $L(\boldsymbol{\theta}) = f(\mathbf{x}; \boldsymbol{\theta})$. If the measurements x_i are statistically independent and each follow the p.d.f. $f(x; \boldsymbol{\theta})$, then the joint p.d.f. for \mathbf{x} factorizes and the likelihood function is

$$L(\boldsymbol{\theta}) = \prod_{i=1}^N f(x_i; \boldsymbol{\theta}) . \quad (32.8)$$

The method of maximum likelihood takes the estimators $\hat{\boldsymbol{\theta}}$ to be those values of $\boldsymbol{\theta}$ that maximize $L(\boldsymbol{\theta})$.

Note that the likelihood function is *not* a p.d.f. for the parameters $\boldsymbol{\theta}$; in frequentist statistics this is not defined. In Bayesian statistics, one can obtain from the likelihood the posterior p.d.f. for $\boldsymbol{\theta}$, but this requires multiplying by a prior p.d.f. (see Sec. 32.3.1).

It is usually easier to work with $\ln L$, and since both are maximized for the same parameter values $\boldsymbol{\theta}$, the maximum likelihood (ML) estimators can be found by solving the *likelihood equations*,

$$\frac{\partial \ln L}{\partial \theta_i} = 0 , \quad i = 1, \dots, n . \quad (32.9)$$

Maximum likelihood estimators are important because they are approximately unbiased and efficient for large data samples, under quite general conditions, and the method has a wide range of applicability.

In evaluating the likelihood function, it is important that any normalization factors in the p.d.f. that involve $\boldsymbol{\theta}$ be included. However, we will only be interested in the maximum of L and in ratios of L at different values of the parameters; hence any multiplicative factors that do not involve the parameters that we want to estimate may be dropped, including factors that depend on the data but not on $\boldsymbol{\theta}$.

Under a one-to-one change of parameters from $\boldsymbol{\theta}$ to $\boldsymbol{\eta}$, the ML estimators $\hat{\boldsymbol{\theta}}$ transform to $\boldsymbol{\eta}(\hat{\boldsymbol{\theta}})$. That is, the ML solution is invariant under change of parameter. However, other properties of ML estimators, in particular the bias, are not invariant under change of parameter.

The inverse V^{-1} of the covariance matrix $V_{ij} = \text{cov}[\hat{\theta}_i, \hat{\theta}_j]$ for a set of ML estimators can be estimated by using

$$(\hat{V}^{-1})_{ij} = - \left. \frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \right|_{\hat{\boldsymbol{\theta}}} . \quad (32.10)$$

For finite samples, however, Eq. (32.10) can result in an underestimate of the variances. In the large sample limit (or in a linear model with Gaussian errors), L has a Gaussian

form and $\ln L$ is (hyper)parabolic. In this case, it can be seen that a numerically equivalent way of determining s -standard-deviation errors is from the contour given by the $\boldsymbol{\theta}'$ such that

$$\ln L(\boldsymbol{\theta}') = \ln L_{\max} - s^2/2, \quad (32.11)$$

where $\ln L_{\max}$ is the value of $\ln L$ at the solution point (compare with Eq. (32.48)). The extreme limits of this contour on the θ_i axis give an approximate s -standard-deviation confidence interval for θ_i (see Section 32.3.2.4).

In the case where the size n of the data sample x_1, \dots, x_n is small, the unbinned maximum likelihood method, *i.e.*, use of equation (32.8), is preferred since binning can only result in a loss of information, and hence larger statistical errors for the parameter estimates. The sample size n can be regarded as fixed, or the user can choose to treat it as a Poisson-distributed variable; this latter option is sometimes called “extended maximum likelihood” (see, *e.g.*, [6–8]). If the sample is large, it can be convenient to bin the values in a histogram, so that one obtains a vector of data $\mathbf{n} = (n_1, \dots, n_N)$ with expectation values $\boldsymbol{\nu} = E[\mathbf{n}]$ and probabilities $f(\mathbf{n}; \boldsymbol{\nu})$. Then one may maximize the likelihood function based on the contents of the bins (so i labels bins). This is equivalent to maximizing the likelihood ratio $\lambda(\boldsymbol{\theta}) = f(\mathbf{n}; \boldsymbol{\nu}(\boldsymbol{\theta})) / f(\mathbf{n}; \mathbf{n})$, or to minimizing the quantity [9]

$$-2 \ln \lambda(\boldsymbol{\theta}) = 2 \sum_{i=1}^N \left[\nu_i(\boldsymbol{\theta}) - n_i + n_i \ln \frac{n_i}{\nu_i(\boldsymbol{\theta})} \right], \quad (32.12)$$

where in bins where $n_i = 0$, the last term in (32.12) is zero. In the limit of zero bin width, maximizing (32.12) is equivalent to maximizing the unbinned likelihood function (32.8).

A benefit of binning is that it allows for a goodness-of-fit test (see Sec. 32.2.2). The minimum of $-2 \ln \lambda$ as defined by Eq. (32.12) follows a χ^2 distribution in the large sample limit. If there are N bins and m fitted parameters, then the number of degrees of freedom for the χ^2 distribution is $N - m$ if the data are treated as Poisson-distributed, and $N - m - 1$ if the n_i are multinomially distributed. If the n_i are Poisson-distributed and the overall normalization $\nu_{\text{tot}} = \sum_i \nu_i$ is taken as an adjustable parameter, then by minimizing Eq. (32.12), one obtains that the area under the fitted function is equal to the sum of the histogram contents, *i.e.*, $\sum_i \nu_i = \sum_i n_i$. This is not the case for parameter estimation methods based on a least-squares procedure with traditional weights (see, *e.g.*, Ref. 8).

6 32. Statistics

32.1.3. The method of least squares :

The *method of least squares* (LS) coincides with the method of maximum likelihood in the following special case. Consider a set of N independent measurements y_i at known points x_i . The measurement y_i is assumed to be Gaussian distributed with mean $F(x_i; \boldsymbol{\theta})$ and known variance σ_i^2 . The goal is to construct estimators for the unknown parameters $\boldsymbol{\theta}$. The likelihood function contains the sum of squares

$$\chi^2(\boldsymbol{\theta}) = -2 \ln L(\boldsymbol{\theta}) + \text{constant} = \sum_{i=1}^N \frac{(y_i - F(x_i; \boldsymbol{\theta}))^2}{\sigma_i^2}. \quad (32.13)$$

The set of parameters $\boldsymbol{\theta}$ which maximize L is the same as those which minimize χ^2 .

The minimum of Equation (32.13) defines the least-squares estimators $\hat{\boldsymbol{\theta}}$ for the more general case where the y_i are not Gaussian distributed as long as they are independent. If they are not independent but rather have a covariance matrix $V_{ij} = \text{cov}[y_i, y_j]$, then the LS estimators are determined by the minimum of

$$\chi^2(\boldsymbol{\theta}) = (\mathbf{y} - \mathbf{F}(\boldsymbol{\theta}))^T V^{-1} (\mathbf{y} - \mathbf{F}(\boldsymbol{\theta})), \quad (32.14)$$

where $\mathbf{y} = (y_1, \dots, y_N)$ is the vector of measurements, $\mathbf{F}(\boldsymbol{\theta})$ is the corresponding vector of predicted values (understood as a column vector in (32.14)), and the superscript T denotes transposed (*i.e.*, row) vector.

In many practical cases, one further restricts the problem to the situation where $F(x_i; \boldsymbol{\theta})$ is a linear function of the parameters, *i.e.*,

$$F(x_i; \boldsymbol{\theta}) = \sum_{j=1}^m \theta_j h_j(x_i). \quad (32.15)$$

Here the $h_j(x)$ are m linearly independent functions, *e.g.*, $1, x, x^2, \dots, x^{m-1}$, or Legendre polynomials. We require $m < N$ and at least m of the x_i must be distinct.

Minimizing χ^2 in this case with m parameters reduces to solving a system of m linear equations. Defining $H_{ij} = h_j(x_i)$ and minimizing χ^2 by setting its derivatives with respect to the θ_i equal to zero gives the LS estimators,

$$\hat{\boldsymbol{\theta}} = (H^T V^{-1} H)^{-1} H^T V^{-1} \mathbf{y} \equiv D \mathbf{y}. \quad (32.16)$$

The covariance matrix for the estimators $U_{ij} = \text{cov}[\hat{\theta}_i, \hat{\theta}_j]$ is given by

$$U = D V D^T = (H^T V^{-1} H)^{-1}, \quad (32.17)$$

or equivalently, its inverse U^{-1} can be found from

$$(U^{-1})_{ij} = \frac{1}{2} \frac{\partial^2 \chi^2}{\partial \theta_i \partial \theta_j} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} = \sum_{k,l=1}^N h_i(x_k) (V^{-1})_{kl} h_j(x_l). \quad (32.18)$$

The LS estimators can also be found from the expression

$$\hat{\boldsymbol{\theta}} = U\mathbf{g}, \quad (32.19)$$

where the vector \mathbf{g} is defined by

$$g_i = \sum_{j,k=1}^N y_j h_i(x_k) (V^{-1})_{jk}. \quad (32.20)$$

For the case of uncorrelated y_i , for example, one can use (32.19) with

$$(U^{-1})_{ij} = \sum_{k=1}^N \frac{h_i(x_k) h_j(x_k)}{\sigma_k^2}, \quad (32.21)$$

$$g_i = \sum_{k=1}^N \frac{y_k h_i(x_k)}{\sigma_k^2}. \quad (32.22)$$

Expanding $\chi^2(\boldsymbol{\theta})$ about $\hat{\boldsymbol{\theta}}$, one finds that the contour in parameter space defined by

$$\chi^2(\boldsymbol{\theta}) = \chi^2(\hat{\boldsymbol{\theta}}) + 1 = \chi_{\min}^2 + 1 \quad (32.23)$$

has tangent planes located at plus-or-minus-one standard deviation $\sigma_{\hat{\boldsymbol{\theta}}}$ from the LS estimates $\hat{\boldsymbol{\theta}}$.

