



# Maximum likelihood estimation

(Redirected from [Maximum likelihood](#))

In [statistics](#), **maximum likelihood estimation** (**MLE**) is a method of [estimating](#) the [parameters](#) of an assumed probability distribution, given some observed data. This is achieved by [maximizing](#) a [likelihood function](#) so that, under the assumed [statistical model](#), the [observed data](#) is most probable. The [point](#) in the [parameter space](#) that maximizes the likelihood function is called the maximum likelihood estimate.<sup>[1]</sup> The logic of maximum likelihood is both intuitive and flexible, and as such the method has become a dominant means of [statistical inference](#).<sup>[2][3][4]</sup>

If the likelihood function is [differentiable](#), the [derivative test](#) for finding maxima can be applied. In some cases, the first-order conditions of the likelihood function can be solved analytically; for instance, the [ordinary least squares](#) estimator for a [linear regression](#) model maximizes the likelihood when the random errors are assumed to have [normal](#) distributions with the same variance.<sup>[5]</sup>

From the perspective of [Bayesian inference](#), MLE is generally equivalent to [maximum a posteriori](#) (MAP) [estimation](#) with [uniform prior distributions](#) (or a [normal prior distribution](#) with a standard deviation of infinity). In [frequentist inference](#), MLE is a special case of an [extremum estimator](#), with the objective function being the likelihood.

## Principles

We model a set of observations as a random sample from an unknown [joint probability distribution](#) which is expressed in terms of a set of [parameters](#). The goal of maximum likelihood estimation is to determine the parameters for which the observed data have the highest joint probability. We write the parameters governing the joint distribution as a vector  $\boldsymbol{\theta} = [\theta_1, \theta_2, \dots, \theta_k]^T$  so that this distribution falls within a [parametric family](#)  $\{f(\cdot; \boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \boldsymbol{\Theta}\}$ , where  $\boldsymbol{\Theta}$  is called the *parameter space*, a finite-dimensional subset of [Euclidean space](#). Evaluating the joint density at the observed data sample  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  gives a real-valued function,

$$\mathcal{L}_n(\boldsymbol{\theta}) = \mathcal{L}_n(\boldsymbol{\theta}; \mathbf{y}) = f_n(\mathbf{y}; \boldsymbol{\theta}) ,$$

which is called the [likelihood function](#). For [independent and identically distributed](#) random variables,  $f_n(\mathbf{y}; \boldsymbol{\theta})$  will be the product of univariate [density functions](#):

$$f_n(\mathbf{y}; \boldsymbol{\theta}) = \prod_{k=1}^n f_k^{\text{univar}}(y_k; \boldsymbol{\theta}) .$$

The goal of maximum likelihood estimation is to find the values of the model parameters that maximize the likelihood function over the parameter space,<sup>[6]</sup> that is

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \mathcal{L}_n(\boldsymbol{\theta}; \mathbf{y}) .$$

Intuitively, this selects the parameter values that make the observed data most probable. The specific value  $\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}_n(\mathbf{y}) \in \boldsymbol{\Theta}$  that maximizes the likelihood function  $\mathcal{L}_n$  is called the maximum likelihood estimate. Further, if the function  $\hat{\boldsymbol{\theta}}_n : \mathbb{R}^n \rightarrow \boldsymbol{\Theta}$  so defined is [measurable](#), then it is called the maximum likelihood [estimator](#). It is

generally a function defined over the sample space, i.e. taking a given sample as its argument. A sufficient but not necessary condition for its existence is for the likelihood function to be continuous over a parameter space  $\Theta$  that is compact.<sup>[7]</sup> For an open  $\Theta$  the likelihood function may increase without ever reaching a supremum value.

In practice, it is often convenient to work with the natural logarithm of the likelihood function, called the log-likelihood:

$$\ell(\theta; \mathbf{y}) = \ln \mathcal{L}_n(\theta; \mathbf{y}) .$$

Since the logarithm is a monotonic function, the maximum of  $\ell(\theta; \mathbf{y})$  occurs at the same value of  $\theta$  as does the maximum of  $\mathcal{L}_n$ .<sup>[8]</sup> If  $\ell(\theta; \mathbf{y})$  is differentiable in  $\Theta$ , sufficient conditions for the occurrence of a maximum (or a minimum) are

$$\frac{\partial \ell}{\partial \theta_1} = 0, \quad \frac{\partial \ell}{\partial \theta_2} = 0, \quad \dots, \quad \frac{\partial \ell}{\partial \theta_k} = 0 ,$$

known as the likelihood equations. For some models, these equations can be explicitly solved for  $\hat{\theta}$ , but in general no closed-form solution to the maximization problem is known or available, and an MLE can only be found via numerical optimization. Another problem is that in finite samples, there may exist multiple roots for the likelihood equations.<sup>[9]</sup> Whether the identified root  $\hat{\theta}$  of the likelihood equations is indeed a (local) maximum depends on whether the matrix of second-order partial and cross-partial derivatives, the so-called Hessian matrix

$$\mathbf{H}(\hat{\theta}) = \begin{bmatrix} \left. \frac{\partial^2 \ell}{\partial \theta_1^2} \right|_{\theta=\hat{\theta}} & \left. \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} \right|_{\theta=\hat{\theta}} & \cdots & \left. \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_k} \right|_{\theta=\hat{\theta}} \\ \left. \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_1} \right|_{\theta=\hat{\theta}} & \left. \frac{\partial^2 \ell}{\partial \theta_2^2} \right|_{\theta=\hat{\theta}} & \cdots & \left. \frac{\partial^2 \ell}{\partial \theta_2 \partial \theta_k} \right|_{\theta=\hat{\theta}} \\ \vdots & \vdots & \ddots & \vdots \\ \left. \frac{\partial^2 \ell}{\partial \theta_k \partial \theta_1} \right|_{\theta=\hat{\theta}} & \left. \frac{\partial^2 \ell}{\partial \theta_k \partial \theta_2} \right|_{\theta=\hat{\theta}} & \cdots & \left. \frac{\partial^2 \ell}{\partial \theta_k^2} \right|_{\theta=\hat{\theta}} \end{bmatrix} ,$$

is negative semi-definite at  $\hat{\theta}$ , as this indicates local concavity. Conveniently, most common probability distributions – in particular the exponential family – are logarithmically concave.<sup>[10][11]</sup>

## Restricted parameter space

While the domain of the likelihood function—the parameter space—is generally a finite-dimensional subset of Euclidean space, additional restrictions sometimes need to be incorporated into the estimation process. The parameter space can be expressed as

$$\Theta = \{ \theta : \theta \in \mathbb{R}^k, h(\theta) = 0 \} ,$$

where  $h(\theta) = [h_1(\theta), h_2(\theta), \dots, h_r(\theta)]$  is a vector-valued function mapping  $\mathbb{R}^k$  into  $\mathbb{R}^r$ . Estimating the true parameter  $\theta$  belonging to  $\Theta$  then, as a practical matter, means to find the maximum of the likelihood function subject to the constraint  $h(\theta) = 0$ .

