

7. The Method of Maximum Likelihood

7.1 Likelihood Ratio: Likelihood Function

In the last chapter we introduced the concept of parameter estimation. We have also described the desirable properties of estimators, though without specifying how such estimators can be constructed in a particular case. We have derived estimators only for the important quantities expectation value and variance. We now take on the general problem.

In order to specify explicitly the parameters

$$\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_p) \quad ,$$

we now write the probability density of the random variables

$$\mathbf{x} = (x_1, x_2, \dots, x_n)$$

in the form

$$f = f(\mathbf{x}; \boldsymbol{\lambda}) \quad . \quad (7.1.1)$$

If we now carry out a certain number of experiments, say N , or we draw a sample of size N from a population, then we can give a number to each experiment j :

$$dP^{(j)} = f(\mathbf{x}^{(j)}; \boldsymbol{\lambda}) d\mathbf{x} \quad . \quad (7.1.2)$$

The number $dP^{(j)}$ has the character of an *a posteriori probability*, i.e., given *after* the experiment, how probable it was to find the result $\mathbf{x}^{(j)}$ (within a small interval). The total probability to find exactly all of the events

$$\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(j)}, \dots, \mathbf{x}^{(N)}$$

is then the product

$$dP = \prod_{j=1}^N f(\mathbf{x}^{(j)}; \boldsymbol{\lambda}) d\mathbf{x} \quad . \quad (7.1.3)$$

This probability still clearly depends on λ . There are cases where the population is determined by only two possible sets of parameters, λ_1 and λ_2 . Such cases occur, for example, in nuclear physics, where the parity of a state is necessarily “even” or “odd”. One can construct the ratio

$$Q = \frac{\prod_{j=1}^N f(\mathbf{x}^{(j)}; \lambda_1)}{\prod_{j=1}^N f(\mathbf{x}^{(j)}; \lambda_2)} \quad (7.1.4)$$

and say that the values λ_1 are “ Q times more probable” than the values λ_2 . This factor is called the *likelihood ratio*.*

A product of the form

$$L = \prod_{j=1}^N f(\mathbf{x}^{(j)}; \lambda) \quad (7.1.5)$$

is called a *likelihood function*. One must clearly distinguish between a probability density and a likelihood function, which is a function of a sample and is hence a random variable. In particular, the a posteriori nature of the probability in (7.1.5) is of significance in many discussions.

Example 7.1: Likelihood ratio

Suppose one wishes to decide whether a coin belongs to type A or B by means of a number of tosses. The coins in question are asymmetric in such a way that A shows heads with a probability of $1/3$, and B shows heads with a probability of $2/3$.

	A	B
Heads	$1/3$	$2/3$
Tails	$2/3$	$1/3$

If an experiment yields heads once and tails four times, then one has

$$L_A = \frac{1}{3} \cdot \left(\frac{2}{3}\right)^4 \text{ and } L_B = \frac{2}{3} \cdot \left(\frac{1}{3}\right)^4,$$

$$Q = \frac{L_A}{L_B} = 8.$$

One would therefore tend towards the position that the coin is of type A . ■

*Although the likelihood ratio Q and the likelihood functions L and ℓ introduced below are random variables, since they are functions of a sample, we do not write them here with a special character type.

7.2 The Method of Maximum Likelihood

The generalization of the likelihood ratio is now clear. One gives the greatest confidence to that choice of the parameters λ for which the likelihood function (7.1.5) is a maximum. Figure 7.1 illustrates the situation for various forms of the likelihood function for the case of a single parameter λ .

The maximum can be located simply by setting the first derivative of the likelihood function with respect to the parameter λ_i equal to zero. The derivative of a product with many factors is, however, unpleasant to deal with. One first constructs therefore the logarithm of the likelihood function,

$$\ell = \ln L = \sum_{j=1}^N \ln f(\mathbf{x}^{(j)}; \lambda). \quad (7.2.1)$$

The function ℓ is also often called the likelihood function. Sometimes one says explicitly “*log-likelihood function*”. Clearly the maxima of (7.2.1) are identical with those of (7.1.5). For the case of a single parameter we now construct

$$\ell' = d\ell/d\lambda = 0. \quad (7.2.2)$$

The problem of estimating a parameter is now reduced to solving this *likelihood equation*. By application of (7.2.1) we can write

$$\ell' = \sum_{j=1}^N \frac{d}{d\lambda} \ln f(\mathbf{x}^{(j)}; \lambda) = \sum_{j=1}^N \frac{f'}{f} = \sum_{j=1}^N \varphi(\mathbf{x}^{(j)}; \lambda), \quad (7.2.3)$$

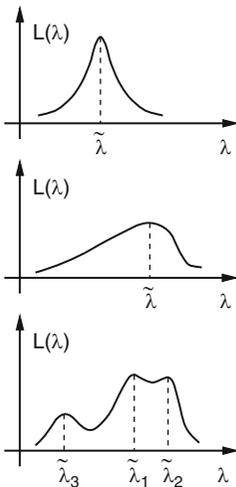


Fig. 7.1: Likelihood functions.

where

$$\varphi(\mathbf{x}^{(j)}; \lambda) = \left(\frac{d}{d\lambda} f(\mathbf{x}^{(j)}; \lambda) \right) / f(\mathbf{x}^{(j)}; \lambda) \quad (7.2.4)$$

is the *logarithmic derivative* of the density f with respect to λ .

In the general case of p parameters the likelihood equation (7.2.2) is replaced by the system of p simultaneous equations,

$$\frac{\partial \ell}{\partial \lambda_i} = 0, \quad i = 1, 2, \dots, p. \quad (7.2.5)$$

Example 7.2: Repeated measurements of differing accuracy

If a quantity is measured with different instruments, then the measurement errors are in general different. The measurements $\mathbf{x}^{(j)}$ are spread about the true value λ . Suppose the errors are normally distributed, so that a measurement corresponds to obtaining a sample from a Gaussian distribution with mean λ and standard deviation σ_j . The a posteriori probability for a measured value is then

$$f(\mathbf{x}^{(j)}; \lambda) dx = \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(\mathbf{x}^{(j)} - \lambda)^2}{2\sigma_j^2}\right) dx.$$

From all N measurements one obtains the likelihood function

$$L = \prod_{j=1}^N \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(\mathbf{x}^{(j)} - \lambda)^2}{2\sigma_j^2}\right) \quad (7.2.6)$$

with the logarithm

$$\ell = -\frac{1}{2} \sum_{j=1}^N \frac{(\mathbf{x}^{(j)} - \lambda)^2}{\sigma_j^2} + \text{const.} \quad (7.2.7)$$

