

1. Introduction

Welcome to the series of E-learning modules on Method of Maximum Likelihood and their properties. In this module, we are going to cover the basic principle of method of maximum likelihood, advantages and properties of maximum likelihood estimators and estimation of certain population parameters by the method of maximum likelihood.

By the end of this session, you will be able to:

- Explain the basic concept of method of maximum likelihood
- Explain the principle of the method and properties of likelihood estimators
- Explain the advantages and disadvantages of the method of maximum likelihood
- Explain the procedure to estimate certain parameters of the population by the method of maximum likelihood

Selection of a random sample is very much essential for drawing valid inferences about the population. Methods that are developed for estimating the population parameters provide a theoretical basis of connection between sample information and population parameters, which assists in drawing efficient estimates from the sample, and draw inference about the population.

In statistics, Maximum-Likelihood Estimation (MLE) is a method of estimating the parameters of a statistical model, which is generally used in estimation. When applied to a data set and given a statistical model, maximum-likelihood estimation provides estimates for the model's parameters.

The method of maximum likelihood corresponds to many well-known estimation methods in statistics. For example, one may be interested in the intelligence of school going children. However, will be unable to measure the intelligence of every single child in a population due to cost or time constraints.

Assuming that the intelligence is normally distributed with some unknown mean and variance, the mean and variance can be estimated with MLE while only knowing the intelligence of some sample of the overall population.

MLE would accomplish this by taking the mean and variance as parameters and finding particular parametric values that make the observed results the most probable.

In general, for a fixed set of data and underlying statistical model, the method of maximum likelihood selects values of the model parameters that produce a distribution that gives the observed data the greatest probability.

2. Applications of the Method of Maximum Likelihood

Applications of the Method of Maximum Likelihood:

Maximum likelihood estimation is used for a wide range of statistical models, including:

Linear Models and generalized linear models

Exploratory and confirmatory factor analysis

Structural Equation Modelling

Many situations in the context of hypothesis testing and confidence interval formation

These uses arise across applications in widespread set of fields, including:

- Communication systems
- Psychometrics
- Econometrics
- Data modelling in nuclear and particle physics
- Magnetic resonance imaging
- Computational Phylogenetics
- Geographical satellite-image classification

On the other hand, MLE is not as widely recognized among modellers in psychology, but it is a standard approach to parameter estimation and inference in statistics. MLE has many optimal properties in estimation, such as sufficiency, consistency, efficiency etc.

Further, many of the inference methods in statistics are developed based on MLE. For example, MLE is a prerequisite for the chi-square test, the Gsquare test, Bayesian methods, inference with missing data, modelling of random effects, and many model selection criteria.

Once data has been collected and the likelihood function of a model is determined, we can make statistical inferences about the population, that is, the probability distribution that underlies the data. Given the different parameter value index and probability distributions, we are interested in finding the parameter value that corresponds to the desired probability distribution.

The principle of maximum likelihood estimation originally developed by R. A. Fisher in the nineteen twenty's. It states that the desired probability distribution is the one that makes the observed data "most likely," which means that one must seek the value of the parameter vector that maximizes the likelihood function. The resulting parameter vector, which is sought by searching the multi-dimensional parameter space, is called the MLE.

3. Likelihood Function

Likelihood Function:

A likelihood function of a number of sample observations is defined to be their joint density function. To use the method of maximum likelihood, we should first specify the joint density function for all observations. For an Independent and Identically Distributed (IID) sample, this joint density function is

$f(x_1, x_2, \dots, x_n, \theta)$ is equal to $f(x_1, \theta) \times f(x_2, \theta) \times \dots \times f(x_n, \theta)$

Now, we look at this function from a different perspective by considering the observed values x_1, x_2, \dots, x_n to be fixed "parameters" of this function, whereas θ will be the function's variable and allowed to vary freely. This function will be called the likelihood.

In case, the sample observations are independent, the likelihood function happens to be the product of the density functions of the random observations. If x_1, x_2, \dots, x_n are n independent and identically distributed observations, from a population with an unknown parameter θ then the likelihood function of the random observations denoted by

L is equal to $L(x_1, x_2, \dots, x_n, \theta)$ **is given by**

L is equal to $L(x_1, x_2, \dots, x_n, \theta)$ **is equal to product of $f(x_i, \theta)$**

Where, $f(x_i, \theta)$ is the p.d.f of the population.

L gives relative likelihood or chance that the random observations assume a particular set of values. For a particular sample (x_1, x_2, \dots, x_n) , L becomes a function of (x_1, x_2, \dots, x_n) and an unknown parameter θ and p.d.f. $f(x, \theta)$. It is desirable to find an estimator $\hat{\theta}$, which would be closest to the true value θ . Both the observed variables x_i and the parameter θ can be vectors.