Theoretically, the most natural approach to this constrained optimization problem is the method of substitution, that is "filling out" the restrictions  $h_1, h_2, \dots, h_r$  to a set  $h_1, h_2, \dots, h_r, h_{r+1}, \dots, h_k$  in such a way that  $h^* = [h_1, h_2, \dots, h_k]$  is a one-to-one function from  $\mathbb{R}^k$  to itself, and reparameterize the likelihood function by setting  $\phi_i = h_i(\theta_1, \theta_2, \dots, \theta_k)$ .<sup>[12]</sup> Because of the equivariance of the maximum likelihood estimator, the

properties of the MLE apply to the restricted estimates also.<sup>[13]</sup> For instance, in a multivariate normal distribution the covariance matrix  $\Sigma$  must be positive-definite; this restriction can be imposed by replacing  $\Sigma = \Gamma^T \Gamma$ , where  $\Gamma$  is a real upper triangular matrix and  $\Gamma^T$  is its transpose.<sup>[14]</sup>

In practice, restrictions are usually imposed using the method of Lagrange which, given the constraints as defined above, leads to the *restricted likelihood equations*

$$\frac{\partial \ell}{\partial \theta} - \frac{\partial h(\theta)^T}{\partial \theta} \lambda = 0 \text{ and } h(\theta) = 0,$$

where  $\lambda = [\lambda_1, \lambda_2, \dots, \lambda_r]^T$  is a column-vector of Lagrange multipliers and  $\frac{\partial h(\theta)^T}{\partial \theta}$  is the  $k \times r$  Jacobian matrix of partial derivatives.<sup>[12]</sup> Naturally, if the constraints are not binding at the maximum, the Lagrange multipliers should be zero.<sup>[15]</sup> This in turn allows for a statistical test of the "validity" of the constraint, known as the Lagrange multiplier test.

## Nonparametric maximum likelihood estimation

Nonparametric maximum likelihood estimation can be performed using the empirical likelihood.

## Properties

---

A maximum likelihood estimator is an extremum estimator obtained by maximizing, as a function of  $\theta$ , the objective function  $\hat{\ell}(\theta; x)$ . If the data are independent and identically distributed, then we have

$$\hat{\ell}(\theta; x) = \frac{1}{n} \sum_{i=1}^n \ln f(x_i | \theta),$$

this being the sample analogue of the expected log-likelihood  $\ell(\theta) = \mathbb{E}[\ln f(x_i | \theta)]$ , where this expectation is taken with respect to the true density.

Maximum-likelihood estimators have no optimum properties for finite samples, in the sense that (when evaluated on finite samples) other estimators may have greater concentration around the true parameter-value.<sup>[16]</sup> However, like other estimation methods, maximum likelihood estimation possesses a number of attractive limiting properties: As the sample size increases to infinity, sequences of maximum likelihood estimators have these properties:

- Consistency: the sequence of MLEs converges in probability to the value being estimated.
- Invariance: If  $\hat{\theta}$  is the maximum likelihood estimator for  $\theta$ , and if  $g(\theta)$  is any transformation of  $\theta$ , then the maximum likelihood estimator for  $\alpha = g(\theta)$  is  $\hat{\alpha} = g(\hat{\theta})$ . This property is less commonly known as *functional equivariance*. The invariance property holds for arbitrary transformation  $g$ , although the proof simplifies if  $g$  is restricted to one-to-one transformations.
- Efficiency, i.e. it achieves the Cramér–Rao lower bound when the sample size tends to infinity. This means that no consistent estimator has lower asymptotic mean squared error than the MLE (or other estimators attaining this bound), which also means that MLE has asymptotic normality.
- Second-order efficiency after correction for bias.

## Consistency

Under the conditions outlined below, the maximum likelihood estimator is consistent. The consistency means that if the data were generated by  $f(\cdot; \theta_0)$  and we have a sufficiently large number of observations  $n$ , then it is possible to find the value of  $\theta_0$  with arbitrary precision. In mathematical terms this means that as  $n$  goes to infinity the estimator  $\hat{\theta}$  converges in probability to its true value:

$$\hat{\theta}_{\text{mle}} \xrightarrow{P} \theta_0.$$

Under slightly stronger conditions, the estimator converges almost surely (or *strongly*):

$$\hat{\theta}_{\text{mle}} \xrightarrow{\text{a.s.}} \theta_0.$$

In practical applications, data is never generated by  $f(\cdot; \theta_0)$ . Rather,  $f(\cdot; \theta_0)$  is a model, often in idealized form, of the process generated by the data. It is a common aphorism in statistics that all models are wrong. Thus, true consistency does not occur in practical applications. Nevertheless, consistency is often considered to be a desirable property for an estimator to have.

To establish consistency, the following conditions are sufficient.<sup>[17]</sup>

1. Identification of the model:

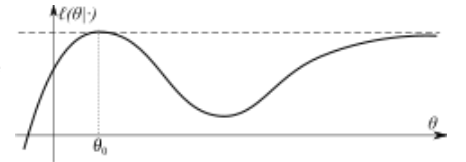
$$\theta \neq \theta_0 \Leftrightarrow f(\cdot | \theta) \neq f(\cdot | \theta_0).$$

In other words, different parameter values  $\theta$  correspond to different distributions within the model. If this condition did not hold, there would be some value  $\theta_1$  such that  $\theta_0$  and  $\theta_1$  generate an identical distribution of the observable data. Then we would not be able to distinguish between these two parameters even with an infinite amount of data—these parameters would have been observationally equivalent.

The identification condition is absolutely necessary for the ML estimator to be consistent. When this condition holds, the limiting likelihood function  $\ell(\theta|\cdot)$  has unique global maximum at  $\theta_0$ .

2. Compactness: the parameter space  $\Theta$  of the model is compact.

The identification condition establishes that the log-likelihood has a unique global maximum. Compactness implies that the likelihood cannot approach the maximum value arbitrarily close at some other point (as demonstrated for example in the picture on the right).



Compactness is only a sufficient condition and not a necessary condition. Compactness can be replaced by some other conditions, such as:

- both concavity of the log-likelihood function and compactness of some (nonempty) upper level sets of the log-likelihood function, or
- existence of a compact neighborhood  $N$  of  $\theta_0$  such that outside of  $N$  the log-likelihood function is less than the maximum by at least some  $\varepsilon > 0$ .

3. Continuity: the function  $\ln f(x | \theta)$  is continuous in  $\theta$  for almost all values of  $x$ :

$$\mathbb{P} \left[ \ln f(x | \theta) \in C^0(\Theta) \right] = 1.$$

The continuity here can be replaced with a slightly weaker condition of upper semi-continuity.

4. Dominance: there exists  $D(x)$  integrable with respect to the distribution  $f(x | \theta_0)$  such that

$$|\ln f(x | \theta)| < D(x) \quad \text{for all } \theta \in \Theta.$$

By the uniform law of large numbers, the dominance condition together with continuity establish the uniform convergence in probability of the log-likelihood:

$$\sup_{\theta \in \Theta} |\hat{\ell}(\theta | \mathbf{x}) - \ell(\theta)| \xrightarrow{P} 0.$$

The dominance condition can be employed in the case of i.i.d. observations. In the non-i.i.d. case, the uniform convergence in probability can be checked by showing that the sequence  $\hat{\ell}(\theta | \mathbf{x})$  is stochastically equicontinuous. If one wants to demonstrate that the ML estimator  $\hat{\theta}$  converges to  $\theta_0$  almost surely, then a stronger condition of uniform convergence almost surely has to be imposed:

$$\sup_{\theta \in \Theta} \|\hat{\ell}(\theta | \mathbf{x}) - \ell(\theta)\| \xrightarrow{\text{a.s.}} 0.$$

Additionally, if (as assumed above) the data were generated by  $f(\cdot; \theta_0)$ , then under certain conditions, it can also be shown that the maximum likelihood estimator converges in distribution to a normal distribution. Specifically,<sup>[18]</sup>

$$\sqrt{n}(\hat{\theta}_{\text{mle}} - \theta_0) \xrightarrow{d} \mathcal{N}(0, I^{-1})$$

where  $I$  is the Fisher information matrix.