The likelihood equation thus becomes

$$\frac{d\ell}{d\lambda} = \sum_{j=1}^N \frac{\mathbf{x}^{(j)} - \lambda}{\sigma_j^2} = 0.$$

It has the solution

$$\tilde{\lambda} = \frac{\sum_{j=1}^N \frac{\mathbf{x}^{(j)}}{\sigma_j^2}}{\sum_{j=1}^N \frac{1}{\sigma_j^2}}. \quad (7.2.8)$$

Since $d^2\ell/d\lambda^2 = -\sum\sigma_j^{-2} < 0$, the solution is, in fact, a maximum. Thus we see that we obtain the maximum likelihood estimator as the mean of the N measurements weighted inversely by the variances of the individual measurements. ■

Example 7.3: Estimation of the parameter N of the hypergeometric distribution

As in the example with coins at the beginning of this chapter, sometimes parameters to be estimated can only take on discrete values. In Example 5.2 we indicated the possibility of estimating zoological population densities by means of tagging and recapture. According to (5.3.1), the probability to catch exactly n fish of which k are tagged out of a pond with an unknown total of N fish, out of which K are tagged, is given by

$$L(k; n, K, N) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}.$$

We must now find the value of N for which the function L is maximum. For this we use the ratio

$$\frac{L(k; n, k, N)}{L(k; n, k, N-1)} = \frac{(N-n)(N-k)}{(N-n-K+k)N} \begin{cases} > 1, & Nk < nK, \\ < 1, & Nk > nK. \end{cases}$$

The function L is thus maximum when N is the integer closest to nK/k . ■

7.3 Information Inequality. Minimum Variance Estimators. Sufficient Estimators

We now want to discuss once more the quality of an estimator. In Sect. 6.1 we called an estimator unbiased if for every sample the *bias* vanished,

$$B(\lambda) = E(\mathbf{S}) - \lambda = 0. \quad (7.3.1)$$

Lack of bias is, however, not the only characteristic required of a “good” estimator. More importantly one should require that the variance

$$\sigma^2(\mathbf{S})$$

is small. Here one must often find a compromise, since there is a connection between B and σ^2 , described by the *information inequality*.[†]

One immediately sees that it is easy to achieve $\sigma^2(\mathbf{S}) = 0$ simply by using a constant for \mathbf{S} . We consider an estimator $\mathbf{S}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})$ that is

[†]This inequality was independently found by H. Cramer, M. Fréchet, and C. R. Rao as well as by other authors. It is also called the Cramer–Rao or Fréchet inequality.

a function of the sample $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}$. According to (6.1.3) and (6.1.4) the joint probability density of the elements of the sample is

$$f(x^{(1)}, x^{(2)}, \dots, x^{(N)}; \lambda) = f(x^{(1)}; \lambda) f(x^{(2)}; \lambda) \cdots f(x^{(N)}; \lambda).$$

The expectation value of \mathbf{S} is thus

$$E(\mathbf{S}) = \int \mathbf{S}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) f(x^{(1)}; \lambda) \cdots f(x^{(N)}; \lambda) \times dx^{(1)} dx^{(2)} \cdots dx^{(N)}. \quad (7.3.2)$$

According to (7.3.1), however, one also has

$$E(\mathbf{S}) = B(\lambda) + \lambda.$$

We now assume that we can differentiate with respect to λ in the integral. We then obtain

$$1 + B'(\lambda) = \int \mathbf{S} \left(\sum_{j=1}^N \frac{f'(x^{(j)}; \lambda)}{f(x^{(j)}; \lambda)} \right) f(x^{(1)}; \lambda) \cdots f(x^{(N)}; \lambda) dx^{(1)} \cdots dx^{(N)},$$

which is equivalent to

$$1 + B'(\lambda) = E \left\{ \mathbf{S} \sum_{j=1}^N \frac{f'(x^{(j)}; \lambda)}{f(x^{(j)}; \lambda)} \right\} = E \left\{ \mathbf{S} \sum_{j=1}^N \varphi(x^{(j)}; \lambda) \right\}.$$

From (7.2.3) we have

$$\ell' = \sum_{j=1}^N \varphi(\mathbf{x}^{(j)}; \lambda)$$

and therefore

$$1 + B'(\lambda) = E\{\mathbf{S}\ell'\}. \quad (7.3.3)$$

One clearly has

$$\int f(x^{(1)}; \lambda) \cdots f(x^{(N)}; \lambda) dx^{(1)} \cdots dx^{(N)} = 1.$$

If we also compute the derivative with respect to λ , we obtain

$$\int \sum_{j=1}^N \frac{f'(x^{(j)}; \lambda)}{f(x^{(j)}; \lambda)} f(x^{(1)}; \lambda) \cdots f(x^{(N)}; \lambda) dx^{(1)} \cdots dx^{(N)} = E(\ell') = 0.$$

By multiplying this equation by $E(\mathbf{S})$ and subtracting the result of (7.3.3) one obtains

$$1 + B'(\lambda) = E\{\mathbf{S}\ell'\} - E(\mathbf{S})E(\ell') = E\{[\mathbf{S} - E(\mathbf{S})]\ell'\}. \quad (7.3.4)$$

In order to see the significance of this expression, we need to use the Cauchy–Schwarz inequality in the following form:

If \mathbf{x} and \mathbf{y} are random variables and if \mathbf{x}^2 and \mathbf{y}^2 have finite expectation values, then

$$\{E(\mathbf{xy})\}^2 \leq E(\mathbf{x}^2)E(\mathbf{y}^2). \quad (7.3.5)$$

To prove this inequality we consider the expression

$$E((a\mathbf{x} + \mathbf{y})^2) = a^2E(\mathbf{x}^2) + 2aE(\mathbf{xy}) + E(\mathbf{y}^2) \geq 0. \quad (7.3.6)$$

This is a non-negative number for all values of a . If we consider for the moment the case of equality, then this is a quadratic equation for a with the solutions

$$a_{1,2} = -\frac{E(\mathbf{xy})}{E(\mathbf{x}^2)} \pm \sqrt{\left(\frac{E(\mathbf{xy})}{E(\mathbf{x}^2)}\right)^2 - \frac{E(\mathbf{y}^2)}{E(\mathbf{x}^2)}}. \quad (7.3.7)$$

The inequality (7.3.6) is then valid for all a if the term under the square root is negative or zero. From this follows the assertion

$$\frac{\{E(\mathbf{xy})\}^2}{\{E(\mathbf{x}^2)\}^2} - \frac{E(\mathbf{y}^2)}{E(\mathbf{x}^2)} \leq 0.$$