In constructing the quantity $\chi^2(\boldsymbol{\theta})$, one requires the variances or, in the case of correlated measurements, the covariance matrix. Often these quantities are not known *a priori* and must be estimated from the data; an important example is where the measured value y_i represents a counted number of events in the bin of a histogram. If, for example, y_i represents a Poisson variable, for which the variance is equal to the mean, then one can either estimate the variance from the predicted value, $F(x_i; \boldsymbol{\theta})$, or from the observed number itself, y_i . In the first option, the variances become functions of the fitted parameters, which may lead to calculational difficulties. The second option can be undefined if y_i is zero, and in both cases for small y_i , the variance will be poorly estimated. In either case, one should constrain the normalization of the fitted curve to the correct value, *i.e.*, one should determine the area under the fitted curve directly from the number of entries in the histogram (see Ref. 8, Section 7.4). A further alternative is to use the method of maximum likelihood; for binned data this can be done by minimizing Eq. (32.12)

As the minimum value of the χ^2 represents the level of agreement between the measurements and the fitted function, it can be used for assessing the goodness-of-fit; this is discussed further in Section 32.2.2.

8 32. Statistics

32.1.4. Propagation of errors :

Consider a set of n quantities $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)$ and a set of m functions $\boldsymbol{\eta}(\boldsymbol{\theta}) = (\eta_1(\boldsymbol{\theta}), \dots, \eta_m(\boldsymbol{\theta}))$. Suppose we have estimated $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_n)$, using, say, maximum-likelihood or least-squares, and we also know or have estimated the covariance matrix $V_{ij} = \text{cov}[\hat{\theta}_i, \hat{\theta}_j]$. The goal of *error propagation* is to determine the covariance matrix for the functions, $U_{ij} = \text{cov}[\hat{\eta}_i, \hat{\eta}_j]$, where $\hat{\boldsymbol{\eta}} = \boldsymbol{\eta}(\hat{\boldsymbol{\theta}})$. In particular, the diagonal elements $U_{ii} = V[\hat{\eta}_i]$ give the variances. The new covariance matrix can be found by expanding the functions $\boldsymbol{\eta}(\boldsymbol{\theta})$ about the estimates $\hat{\boldsymbol{\theta}}$ to first order in a Taylor series. Using this one finds

$$U_{ij} \approx \sum_{k,l} \left. \frac{\partial \eta_i}{\partial \theta_k} \frac{\partial \eta_j}{\partial \theta_l} \right|_{\hat{\boldsymbol{\theta}}} V_{kl}. \quad (32.24)$$

This can be written in matrix notation as $U \approx AVA^T$ where the matrix of derivatives A is

$$A_{ij} = \left. \frac{\partial \eta_i}{\partial \theta_j} \right|_{\hat{\boldsymbol{\theta}}}, \quad (32.25)$$

and A^T is its transpose. The approximation is exact if $\boldsymbol{\eta}(\boldsymbol{\theta})$ is linear (it holds, for example, in equation (32.17)). If this is not the case, the approximation can break down if, for example, $\boldsymbol{\eta}(\boldsymbol{\theta})$ is significantly nonlinear close to $\hat{\boldsymbol{\theta}}$ in a region of a size comparable to the standard deviations of $\hat{\boldsymbol{\theta}}$.

32.2. Statistical tests

In addition to estimating parameters, one often wants to assess the validity of certain statements concerning the data's underlying distribution. *Hypothesis tests* provide a rule for accepting or rejecting hypotheses depending on the outcome of a measurement. In *significance tests*, one gives the probability to obtain a level of incompatibility with a certain hypothesis that is greater than or equal to the level observed with the actual data.

32.2.1. Hypothesis tests :

Consider an experiment whose outcome is characterized by a vector of data \boldsymbol{x} . A *hypothesis* is a statement about the distribution of \boldsymbol{x} . It could, for example, define completely the p.d.f. for the data (a simple hypothesis), or it could specify only the functional form of the p.d.f., with the values of one or more parameters left open (a composite hypothesis).

A *statistical test* is a rule that states for which values of \boldsymbol{x} a given hypothesis (often called the null hypothesis, H_0) should be rejected in favor of its complementary alternative H_1 . This is done by defining a region of \boldsymbol{x} -space called the critical region; if the outcome of the experiment lands in this region, H_0 is rejected, otherwise it is accepted.

Rejecting H_0 if it is true is called an error of the first kind. The probability for this to occur is called the *size* or *significance level* of the test, α , which is chosen to be equal to

some pre-specified value. It can also happen that H_0 is false and the true hypothesis is the alternative, H_1 . If H_0 is accepted in such a case, this is called an error of the second kind, which will have some probability β . The quantity $1 - \beta$ is called the *power* of the test to reject H_1 .

In high-energy physics, the components of \mathbf{x} might represent the measured properties of candidate events, and the acceptance region is defined by the cuts that one imposes in order to select events of a certain desired type. That is, H_0 could represent the signal hypothesis, and various alternatives, H_1 , H_2 , *etc.*, could represent background processes.

Often rather than using the full set of quantities \mathbf{x} , it is convenient to define a *test statistic*, t , which can be a single number, or in any case a vector with fewer components than \mathbf{x} . Each hypothesis for the distribution of \mathbf{x} will determine a distribution for t , and the acceptance region in \mathbf{x} -space will correspond to a specific range of values of t . In constructing t , one attempts to reduce the volume of data without losing the ability to discriminate between different hypotheses.

In particle physics terminology, the probability to accept the signal hypothesis, H_0 , is the selection efficiency, *i.e.*, one minus the significance level. The efficiencies for the various background processes are given by one minus the power. Often one tries to construct a test to minimize the background efficiency for a given signal efficiency. The *Neyman–Pearson lemma* states that this is done by defining the acceptance region such that, for \mathbf{x} in that region, the ratio of p.d.f.s for the hypotheses H_0 and H_1 ,

$$\lambda(\mathbf{x}) = \frac{f(\mathbf{x}|H_0)}{f(\mathbf{x}|H_1)}, \quad (32.26)$$

is greater than a given constant, the value of which is chosen to give the desired signal efficiency. This is equivalent to the statement that (32.26) represents the test statistic with which one may obtain the highest purity sample for a given signal efficiency. It can be difficult in practice, however, to determine $\lambda(\mathbf{x})$, since this requires knowledge of the joint p.d.f.s $f(\mathbf{x}|H_0)$ and $f(\mathbf{x}|H_1)$.

In the usual case where the likelihood ratio (32.26) cannot be used explicitly, there exist a variety of other multivariate classifiers that effectively separate different types of events. Methods often used in HEP include *neural networks* or *Fisher discriminants* (see Ref. 10). Recently, further classification methods from machine-learning have been applied in HEP analyses; these include *probability density estimation (PDE)* techniques, *kernel-based PDE (KDE or Parzen window)*, *support vector machines*, and *decision trees*. Techniques such as “boosting” and “bagging” can be applied to combine a number of classifiers into a stronger one with greater stability with respect to fluctuations in the training data. Descriptions of these methods can be found in [11–13], and Proceedings of the PHYSTAT conference series [14]. Software for HEP includes the `TMVA` [15] and `StatPatternRecognition` [16] packages.

10 32. Statistics

32.2.2. Significance tests :

Often one wants to quantify the level of agreement between the data and a hypothesis without explicit reference to alternative hypotheses. This can be done by defining a statistic t , which is a function of the data whose value reflects in some way the level of agreement between the data and the hypothesis. The user must decide what values of the statistic correspond to better or worse levels of agreement with the hypothesis in question; for many goodness-of-fit statistics, there is an obvious choice.

The hypothesis in question, say, H_0 , will determine the p.d.f. $g(t|H_0)$ for the statistic. The significance of a discrepancy between the data and what one expects under the assumption of H_0 is quantified by giving the p -value, defined as the probability to find t in the region of equal or lesser compatibility with H_0 than the level of compatibility observed with the actual data. For example, if t is defined such that large values correspond to poor agreement with the hypothesis, then the p -value would be

$$p = \int_{t_{\text{obs}}}^{\infty} g(t|H_0) dt , \quad (32.27)$$

where t_{obs} is the value of the statistic obtained in the actual experiment. The p -value should not be confused with the size (significance level) of a test, or the confidence level of a confidence interval (Section 32.3), both of which are pre-specified constants.