## Functional invariance

The maximum likelihood estimator selects the parameter value which gives the observed data the largest possible probability (or probability density, in the continuous case). If the parameter consists of a number of components, then we define their separate maximum likelihood estimators, as the corresponding component of the MLE of the complete parameter. Consistent with this, if  $\hat{\theta}$  is the MLE for  $\theta$ , and if  $g(\theta)$  is any transformation of  $\theta$ , then the MLE for  $\alpha = g(\theta)$  is by definition<sup>[19]</sup>

$$\hat{\alpha} = g(\hat{\theta}).$$

It maximizes the so-called profile likelihood:

$$\bar{L}(\alpha) = \sup_{\theta: \alpha = g(\theta)} L(\theta).$$

The MLE is also equivariant with respect to certain transformations of the data. If  $\mathbf{y} = g(\mathbf{x})$  where  $g$  is one to one and does not depend on the parameters to be estimated, then the density functions satisfy

$$f_Y(\mathbf{y}) = \frac{f_X(\mathbf{x})}{|g'(\mathbf{x})|}$$

and hence the likelihood functions for  $\mathbf{X}$  and  $\mathbf{Y}$  differ only by a factor that does not depend on the model parameters.

For example, the MLE parameters of the log-normal distribution are the same as those of the normal distribution fitted to the logarithm of the data.

## Efficiency

As assumed above, if the data were generated by  $f(\cdot; \theta_0)$ , then under certain conditions, it can also be shown that the maximum likelihood estimator converges in distribution to a normal distribution. It is  $\sqrt{n}$ -consistent and asymptotically efficient, meaning that it reaches the Cramér–Rao bound. Specifically,<sup>[18]</sup>

$$\sqrt{n} \left( \hat{\theta}_{\text{mle}} - \theta_0 \right) \xrightarrow{d} \mathcal{N}(0, \mathcal{I}^{-1}),$$

where  $\mathcal{I}$  is the Fisher information matrix:

$$\mathcal{I}_{jk} = \mathbb{E} \left[ - \frac{\partial^2 \ln f_{\theta_0}(X_t)}{\partial \theta_j \partial \theta_k} \right].$$

In particular, it means that the bias of the maximum likelihood estimator is equal to zero up to the order  $\frac{1}{\sqrt{n}}$ .

## Second-order efficiency after correction for bias

However, when we consider the higher-order terms in the expansion of the distribution of this estimator, it turns out that  $\theta_{\text{mle}}$  has bias of order  $\frac{1}{n}$ . This bias is equal to (componentwise)<sup>[20]</sup>

$$b_h \equiv \mathbb{E} \left[ \left( \hat{\theta}_{\text{mle}} - \theta_0 \right)_h \right] = \frac{1}{n} \sum_{i,j,k=1}^m \mathcal{I}^{hi} \mathcal{I}^{jk} \left( \frac{1}{2} K_{ijk} + J_{j,ik} \right)$$

where  $\mathcal{I}^{jk}$  (with superscripts) denotes the  $(j,k)$ -th component of the *inverse* Fisher information matrix  $\mathcal{I}^{-1}$ , and

$$\frac{1}{2} K_{ijk} + J_{j,ik} = \mathbb{E} \left[ \frac{1}{2} \frac{\partial^3 \ln f_{\theta_0}(X_t)}{\partial \theta_i \partial \theta_j \partial \theta_k} + \frac{\partial \ln f_{\theta_0}(X_t)}{\partial \theta_j} \frac{\partial^2 \ln f_{\theta_0}(X_t)}{\partial \theta_i \partial \theta_k} \right].$$

Using these formulae it is possible to estimate the second-order bias of the maximum likelihood estimator, and *correct* for that bias by subtracting it:

$$\hat{\theta}_{\text{mle}}^* = \hat{\theta}_{\text{mle}} - \hat{b}.$$

This estimator is unbiased up to the terms of order  $\frac{1}{n}$ , and is called the **bias-corrected maximum likelihood estimator**.

This bias-corrected estimator is *second-order efficient* (at least within the curved exponential family), meaning that it has minimal mean squared error among all second-order bias-corrected estimators, up to the terms of the order  $\frac{1}{n^2}$ . It is possible to continue this process, that is to derive the third-order bias-correction term, and so on. However, the maximum likelihood estimator is *not* third-order efficient.<sup>[21]</sup>

## Relation to Bayesian inference

A maximum likelihood estimator coincides with the most probable Bayesian estimator given a uniform prior distribution on the parameters. Indeed, the maximum a posteriori estimate is the parameter  $\theta$  that maximizes the probability of  $\theta$  given the data, given by Bayes' theorem:

$$\mathbb{P}(\theta \mid x_1, x_2, \dots, x_n) = \frac{f(x_1, x_2, \dots, x_n \mid \theta) \mathbb{P}(\theta)}{\mathbb{P}(x_1, x_2, \dots, x_n)}$$

where  $\mathbb{P}(\theta)$  is the prior distribution for the parameter  $\theta$  and where  $\mathbb{P}(x_1, x_2, \dots, x_n)$  is the probability of the data averaged over all parameters. Since the denominator is independent of  $\theta$ , the Bayesian estimator is obtained by maximizing  $f(x_1, x_2, \dots, x_n | \theta) \mathbb{P}(\theta)$  with respect to  $\theta$ . If we further assume that the prior  $\mathbb{P}(\theta)$  is a uniform distribution, the Bayesian estimator is obtained by maximizing the likelihood function  $f(x_1, x_2, \dots, x_n | \theta)$ . Thus the Bayesian estimator coincides with the maximum likelihood estimator for a uniform prior distribution  $\mathbb{P}(\theta)$ .

### Application of maximum-likelihood estimation in Bayes decision theory

In many practical applications in machine learning, maximum-likelihood estimation is used as the model for parameter estimation.

The Bayesian Decision theory is about designing a classifier that minimizes total expected risk, especially, when the costs (the loss function) associated with different decisions are equal, the classifier is minimizing the error over the whole distribution.<sup>[22]</sup>

Thus, the Bayes Decision Rule is stated as

"decide  $w_1$  if  $\mathbb{P}(w_1|x) > \mathbb{P}(w_2|x)$ ; otherwise decide  $w_2$  "

where  $w_1, w_2$  are predictions of different classes. From a perspective of minimizing error, it can also be stated as

$$w = \arg \max_w \int_{-\infty}^{\infty} \mathbb{P}(\text{error} | x) \mathbb{P}(x) dx$$

where

$$\mathbb{P}(\text{error} | x) = \mathbb{P}(w_1 | x)$$

if we decide  $w_2$  and  $\mathbb{P}(\text{error} | x) = \mathbb{P}(w_2 | x)$  if we decide  $w_1$  .

By applying Bayes' theorem

$$\mathbb{P}(w_i | x) = \frac{\mathbb{P}(x | w_i) \mathbb{P}(w_i)}{\mathbb{P}(x)},$$

and if we further assume the zero-or-one loss function, which is a same loss for all errors, the Bayes Decision rule can be reformulated as:

$$h_{\text{Bayes}} = \arg \max_w [\mathbb{P}(x | w) \mathbb{P}(w)] ,$$

where  $h_{\text{Bayes}}$  is the prediction and  $\mathbb{P}(w)$  is the prior probability.