If we now apply the inequality (7.3.5) to (7.3.4), one obtains

$$\{1 + B'(\lambda)\}^2 \leq E\{[\mathbf{S} - E(\mathbf{S})]^2\}E(\ell'^2). \quad (7.3.8)$$

We now use (7.2.3) in order to rewrite the expression for $E(\ell'^2)$,

$$\begin{aligned} E(\ell'^2) &= E\left\{\left(\sum_{j=1}^N \varphi(\mathbf{x}^{(j)}; \lambda)\right)^2\right\} \\ &= E\left\{\sum_{j=1}^N (\varphi(\mathbf{x}^{(j)}; \lambda))^2\right\} + E\left\{\sum_{i \neq j} \varphi(\mathbf{x}^{(i)}; \lambda)\varphi(\mathbf{x}^{(j)}; \lambda)\right\}. \end{aligned}$$

All terms on the right-hand side vanish, since for $i \neq j$

$$\begin{aligned} E\{\varphi(\mathbf{x}^{(i)}; \lambda)\varphi(\mathbf{x}^{(j)}; \lambda)\} &= E\{\varphi(\mathbf{x}^{(i)}; \lambda)\}E\{\varphi(\mathbf{x}^{(j)}; \lambda)\}, \\ E\{\varphi(\mathbf{x}; \lambda)\} &= \int_{-\infty}^{\infty} \frac{f'(x; \lambda)}{f(x; \lambda)} f(x; \lambda) dx = \int f'(x; \lambda) dx, \end{aligned}$$

and

$$\int_{-\infty}^{\infty} f(x; \lambda) dx = 1.$$

By differentiating the last line with respect to λ one obtains

$$\int_{-\infty}^{\infty} f'(x; \lambda) dx = 0.$$

Thus one has simply

$$E(\ell'^2) = E \left\{ \sum_{j=1}^N (\varphi(\mathbf{x}^{(j)}; \lambda))^2 \right\} = E \left\{ \sum_{j=1}^N \left(\frac{f'(\mathbf{x}^{(j)}; \lambda)}{f(\mathbf{x}^{(j)}; \lambda)} \right)^2 \right\}.$$

Since the individual terms of the sum are independent, the expectation value of the sum is simply the sum of the expectation values. The individual expectation values do not depend on the elements of the sample. Therefore one has

$$I(\lambda) = E(\ell'^2) = NE \left\{ \left(\frac{f'(x; \lambda)}{f(x; \lambda)} \right)^2 \right\}.$$

This expression is called the *information of the sample with respect to λ* . It is a non-negative number, which vanishes if the likelihood function does not depend on the parameter λ .

It is sometimes useful to write the information in a somewhat different form. To do this we differentiate the expression

$$E \left(\frac{f'(x; \lambda)}{f(x; \lambda)} \right) = \int_{-\infty}^{\infty} \frac{f'(x; \lambda)}{f(x; \lambda)} f(x; \lambda) dx = 0$$

once more with respect to λ and obtain

$$\begin{aligned} 0 &= \int_{-\infty}^{\infty} \left\{ \frac{f'^2}{f} + f \left(\frac{f'}{f} \right)' \right\} dx = \int_{-\infty}^{\infty} \left\{ \left(\frac{f'}{f} \right)^2 + \left(\frac{f'}{f} \right)' \right\} f dx \\ &= E \left\{ \left(\frac{f'}{f} \right)^2 \right\} + E \left\{ \left(\frac{f'}{f} \right)' \right\}. \end{aligned}$$

The information can then be written as

$$I(\lambda) = NE \left\{ \left(\frac{f'(x; \lambda)}{f(x; \lambda)} \right)^2 \right\} = -NE \left\{ \left(\frac{f'(x; \lambda)}{f(x; \lambda)} \right)' \right\}$$

or

$$I(\lambda) = E(\ell'^2) = -E(\ell''). \quad (7.3.9)$$

The inequality (7.3.8) can now be written in the following way:

$$\{1 + B'(\lambda)\}^2 \leq \sigma^2(\mathbf{S})I(\lambda)$$

or

$$\sigma^2(\mathbf{S}) \geq \frac{\{1 + B'(\lambda)\}^2}{I(\lambda)}. \quad (7.3.10)$$

This is the *information inequality*. It gives the connection between the bias and the variance of an estimator and the information of a sample. It should be noted that in its derivation no assumption about the estimator was made. The right-hand side of the inequality (7.3.10) is therefore a lower bound for the variance of an estimator. It is called the *minimum variance bound* or *Cramer–Rao bound*. In cases where the bias does not depend on λ , i.e., particularly in cases of vanishing bias, the inequality (7.3.10) simplifies to

$$\sigma^2(\mathbf{S}) \geq 1/I(\lambda). \quad (7.3.11)$$

This relation justifies using the name information. As the information of a sample increases, the variance of an estimator can be made smaller.

We now ask under which circumstances the minimum variance bound is attained, or explicitly, when the equals sign in the relation (7.3.10) holds. In the inequality (7.3.6), one has equality if $(a\mathbf{x} + \mathbf{y})$ vanishes, since only then does one have $E\{(a\mathbf{x} + \mathbf{y})^2\} = 0$ for all values of a , \mathbf{x} , and \mathbf{y} . Applied to (7.3.8), this means that

$$\ell' + a(\mathbf{S} - E(\mathbf{S})) = 0$$

or

$$\ell' = A(\lambda)(\mathbf{S} - E(\mathbf{S})). \quad (7.3.12)$$

Here A means an arbitrary quantity that does not depend on the sample $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}$, but may be, however, a function of λ . By integration we obtain

$$\ell = \int \ell' d\lambda = B(\lambda)\mathbf{S} + C(\lambda) + D \quad (7.3.13)$$

and finally

$$L = d \exp\{B(\lambda)\mathbf{S} + C(\lambda)\}. \quad (7.3.14)$$

The quantities d and D do not depend on λ .

We thus see that estimators attain the minimum variance bound when the likelihood function is of the special form (7.3.14). They are therefore called *minimum variance estimators*.

For the case of an unbiased minimum variance estimator we obtain from (7.3.11)

$$\sigma^2(\mathbf{S}) = \frac{1}{I(\lambda)} = \frac{1}{E(\ell'^2)}. \quad (7.3.15)$$

By substituting (7.3.12) one obtains

$$\sigma^2(\mathbf{S}) = \frac{1}{(A(\lambda))^2 E\{(\mathbf{S} - E(\mathbf{S}))^2\}} = \frac{1}{(A(\lambda))^2 \sigma^2(\mathbf{S})}$$

or

$$\sigma^2(\mathbf{S}) = \frac{1}{|A(\lambda)|}. \quad (7.3.16)$$

If instead of (7.3.14) only the weaker requirement

$$L = g(\mathbf{S}, \lambda) c(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}) \quad (7.3.17)$$

holds, then the estimator \mathbf{S} is said to be *sufficient* for λ . One can show [see, e.g., Kendall and Stuart, Vol. 2 (1967)], that no other estimator can contribute information about λ that is not already contained in S , if the requirement (7.3.17) is fulfilled. Hence the name “sufficient estimator” (or statistic).