The p -value is a function of the data, and is therefore itself a random variable. If the hypothesis used to compute the p -value is true, then for continuous data, p will be uniformly distributed between zero and one. Note that the p -value is not the probability for the hypothesis; in frequentist statistics, this is not defined. Rather, the p -value is the probability, under the assumption of a hypothesis H_0 , of obtaining data at least as incompatible with H_0 as the data actually observed.

When estimating parameters using the method of least squares, one obtains the minimum value of the quantity χ^2 (32.13). This statistic can be used to test the *goodness-of-fit*, *i.e.*, the test provides a measure of the significance of a discrepancy between the data and the hypothesized functional form used in the fit. It may also happen that no parameters are estimated from the data, but that one simply wants to compare a histogram, *e.g.*, a vector of Poisson distributed numbers $\mathbf{n} = (n_1, \dots, n_N)$, with a hypothesis for their expectation values $\nu_i = E[n_i]$. As the distribution is Poisson with variances $\sigma_i^2 = \nu_i$, the χ^2 (32.13) becomes *Pearson's χ^2 statistic*,

$$\chi^2 = \sum_{i=1}^N \frac{(n_i - \nu_i)^2}{\nu_i} . \quad (32.28)$$

If the hypothesis $\boldsymbol{\nu} = (\nu_1, \dots, \nu_N)$ is correct, and if the measured values n_i in (32.28) are sufficiently large (in practice, this will be a good approximation if all $n_i > 5$), then the χ^2 statistic will follow the χ^2 p.d.f. with the number of degrees of freedom equal to the number of measurements N minus the number of fitted parameters. The same holds for the minimized χ^2 from Eq. (32.13) if the y_i are Gaussian.

Alternatively, one may fit parameters and evaluate goodness-of-fit by minimizing $-2 \ln \lambda$ from Eq. (32.12). One finds that the distribution of this statistic approaches the asymptotic limit faster than does Pearson's χ^2 , and thus computing the p -value with the χ^2 p.d.f. will in general be better justified (see Ref. 9 and references therein).

Assuming the goodness-of-fit statistic follows a χ^2 p.d.f., the p -value for the hypothesis is then

$$p = \int_{\chi^2}^{\infty} f(z; n_d) dz, \quad (32.29)$$

where $f(z; n_d)$ is the χ^2 p.d.f. and n_d is the appropriate number of degrees of freedom. Values can be obtained from Fig. 32.1 or from the CERNLIB routine PROB or the ROOT function TMath::Prob. If the conditions for using the χ^2 p.d.f. do not hold, the statistic can still be defined as before, but its p.d.f. must be determined by other means in order to obtain the p -value, *e.g.*, using a Monte Carlo calculation.

If one finds a χ^2 value much greater than n_d , and a correspondingly small p -value, one may be tempted to expect a high degree of uncertainty for any fitted parameters. Although this may be true for systematic errors in the parameters, it is not in general the case for statistical uncertainties. If, for example, the error bars (or covariance matrix) used in constructing the χ^2 are underestimated, then this will lead to underestimated statistical errors for the fitted parameters. But in such a case, an estimate $\hat{\theta}$ can differ from the true value θ by an amount much greater than its estimated statistical error. The standard deviations of estimators that one finds from, say, Eq. (32.11) reflect how widely the estimates would be distributed if one were to repeat the measurement many times, assuming that the measurement errors used in the χ^2 are also correct. They do not include the systematic error which may result from an incorrect hypothesis or incorrectly estimated measurement errors in the χ^2 .

Since the mean of the χ^2 distribution is equal to n_d , one expects in a “reasonable” experiment to obtain $\chi^2 \approx n_d$. Hence the quantity χ^2/n_d is sometimes reported. Since the p.d.f. of χ^2/n_d depends on n_d , however, one must report n_d as well in order to make a meaningful statement. The p -values obtained for different values of χ^2/n_d are shown in Fig. 32.2.

32.3. Confidence intervals and limits

When the goal of an experiment is to determine a parameter θ , the result is usually expressed by quoting, in addition to the point estimate, some sort of interval which reflects the statistical precision of the measurement. In the simplest case, this can be given by the parameter's estimated value $\hat{\theta}$ plus or minus an estimate of the standard deviation of θ , $\sigma_{\hat{\theta}}$. If, however, the p.d.f. of the estimator is not Gaussian or if there are physical boundaries on the possible values of the parameter, then one usually quotes instead an interval according to one of the procedures described below.

In reporting an interval or limit, the experimenter may wish to

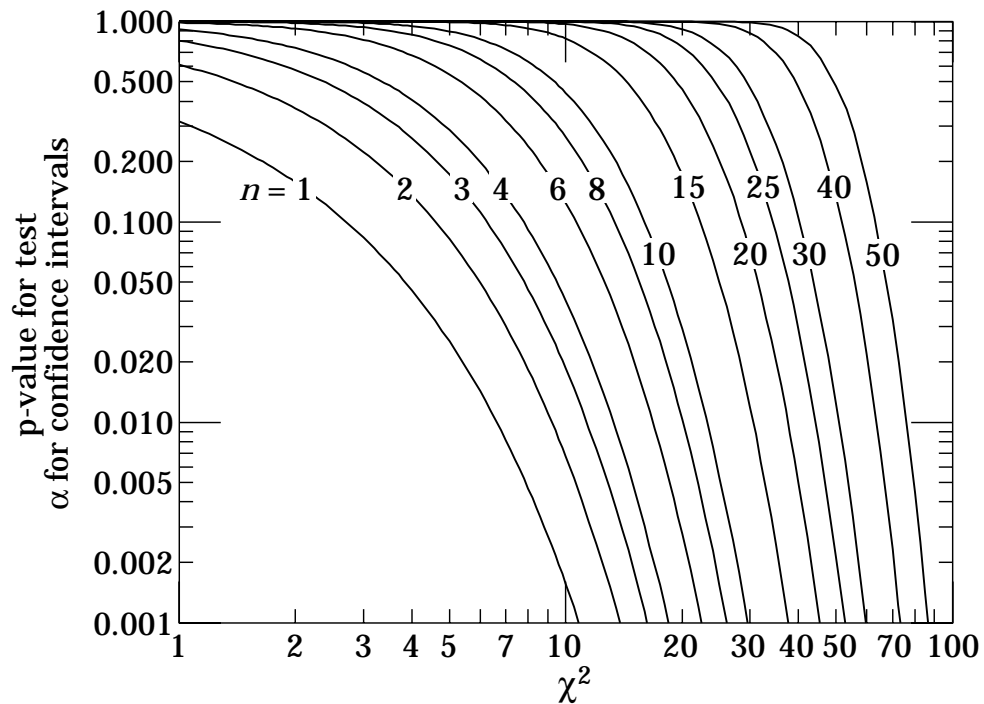


Figure 32.1: One minus the χ^2 cumulative distribution, $1 - F(\chi^2; n)$, for n degrees of freedom. This gives the p -value for the χ^2 goodness-of-fit test as well as one minus the coverage probability for confidence regions (see Sec. 32.3.2.4).

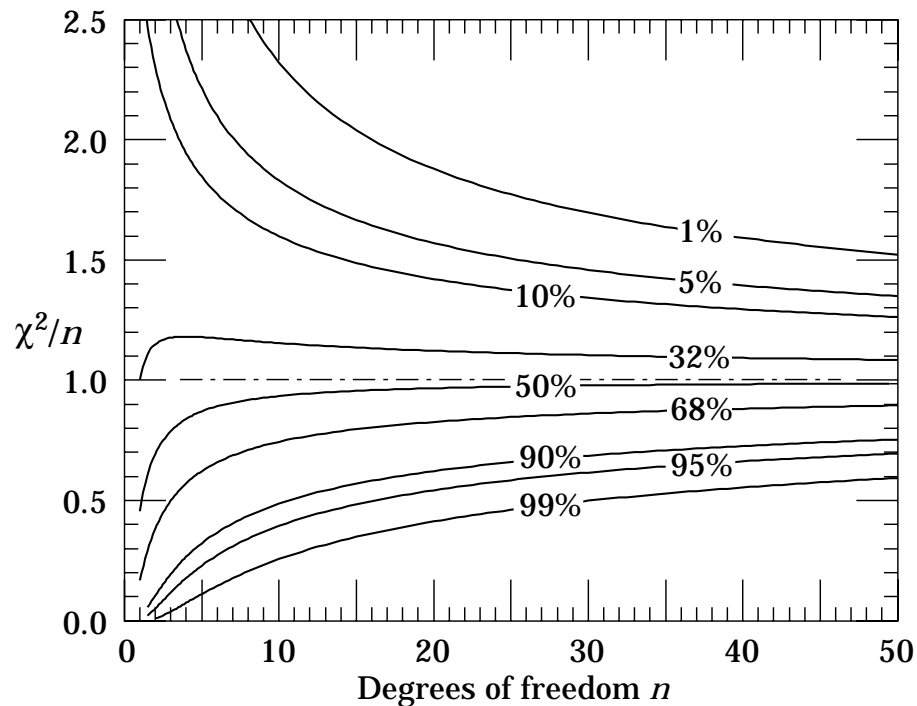


Figure 32.2: The ‘reduced’ χ^2 , equal to χ^2/n , for n degrees of freedom. The curves show as a function of n the χ^2/n that corresponds to a given p -value.