### Relation to minimizing Kullback–Leibler divergence and cross entropy

Finding  $\hat{\theta}$  that maximizes the likelihood is asymptotically equivalent to finding the  $\hat{\theta}$  that defines a probability distribution ( $Q_{\hat{\theta}}$ ) that has a minimal distance, in terms of Kullback–Leibler divergence, to the real probability distribution from which our data were generated (i.e., generated by  $P_{\theta_0}$ ).<sup>[23]</sup> In an ideal world, P and Q are the same (and the only thing unknown is  $\theta$  that defines P), but even if they are not and the model we use is misspecified, still the MLE will give us the "closest" distribution (within the restriction of a model Q that depends on  $\hat{\theta}$ ) to the real distribution  $P_{\theta_0}$ .<sup>[24]</sup>

For simplicity of notation, let's assume that  $P=Q$ . Let there be  $n$  i.i.d data samples  $\mathbf{y} = (y_1, y_2, \dots, y_n)$  from some probability  $y \sim P_{\theta_0}$ , that we try to estimate by finding  $\hat{\theta}$  that will maximize the likelihood using  $P_{\theta}$ , then:

$$\begin{aligned}
 \hat{\theta} &= \arg \max_{\theta} L_{P_{\theta}}(\mathbf{y}) = \arg \max_{\theta} P_{\theta}(\mathbf{y}) = \arg \max_{\theta} P(\mathbf{y} | \theta) \\
 &= \arg \max_{\theta} \prod_{i=1}^n P(y_i | \theta) = \arg \max_{\theta} \sum_{i=1}^n \log P(y_i | \theta) \\
 &= \arg \max_{\theta} \left( \sum_{i=1}^n \log P(y_i | \theta) - \sum_{i=1}^n \log P(y_i | \theta_0) \right) = \arg \max_{\theta} \sum_{i=1}^n (\log P(y_i | \theta) - \log P(y_i | \theta_0)) \\
 &= \arg \max_{\theta} \sum_{i=1}^n \log \frac{P(y_i | \theta)}{P(y_i | \theta_0)} = \arg \min_{\theta} \sum_{i=1}^n \log \frac{P(y_i | \theta_0)}{P(y_i | \theta)} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \log \frac{P(y_i | \theta_0)}{P(y_i | \theta)} \\
 &= \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n h_{\theta}(y_i) \xrightarrow{n \rightarrow \infty} \arg \min_{\theta} E[h_{\theta}(y)] \\
 &= \arg \min_{\theta} \int P_{\theta_0}(y) h_{\theta}(y) dy = \arg \min_{\theta} \int P_{\theta_0}(y) \log \frac{P(y | \theta_0)}{P(y | \theta)} dy \\
 &= \arg \min_{\theta} D_{\text{KL}}(P_{\theta_0} \parallel P_{\theta})
 \end{aligned}$$

Where  $h_{\theta}(x) = \log \frac{P(x | \theta_0)}{P(x | \theta)}$ . Using  $h$  helps see how we are using the law of large numbers to move from the average of  $h(x)$  to the expectancy of it using the law of the unconscious statistician. The first several transitions have to do with laws of logarithm and that finding  $\hat{\theta}$  that maximizes some function will also be the one that maximizes some monotonic transformation of that function (i.e.: adding/multiplying by a constant).

Since cross entropy is just Shannon's entropy plus KL divergence, and since the entropy of  $P_{\theta_0}$  is constant, then the MLE is also asymptotically minimizing cross entropy.<sup>[25]</sup>

## Examples

### Discrete uniform distribution

Consider a case where  $n$  tickets numbered from 1 to  $n$  are placed in a box and one is selected at random (see uniform distribution); thus, the sample size is 1. If  $n$  is unknown, then the maximum likelihood estimator  $\hat{n}$  of  $n$  is the number  $m$  on the drawn ticket. (The likelihood is 0 for  $n < m$ ,  $1/n$  for  $n \geq m$ , and this is greatest when  $n = m$ . Note that the maximum likelihood estimate of  $n$  occurs at the lower extreme of possible values  $\{m, m + 1, \dots\}$ , rather than somewhere in the "middle" of the range of possible values, which would result in less bias.) The expected value of the number  $m$  on the drawn ticket, and therefore the expected value of  $\hat{n}$ , is  $(n + 1)/2$ . As a result, with a sample size of 1, the maximum likelihood estimator for  $n$  will systematically underestimate  $n$  by  $(n - 1)/2$ .

### Discrete distribution, finite parameter space

Suppose one wishes to determine just how biased an unfair coin is. Call the probability of tossing a 'head'  $p$ . The goal then becomes to determine  $p$ .

Suppose the coin is tossed 80 times: i.e. the sample might be something like  $x_1 = H, x_2 = T, \dots, x_{80} = T$ , and the count of the number of heads "H" is observed.



The probability of tossing tails is  $1 - p$  (so here  $p$  is  $\theta$  above). Suppose the outcome is 49 heads and 31 tails, and suppose the coin was taken from a box containing three coins: one which gives heads with probability  $p = \frac{1}{3}$ , one which gives heads with probability  $p = \frac{1}{2}$  and another which gives heads with probability  $p = \frac{2}{3}$ . The coins have lost their labels, so which one it was is unknown. Using maximum likelihood estimation, the coin that has the largest likelihood can be found, given the data that were observed. By using the probability mass function of the binomial distribution with sample size equal to 80, number successes equal to 49 but for different values of  $p$  (the "probability of success"), the likelihood function (defined below) takes one of three values:

$$\mathbb{P} [ H = 49 \mid p = \frac{1}{3} ] = \binom{80}{49} \left(\frac{1}{3}\right)^{49} \left(1 - \frac{1}{3}\right)^{31} \approx 0.000,$$

$$\mathbb{P} [ H = 49 \mid p = \frac{1}{2} ] = \binom{80}{49} \left(\frac{1}{2}\right)^{49} \left(1 - \frac{1}{2}\right)^{31} \approx 0.012,$$

$$\mathbb{P} [ H = 49 \mid p = \frac{2}{3} ] = \binom{80}{49} \left(\frac{2}{3}\right)^{49} \left(1 - \frac{2}{3}\right)^{31} \approx 0.054 .$$

The likelihood is maximized when  $p = \frac{2}{3}$ , and so this is the *maximum likelihood estimate* for  $p$ .

### Discrete distribution, continuous parameter space

Now suppose that there was only one coin but its  $p$  could have been any value  $0 \leq p \leq 1$ . The likelihood function to be maximised is

$$L(p) = f_D(H = 49 \mid p) = \binom{80}{49} p^{49} (1 - p)^{31} ,$$

and the maximisation is over all possible values  $0 \leq p \leq 1$ .

One way to maximize this function is by differentiating with respect to  $p$  and setting to zero:

$$0 = \frac{\partial}{\partial p} \left( \binom{80}{49} p^{49} (1 - p)^{31} \right) ,$$

$$0 = 49p^{48} (1 - p)^{31} - 31p^{49} (1 - p)^{30}$$

$$= p^{48} (1 - p)^{30} [49(1 - p) - 31p]$$

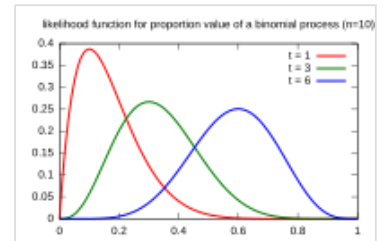
$$= p^{48} (1 - p)^{30} [49 - 80p] .$$

This is a product of three terms. The first term is 0 when  $p = 0$ . The second is 0 when  $p = 1$ . The third is zero when  $p = \frac{49}{80}$ . The solution that maximizes the likelihood is clearly  $p = \frac{49}{80}$  (since  $p = 0$  and  $p = 1$  result in a likelihood of 0). Thus the *maximum likelihood estimator* for  $p$  is  $\frac{49}{80}$ .

This result is easily generalized by substituting a letter such as  $s$  in the place of 49 to represent the observed number of 'successes' of our Bernoulli trials, and a letter such as  $n$  in the place of 80 to represent the number of Bernoulli trials. Exactly the same calculation yields  $\frac{s}{n}$  which is the maximum likelihood estimator for any sequence of  $n$  Bernoulli trials resulting in  $s$  'successes'.

### Continuous distribution, continuous parameter space

For the normal distribution  $\mathcal{N}(\mu, \sigma^2)$  which has probability density function



Likelihood function for proportion value of a binomial process  
( $n = 10$ )

$$f(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right),$$

the corresponding probability density function for a sample of  $n$  independent identically distributed normal random variables (the likelihood) is

$$f(x_1, \dots, x_n | \mu, \sigma^2) = \prod_{i=1}^n f(x_i | \mu, \sigma^2) = \left(\frac{1}{2\pi\sigma^2}\right)^{n/2} \exp\left(-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right).$$

This family of distributions has two parameters:  $\theta = (\mu, \sigma)$ ; so we maximize the likelihood,  $\mathcal{L}(\mu, \sigma^2) = f(x_1, \dots, x_n | \mu, \sigma^2)$ , over both parameters simultaneously, or if possible, individually.