The method of maximum likelihood is the one in which for a given set of values x_1, x_2, \dots, x_n , an estimator of θ is found that maximizes L . Thus, if there exists $\hat{\theta}$, a function of x_1, x_2, \dots, x_n for which L is maximum for variations in θ . Then, $\hat{\theta}$ is the M.L.E of θ if

$\frac{\partial L}{\partial \theta} = 0$ and $\frac{\partial^2 L}{\partial \theta^2} < 0$

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

Since $\log L$ attains maximum, when L attains maximum the solution of the equation.

$\log L(x_1, x_2, \dots, x_n, \theta)$ is equal to $\sum \log f(x_i, \theta)$

$\frac{\partial \log L}{\partial \theta} = 0$ and $\frac{\partial^2 \log L}{\partial \theta^2} < 0$ also gives the M.L.E of θ

An MLE estimate is the same regardless of whether we maximize the likelihood or the log-likelihood function, since \log is a monotone transformation.

4. Principle, Properties and Advantages of Maximum Likelihood Function

Principle:

The likelihood function L of $(x_1, x_2, \dots, x_n, \theta)$ in a way gives the probability of the random sample x_1, x_2, \dots, x_n , when the parameter of the distribution is θ . For different values of θ , L gives the different probabilities. Since the event of getting a sample x_1, x_2, \dots, x_n , has occurred, the probability of the event must be high or else the value of L should be large.

So the value of θ , which gives the maximum probability for the sample is found and it is taken as the maximum likelihood estimate of θ . In fact, the function of the sample observations, which maximizes L , is determined and that is the MLE of θ .

For many models, a maximum likelihood estimator can be found as an explicit function of the observed data x_1, x_2, \dots, x_n . For many other models, no closed-form solution to the maximization problem is known or available and an MLE has to be found numerically using optimization methods.

For some problems, there may be multiple estimates that maximize the likelihood. For other problems, no maximum likelihood estimate exists (meaning that the log-likelihood function increases without attaining the supremum value).

In the exposition above, it is assumed that the data are independent and identically distributed. The method can be applied however to a broader setting, as long as it is possible to write the joint density function f of $(x_1, x_2, \dots, x_n, \theta)$.

In a simpler extension, an allowance can be made for data heterogeneity, so that the joint density is equal to f_1 of (x_1, θ) into f_2 of (x_2, θ) up to f_n of (x_n, θ) . In the more complicated case of time series models, the independence assumption may have to be dropped as well.

Properties

A maximum-likelihood estimator is an extremum estimator obtained by maximizing, as a function of θ .

Maximum-likelihood estimators have no optimum properties for finite samples.

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: A subsequence of the sequence of MLEs converges in probability to the value being estimated.

Asymptotic normality: As the sample size increases, the distribution of the MLE tends to the Gaussian distribution.

Efficiency: MLEs are most efficient among the class of all consistent estimators. That is among the class of consistent estimators MLE has the minimum variance.

Sufficiency: Complete information about the parameter of interest is contained in its MLE.

Maximum-likelihood estimators can lack asymptotic normality and can be inconsistent if there is a failure of one (or more) of the below regularity conditions:

Estimate on boundary

Sometimes the maximum likelihood estimate lies on the boundary of the set of possible parameters. Standard asymptotic theory needs the assumption that the true parameter value lies away from the boundary.

Data boundary parameter-dependent

For the theory to apply in a simple way, the set of data values, which has positive probability (or positive probability density), should not depend on the unknown parameter.

A simple example where such parameter-dependence does hold is the case of estimating θ from a set of IID variables is when the common distribution is uniform on the range $(0, \theta)$. For estimation purposes, the relevant range of θ is such that θ cannot be less than the largest observation. Because the interval $(0, \theta)$ is not compact, there exists no maximum for the likelihood function.

Nuisance parameters

For maximum likelihood estimations, a model may have a number of [nuisance parameters](#).

Increasing information

For the asymptotic to hold in cases where the assumption of [IID](#) observations does not hold, a basic requirement is that the amount of information in the data increases indefinitely as the sample size increases. Either which may not be met because of too much dependence or new independent observations are subject to an increasing observation error.

Advantages:

- i. The method has a good intuitive foundation. The underlying concept is that the best estimate of a parameter is giving the highest probability that the observed set of measurements will be obtained.
- ii. The least-squares method and various approaches for combining errors or calculating weighted averages, etc. can be derived or justified in terms of the maximum likelihood approach.
- iii. The method is of sufficient generality that most problems are amenable to a straightforward application of this method, even in cases where other techniques become difficult. Inelegant but conceptually simple approaches often provide useful results where there is no easy alternative.