- communicate as objectively as possible the result of the experiment;
- provide an interval that is constructed to cover the true value of the parameter with a specified probability;
- provide the information needed by the consumer of the result to draw conclusions about the parameter or to make a particular decision;
- draw conclusions about the parameter that incorporate stated prior beliefs.

With a sufficiently large data sample, the point estimate and standard deviation (or for the multiparameter case, the parameter estimates and covariance matrix) satisfy essentially all of these goals. For finite data samples, no single method for quoting an interval will achieve all of them.

In addition to the goals listed above, the choice of method may be influenced by practical considerations such as ease of producing an interval from the results of several measurements. Of course the experimenter is not restricted to quoting a single interval or limit; one may choose, for example, first to communicate the result with a confidence interval having certain frequentist properties, and then in addition to draw conclusions about a parameter using Bayesian statistics. It is recommended, however, that there be a clear separation between these two aspects of reporting a result. In the remainder of this section, we assess the extent to which various types of intervals achieve the goals stated here.

32.3.1. *The Bayesian approach :*

Suppose the outcome of the experiment is characterized by a vector of data \mathbf{x} , whose probability distribution depends on an unknown parameter (or parameters) $\boldsymbol{\theta}$ that we wish to determine. In Bayesian statistics, all knowledge about $\boldsymbol{\theta}$ is summarized by the posterior p.d.f. $p(\boldsymbol{\theta}|\mathbf{x})$, which gives the degree of belief for $\boldsymbol{\theta}$ to take on values in a certain region given the data \mathbf{x} . It is obtained by using Bayes' theorem,

$$p(\boldsymbol{\theta}|\mathbf{x}) = \frac{L(\mathbf{x}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\int L(\mathbf{x}|\boldsymbol{\theta}')\pi(\boldsymbol{\theta}') d\boldsymbol{\theta}'}, \quad (32.30)$$

where $L(\mathbf{x}|\boldsymbol{\theta})$ is the likelihood function, *i.e.*, the joint p.d.f. for the data given a certain value of $\boldsymbol{\theta}$, evaluated with the data actually obtained in the experiment, and $\pi(\boldsymbol{\theta})$ is the prior p.d.f. for $\boldsymbol{\theta}$. Note that the denominator in Eq. (32.30) serves simply to normalize the posterior p.d.f. to unity.

Bayesian statistics supplies no unique rule for determining $\pi(\boldsymbol{\theta})$; this reflects the experimenter's subjective degree of belief about $\boldsymbol{\theta}$ before the measurement was carried out. By itself, therefore, the posterior p.d.f. is not a good way to report the result of an observation objectively, since it contains both the result (through the likelihood function) and the experimenter's prior beliefs. Without the likelihood function, someone with different prior beliefs would be unable to substitute these to determine his or her own posterior p.d.f. This is an important reason, therefore, to publish wherever possible the likelihood function or an appropriate summary of it. Often this can be achieved by reporting the ML estimate and one or several low order derivatives of L evaluated at the estimate.

14 32. Statistics

In the single parameter case, for example, an interval (called a Bayesian or credible interval) $[\theta_{\text{lo}}, \theta_{\text{up}}]$ can be determined which contains a given fraction $1 - \alpha$ of the posterior probability, *i.e.*,

$$1 - \alpha = \int_{\theta_{\text{lo}}}^{\theta_{\text{up}}} p(\theta|\mathbf{x}) d\theta . \quad (32.31)$$

Sometimes an upper or lower limit is desired, *i.e.*, θ_{lo} can be set to zero or θ_{up} to infinity. In other cases, one might choose θ_{lo} and θ_{up} such that $p(\theta|\mathbf{x})$ is higher everywhere inside the interval than outside; these are called *highest posterior density* (HPD) intervals. Note that HPD intervals are not invariant under a nonlinear transformation of the parameter.

The main difficulty with Bayesian intervals is in quantifying the prior beliefs. Sometimes one attempts to construct $\pi(\boldsymbol{\theta})$ to represent complete ignorance about the parameters by setting it equal to a constant. A problem here is that if the prior p.d.f. is flat in $\boldsymbol{\theta}$, then it is not flat for a nonlinear function of $\boldsymbol{\theta}$, and so a different parametrization of the problem would lead in general to a different posterior p.d.f. In practice, one does not choose a flat prior as a true expression of degree of belief about a parameter; rather, it is used as a recipe to construct an interval, which in the end will have certain frequentist properties.

If a parameter is constrained to be non-negative, then the prior p.d.f. can simply be set to zero for negative values. An important example is the case of a Poisson variable n , which counts signal events with unknown mean s , as well as background with mean b , assumed known. For the signal mean s , one often uses the prior

$$\pi(s) = \begin{cases} 0 & s < 0 \\ 1 & s \geq 0 \end{cases} . \quad (32.32)$$

As mentioned above, this is regarded as providing an interval whose frequentist properties can be studied, rather than as representing a degree of belief. In the absence of a clear discovery, (*e.g.*, if $n = 0$ or if in any case n is compatible with the expected background), one usually wishes to place an upper limit on s . Using the likelihood function for Poisson distributed n ,

$$L(n|s) = \frac{(s+b)^n}{n!} e^{-(s+b)} , \quad (32.33)$$

along with the prior (32.32) in (32.30) gives the posterior density for s . An upper limit s_{up} at confidence level (or here, rather, credibility level) $1 - \alpha$ can be obtained by requiring

$$1 - \alpha = \int_{-\infty}^{s_{\text{up}}} p(s|n) ds = \frac{\int_{-\infty}^{s_{\text{up}}} L(n|s) \pi(s) ds}{\int_{-\infty}^{\infty} L(n|s) \pi(s) ds} , \quad (32.34)$$

where the lower limit of integration is effectively zero because of the cut-off in $\pi(s)$. By relating the integrals in Eq. (32.34) to incomplete gamma functions, the equation reduces to

$$\alpha = e^{-s_{\text{up}}} \frac{\sum_{m=0}^n (s_{\text{up}} + b)^m / m!}{\sum_{m=0}^n b^m / m!} . \quad (32.35)$$

This must be solved numerically for the limit s_{up} . For the special case of $b = 0$, the sums can be related to the *quantile* $F_{\chi^2}^{-1}$ of the χ^2 distribution (inverse of the cumulative distribution) to give

$$s_{\text{up}} = \frac{1}{2} F_{\chi^2}^{-1}(1 - \alpha; n_{\text{d}}) , \quad (32.36)$$

where the number of degrees of freedom is $n_{\text{d}} = 2(n + 1)$. The quantile of the χ^2 distribution can be obtained using the CERNLIB routine `CHISIN`, or the ROOT function `TMath::ChisquareQuantile`. It so happens that for the case of $b = 0$, the upper limits from Eq. (32.36) coincide numerically with the values of the frequentist upper limits discussed in Section 32.3.2.5. Values for $1 - \alpha = 0.9$ and 0.95 are given by the values ν_{up} in Table 32.3. The frequentist properties of confidence intervals for the Poisson mean obtained in this way are discussed in Refs. [2] and [17].

Bayesian statistics provides a framework for incorporating systematic uncertainties into a result. Suppose, for example, that a model depends not only on parameters of interest $\boldsymbol{\theta}$, but on *nuisance parameters* $\boldsymbol{\nu}$, whose values are known with some limited accuracy. For a single nuisance parameter ν , for example, one might have a p.d.f. centered about its nominal value with a certain standard deviation σ_{ν} . Often a Gaussian p.d.f. provides a reasonable model for one's degree of belief about a nuisance parameter; in other cases, more complicated shapes may be appropriate. The likelihood function, prior, and posterior p.d.f.s then all depend on both $\boldsymbol{\theta}$ and $\boldsymbol{\nu}$, and are related by Bayes' theorem, as usual. One can obtain the posterior p.d.f. for $\boldsymbol{\theta}$ alone by integrating over the nuisance parameters, *i.e.*,

$$p(\boldsymbol{\theta}|\mathbf{x}) = \int p(\boldsymbol{\theta}, \boldsymbol{\nu}|\mathbf{x}) d\boldsymbol{\nu} . \quad (32.37)$$

If the prior joint p.d.f. for $\boldsymbol{\theta}$ and $\boldsymbol{\nu}$ factorizes, then integrating the posterior p.d.f. over $\boldsymbol{\nu}$ is equivalent to replacing the likelihood function by (see Ref. 18),

$$L'(\mathbf{x}|\boldsymbol{\theta}) = \int L(\mathbf{x}|\boldsymbol{\theta}, \boldsymbol{\nu}) \pi(\boldsymbol{\nu}) d\boldsymbol{\nu} . \quad (32.38)$$

The function $L'(\mathbf{x}|\boldsymbol{\theta})$ can also be used together with frequentist methods that employ the likelihood function such as ML estimation of parameters. The results then have a mixed frequentist/Bayesian character, where the systematic uncertainty due to limited knowledge of the nuisance parameters is built in. Although this may make it more difficult to disentangle statistical from systematic effects, such a hybrid approach may satisfy the objective of reporting the result in a convenient way.