Since the logarithm function itself is a continuous strictly increasing function over the range of the likelihood, the values which maximize the likelihood will also maximize its logarithm (the log-likelihood itself is not necessarily strictly increasing). The log-likelihood can be written as follows:

$$\log(\mathcal{L}(\mu, \sigma^2)) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

(Note: the log-likelihood is closely related to information entropy and Fisher information.)

We now compute the derivatives of this log-likelihood as follows.

$$0 = \frac{\partial}{\partial \mu} \log(\mathcal{L}(\mu, \sigma^2)) = 0 - \frac{-2n(\bar{x} - \mu)}{2\sigma^2}.$$

where  $\bar{x}$  is the sample mean. This is solved by

$$\hat{\mu} = \bar{x} = \sum_{i=1}^n \frac{x_i}{n}.$$

This is indeed the maximum of the function, since it is the only turning point in  $\mu$  and the second derivative is strictly less than zero. Its expected value is equal to the parameter  $\mu$  of the given distribution,

$$\mathbb{E}[\hat{\mu}] = \mu,$$

which means that the maximum likelihood estimator  $\hat{\mu}$  is unbiased.

Similarly we differentiate the log-likelihood with respect to  $\sigma$  and equate to zero:

$$0 = \frac{\partial}{\partial \sigma} \log(\mathcal{L}(\mu, \sigma^2)) = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2.$$

which is solved by

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2.$$

Inserting the estimate  $\mu = \hat{\mu}$  we obtain

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n x_i x_j.$$

To calculate its expected value, it is convenient to rewrite the expression in terms of zero-mean random variables (statistical error)  $\delta_i \equiv \mu - x_i$ . Expressing the estimate in these variables yields

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (\mu - \delta_i)^2 - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (\mu - \delta_i)(\mu - \delta_j).$$

Simplifying the expression above, utilizing the facts that  $\mathbb{E} [\delta_i] = 0$  and  $\mathbb{E} [\delta_i^2] = \sigma^2$ , allows us to obtain

$$\mathbb{E} [\hat{\sigma}^2] = \frac{n-1}{n} \sigma^2.$$

This means that the estimator  $\hat{\sigma}^2$  is biased for  $\sigma^2$ . It can also be shown that  $\hat{\sigma}$  is biased for  $\sigma$ , but that both  $\hat{\sigma}^2$  and  $\hat{\sigma}$  are consistent.

Formally we say that the *maximum likelihood estimator* for  $\theta = (\mu, \sigma^2)$  is

$$\hat{\theta} = (\hat{\mu}, \hat{\sigma}^2).$$

In this case the MLEs could be obtained individually. In general this may not be the case, and the MLEs would have to be obtained simultaneously.

The normal log-likelihood at its maximum takes a particularly simple form:

$$\log (\mathcal{L}(\hat{\mu}, \hat{\sigma})) = \frac{-n}{2} (\log(2\pi\hat{\sigma}^2) + 1)$$

This maximum log-likelihood can be shown to be the same for more general least squares, even for non-linear least squares. This is often used in determining likelihood-based approximate confidence intervals and confidence regions, which are generally more accurate than those using the asymptotic normality discussed above.

## Non-independent variables

---

It may be the case that variables are correlated, that is, not independent. Two random variables  $y_1$  and  $y_2$  are independent only if their joint probability density function is the product of the individual probability density functions, i.e.

$$f(y_1, y_2) = f(y_1)f(y_2)$$

Suppose one constructs an order- $n$  Gaussian vector out of random variables  $(y_1, \dots, y_n)$ , where each variable has means given by  $(\mu_1, \dots, \mu_n)$ . Furthermore, let the covariance matrix be denoted by  $\Sigma$ . The joint probability density function of these  $n$  random variables then follows a multivariate normal distribution given by:

$$f(y_1, \dots, y_n) = \frac{1}{(2\pi)^{n/2} \sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2} [y_1 - \mu_1, \dots, y_n - \mu_n] \Sigma^{-1} [y_1 - \mu_1, \dots, y_n - \mu_n]^T\right)$$

In the bivariate case, the joint probability density function is given by:

$$f(y_1, y_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left[ -\frac{1}{2(1-\rho^2)} \left( \frac{(y_1 - \mu_1)^2}{\sigma_1^2} - \frac{2\rho(y_1 - \mu_1)(y_2 - \mu_2)}{\sigma_1\sigma_2} + \frac{(y_2 - \mu_2)^2}{\sigma_2^2} \right) \right]$$

In this and other cases where a joint density function exists, the likelihood function is defined as above, in the section "principles," using this density.

## Example

$X_1, X_2, \dots, X_m$  are counts in cells / boxes 1 up to  $m$ ; each box has a different probability (think of the boxes being bigger or smaller) and we fix the number of balls that fall to be  $n: x_1 + x_2 + \dots + x_m = n$ . The probability of each box is  $p_i$ , with a constraint:  $p_1 + p_2 + \dots + p_m = 1$ . This is a case in which the  $X_i$  s are not independent, the joint probability of a vector  $x_1, x_2, \dots, x_m$  is called the multinomial and has the form:

$$f(x_1, x_2, \dots, x_m \mid p_1, p_2, \dots, p_m) = \frac{n!}{\prod x_i!} \prod p_i^{x_i} = \binom{n}{x_1, x_2, \dots, x_m} p_1^{x_1} p_2^{x_2} \dots p_m^{x_m}$$

Each box taken separately against all the other boxes is a binomial and this is an extension thereof.

The log-likelihood of this is:

$$\ell(p_1, p_2, \dots, p_m) = \log n! - \sum_{i=1}^m \log x_i! + \sum_{i=1}^m x_i \log p_i$$

The constraint has to be taken into account and use the Lagrange multipliers:

$$L(p_1, p_2, \dots, p_m, \lambda) = \ell(p_1, p_2, \dots, p_m) + \lambda \left( 1 - \sum_{i=1}^m p_i \right)$$

By posing all the derivatives to be 0, the most natural estimate is derived

$$\hat{p}_i = \frac{x_i}{n}$$

Maximizing log likelihood, with and without constraints, can be an unsolvable problem in closed form, then we have to use iterative procedures.

## Iterative procedures

---

Except for special cases, the likelihood equations

$$\frac{\partial \ell(\theta; \mathbf{y})}{\partial \theta} = 0$$

cannot be solved explicitly for an estimator  $\hat{\theta} = \hat{\theta}(\mathbf{y})$ . Instead, they need to be solved iteratively: starting from an initial guess of  $\theta$  (say  $\hat{\theta}_1$ ), one seeks to obtain a convergent sequence  $\{\hat{\theta}_r\}$ . Many methods for this kind of optimization problem are available,<sup>[26][27]</sup> but the most commonly used ones are algorithms based on an updating formula of the form

$$\hat{\theta}_{r+1} = \hat{\theta}_r + \eta_r \mathbf{d}_r(\hat{\theta})$$

where the vector  $\mathbf{d}_r(\hat{\theta})$  indicates the descent direction of the  $r$ th "step," and the scalar  $\eta_r$  captures the "step length,"<sup>[28][29]</sup> also known as the learning rate.<sup>[30]</sup>

## Gradient descent method

(Note: here it is a maximization problem, so the sign before gradient is flipped)

$$\eta_r \in \mathbb{R}^+ \text{ that is small enough for convergence and } \mathbf{d}_r(\hat{\theta}) = \nabla \ell(\hat{\theta}_r; \mathbf{y})$$

Gradient descent method requires to calculate the gradient at the  $r$ th iteration, but no need to calculate the inverse of second-order derivative, i.e., the Hessian matrix. Therefore, it is computationally faster than Newton-Raphson method.