Example 7.4: Estimator for the parameter of the Poisson distribution

Consider the *Poisson distribution* (5.4.1)

$$f(k) = \frac{\lambda^k}{k!} e^{-\lambda}.$$

The likelihood function of a sample $\mathbf{k}^{(1)}, \mathbf{k}^{(2)}, \dots, \mathbf{k}^{(N)}$ is

$$\ell = \sum_{j=1}^N \{\mathbf{k}^{(j)} \ln \lambda - \ln(\mathbf{k}^{(j)}!) - \lambda\}$$

and its derivative with respect to λ is

$$\begin{aligned} \frac{d\ell}{d\lambda} = \ell' &= \sum_{j=1}^N \left\{ \frac{\mathbf{k}^{(j)}}{\lambda} - 1 \right\} = \frac{1}{\lambda} \sum_{j=1}^N \{\mathbf{k}^{(j)} - \lambda\}, \\ \ell' &= \frac{N}{\lambda} (\bar{\mathbf{k}} - \lambda). \end{aligned} \quad (7.3.18)$$

Comparing with (7.3.12) and (7.3.16) shows that the arithmetic mean \bar{k} is an unbiased minimum variance estimator with variance λ/N . ■

Example 7.5: Estimator for the parameter of the binomial distribution

The likelihood function of a sample from a *binomial distribution* with the parameters $p = \lambda$, $q = 1 - \lambda$ is given directly by (5.1.3),

$$L(\mathbf{k}, \lambda) = \binom{n}{\mathbf{k}} \lambda^{\mathbf{k}} (1 - \lambda)^{n - \mathbf{k}}.$$

(The result of the sample can be summarized by the statement that in n experiments, the event A occurred \mathbf{k} times; see Sect. 5.1.) One then has

$$\begin{aligned} \ell &= \ln L = \mathbf{k} \ln \lambda + (n - \mathbf{k}) \ln(1 - \lambda) + \ln \binom{n}{\mathbf{k}}, \\ \ell' &= \frac{\mathbf{k}}{\lambda} - \frac{n - \mathbf{k}}{1 - \lambda} = \frac{n}{\lambda(1 - \lambda)} \left(\frac{\mathbf{k}}{n} - \lambda \right). \end{aligned}$$

By comparing with (7.3.12) and (7.3.16) one finds \mathbf{k}/n to be a minimum variance estimator with variance $\lambda(1 - \lambda)/n$. ■

Example 7.6: Law of error combination (“Quadratic averaging of individual errors”)

We now return to the problem of Example 7.2 of repeated measurements of the same quantity with varying uncertainties, or expressed in another way, to the problem of drawing a sample from normal distributions with the same mean λ and different but known variances σ_j . From (7.2.7) we obtain

$$\frac{d\ell}{d\lambda} = \ell' = \sum_{j=1}^N \frac{\mathbf{x}^{(j)} - \lambda}{\sigma_j^2}.$$

We can rewrite this expression as

$$\begin{aligned} \ell' &= \sum \frac{\mathbf{x}^{(j)}}{\sigma_j^2} - \sum \frac{\lambda}{\sigma_j^2} \\ &= \sum_{j=1}^N \frac{1}{\sigma_j^2} \left\{ \frac{\sum \frac{\mathbf{x}^{(j)}}{\sigma_j^2}}{\sum \frac{1}{\sigma_j^2}} - \lambda \right\}. \end{aligned}$$

As in Example 7.2 we recognize

$$\mathbf{S} = \tilde{\lambda} = \frac{\sum \frac{\mathbf{x}^{(j)}}{\sigma_j^2}}{\sum \frac{1}{\sigma_j^2}} \quad (7.3.19)$$

as an unbiased estimator for λ . Comparing with (7.3.12) shows that it is also a minimum variance estimator. From (7.3.16) one determines that its variance is

$$\sigma^2(\tilde{\lambda}) = \left(\sum_{j=1}^N \frac{1}{\sigma_j^2} \right)^{-1}. \quad (7.3.20)$$

The relation (7.3.20) often goes by the name of the *law of error combination* or *quadratic averaging of individual errors*. It could have been obtained by application of the rule of error propagation (3.8.7) to (7.3.19). If we identify $\sigma(\tilde{\lambda})$ as the error of the estimator $\tilde{\lambda}$ and σ_j as the error of the j th measurement, then we can write it in its usual form

$$\Delta\tilde{\lambda} = \left(\frac{1}{(\Delta x_1)^2} + \frac{1}{(\Delta x_2)^2} + \cdots + \frac{1}{(\Delta x_n)^2} \right)^{-\frac{1}{2}}. \quad (7.3.21)$$

If all of the measurements have the same error $\sigma = \sigma_j$, Eqs. (7.3.19), (7.3.20) simplify to

$$\tilde{\lambda} = \bar{\mathbf{x}}, \quad \sigma^2(\tilde{\lambda}) = \sigma^2/n,$$

which we have already found in Sect. 6.2. ■

7.4 Asymptotic Properties of the Likelihood Function and Maximum-Likelihood Estimators

We can now show heuristically several important properties of the likelihood function and maximum-likelihood estimators for very large data samples, that is, for the limit $N \rightarrow \infty$. The estimator $\mathbf{S} = \tilde{\lambda}$ was defined as the solution to the likelihood equation

$$\ell'(\lambda) = \sum_{j=1}^N \left(\frac{f'(\mathbf{x}^{(j)}; \lambda)}{f(\mathbf{x}^{(j)}; \lambda)} \right)_{\tilde{\lambda}} = 0. \quad (7.4.1)$$

Let us assume that $\ell'(\lambda)$ can be differentiated with respect to λ one more time, so that we can expand it in a series around the point $\lambda = \tilde{\lambda}$,

$$\ell'(\lambda) = \ell'(\tilde{\lambda}) + (\lambda - \tilde{\lambda})\ell''(\tilde{\lambda}) + \cdots. \quad (7.4.2)$$

The first term on the right side vanishes because of Eq. (7.4.1). In the second term one can write explicitly

$$\ell''(\tilde{\lambda}) = \sum_{j=1}^N \left(\frac{f'(\mathbf{x}^{(j)}; \lambda)}{f(\mathbf{x}^{(j)}; \lambda)} \right)'_{\tilde{\lambda}}.$$