Even if the subjective Bayesian approach is not used explicitly, Bayes' theorem represents the way that people evaluate the impact of a new result on their beliefs. One

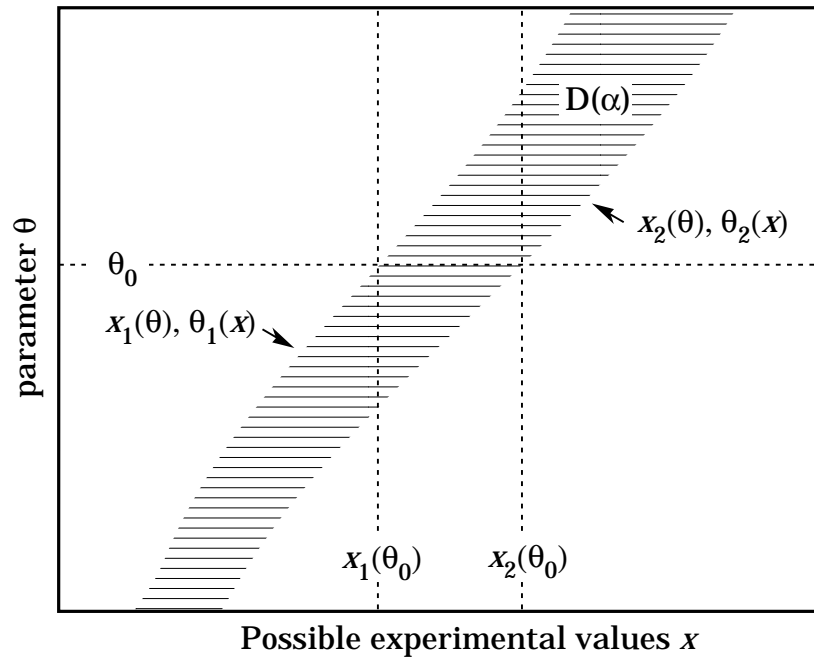


Figure 32.3: Construction of the confidence belt (see text).

of the criteria in choosing a method for reporting a measurement, therefore, should be the ease and convenience with which the consumer of the result can carry out this exercise.

32.3.2. Frequentist confidence intervals :

The unqualified phrase “confidence intervals” refers to frequentist intervals obtained with a procedure due to Neyman [19], described below. These are intervals (or in the multiparameter case, regions) constructed so as to include the true value of the parameter with a probability greater than or equal to a specified level, called the *coverage probability*. In this section, we discuss several techniques for producing intervals that have, at least approximately, this property.

32.3.2.1. The Neyman construction for confidence intervals:

Consider a p.d.f. $f(x; \theta)$ where x represents the outcome of the experiment and θ is the unknown parameter for which we want to construct a confidence interval. The variable x could (and often does) represent an estimator for θ . Using $f(x; \theta)$, we can find for a pre-specified probability $1 - \alpha$, and for every value of θ , a set of values $x_1(\theta, \alpha)$ and $x_2(\theta, \alpha)$ such that

$$P(x_1 < x < x_2; \theta) = 1 - \alpha = \int_{x_1}^{x_2} f(x; \theta) dx . \quad (32.39)$$

This is illustrated in Fig. 32.3: a horizontal line segment $[x_1(\theta, \alpha), x_2(\theta, \alpha)]$ is drawn for representative values of θ . The union of such intervals for all values of θ , designated in the figure as $D(\alpha)$, is known as the *confidence belt*. Typically the curves $x_1(\theta, \alpha)$ and $x_2(\theta, \alpha)$ are monotonic functions of θ , which we assume for this discussion.

Upon performing an experiment to measure x and obtaining a value x_0 , one draws a vertical line through x_0 . The confidence interval for θ is the set of all values of θ for which the corresponding line segment $[x_1(\theta, \alpha), x_2(\theta, \alpha)]$ is intercepted by this vertical line. Such confidence intervals are said to have a *confidence level* (CL) equal to $1 - \alpha$.

Now suppose that the true value of θ is θ_0 , indicated in the figure. We see from the figure that θ_0 lies between $\theta_1(x)$ and $\theta_2(x)$ if and only if x lies between $x_1(\theta_0)$ and $x_2(\theta_0)$. The two events thus have the same probability, and since this is true for any value θ_0 , we can drop the subscript 0 and obtain

$$1 - \alpha = P(x_1(\theta) < x < x_2(\theta)) = P(\theta_2(x) < \theta < \theta_1(x)). \quad (32.40)$$

In this probability statement, $\theta_1(x)$ and $\theta_2(x)$, *i.e.*, the endpoints of the interval, are the random variables and θ is an unknown constant. If the experiment were to be repeated a large number of times, the interval $[\theta_1, \theta_2]$ would vary, covering the fixed value θ in a fraction $1 - \alpha$ of the experiments.

The condition of coverage in Eq. (32.39) does not determine x_1 and x_2 uniquely, and additional criteria are needed. The most common criterion is to choose *central intervals* such that the probabilities excluded below x_1 and above x_2 are each $\alpha/2$. In other cases, one may want to report only an upper or lower limit, in which case the probability excluded below x_1 or above x_2 can be set to zero. Another principle based on *likelihood ratio ordering* for determining which values of x should be included in the confidence belt is discussed in Sec. 32.3.2.2

When the observed random variable x is continuous, the coverage probability obtained with the Neyman construction is $1 - \alpha$, regardless of the true value of the parameter. If x is discrete, however, it is not possible to find segments $[x_1(\theta, \alpha), x_2(\theta, \alpha)]$ that satisfy Eq. (32.39) exactly for all values of θ . By convention, one constructs the confidence belt requiring the probability $P(x_1 < x < x_2)$ to be *greater than or equal to* $1 - \alpha$. This gives confidence intervals that include the true parameter with a probability greater than or equal to $1 - \alpha$.

32.3.2.2. Relationship between intervals and tests:

An equivalent method of constructing confidence intervals is to consider a test (see Sec. 32.2) of the hypothesis that the parameter's true value is θ (assume one constructs a test for all physical values of θ). One then excludes all values of θ where the hypothesis would be rejected at a significance level less than α . The remaining values constitute the confidence interval at confidence level $1 - \alpha$.

In this procedure, one is still free to choose the test to be used; this corresponds to the freedom in the Neyman construction as to which values of the data are included in the confidence belt. One possibility is use a test statistic based on the *likelihood ratio*,

$$\lambda = \frac{f(x; \theta)}{f(x; \hat{\theta})}, \quad (32.41)$$

where $\hat{\theta}$ is the value of the parameter which, out of all allowed values, maximizes $f(x; \theta)$. This results in the intervals described in Ref. 20 by Feldman and Cousins. The same

18 32. Statistics

intervals can be obtained from the Neyman construction described in the previous section by including in the confidence belt those values of x which give the greatest values of λ .

Another technique that can be formulated in the language of statistical tests has been used to set limits on the Higgs mass from measurements at LEP [21,22]. For each value of the Higgs mass, a statistic called CL_s is determined from the ratio

$$CL_s = \frac{p\text{-value of signal plus background hypothesis}}{1 - p\text{-value of hypothesis of background only}} . \quad (32.42)$$

The p -values in Eq. (32.42) are themselves based on a test statistic which depends in general on the signal being tested, *i.e.*, on the hypothesized Higgs mass. Smaller CL_s corresponds to a lesser level of agreement with the signal hypothesis.

In the usual procedure for constructing confidence intervals, one would exclude the signal hypothesis if the probability to obtain a value of CL_s less than the one actually observed is less than α . The LEP Higgs group has in fact followed a more conservative approach, and excludes the signal at a confidence level $1 - \alpha$ if CL_s itself (not the probability to obtain a lower CL_s value) is less than α . This results in a coverage probability that is in general greater than $1 - \alpha$. The interpretation of such intervals is discussed in Refs.[21,22].