## Newton–Raphson method

$$\eta_r = 1 \text{ and } \mathbf{d}_r(\hat{\theta}) = -\mathbf{H}_r^{-1}(\hat{\theta}) \mathbf{s}_r(\hat{\theta})$$

where  $\mathbf{s}_r(\hat{\theta})$  is the score and  $\mathbf{H}_r^{-1}(\hat{\theta})$  is the inverse of the Hessian matrix of the log-likelihood function, both evaluated the  $r$ th iteration.<sup>[31][32]</sup> But because the calculation of the Hessian matrix is computationally costly, numerous alternatives have been proposed. The popular Berndt–Hall–Hall–Hausman algorithm approximates the Hessian with the outer product of the expected gradient, such that

$$\mathbf{d}_r(\hat{\theta}) = -\left[ \frac{1}{n} \sum_{t=1}^n \frac{\partial \ell(\theta; \mathbf{y})}{\partial \theta} \left( \frac{\partial \ell(\theta; \mathbf{y})}{\partial \theta} \right)^T \right]^{-1} \mathbf{s}_r(\hat{\theta})$$

## Quasi-Newton methods

Other quasi-Newton methods use more elaborate secant updates to give approximation of Hessian matrix.

### Davidon–Fletcher–Powell formula

DFP formula finds a solution that is symmetric, positive-definite and closest to the current approximate value of second-order derivative:

$$\mathbf{H}_{k+1} = (I - \gamma_k \mathbf{y}_k \mathbf{s}_k^T) \mathbf{H}_k (I - \gamma_k \mathbf{s}_k \mathbf{y}_k^T) + \gamma_k \mathbf{y}_k \mathbf{y}_k^T,$$

where

$$\begin{aligned} \mathbf{y}_k &= \nabla \ell(\mathbf{x}_k + \mathbf{s}_k) - \nabla \ell(\mathbf{x}_k), \\ \gamma_k &= \frac{1}{\mathbf{y}_k^T \mathbf{s}_k}, \\ \mathbf{s}_k &= \mathbf{x}_{k+1} - \mathbf{x}_k. \end{aligned}$$

### Broyden–Fletcher–Goldfarb–Shanno algorithm

BFGS also gives a solution that is symmetric and positive-definite:

$$B_{k+1} = B_k + \frac{y_k y_k^\top}{y_k^\top s_k} - \frac{B_k s_k s_k^\top B_k^\top}{s_k^\top B_k s_k},$$

where

$$\begin{aligned} y_k &= \nabla \ell(x_k + s_k) - \nabla \ell(x_k), \\ s_k &= x_{k+1} - x_k. \end{aligned}$$

BFGS method is not guaranteed to converge unless the function has a quadratic Taylor expansion near an optimum. However, BFGS can have acceptable performance even for non-smooth optimization instances

### **Fisher's scoring**

Another popular method is to replace the Hessian with the Fisher information matrix,  $\mathcal{I}(\theta) = \mathbb{E}[\mathbf{H}_r(\hat{\theta})]$ , giving us the Fisher scoring algorithm. This procedure is standard in the estimation of many methods, such as generalized linear models.

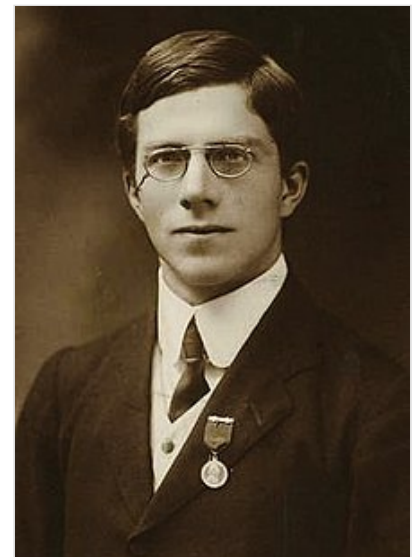
Although popular, quasi-Newton methods may converge to a stationary point that is not necessarily a local or global maximum,<sup>[33]</sup> but rather a local minimum or a saddle point. Therefore, it is important to assess the validity of the obtained solution to the likelihood equations, by verifying that the Hessian, evaluated at the solution, is both negative definite and well-conditioned.<sup>[34]</sup>

## **History**

---

Early users of maximum likelihood include Carl Friedrich Gauss, Pierre-Simon Laplace, Thorvald N. Thiele, and Francis Ysidro Edgeworth.<sup>[35][36]</sup> It was Ronald Fisher however, between 1912 and 1922, who singlehandedly created the modern version of the method.<sup>[37][38]</sup>

Maximum-likelihood estimation finally transcended heuristic justification in a proof published by Samuel S. Wilks in 1938, now called Wilks' theorem.<sup>[39]</sup> The theorem shows that the error in the logarithm of likelihood values for estimates from multiple independent observations is asymptotically  $\chi^2$ -distributed, which enables convenient determination of a confidence region around any estimate of the parameters. The only difficult part of Wilks' proof depends on the expected value of the Fisher information matrix, which is provided by a theorem proven by Fisher.<sup>[40]</sup> Wilks continued to improve on the generality of the theorem throughout his life, with his most general proof published in 1962.<sup>[41]</sup>



Ronald Fisher in 1913

Reviews of the development of maximum likelihood estimation have been provided by a number of authors.<sup>[42][43][44][45][46][47][48][49]</sup>

## **See also**

---

### **Related concepts**

- Akaike information criterion: a criterion to compare statistical models, based on MLE

- Extremum estimator: a more general class of estimators to which MLE belongs
- Fisher information: information matrix, its relationship to covariance matrix of ML estimates
- Mean squared error: a measure of how 'good' an estimator of a distributional parameter is (be it the maximum likelihood estimator or some other estimator)
- RANSAC: a method to estimate parameters of a mathematical model given data that contains outliers
- Rao–Blackwell theorem: yields a process for finding the best possible unbiased estimator (in the sense of having minimal mean squared error); the MLE is often a good starting place for the process
- Wilks' theorem: provides a means of estimating the size and shape of the region of roughly equally-probable estimates for the population's parameter values, using the information from a single sample, using a chi-squared distribution

## Other estimation methods

- Generalized method of moments: methods related to the likelihood equation in maximum likelihood estimation
- M-estimator: an approach used in robust statistics
- Maximum a posteriori (MAP) estimator: for a contrast in the way to calculate estimators when prior knowledge is postulated
- Maximum spacing estimation: a related method that is more robust in many situations
- Maximum entropy estimation
- Method of moments (statistics): another popular method for finding parameters of distributions
- Method of support, a variation of the maximum likelihood technique
- Minimum-distance estimation
- Partial likelihood methods for panel data
- Quasi-maximum likelihood estimator: an MLE estimator that is misspecified, but still consistent
- Restricted maximum likelihood: a variation using a likelihood function calculated from a transformed set of data