This expression has the form of the mean value of a sample. For very large N it can be replaced by the expectation value of the population (Sect. 6.2),

$$\ell''(\tilde{\lambda}) = NE \left\{ \left(\frac{f'(x; \lambda)}{f(x; \lambda)} \right)'_{\tilde{\lambda}} \right\}. \quad (7.4.3)$$

Using Eq. (7.3.9) we can now write

$$\ell''(\tilde{\lambda}) = E(\ell''(\tilde{\lambda})) = -E(\ell'^2(\tilde{\lambda})) = -I(\tilde{\lambda}) = -1/b^2. \quad (7.4.4)$$

In this way we can replace the expression for $\ell''(\tilde{\lambda})$, which is a function the sample $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}$, by the quantity $-1/b^2$, which only depends on the probability density f and the estimator $\tilde{\lambda}$. If one neglects higher-order terms, Eq. (7.4.2) can be expressed as

$$\ell'(\lambda) = -\frac{1}{b^2}(\lambda - \tilde{\lambda}). \quad (7.4.5)$$

By integration one obtains

$$\ell(\lambda) = -\frac{1}{2b^2}(\lambda - \tilde{\lambda})^2 + c.$$

Inserting $\lambda = \tilde{\lambda}$ gives $c = \ell(\tilde{\lambda})$, or

$$\ell(\lambda) - \ell(\tilde{\lambda}) = -\frac{1}{2} \frac{(\lambda - \tilde{\lambda})^2}{b^2}. \quad (7.4.6)$$

By exponentiation one obtains

$$L(\lambda) = k \exp\{-(\lambda - \tilde{\lambda})^2/2b^2\}, \quad (7.4.7)$$

where k is a constant. The likelihood function $L(\lambda)$ has the form of a normal distribution with mean $\tilde{\lambda}$ and variance b^2 . At the values $\lambda = \tilde{\lambda} \pm b$, where λ is one standard deviation from $\tilde{\lambda}$, one has

$$-(\ell(\lambda) - \ell(\tilde{\lambda})) = \frac{1}{2}. \quad (7.4.8)$$

We can now compare (7.4.7) with Eqs. (7.3.12) and (7.3.16). Since we are estimating the parameter λ , we must write $\mathbf{S} = \tilde{\lambda}$ and thus $E(\mathbf{S}) = \lambda$. The estimator $\tilde{\lambda}$ is therefore an unbiased minimum variance estimator with variance

$$\sigma^2(\tilde{\lambda}) = b^2 = \frac{1}{I(\tilde{\lambda})} = \frac{1}{E(\ell'^2(\tilde{\lambda}))} = -\frac{1}{E(\ell''(\tilde{\lambda}))}. \quad (7.4.9)$$

Since the estimator $\tilde{\lambda}$ only possesses this property for the limiting case $N \rightarrow \infty$, we call it *asymptotically unbiased*. This is equivalent to the statement

that the maximum likelihood estimator is consistent (Sect. 6.1). For the same reason the likelihood function is called *asymptotically normal*.

In Sect. 7.2 we interpreted the likelihood function $L(\lambda)$ as a measure of the probability that the true value λ_0 of a parameter is equal to λ . The result of an estimator is often represented in abbreviated form,

$$\lambda = \tilde{\lambda} \pm \sigma(\tilde{\lambda}) = \tilde{\lambda} \pm \Delta\tilde{\lambda}.$$

Since the likelihood function is asymptotically normal, at least in the case of large samples, i.e., many measurements, this can be interpreted by saying that the probability that the true value λ_0 lies in the interval

$$\tilde{\lambda} - \Delta\tilde{\lambda} < \lambda_0 < \tilde{\lambda} + \Delta\tilde{\lambda}$$

is 68.3% (Sect. 5.8). In practice the relation above is used for large but finite samples. Unfortunately one cannot construct any general rule for determining when a sample is large enough for this procedure to be reliable. Clearly if N is finite, (7.4.3) is only an approximation, whose accuracy depends not only on N , but also on the particular probability density $f(x; \lambda)$.

Example 7.7: Determination of the mean lifetime from a small number of decays

The probability that a radioactive nucleus, which exists at time $t = 0$, decays in the time interval between t and $t + dt$ is

$$f(t) dt = \frac{1}{\tau} \exp(-t/\tau) dt.$$

For observed decay times t_1, t_2, \dots, t_N the likelihood function is

$$L = \frac{1}{\tau^N} \exp \left\{ -\frac{1}{\tau} \sum_{i=1}^N t_i \right\} = \frac{1}{\tau^N} \exp \left\{ -\frac{N}{\tau} \bar{t} \right\}.$$

Its logarithm is

$$\ell = \ln L = -\frac{N}{\tau} \bar{t} - N \ln \tau$$

with the derivative

$$\ell' = \frac{N}{\tau} \left(\frac{\bar{t}}{\tau} - 1 \right) = \frac{N}{\tau^2} (\bar{t} - \tau).$$

Comparing with (7.3.12) we see that $\tilde{\tau} = \bar{t}$ is the maximum likelihood solution, which has a variance of $\sigma^2(\tau) = \tau^2/N$. For $\tau = \tilde{\tau} = \bar{t}$ one obtains $\Delta\tilde{\tau} = \bar{t}/\sqrt{N}$.

For $\tilde{\tau} = \bar{t}$ one has

$$\ell(\tilde{\tau}) = \ell_{\max} = -N(1 + \ln \bar{t}).$$

We can write

$$-(\ell(\tau) - \ell(\tilde{\tau})) = N \left(\frac{\tilde{t}}{\tau} + \ln \frac{\tau}{\tilde{t}} - 1 \right).$$

From this expression for the log-likelihood function one cannot so easily recognize the asymptotic form (7.4.6) for $N \rightarrow \infty$. For small values of N it clearly does not have this form. Corresponding to (7.4.8), we want to use the values $\tau_+ = \tilde{\tau} + \Delta_+$ and $\tau_- = \tilde{\tau} - \Delta_-$, where one has

$$-(\ell(\tau_{\pm}) - \ell(\tilde{\tau})) = \frac{1}{2}$$

for the *asymmetric errors* Δ_+ , Δ_- . Clearly we expect for $N \rightarrow \infty$ that Δ_+ , $\Delta_- \rightarrow \Delta\tilde{\tau} = \sigma(\tilde{\tau})$.