32.3.2.3. Profile likelihood and treatment of nuisance parameters:

As mentioned in Section 32.3.1, one may have a model containing parameters that must be determined from data, but which are not of any interest in the final result (nuisance parameters). Suppose the likelihood $L(\boldsymbol{\theta}, \boldsymbol{\nu})$ depends on parameters of interest $\boldsymbol{\theta}$ and nuisance parameters $\boldsymbol{\nu}$. The nuisance parameters can be effectively removed from the problem by constructing the *profile likelihood*, defined by

$$L_p(\boldsymbol{\theta}) = L(\boldsymbol{\theta}, \widehat{\boldsymbol{\nu}}(\boldsymbol{\theta})) , \quad (32.43)$$

where $\widehat{\boldsymbol{\nu}}(\boldsymbol{\theta})$ is given by the $\boldsymbol{\nu}$ that maximizes the likelihood for fixed $\boldsymbol{\theta}$. The profile likelihood may then be used to construct tests of or intervals for the parameters of interest. This is analogous to use of the integrated likelihood (32.38) used in the Bayesian approach. For example, one may construct the profile likelihood ratio,

$$\lambda_p(\boldsymbol{\theta}) = \frac{L_p(\boldsymbol{\theta})}{L(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\nu}})} , \quad (32.44)$$

where $\widehat{\boldsymbol{\theta}}$ and $\widehat{\boldsymbol{\nu}}$ are the ML estimators. The ratio λ_p can be used in place of the likelihood ratio (32.41) for inference about $\boldsymbol{\theta}$. The resulting intervals for the parameters of interest are not guaranteed to have the exact coverage probability for all values of the nuisance parameters, but in cases of practical interest the approximation is found to be very good. Further discussion on use of the profile likelihood can be found in, *e.g.*, Refs.[25,26] and other contributions to the PHYSTAT conferences [14].

32.3.2.4. Gaussian distributed measurements:

An important example of constructing a confidence interval is when the data consists of a single random variable x that follows a Gaussian distribution; this is often the case when x represents an estimator for a parameter and one has a sufficiently large data sample. If there is more than one parameter being estimated, the multivariate Gaussian is used. For the univariate case with known σ ,

$$1 - \alpha = \frac{1}{\sqrt{2\pi}\sigma} \int_{\mu-\delta}^{\mu+\delta} e^{-(x-\mu)^2/2\sigma^2} dx = \operatorname{erf}\left(\frac{\delta}{\sqrt{2}\sigma}\right) \quad (32.45)$$

is the probability that the measured value x will fall within $\pm\delta$ of the true value μ . From the symmetry of the Gaussian with respect to x and μ , this is also the probability for the interval $x \pm \delta$ to include μ . Fig. 32.4 shows a $\delta = 1.64\sigma$ confidence interval unshaded. The choice $\delta = \sigma$ gives an interval called the *standard error* which has $1 - \alpha = 68.27\%$ if σ is known. Values of α for other frequently used choices of δ are given in Table 32.1.

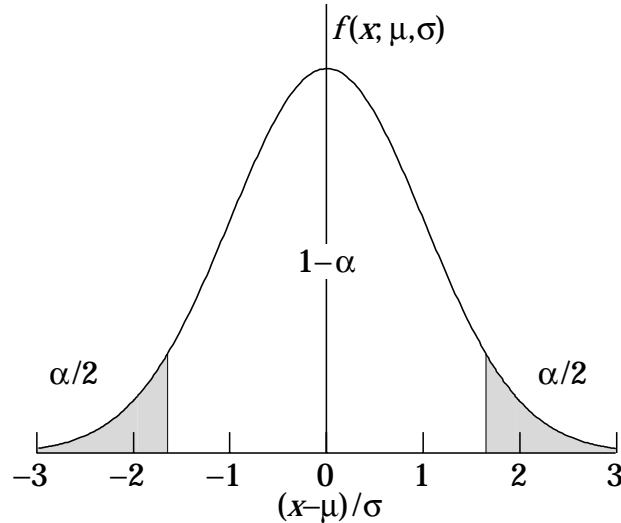


Figure 32.4: Illustration of a symmetric 90% confidence interval (unshaded) for a measurement of a single quantity with Gaussian errors. Integrated probabilities, defined by α , are as shown.

We can set a one-sided (upper or lower) limit by excluding above $x + \delta$ (or below $x - \delta$). The values of α for such limits are half the values in Table 32.1.

Table 32.1: Area of the tails α outside $\pm\delta$ from the mean of a Gaussian distribution.

α	δ	α	δ
0.3173	1σ	0.2	1.28σ
4.55×10^{-2}	2σ	0.1	1.64σ
2.7×10^{-3}	3σ	0.05	1.96σ
6.3×10^{-5}	4σ	0.01	2.58σ
5.7×10^{-7}	5σ	0.001	3.29σ
2.0×10^{-9}	6σ	10^{-4}	3.89σ

In addition to Eq. (32.45), α and δ are also related by the cumulative distribution function for the χ^2 distribution,

$$\alpha = 1 - F(\chi^2; n), \quad (32.46)$$

for $\chi^2 = (\delta/\sigma)^2$ and $n = 1$ degree of freedom. This can be obtained from Fig. 32.1 on the $n = 1$ curve or by using the CERNLIB routine `PROB` or the ROOT function `TMath::Prob`.

For multivariate measurements of, say, n parameter estimates $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_n)$, one requires the full covariance matrix $V_{ij} = \text{cov}[\hat{\theta}_i, \hat{\theta}_j]$, which can be estimated as described in Sections 32.1.2 and 32.1.3. Under fairly general conditions with the methods of maximum-likelihood or least-squares in the large sample limit, the estimators will be distributed according to a multivariate Gaussian centered about the true (unknown) values $\boldsymbol{\theta}$, and furthermore, the likelihood function itself takes on a Gaussian shape.

The standard error ellipse for the pair $(\hat{\theta}_i, \hat{\theta}_j)$ is shown in Fig. 32.5, corresponding to a contour $\chi^2 = \chi_{\min}^2 + 1$ or $\ln L = \ln L_{\max} - 1/2$. The ellipse is centered about the estimated values $\hat{\boldsymbol{\theta}}$, and the tangents to the ellipse give the standard deviations of the estimators, σ_i and σ_j . The angle of the major axis of the ellipse is given by

$$\tan 2\phi = \frac{2\rho_{ij}\sigma_i\sigma_j}{\sigma_j^2 - \sigma_i^2}, \quad (32.47)$$

where $\rho_{ij} = \text{cov}[\hat{\theta}_i, \hat{\theta}_j]/\sigma_i\sigma_j$ is the correlation coefficient.

The correlation coefficient can be visualized as the fraction of the distance σ_i from the ellipse's horizontal centerline at which the ellipse becomes tangent to vertical, *i.e.*, at the distance $\rho_{ij}\sigma_i$ below the centerline as shown. As ρ_{ij} goes to $+1$ or -1 , the ellipse thins to a diagonal line.

It could happen that one of the parameters, say, θ_j , is known from previous measurements to a precision much better than σ_j , so that the current measurement contributes almost nothing to the knowledge of θ_j . However, the current measurement of θ_i and its dependence on θ_j may still be important. In this case, instead of quoting both parameter estimates and their correlation, one sometimes reports the value of θ_i , which

Table 32.2: $\Delta\chi^2$ or $2\Delta\ln L$ corresponding to a coverage probability $1 - \alpha$ in the large data sample limit, for joint estimation of m parameters.

$(1 - \alpha)$ (%)	$m = 1$	$m = 2$	$m = 3$
68.27	1.00	2.30	3.53
90.	2.71	4.61	6.25
95.	3.84	5.99	7.82
95.45	4.00	6.18	8.03
99.	6.63	9.21	11.34
99.73	9.00	11.83	14.16

minimizes χ^2 at a fixed value of θ_j , such as the PDG best value. This θ_i value lies along the dotted line between the points where the ellipse becomes tangent to vertical, and has statistical error σ_{inner} as shown on the figure, where $\sigma_{\text{inner}} = (1 - \rho_{ij}^2)^{1/2}\sigma_i$. Instead of the correlation ρ_{ij} , one reports the dependency $d\hat{\theta}_i/d\theta_j$ which is the slope of the dotted line. This slope is related to the correlation coefficient by $d\hat{\theta}_i/d\theta_j = \rho_{ij} \times \frac{\sigma_i}{\sigma_j}$.

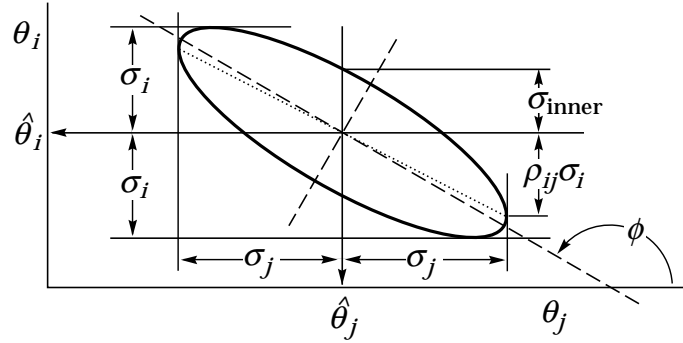


Figure 32.5: Standard error ellipse for the estimators $\hat{\theta}_i$ and $\hat{\theta}_j$. In this case the correlation is negative.