## References

---

1. Rossi, Richard J. (2018). *Mathematical Statistics: An Introduction to Likelihood Based Inference*. New York: John Wiley & Sons. p. 227. ISBN 978-1-118-77104-4.
2. Hendry, David F.; Nielsen, Bent (2007). *Econometric Modeling: A Likelihood Approach*. Princeton: Princeton University Press. ISBN 978-0-691-13128-3.
3. Chambers, Raymond L.; Steel, David G.; Wang, Suojin; Welsh, Alan (2012). *Maximum Likelihood Estimation for Sample Surveys*. Boca Raton: CRC Press. ISBN 978-1-58488-632-7.
4. Ward, Michael Don; Ahlquist, John S. (2018). *Maximum Likelihood for Social Science: Strategies for Analysis*. New York: Cambridge University Press. ISBN 978-1-107-18582-1.
5. Press, W.H.; Flannery, B.P.; Teukolsky, S.A.; Vetterling, W.T. (1992). "Least Squares as a Maximum Likelihood Estimator" ([https://books.google.com/books?id=gn\\_4mpdN9WkC&pg=PA651](https://books.google.com/books?id=gn_4mpdN9WkC&pg=PA651)). *Numerical Recipes in FORTRAN: The Art of Scientific Computing* (2nd ed.). Cambridge: Cambridge University Press. pp. 651–655. ISBN 0-521-43064-X.
6. Myung, I.J. (2003). "Tutorial on maximum likelihood Estimation". *Journal of Mathematical Psychology*. 47 (1): 90–100. doi:10.1016/S0022-2496(02)00028-7 ([https://doi.org/10.1016/S0022-2496\(02\)00028-7](https://doi.org/10.1016/S0022-2496(02)00028-7)).
7. Gourieroux, Christian; Monfort, Alain (1995). *Statistics and Econometrics Models* ([https://archive.org/details/statisticseconom00gour\\_434](https://archive.org/details/statisticseconom00gour_434)). Cambridge University Press. p. 161 ([https://archive.org/details/statisticseconom00gour\\_434/page/n172](https://archive.org/details/statisticseconom00gour_434/page/n172)). ISBN 0-521-40551-3.
8. Kane, Edward J. (1968). *Economic Statistics and Econometrics* (<https://archive.org/details/economicstatisti00kane>). New York, NY: Harper & Row. p. 179 (<https://archive.org/details/economicstatisti00kane/page/n200>).
9. Small, Christopher G.; Wang, Jinfang (2003). "Working with roots" (<https://books.google.com/books?id=hMrwQViiY5AC&pg=PA74>). *Numerical Methods for Nonlinear Estimating Equations*. Oxford University Press. pp. 74–124. ISBN 0-19-850688-0.

10. Kass, Robert E.; Vos, Paul W. (1997). *Geometrical Foundations of Asymptotic Inference* (<https://books.google.com/books?id=e43EAlfUPCwC&pg=PA14>). New York, NY: John Wiley & Sons. p. 14. ISBN 0-471-82668-5.
11. Papadopoulos, Alecos (25 September 2013). "Why we always put  $\log()$  before the joint pdf when we use MLE (Maximum likelihood Estimation)?" (<https://stats.stackexchange.com/q/70975>). *Stack Exchange*.
12. Silvey, S. D. (1975). *Statistical Inference* (<https://books.google.com/books?id=qIKLejbVMf4C&pg=PA79>). London, UK: Chapman and Hall. p. 79. ISBN 0-412-13820-4.
13. Olive, David (2004). "Does the MLE maximize the likelihood?" (<http://lagrange.math.siu.edu/Olive/simle.pdf>) (PDF). *Southern Illinois University*.
14. Schwallie, Daniel P. (1985). "Positive definite maximum likelihood covariance estimators". *Economics Letters*. **17** (1–2): 115–117. doi:[10.1016/0165-1765\(85\)90139-9](https://doi.org/10.1016/0165-1765(85)90139-9) (<https://doi.org/10.1016%2F0165-1765%2885%2990139-9>).
15. Magnus, Jan R. (2017). *Introduction to the Theory of Econometrics*. Amsterdam: VU University Press. pp. 64–65. ISBN 978-90-8659-766-6.
16. Pfanzagl (1994, p. 206)
17. By Theorem 2.5 in Newey, Whitney K.; McFadden, Daniel (1994). "Chapter 36: Large sample estimation and hypothesis testing". In Engle, Robert; McFadden, Dan (eds.). *Handbook of Econometrics, Vol. 4*. Elsevier Science. pp. 2111–2245. ISBN 978-0-444-88766-5.
18. By Theorem 3.3 in Newey, Whitney K.; McFadden, Daniel (1994). "Chapter 36: Large sample estimation and hypothesis testing". In Engle, Robert; McFadden, Dan (eds.). *Handbook of Econometrics, Vol. 4*. Elsevier Science. pp. 2111–2245. ISBN 978-0-444-88766-5.
19. Zacks, Shelemyahu (1971). *The Theory of Statistical Inference*. New York: John Wiley & Sons. p. 223. ISBN 0-471-98103-6.
20. See formula 20 in Cox, David R.; Snell, E. Joyce (1968). "A general definition of residuals". *Journal of the Royal Statistical Society, Series B*. **30** (2): 248–275. JSTOR 2984505 (<https://www.jstor.org/stable/2984505>).
21. Kano, Yutaka (1996). "Third-order efficiency implies fourth-order efficiency" (<https://doi.org/10.14490%2Fjjss1995.26.101>). *Journal of the Japan Statistical Society*. **26**: 101–117. doi:[10.14490/jjss1995.26.101](https://doi.org/10.14490/jjss1995.26.101) (<https://doi.org/10.14490%2Fjjss1995.26.101>).
22. Christensen, Henrik I. "Pattern Recognition" (<https://www.cc.gatech.edu/~hic/CS7616/pdf/lecture2.pdf>) (PDF) (lecture). Bayesian Decision Theory - CS 7616. Georgia Tech.
23. cmplx96 (<https://stats.stackexchange.com/users/177679/cmplx96>), Kullback–Leibler divergence, URL (version: 2017-11-18): <https://stats.stackexchange.com/q/314472> (at the youtube video, look at minutes 13 to 25)
24. Introduction to Statistical Inference | Stanford (Lecture 16 — MLE under model misspecification) (<https://web.stanford.edu/class/stats200/Lecture16.pdf>)
25. Sycorax says Reinstate Monica (<https://stats.stackexchange.com/users/22311/sycorax-says-reinstate-monica>), the relationship between maximizing the likelihood and minimizing the cross-entropy, URL (version: 2019-11-06): <https://stats.stackexchange.com/q/364237>
26. Fletcher, R. (1987). *Practical Methods of Optimization* (<https://archive.org/details/practicalmethods0000flet>) (Second ed.). New York, NY: John Wiley & Sons. ISBN 0-471-91547-5.
27. Nocedal, Jorge; Wright, Stephen J. (2006). *Numerical Optimization* (Second ed.). New York, NY: Springer. ISBN 0-387-30303-0.
28. Daganzo, Carlos (1979). *Multinomial Probit: The Theory and its Application to Demand Forecasting*. New York: Academic Press. pp. 61–78. ISBN 0-12-201150-3.
29. Gould, William; Pitblado, Jeffrey; Poi, Brian (2010). *Maximum Likelihood Estimation with Stata* (Fourth ed.). College Station: Stata Press. pp. 13–20. ISBN 978-1-59718-078-8.
30. Murphy, Kevin P. (2012). *Machine Learning: A Probabilistic Perspective* (<https://books.google.com/books?id=NZP6AQAAQBAJ&pg=PA247>). Cambridge: MIT Press. p. 247. ISBN 978-0-262-01802-9.
31. Amemiya, Takeshi (1985). *Advanced Econometrics* (<https://archive.org/details/advancedeconomet00amem/page/137>). Cambridge: Harvard University Press. pp. 137–138 (<https://archive.org/details/advancedeconomet00amem/page/137>). ISBN 0-674-00560-0.
32. Sargan, Denis (1988). "Methods of Numerical Optimization". *Lecture Notes on Advanced Econometric Theory*. Oxford: Basil Blackwell. pp. 161–169. ISBN 0-631-14956-2.