5. Illustrations

Illustrations:

For a random sampling from a normal population $N(\theta, \sigma^2)$, find the MLE for

- i. θ when σ^2 is known
- ii. σ^2 when θ is known
- iii. θ and σ^2 when both are unknown

Let x_1, x_2, \dots, x_n be a random sample of size n drawn from a given population

L is equal to $L(x_1, x_2, \dots, x_n, \theta)$ is equal to product of $f(x_i, \theta)$

Which is equal to product of $\frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma^2}(x_i - \theta)^2\right)$

Which is equal to $\left(\frac{1}{\sigma \sqrt{2\pi}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2\right)$

$\ln L$ is equal to $-\frac{n}{2} \ln \sigma^2 - \frac{n}{2} \ln 2\pi - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2$

i) Maximum likelihood estimate of θ when σ^2 is known

$\frac{\partial \ln L}{\partial \theta}$ is equal to zero implies $\sum_{i=1}^n (x_i - \theta) = 0$

which implies $\sum_{i=1}^n x_i - n\theta = 0$, which implies $\hat{\theta} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$

Therefore, MLE of θ is $\hat{\theta} = \bar{x}$ which is equal to $\sum_{i=1}^n x_i / n$ which is equal to \bar{x} with $\frac{\partial^2 \ln L}{\partial \theta^2} < 0$

ii) Maximum likelihood estimate of σ^2 when θ is known

$\frac{\partial \ln L}{\partial \sigma^2}$ is equal to zero implies $-\frac{n}{2\sigma^4} + \frac{1}{2\sigma^6} \sum_{i=1}^n (x_i - \theta)^2 = 0$ implies $-\frac{n}{2} + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2 = 0$

which implies $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \theta)^2$

Therefore, MLE of σ^2 when θ is known is $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \theta)^2$

iii) Maximum likelihood estimate of θ and σ^2 when both are unknown

$\frac{\partial \ln L}{\partial \theta}$ is equal to zero implies $\sum_{i=1}^n (x_i - \theta) = 0$ implies $\hat{\theta} = \bar{x}$. Call this as equation 1

$\frac{\partial \ln L}{\partial \sigma^2}$ is equal to zero implies $-\frac{n}{2\sigma^4} + \frac{1}{2\sigma^6} \sum_{i=1}^n (x_i - \theta)^2 = 0$

implies $-\frac{n}{2} + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \hat{\theta})^2 = 0$ implies $-\frac{n}{2} + \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \bar{x})^2 = 0$ from equation (1)

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

Therefore, MLE of θ is $\hat{\theta} = \bar{x}$ and the MLE of σ^2 when θ is unknown is $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

Where we compare the above example with the method of moments, we see that the method

of moment estimators and the maximum likelihood estimators for θ and σ^2 are the same. However, it happens rarely, if they are not the same.

Due to the fact that the maximum likelihood estimator of \bar{x} has an approximate normal distribution with mean θ and a variance σ^2/n that is equal to a certain lower bound, thus at least approximately, it is unbiased minimum variance estimator. Accordingly, most statisticians prefer the maximum likelihood estimators than estimators found using the method of moments.

2) Find the MLE of θ in the following uniform distribution $U(0, \theta)$

Let x_1, x_2, \dots, x_n be a random sample of size n drawn from a uniform population with probability density function

$f(x, \theta) = \frac{1}{\theta}$ for $0 < x < \theta$

$L = L(x_1, x_2, \dots, x_n, \theta) = \prod_{i=1}^n f(x_i, \theta)$

$L = \left(\frac{1}{\theta}\right)^n$

$\log L = -n \log \theta$

$\frac{\partial \log L}{\partial \theta} = 0$ implies $-\frac{n}{\theta} = 0$ which implies $\hat{\theta} = \infty$ or $n = 0$, which is meaningless.

Hence, we use the basic principle of the maximum likelihood estimation, where the value of θ for which L is maximum

L is maximum when $\frac{1}{\theta^n}$ is maximum. That is when θ^n is minimum, when θ is minimum

But $0 < x < \theta$

Implies, $0 < (x_1) < (x_2) < \dots < (x_n) < \theta$

Where, $(x_1), (x_2), \dots, (x_n)$ are order statistics.

Minimum value of θ is the maximum of the observations

Therefore, MLE of θ is $\hat{\theta} = x_n$, the maximum of the observations

Here's a summary of our learning in this session, where we understood:

- The principle of method of maximum likelihood estimation
- The applications and advantages of the method
- The properties of the maximum likelihood estimators
- The illustrative examples to obtain MLE