As in the single-variable case, because of the symmetry of the Gaussian function between $\boldsymbol{\theta}$ and $\hat{\boldsymbol{\theta}}$, one finds that contours of constant $\ln L$ or χ^2 cover the true values with a certain, fixed probability. That is, the confidence region is determined by

$$\ln L(\boldsymbol{\theta}) \geq \ln L_{\text{max}} - \Delta \ln L, \quad (32.48)$$

or where a χ^2 has been defined for use with the method of least-squares,

$$\chi^2(\boldsymbol{\theta}) \leq \chi_{\text{min}}^2 + \Delta\chi^2. \quad (32.49)$$

Values of $\Delta\chi^2$ or $2\Delta\ln L$ are given in Table 32.2 for several values of the coverage probability and number of fitted parameters.

22 32. Statistics

For finite data samples, the probability for the regions determined by equations (32.48) or (32.49) to cover the true value of θ will depend on θ , so these are not exact confidence regions according to our previous definition. Nevertheless, they can still have a coverage probability only weakly dependent on the true parameter, and approximately as given in Table 32.2. In any case, the coverage probability of the intervals or regions obtained according to this procedure can in principle be determined as a function of the true parameter(s), for example, using a Monte Carlo calculation.

One of the practical advantages of intervals that can be constructed from the log-likelihood function or χ^2 is that it is relatively simple to produce the interval for the combination of several experiments. If N independent measurements result in log-likelihood functions $\ln L_i(\theta)$, then the combined log-likelihood function is simply the sum,

$$\ln L(\theta) = \sum_{i=1}^N \ln L_i(\theta). \quad (32.50)$$

This can then be used to determine an approximate confidence interval or region with Eq. (32.48), just as with a single experiment.

32.3.2.5. Poisson or binomial data:

Another important class of measurements consists of counting a certain number of events, n . In this section, we will assume these are all events of the desired type, *i.e.*, there is no background. If n represents the number of events produced in a reaction with cross section σ , say, in a fixed integrated luminosity \mathcal{L} , then it follows a Poisson distribution with mean $\nu = \sigma\mathcal{L}$. If, on the other hand, one has selected a larger sample of N events and found n of them to have a particular property, then n follows a binomial distribution where the parameter p gives the probability for the event to possess the property in question. This is appropriate, *e.g.*, for estimates of branching ratios or selection efficiencies based on a given total number of events.

For the case of Poisson distributed n , the upper and lower limits on the mean value ν can be found from the Neyman procedure to be

$$\nu_{\text{lo}} = \frac{1}{2} F_{\chi^2}^{-1}(\alpha_{\text{lo}}; 2n), \quad (32.51a)$$

$$\nu_{\text{up}} = \frac{1}{2} F_{\chi^2}^{-1}(1 - \alpha_{\text{up}}; 2(n + 1)), \quad (32.51b)$$

where the upper and lower limits are at confidence levels of $1 - \alpha_{\text{lo}}$ and $1 - \alpha_{\text{up}}$, respectively, and $F_{\chi^2}^{-1}$ is the *quantile* of the χ^2 distribution (inverse of the cumulative distribution). The quantiles $F_{\chi^2}^{-1}$ can be obtained from standard tables or from the CERNLIB routine CHISIN. For central confidence intervals at confidence level $1 - \alpha$, set $\alpha_{\text{lo}} = \alpha_{\text{up}} = \alpha/2$.

It happens that the upper limit from Eq. (32.51a) coincides numerically with the Bayesian upper limit for a Poisson parameter, using a uniform prior p.d.f. for ν . Values for confidence levels of 90% and 95% are shown in Table 32.3.

Table 32.3: Lower and upper (one-sided) limits for the mean ν of a Poisson variable given n observed events in the absence of background, for confidence levels of 90% and 95%.

n	$1 - \alpha = 90\%$		$1 - \alpha = 95\%$	
	ν_{lo}	ν_{up}	ν_{lo}	ν_{up}
0	–	2.30	–	3.00
1	0.105	3.89	0.051	4.74
2	0.532	5.32	0.355	6.30
3	1.10	6.68	0.818	7.75
4	1.74	7.99	1.37	9.15
5	2.43	9.27	1.97	10.51
6	3.15	10.53	2.61	11.84
7	3.89	11.77	3.29	13.15
8	4.66	12.99	3.98	14.43
9	5.43	14.21	4.70	15.71
10	6.22	15.41	5.43	16.96

For the case of binomially distributed n successes out of N trials with probability of success p , the upper and lower limits on p are found to be

$$p_{\text{lo}} = \frac{nF_F^{-1}[\alpha_{\text{lo}}; 2n, 2(N - n + 1)]}{N - n + 1 + nF_F^{-1}[\alpha_{\text{lo}}; 2n, 2(N - n + 1)]}, \quad (32.52a)$$

$$p_{\text{up}} = \frac{(n + 1)F_F^{-1}[1 - \alpha_{\text{up}}; 2(n + 1), 2(N - n)]}{(N - n) + (n + 1)F_F^{-1}[1 - \alpha_{\text{up}}; 2(n + 1), 2(N - n)]}. \quad (32.52b)$$

Here F_F^{-1} is the quantile of the F distribution (also called the Fisher–Snedecor distribution; see Ref. 4).

32.3.2.6. *Difficulties with intervals near a boundary:*

A number of issues arise in the construction and interpretation of confidence intervals when the parameter can only take on values in a restricted range. An important example is where the mean of a Gaussian variable is constrained on physical grounds to be non-negative. This arises, for example, when the square of the neutrino mass is estimated from $\hat{m}^2 = \hat{E}^2 - \hat{p}^2$, where \hat{E} and \hat{p} are independent, Gaussian-distributed estimates of the energy and momentum. Although the true m^2 is constrained to be positive, random errors in \hat{E} and \hat{p} can easily lead to negative values for the estimate \hat{m}^2 .

If one uses the prescription given above for Gaussian distributed measurements, which says to construct the interval by taking the estimate plus-or-minus-one standard deviation, then this can give intervals that are partially or entirely in the unphysical region. In fact, by following strictly the Neyman construction for the central confidence interval, one finds that the interval is truncated below zero; nevertheless an extremely small or even a zero-length interval can result.

An additional important example is where the experiment consists of counting a certain number of events, n , which is assumed to be Poisson-distributed. Suppose the expectation value $E[n] = \nu$ is equal to $s + b$, where s and b are the means for signal and background processes, and assume further that b is a known constant. Then $\hat{s} = n - b$ is an unbiased estimator for s . Depending on true magnitudes of s and b , the estimate \hat{s} can easily fall in the negative region. Similar to the Gaussian case with the positive mean, the central confidence interval or even the upper limit for s may be of zero length.

The confidence interval is in fact designed not to cover the parameter with a probability of at most α , and if a zero-length interval results, then this is evidently one of those experiments. So although the construction is behaving as it should, a null interval is an unsatisfying result to report, and several solutions to this type of problem are possible.

An additional difficulty arises when a parameter estimate is not significantly far away from the boundary, in which case it is natural to report a one-sided confidence interval (often an upper limit). It is straightforward to force the Neyman prescription to produce only an upper limit by setting $x_2 = \infty$ in Eq. 32.39. Then x_1 is uniquely determined and the upper limit can be obtained. If, however, the data come out such that the parameter estimate is not so close to the boundary, one might wish to report a central (*i.e.*, two-sided) confidence interval. As pointed out by Feldman and Cousins [20], however, if the decision to report an upper limit or two-sided interval is made by looking at the data (“flip-flopping”), then the resulting intervals will not in general cover the parameter with the probability $1 - \alpha$.

With the confidence intervals suggested in Ref. 20, the prescription determines whether the interval is one- or two-sided in a way which preserves the coverage probability. Interval constructions that have this property and avoid the problem of null intervals are said to be unified. The intervals based on the Feldman-Cousins prescription are of this type. For a given choice of $1 - \alpha$, if the parameter estimate is sufficiently close to the boundary, the method gives a one-sided limit. In the case of a Poisson variable in the presence of background, for example, this would occur if the number of observed events is compatible with the expected background. For parameter estimates increasingly far away from the boundary, *i.e.*, for increasing signal significance, the interval makes a smooth transition from one- to two-sided, and far away from the boundary, one obtains a central interval.

The intervals according to this method for the mean of Poisson variable in the absence of background are given in Table 32.4. (Note that α in Ref. 20 is defined following Neyman [19] as the coverage probability; this is opposite the modern convention used here in which the coverage probability is $1 - \alpha$.) The values of $1 - \alpha$ given here refer to the coverage of the true parameter by the whole interval $[\nu_1, \nu_2]$. In Table 32.3 for the one-sided upper and lower limits, however, $1 - \alpha$ refers to the probability to have

individually $\nu_{\text{up}} \geq \nu$ or $\nu_{\text{lo}} \leq \nu$.