33. See theorem 10.1 in Avriel, Mordecai (1976). *Nonlinear Programming: Analysis and Methods* (<https://books.google.com/books?id=byF4Xb1QbvMC&pg=PA293>). Englewood Cliffs, NJ: Prentice-Hall. pp. 293–294. ISBN 978-0-486-43227-4.
34. Gill, Philip E.; Murray, Walter; Wright, Margaret H. (1981). *Practical Optimization* (<https://archive.org/details/practicaloptimiz00gill>). London, UK: Academic Press. pp. 312 (<https://archive.org/details/practicaloptimiz00gill/page/n329>)–313. ISBN 0-12-283950-1.
35. Edgeworth, Francis Y. (Sep 1908). "On the probable errors of frequency-constants" (<https://zenodo.org/record/1449468>). *Journal of the Royal Statistical Society*. **71** (3): 499–512. doi:10.2307/2339293 (<https://doi.org/10.2307%2F2339293>). JSTOR 2339293 (<https://www.jstor.org/stable/2339293>).
36. Edgeworth, Francis Y. (Dec 1908). "On the probable errors of frequency-constants" (<https://zenodo.org/record/1449468>). *Journal of the Royal Statistical Society*. **71** (4): 651–678. doi:10.2307/2339378 (<https://doi.org/10.2307%2F2339378>). JSTOR 2339378 (<https://www.jstor.org/stable/2339378>).
37. Pfanzagl, Johann (1994). *Parametric Statistical Theory*. Walter de Gruyter. pp. 207–208. doi:10.1515/9783110889765 (<https://doi.org/10.1515%2F9783110889765>). ISBN 978-3-11-013863-4. MR 1291393 (<https://mathscinet.ams.org/mathscinet-getitem?mr=1291393>).
38. Hald, Anders (1999). "On the History of Maximum Likelihood in Relation to Inverse Probability and Least Squares" (<https://www.jstor.org/stable/2676741>). *Statistical Science*. **14** (2): 214–222. ISSN 0883-4237 (<http://www.worldcat.org/issn/0883-4237>).
39. Wilks, S.S. (1938). "The large-sample distribution of the likelihood ratio for testing composite hypotheses" (<https://doi.org/10.1214%2Faoms%2F1177732360>). *Annals of Mathematical Statistics*. **9**: 60–62. doi:10.1214/aoms/1177732360 (<https://doi.org/10.1214%2Faoms%2F1177732360>).
40. Owen, Art B. (2001). *Empirical Likelihood*. London, UK; Boca Raton, FL: Chapman & Hall; CRC Press. ISBN 978-1-58488-071-4.
41. Wilks, Samuel S. (1962). *Mathematical Statistics*. New York, NY: John Wiley & Sons. ISBN 978-0-471-94650-2.
42. Savage, Leonard J. (1976). "On rereading R.A. Fisher" (<https://doi.org/10.1214%2Faos%2F1176343456>). *The Annals of Statistics*. **4** (3): 441–500. doi:10.1214/aos/1176343456 (<https://doi.org/10.1214%2Faos%2F1176343456>). JSTOR 2958221 (<https://www.jstor.org/stable/2958221>).
43. Pratt, John W. (1976). "F. Y. Edgeworth and R. A. Fisher on the efficiency of maximum likelihood estimation" (<https://doi.org/10.1214%2Faos%2F1176343457>). *The Annals of Statistics*. **4** (3): 501–514. doi:10.1214/aos/1176343457 (<https://doi.org/10.1214%2Faos%2F1176343457>). JSTOR 2958222 (<https://www.jstor.org/stable/2958222>).
44. Stigler, Stephen M. (1978). "Francis Ysidro Edgeworth, statistician". *Journal of the Royal Statistical Society, Series A*. **141** (3): 287–322. doi:10.2307/2344804 (<https://doi.org/10.2307%2F2344804>). JSTOR 2344804 (<https://www.jstor.org/stable/2344804>).
45. Stigler, Stephen M. (1986). *The history of statistics: the measurement of uncertainty before 1900* (<https://archive.org/details/historyofstatist00stig>). Harvard University Press. ISBN 978-0-674-40340-6.
46. Stigler, Stephen M. (1999). *Statistics on the table: the history of statistical concepts and methods*. Harvard University Press. ISBN 978-0-674-83601-3.
47. Hald, Anders (1998). *A history of mathematical statistics from 1750 to 1930*. New York, NY: Wiley. ISBN 978-0-471-17912-2.
48. Hald, Anders (1999). "On the history of maximum likelihood in relation to inverse probability and least squares" ([http://projecteuclid.org/download/pdf\\_1/euclid.ss/1009212248](http://projecteuclid.org/download/pdf_1/euclid.ss/1009212248)). *Statistical Science*. **14** (2): 214–222. doi:10.1214/ss/1009212248 (<https://doi.org/10.1214%2Fss%2F1009212248>). JSTOR 2676741 (<https://www.jstor.org/stable/2676741>).
49. Aldrich, John (1997). "R.A. Fisher and the making of maximum likelihood 1912–1922" (<https://doi.org/10.1214%2Fss%2F1030037906>). *Statistical Science*. **12** (3): 162–176. doi:10.1214/ss/1030037906 (<https://doi.org/10.1214%2Fss%2F1030037906>). MR 1617519 (<https://mathscinet.ams.org/mathscinet-getitem?mr=1617519>).

## Further reading

- Cramer, J.S. (1986). *Econometric Applications of Maximum Likelihood Methods* (<https://archive.org/details/econometricappli0000cram>). New York, NY: Cambridge University Press. ISBN 0-521-25317-9.

- Eliason, Scott R. (1993). *Maximum Likelihood Estimation: Logic and Practice*. Newbury Park: Sage. ISBN 0-8039-4107-2.
- King, Gary (1989). *Unifying Political Methodology: the Likelihood Theory of Statistical Inference*. Cambridge University Press. ISBN 0-521-36697-6.
- Le Cam, Lucien (1990). "Maximum likelihood: An Introduction". *ISI Review*. **58** (2): 153–171. doi:10.2307/1403464 (<https://doi.org/10.2307%2F1403464>). JSTOR 1403464 (<https://www.jstor.org/stable/1403464>).
- Magnus, Jan R. (2017). "Maximum Likelihood". *Introduction to the Theory of Econometrics*. Amsterdam, NL: VU University Press. pp. 53–68. ISBN 978-90-8659-766-6.
- Millar, Russell B. (2011). *Maximum Likelihood Estimation and Inference*. Hoboken, NJ: Wiley. ISBN 978-0-470-09482-2.
- Pickles, Andrew (1986). *An Introduction to Likelihood Analysis* (<https://archive.org/details/introductiontoli0000pick>). Norwich: W. H. Hutchins & Sons. ISBN 0-86094-190-6.
- Severini, Thomas A. (2000). *Likelihood Methods in Statistics*. New York, NY: Oxford University Press. ISBN 0-19-850650-3.
- Ward, Michael D.; Ahlquist, John S. (2018). *Maximum Likelihood for Social Science: Strategies for Analysis*. Cambridge University Press. ISBN 978-1-316-63682-4.

## External links

---

- Tilevik, Andreas (2022). Maximum likelihood vs least squares in linear regression (<https://www.youtube.com/watch?v=bhTlpGtWtzQ>) (video)
- "Maximum-likelihood method" ([https://www.encyclopediaofmath.org/index.php?title=Maximum-likelihood\\_method](https://www.encyclopediaofmath.org/index.php?title=Maximum-likelihood_method)), *Encyclopedia of Mathematics*, EMS Press, 2001 [1994]
- Purcell, S. "Maximum Likelihood Estimation" ([http://statgen.iop.kcl.ac.uk/bgim/mle/sslike\\_1.html](http://statgen.iop.kcl.ac.uk/bgim/mle/sslike_1.html)).
- Sargent, Thomas; Stachurski, John. "Maximum Likelihood Estimation" (<https://intro.quantecon.org/mle.html>). *Quantitative Economics with Python*.
- Toomet, Ott; Henningsen, Arne (2019-05-19). "maxLik: A package for maximum likelihood estimation in R" (<https://cran.r-project.org/package=maxLik>).
- Lesser, Lawrence M. (2007). "'MLE' song lyrics" (<http://www.math.utep.edu/Faculty/lesser/MLE.html>). Mathematical Sciences / College of Science. *University of Texas*. El Paso, TX. Retrieved 2021-03-06.

---

Retrieved from "[https://en.wikipedia.org/w/index.php?title=Maximum\\_likelihood\\_estimation&oldid=1233238324](https://en.wikipedia.org/w/index.php?title=Maximum_likelihood_estimation&oldid=1233238324)"