In Fig. 7.2 the N observed decay times t_i are marked as vertical tick marks on the abscissa for various small values of N . In addition the function $-(\ell - \ell_{\max}) = -(\ell(\tau) - \ell(\tilde{\tau}))$ is plotted. The points τ_+ and τ_- are found where a horizontal line intersects $-(\ell - \ell_{\max}) = 1/2$. The point $\tilde{\tau}$ is indicated by an additional mark on the horizontal line. One sees that with increasing N the function $-(\ell - \ell_{\max})$ approaches more and more the symmetric parabolic form and that the errors Δ_+ , Δ_- , and $\Delta\tilde{\tau}$ become closer to each other. ■

7.5 Simultaneous Estimation of Several Parameters: Confidence Intervals

We have already given a system of equations (7.2.5) allowing the simultaneous determination of p parameters $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_p)$. It turns out that it is not the parameter determination but rather the estimation of their errors that becomes significantly more complicated in the case of several parameters. In particular we will need to consider correlations as well as errors of the parameters.

We extend our considerations from Sect. 7.4 on the properties of the likelihood function to the case of several parameters. The log-likelihood function

$$\ell(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}; \boldsymbol{\lambda}) = \sum_{j=1}^N \ln f(\mathbf{x}^{(j)}; \boldsymbol{\lambda}) \quad (7.5.1)$$

can be expanded in a series about the point

$$\tilde{\boldsymbol{\lambda}} = (\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_p) \quad (7.5.2)$$

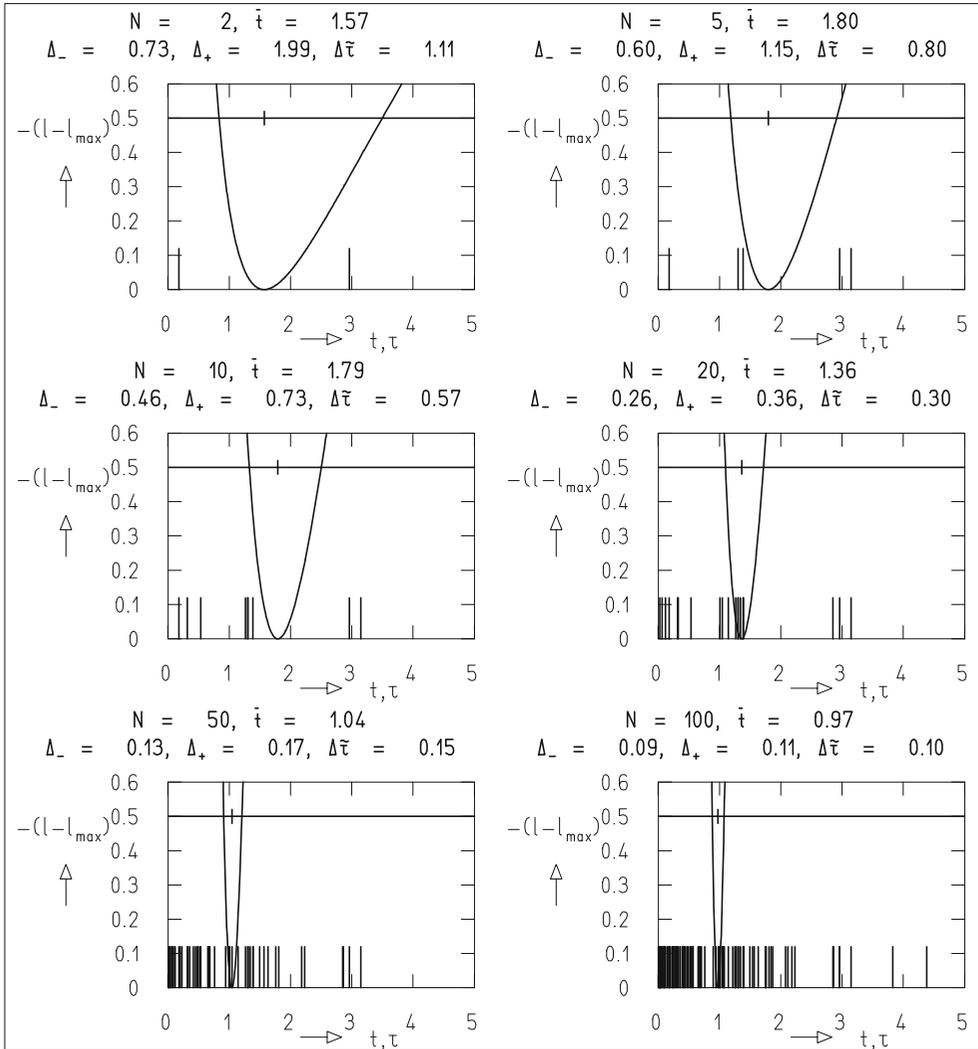


Fig. 7.2: Data and log-likelihood function for Example 7.7.

to give

$$\begin{aligned}
 \ell(\boldsymbol{\lambda}) &= \ell(\tilde{\boldsymbol{\lambda}}) + \sum_{k=1}^p \left(\frac{\partial \ell}{\partial \lambda_k} \right)_{\tilde{\boldsymbol{\lambda}}} (\lambda_k - \tilde{\lambda}_k) \\
 &\quad + \frac{1}{2} \sum_{\ell=1}^p \sum_{m=1}^p \left(\frac{\partial^2 \ell}{\partial \lambda_\ell \partial \lambda_m} \right)_{\tilde{\boldsymbol{\lambda}}} (\lambda_\ell - \tilde{\lambda}_\ell) (\lambda_m - \tilde{\lambda}_m) + \dots \quad (7.5.3)
 \end{aligned}$$

Since by the definition of $\tilde{\boldsymbol{\lambda}}$ one has

$$\left(\frac{\partial \ell}{\partial \lambda_k} \right)_{\tilde{\boldsymbol{\lambda}}} = 0, \quad k = 1, 2, \dots, p, \quad (7.5.4)$$

which holds for all k , the series simplifies to

$$-\left(\ell(\boldsymbol{\lambda}) - \ell(\tilde{\boldsymbol{\lambda}})\right) = \frac{1}{2}(\boldsymbol{\lambda} - \tilde{\boldsymbol{\lambda}})^T A(\boldsymbol{\lambda} - \tilde{\boldsymbol{\lambda}}) + \dots \quad (7.5.5)$$

with

$$-A = \begin{pmatrix} \frac{\partial^2 \ell}{\partial \lambda_1^2} & \frac{\partial^2 \ell}{\partial \lambda_1 \partial \lambda_2} & \cdots & \frac{\partial^2 \ell}{\partial \lambda_1 \partial \lambda_p} \\ \frac{\partial^2 \ell}{\partial \lambda_1 \partial \lambda_2} & \frac{\partial^2 \ell}{\partial \lambda_2^2} & \cdots & \frac{\partial^2 \ell}{\partial \lambda_2 \partial \lambda_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \ell}{\partial \lambda_1 \partial \lambda_p} & \frac{\partial^2 \ell}{\partial \lambda_2 \partial \lambda_p} & \cdots & \frac{\partial^2 \ell}{\partial \lambda_p^2} \end{pmatrix}_{\boldsymbol{\lambda}=\tilde{\boldsymbol{\lambda}}} \quad (7.5.6)$$