Table 32.4: Unified confidence intervals $[\nu_1, \nu_2]$ for a the mean of a Poisson variable given n observed events in the absence of background, for confidence levels of 90% and 95%.

n	$1 - \alpha = 90\%$		$1 - \alpha = 95\%$	
	ν_1	ν_2	ν_1	ν_2
0	0.00	2.44	0.00	3.09
1	0.11	4.36	0.05	5.14
2	0.53	5.91	0.36	6.72
3	1.10	7.42	0.82	8.25
4	1.47	8.60	1.37	9.76
5	1.84	9.99	1.84	11.26
6	2.21	11.47	2.21	12.75
7	3.56	12.53	2.58	13.81
8	3.96	13.99	2.94	15.29
9	4.36	15.30	4.36	16.77
10	5.50	16.50	4.75	17.82

A potential difficulty with unified intervals arises if, for example, one constructs such an interval for a Poisson parameter s of some yet to be discovered signal process with, say, $1 - \alpha = 0.9$. If the true signal parameter is zero, or in any case much less than the expected background, one will usually obtain a one-sided upper limit on s . In a certain fraction of the experiments, however, a two-sided interval for s will result. Since, however, one typically chooses $1 - \alpha$ to be only 0.9 or 0.95 when searching for a new effect, the value $s = 0$ may be excluded from the interval before the existence of the effect is well established. It must then be communicated carefully that in excluding $s = 0$ from the interval, one is not necessarily claiming to have discovered the effect.

The intervals constructed according to the unified procedure in Ref. 20 for a Poisson variable n consisting of signal and background have the property that for $n = 0$ observed events, the upper limit decreases for increasing expected background. This is counter-intuitive, since it is known that if $n = 0$ for the experiment in question, then no background was observed, and therefore one may argue that the expected background should not be relevant. The extent to which one should regard this feature as a drawback is a subject of some controversy (see, *e.g.*, Ref. 24).

Another possibility is to construct a Bayesian interval as described in Section 32.3.1. The presence of the boundary can be incorporated simply by setting the prior density

to zero in the unphysical region. Priors based on invariance principles (rather than subjective degree of belief) for the Poisson mean are rarely used in high-energy physics. An example one may consider for the Poisson problem is a prior inversely proportional to the mean; here one obtains a posterior that diverges for the case of zero events observed, and finds upper limits which undercover when evaluated by the frequentist definition of coverage [2]. Rather, priors uniform in the Poisson mean have been used, although as previously mentioned, this is generally not done to reflect the experimenter's degree of belief, but rather as a procedure for obtaining an interval with certain frequentist properties. The resulting upper limits have a coverage probability that depends on the true value of the Poisson parameter, and is everywhere greater than the stated probability content. Lower limits and two-sided intervals for the Poisson mean based on flat priors undercover, however, for some values of the parameter, although to an extent that in practical cases may not be too severe [2, 17]. Intervals constructed in this way have the advantage of being easy to derive; if several independent measurements are to be combined then one simply multiplies the likelihood functions (cf. Eq. (32.50)).

An additional alternative is presented by the intervals found from the likelihood function or χ^2 using the prescription of Equations (32.48) or (32.49). As in the case of the Bayesian intervals, the coverage probability is not, in general, independent of the true parameter. Furthermore, these intervals can for some parameter values undercover. The coverage probability can, of course, be determined with some extra effort and reported with the result.

Also as in the Bayesian case, intervals derived from the value of the likelihood function from a combination of independent experiments can be determined simply by multiplying the likelihood functions. These intervals are also invariant under transformation of the parameter; this is not true for Bayesian intervals with a conventional flat prior, because a uniform distribution in, say, θ will not be uniform if transformed to θ^2 . Use of the likelihood function to determine approximate confidence intervals is discussed further in Ref. 23.

In any case, it is important to always report sufficient information so that the result can be combined with other measurements. Often this means giving an unbiased estimator and its standard deviation, even if the estimated value is in the unphysical region.

Regardless of the type of interval reported, the consumer of that result will almost certainly use it to derive some impression about the value of the parameter. This will inevitably be done, either explicitly or intuitively, with Bayes' theorem,

$$p(\theta|\text{result}) \propto L(\text{result}|\theta)\pi(\theta), \quad (32.53)$$

where the reader supplies his or her own prior beliefs $\pi(\theta)$ about the parameter, and the 'result' is whatever sort of interval or other information the author has reported. For all of the intervals discussed, therefore, it is not sufficient to know the result; one must also know the probability to have obtained this result as a function of the parameter, *i.e.*, the likelihood. Contours of constant likelihood, for example, provide this information, and so an interval obtained from $\ln L = \ln L_{\max} - \Delta \ln L$ already takes one step in this direction.

It can also be useful with a frequentist interval to calculate its subjective probability content using the posterior p.d.f. based on one or several reasonable guesses for the prior p.d.f. If it turns out to be significantly less than the stated confidence level, this warns that it would be particularly misleading to draw conclusions about the parameter's value from the interval alone.

References:

1. B. Efron, *Am. Stat.* **40**, 11 (1986).
2. R.D. Cousins, *Am. J. Phys.* **63**, 398 (1995).
3. A. Stuart, J.K. Ord, and S. Arnold, *Kendall's Advanced Theory of Statistics*, Vol. 2A: *Classical Inference and the Linear Model*, 6th Ed., Oxford Univ. Press (1999), and earlier editions by Kendall and Stuart. The likelihood-ratio ordering principle is described at the beginning of Ch. 23. Chapter 26 compares different schools of statistical inference.
4. F.E. James, *Statistical Methods in Experimental Physics*, 2nd ed., (World Scientific, Singapore, 2007).
5. H. Cramér, *Mathematical Methods of Statistics*, Princeton Univ. Press, New Jersey (1958).
6. L. Lyons, *Statistics for Nuclear and Particle Physicists*, (Cambridge University Press, New York, 1986).
7. R. Barlow, *Nucl. Instrum. Methods* **A297**, 496 (1990).
8. G. Cowan, *Statistical Data Analysis*, (Oxford University Press, Oxford, 1998).
9. For a review, see S. Baker and R. Cousins, *Nucl. Instrum. Methods* **221**, 437 (1984).
10. For information on neural networks and related topics, see *e.g.*, C.M. Bishop, *Neural Networks for Pattern Recognition*, Clarendon Press, Oxford (1995); C. Peterson and T. Rönngvaldsson, An Intro. to Artificial Neural Networks, in *Proc. of the 1991 CERN School of Computing*, C. Verkerk (ed.), CERN 92-02 (1992).
11. T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning* (Springer, New York, 2001).
12. A. Webb, *Statistical Pattern Recognition*, 2nd ed., (Wiley, New York, 2002).
13. L.I. Kuncheva, *Combining Pattern Classifiers*, (Wiley, New York, 2004).
14. Links to the Proceedings of the PHYSTAT conference series (Durham 2002, Stanford 2003, Oxford 2005, and Geneva 2007) can be found at phystat.org.
15. A. Höcker *et al.*, *TMVA Users Guide*, [physics/0703039\(2007\)](http://physics/0703039(2007)); software available from tmva.sf.net.
16. I. Narsky, *StatPatternRecognition: A C++ Package for Statistical Analysis of High Energy Physics Data*, [physics/0507143\(2005\)](http://physics/0507143(2005)); software avail. from sourceforge.net/projects/statpatrec.
17. B.P. Roe and M.B. Woodroffe, *Phys. Rev.* **D63**, 13009 (2001).
18. P.H. Garthwaite, I.T. Jolliffe, and B. Jones, *Statistical Inference*, (Prentice Hall, 1995).
19. J. Neyman, *Phil. Trans. Royal Soc. London, Series A*, **236**, 333 (1937), reprinted in *A Selection of Early Statistical Papers on J. Neyman*, (University of California Press, Berkeley, 1967).

28 *32. Statistics*

20. G.J. Feldman and R.D. Cousins, Phys. Rev. **D57**, 3873 (1998). This paper does not specify what to do if the ordering principle gives equal rank to some values of x . Eq. 21.6 of Ref. 3 gives the rule: all such points are included in the acceptance region (the domain $D(\alpha)$). Some authors have assumed the contrary, and shown that one can then obtain null intervals.
21. T. Junk, Nucl. Instrum. Methods **A434**, 435 (1999).
22. A.L. Read, *Modified frequentist analysis of search results (the CL_s method)*, in F. James, L. Lyons, and Y. Perrin (eds.), *Workshop on Confidence Limits*, CERN Yellow Report 2000-005, available through cdsweb.cern.ch.
23. F. Porter, Nucl. Instrum. Methods **A368**, 793 (1996).
24. Workshop on Confidence Limits, CERN, 17-18 Jan. 2000, www.cern.ch/CERN/Divisions/EP/Events/CLW/. The proceedings, F. James, L. Lyons, and Y. Perrin (eds.), CERN Yellow Report 2000-005, are available through cdsweb.cern.ch. See also the later Fermilab workshop linked to the CERN web page.
25. N. Reid, *Likelihood Inference in the Presence of Nuisance Parameters*, Proceedings of PHYSTAT2003, L. Lyons, R. Mount, and R. Reitmeyer, eds., eConf C030908, Stanford, 2003.
26. W.A. Rolke, A.M. Lopez, and J. Conrad, Nucl. Instrum. Methods **A551**, 493 (2005); physics/0403059.