In the limit $N \rightarrow \infty$ we can replace the elements of A , which still depend on the specific sample, by the corresponding expectation values,

$$B = E(A) = \begin{pmatrix} E\left(\frac{\partial^2 \ell}{\partial \lambda_1^2}\right) & E\left(\frac{\partial^2 \ell}{\partial \lambda_1 \partial \lambda_2}\right) & \cdots & E\left(\frac{\partial^2 \ell}{\partial \lambda_1 \partial \lambda_p}\right) \\ E\left(\frac{\partial^2 \ell}{\partial \lambda_1 \partial \lambda_2}\right) & E\left(\frac{\partial^2 \ell}{\partial \lambda_2^2}\right) & \cdots & E\left(\frac{\partial^2 \ell}{\partial \lambda_2 \partial \lambda_p}\right) \\ \vdots & \vdots & \ddots & \vdots \\ E\left(\frac{\partial^2 \ell}{\partial \lambda_1 \partial \lambda_p}\right) & E\left(\frac{\partial^2 \ell}{\partial \lambda_2 \partial \lambda_p}\right) & \cdots & E\left(\frac{\partial^2 \ell}{\partial \lambda_p^2}\right) \end{pmatrix}_{\boldsymbol{\lambda}=\tilde{\boldsymbol{\lambda}}} \quad (7.5.7)$$

If we neglect higher-order terms, we can give the likelihood function as

$$L = k \exp\left\{-\frac{1}{2}(\boldsymbol{\lambda} - \tilde{\boldsymbol{\lambda}})^T B(\boldsymbol{\lambda} - \tilde{\boldsymbol{\lambda}})\right\}. \quad (7.5.8)$$

Comparing with (5.10.1) shows that this is a p -dimensional normal distribution with mean $\tilde{\boldsymbol{\lambda}}$ and covariance matrix

$$C = B^{-1}. \quad (7.5.9)$$

The variances of the maximum likelihood estimators $\tilde{\lambda}_1, \tilde{\lambda}_2, \dots, \tilde{\lambda}_p$ are given by the diagonal elements of the matrix (7.5.9). The off-diagonal elements are the covariances between all possible pairs of estimators,

$$\sigma^2(\tilde{\lambda}_i) = c_{ii}, \quad (7.5.10)$$

$$\text{cov}(\tilde{\lambda}_j, \tilde{\lambda}_k) = c_{jk}. \quad (7.5.11)$$

For the correlation coefficient between the estimators $\tilde{\lambda}_j, \tilde{\lambda}_k$ we can define

$$\rho(\tilde{\lambda}_j, \tilde{\lambda}_k) = \frac{\text{cov}(\tilde{\lambda}_j, \tilde{\lambda}_k)}{\sigma(\tilde{\lambda}_j)\sigma(\tilde{\lambda}_k)}. \quad (7.5.12)$$

As in the case of a single parameter, the square roots of the variances are given as the error or standard deviations of the estimators,

$$\Delta\tilde{\lambda}_i = \sigma(\tilde{\lambda}_i) = \sqrt{c_{ii}}. \quad (7.5.13)$$

In Sect. 7.4 we determined that by giving the maximum-likelihood estimator and its error one defines a region that contains the true value of the parameter with a probability of 68.3%. Since the likelihood function in the several parameter case is asymptotically a Gaussian distribution of several variables, this region is not determined only by the errors, but rather by the entire covariance matrix. In the special case of two parameters this is the covariance ellipse, which we introduced in Sect. 5.10.

The expression (7.5.8) has (with the replacement $\mathbf{x} = \boldsymbol{\lambda}$) exactly the form of (5.10.1). We can therefore use it for all of the results of Sect. 5.10. For the exponent one has

$$-\frac{1}{2}(\boldsymbol{\lambda} - \tilde{\boldsymbol{\lambda}})^T B(\boldsymbol{\lambda} - \tilde{\boldsymbol{\lambda}}) = -\frac{1}{2}g(\boldsymbol{\lambda}) = -\left\{\ell(\boldsymbol{\lambda}) - \ell(\tilde{\boldsymbol{\lambda}})\right\}. \quad (7.5.14)$$

In the parameter space spanned by $\lambda_1, \dots, \lambda_p$, the covariance ellipsoid of the distribution (7.5.8) is then determined by the condition

$$g(\boldsymbol{\lambda}) = 1 = 2\left\{\ell(\boldsymbol{\lambda}) - \ell(\tilde{\boldsymbol{\lambda}})\right\}. \quad (7.5.15)$$

For other values of $g(\boldsymbol{\lambda})$ one obtains the confidence ellipsoids introduced in Sect. 5.10. For smaller values of N , the series (7.5.3) cannot be truncated and the approximation (7.5.7) is not valid. Nevertheless, the solution (7.5.4) can clearly still be computed. For a given probability W one obtains instead of a confidence ellipsoid a *confidence region*, contained within the hypersurface

$$g(\boldsymbol{\lambda}) = 2\left\{\ell(\boldsymbol{\lambda}) - \ell(\tilde{\boldsymbol{\lambda}})\right\} = \text{const.} \quad (7.5.16)$$

The value of g is determined in the same way as for the confidence ellipsoid as in Sect. 5.10.

In Example 7.7 we computed the region $\tilde{\lambda} - \Delta_- < \lambda < \tilde{\lambda} + \Delta_+$ for the case of a single variable. This corresponds to a confidence region with the probability content 68.3%.

Example 7.8: Estimation of the mean and variance of a normal distribution

We want to determine the mean λ_1 and standard deviation λ_2 of a normal distribution using a sample of size N . This problem occurs, for example, in the measurement of the range of α -particles in matter. Because of the statistical nature of the energy loss through a large number of independent individual collisions, the range of the individual particles is Gaussian distributed about some mean value. By measuring the range $\mathbf{x}^{(j)}$ of N different particles, the mean λ_1 and “straggling constant” $\lambda_2 = \sigma$ can be estimated. We obtain the likelihood function

$$L = \prod_{j=1}^N \frac{1}{\lambda_2 \sqrt{2\pi}} \exp\left(-\frac{(\mathbf{x}^{(j)} - \lambda_1)^2}{2\lambda_2^2}\right)$$

and

$$\ell = -\frac{1}{2} \sum_{j=1}^N \frac{(\mathbf{x}^{(j)} - \lambda_1)^2}{\lambda_2^2} - N \ln \lambda_2 - \text{const.}$$

The system of likelihood equations is

$$\begin{aligned} \frac{\partial \ell}{\partial \lambda_1} &= \sum_{j=1}^N \frac{\mathbf{x}^{(j)} - \lambda_1}{\lambda_2^2} = 0, \\ \frac{\partial \ell}{\partial \lambda_2} &= \frac{1}{\lambda_2^3} \sum_{j=1}^N (\mathbf{x}^{(j)} - \lambda_1)^2 - \frac{N}{\lambda_2} = 0. \end{aligned}$$

Its solution is

$$\begin{aligned} \tilde{\lambda}_1 &= \frac{1}{N} \sum_{j=1}^N \mathbf{x}^{(j)}, \\ \tilde{\lambda}_2 &= \sqrt{\frac{\sum_{j=1}^N (\mathbf{x}^{(j)} - \tilde{\lambda}_1)^2}{N}}. \end{aligned}$$

For the estimator of the mean, the maximum-likelihood method leads to the arithmetic mean of the individual measurements. For the variance it gives the quantity \mathbf{s}^2 (6.2.4), which has a small bias, and not \mathbf{s}^2 , the unbiased estimator (6.2.6).

Let us now determine the matrix B . The second derivatives are

$$\frac{\partial^2 \ell}{\partial \lambda_1^2} = -\frac{N}{\lambda_2^2},$$

$$\begin{aligned}\frac{\partial^2 \ell}{\partial \lambda_1 \partial \lambda_2} &= -\frac{2 \sum (\mathbf{x}^{(j)} - \lambda_1)}{\lambda_2^3}, \\ \frac{\partial^2 \ell}{\partial \lambda_2^2} &= -\frac{3 \sum (\mathbf{x}^{(j)} - \lambda_1)^2}{\lambda_2^4} + \frac{N}{\lambda_2^2}.\end{aligned}$$

We use the procedure of (7.5.7), substitute λ_1, λ_2 by $\tilde{\lambda}_1, \tilde{\lambda}_2$ and find

$$B = \begin{pmatrix} N/\tilde{\lambda}_2^2 & 0 \\ 0 & 2N/\tilde{\lambda}_2^2 \end{pmatrix}$$

or for the covariance matrix

$$C = B^{-1} = \begin{pmatrix} \tilde{\lambda}_2^2/N & 0 \\ 0 & \tilde{\lambda}_2^2/2N \end{pmatrix}.$$

We interpret the diagonal elements as the errors of the corresponding parameters, i.e.,

$$\Delta \tilde{\lambda}_1 = \tilde{\lambda}_2/\sqrt{N}, \quad \Delta \tilde{\lambda}_2 = \tilde{\lambda}_2/\sqrt{2N}.$$

The estimators for λ_1 and λ_2 are not correlated. ■

Example 7.9: Estimators for the parameters of a two-dimensional normal distribution

To conclude we consider a population described by a two-dimensional normal distribution (Sect. 5.10)

$$\begin{aligned}f(x_1, x_2) &= \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2(1-\rho^2)}\right. \\ &\quad \left. \times \left\{ \frac{(x_1 - a_1)^2}{\sigma_1^2} + \frac{(x_2 - a_2)^2}{\sigma_2^2} - 2\rho \frac{(x_1 - a_1)(x_2 - a_2)}{\sigma_1\sigma_2} \right\} \right].\end{aligned}$$

By constructing and solving a system of five simultaneous likelihood equations for the five parameters $a_1, a_2, \sigma_1^2, \sigma_2^2, \rho$ we obtain for the maximum-likelihood estimators

$$\begin{aligned}\bar{x}_1 &= \frac{1}{N} \sum_{j=1}^N x_1^{(j)}, & \bar{x}_2 &= \frac{1}{N} \sum_{j=1}^N x_2^{(j)}, \\ s_1'^2 &= \frac{1}{N} \sum_{j=1}^N (x_1^{(j)} - \bar{x}_1)^2, & s_2'^2 &= \frac{1}{N} \sum_{j=1}^N (x_2^{(j)} - \bar{x}_2)^2,\end{aligned}$$

$$r = \frac{\sum_{j=1}^N (\mathbf{x}_1^{(j)} - \bar{\mathbf{x}}_1)(\mathbf{x}_2^{(j)} - \bar{\mathbf{x}}_2)}{N \mathbf{s}'_1 \mathbf{s}'_2}. \quad (7.5.17)$$

Exactly as in Example 7.8, the estimators for the variances \mathbf{s}'_1 and \mathbf{s}'_2 are biased. This also holds for the expression (7.5.17), the *sample correlation coefficient* r . Like all maximum likelihood estimators, r is consistent, i.e., it provides a good estimation of ϱ for very large samples. For $N \rightarrow \infty$ the probability density of the random variable r becomes a normal distribution with mean ϱ and variance

$$\sigma^2(r) = (1 - \varrho^2)^2 / N. \quad (7.5.18)$$

For finite samples the distribution is asymmetric. It is therefore important to have a sufficiently large sample before applying Eq. (7.5.17). As a rule of thumb, $N \geq 500$ is usually recommended. ■

7.6 Example Programs

Example Program 7.1: The class E1MaxLife computes the mean lifetime and its asymmetric errors from a small number of radioactive decays

The program performs the computations and the graphical display for the problem described in Example 7.7. First by Monte Carlo method a total of N decay times t_i of radioactive nuclei with a mean lifetime of $\tau = 1$ are simulated. The number N of decays and also the seeds for the random number generator are entered interactively.

Example Program 7.2: The class E2MaxLife computes the maximum-likelihood estimates of the parameters of a bivariate normal distribution from a simulated sample

The program asks interactively for the number n_{exp} of experiments to simulate (i.e., of the samples to be treated consecutively), for the size n_{pt} of each sample and for the means a_1, a_2 , the standard deviations σ_1, σ_2 , and the correlation coefficient ρ of a bivariate Gaussian distribution.

The covariance matrix C of the normal distribution is calculated and the generator of random numbers from a multivariate Gaussian distribution is initialized. Each sample is generated and then analyzed, i.e., the quantities $\bar{x}_1, \bar{x}_2, s'_1, s'_2$, and r are computed, which are estimates of $a_1, a_2, \sigma_1, \sigma_2$, and ρ [cf. (7.5.17)]. The quantities are displayed for each sample.

Suggestions: Choose $n_{\text{exp}} = 20$, keep all other parameters fixed, and study the statistical fluctuations of r for $n_{\text{pt}} = 5, 50, 500$. Use the values $\rho = 0, 0.5, 0.95$.