

MAXIMUM LIKELIHOOD ESTIMATION THEORY  
SUMMARY OF NOTATION AND MAIN RESULTS

**Definition**

Consider a random sample  $\mathbf{X} = (X_1, \dots, X_n)$  of independent identically distributed random variables (or vectors), each with density  $f(x|\theta_X)$  with respect to a  $\sigma$ -finite measure  $\mu$ . We assume that  $f(x|\theta_X) \in \mathcal{F}$ , where

$$\mathcal{F} = \{\text{distributions with density } f(x|\theta), \theta \in \Theta \subseteq \mathbb{R}^d\}$$

represents a parametric model for the distribution of the data.

The model  $\mathcal{F}$  must satisfy the model identifiability condition: For any  $\theta_1 \neq \theta_2$  it holds  $f(x|\theta_1) \neq f(x|\theta_2)$ . In other words, no distribution can be parametrized by several different parameter vectors. Because of independence, the joint density of the random sample  $X_1, \dots, X_n$  is  $\prod_{i=1}^n f(x_i|\theta_X)$ . The maximum likelihood estimator  $\hat{\theta}$  of the parameter  $\theta_X$  is the point from  $\Theta$  that maximizes the joint density evaluated at the observed values of  $X_1, \dots, X_n$ .

**Definition 1** (likelihood, log-likelihood).

- The random function

$$L_n(\theta) \stackrel{\text{df}}{=} \prod_{i=1}^n f(X_i|\theta)$$

is called *the likelihood function* for the parameter  $\theta$  in the model  $\mathcal{F}$ .

- The random function

$$\ell_n(\theta) \stackrel{\text{df}}{=} \log L_n(\theta) = \sum_{i=1}^n \log f(X_i|\theta)$$

is called *the log-likelihood function*.

**Definition 2** (maximum likelihood estimator). *The maximum likelihood estimator* (MLE) of the parameter  $\theta_X$  in the model  $\mathcal{F}$  is defined as

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} L_n(\theta).$$

**Note.** Since the logarithm is strictly increasing,  $L_n(\theta)$  and  $\ell_n(\theta)$  attain the maximum at the same point.

**Definition 3.** Let  $P$  and  $Q$  be probability measures on the same probability space with densities  $p$  and  $q$  with respect to the same  $\sigma$ -finite measure  $\mu$  (for example,  $\mu = P + Q$ ). Define

$$K(P, Q) = \begin{cases} E_P \log \frac{p(X)}{q(X)} = \int_{\{x:p(x)>0\}} \log \frac{p(x)}{q(x)} p(x) d\mu(x) & \text{if } P[q(X) = 0] = 0 \\ +\infty & \text{otherwise.} \end{cases}$$

$K(P, Q)$  is called the *Kullback-Leibler distance (divergence)*.

**Note.** In fact,  $K(P, Q)$  is a pseudo-distance: it holds  $K(P, Q) \geq 0$ , and  $K(P, Q) = 0$  if and only if  $P = Q$ , but it is not symmetric:  $K(P, Q) \neq K(Q, P)$ .

**Theorem 1.** Suppose the support set  $S = \{x \in \mathbb{R} : f(x|\boldsymbol{\theta}) > 0\}$  does not depend on the parameter  $\boldsymbol{\theta}$ . Denote  $P_X$  the induced probability measure of the random variable  $X_i$  and  $P_\theta$  the probability measure associated with the density  $f(x|\boldsymbol{\theta})$ . Then for any  $\boldsymbol{\theta} \neq \boldsymbol{\theta}_X$

$$\frac{1}{n} \log \frac{L_n(\boldsymbol{\theta}_X)}{L_n(\boldsymbol{\theta})} = \frac{1}{n} \sum_{i=1}^n \log \frac{f(X_i|\boldsymbol{\theta}_X)}{f(X_i|\boldsymbol{\theta})} \rightarrow K(P_X, P_\theta) \quad P_X - \text{almost surely,}$$

and hence

$$P[\ell_n(\boldsymbol{\theta}_X) > \ell_n(\boldsymbol{\theta})] \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

**Note.** When the number of observations increases to infinity, the (log-)likelihood function at the true parameter will be with a large probability larger than the (log-)likelihood function at any other parameter. This observation justifies the idea of estimating the parameters by maximizing the log-likelihood over all possible parameter vectors.

## The calculation of the maximum likelihood estimator

The maximum likelihood estimator is usually determined by differentiation of the log-likelihood. The first derivative is set to zero and it is verified that the second derivative is negative definite.

**Definition 4** (score, information).

- The random vector

$$\mathbf{U}(\boldsymbol{\theta}|X_i) \stackrel{\text{df}}{=} \frac{\partial}{\partial \boldsymbol{\theta}} \log f(X_i|\boldsymbol{\theta})$$

is called *the score function* for the parameter  $\boldsymbol{\theta}$  in the model  $\mathcal{F}$ .

- The random vector

$$\mathbf{u}_n(\boldsymbol{\theta}|\mathbf{X}) \stackrel{\text{df}}{=} \sum_{i=1}^n \mathbf{U}(\boldsymbol{\theta}|X_i) = \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} \log f(X_i|\boldsymbol{\theta})$$

is called *the score statistic*.

- The random matrix

$$I(\boldsymbol{\theta}|X_i) \stackrel{\text{df}}{=} -\frac{\partial}{\partial \boldsymbol{\theta}^\top} \mathbf{U}(\boldsymbol{\theta}|X_i) = -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \log f(X_i|\boldsymbol{\theta})$$

is called the contribution of the  $i$ -th observation to the information matrix.

- The random matrix

$$I_n(\boldsymbol{\theta}|\mathbf{X}) \stackrel{\text{df}}{=} -\frac{1}{n} \frac{\partial}{\partial \boldsymbol{\theta}^\top} \mathbf{u}_n(\boldsymbol{\theta}|\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n I(\boldsymbol{\theta}|X_i)$$

is called *the observed information matrix*.

- The matrix

$$I(\boldsymbol{\theta}) \stackrel{\text{df}}{=} E I(\boldsymbol{\theta}|X_i) = -E \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} \log f(X_i|\boldsymbol{\theta})$$

is called *the expected (Fisher) information matrix*.

If the set  $\Theta$  is open, the MLE  $\hat{\theta}_n$  solves the system of equations  $\mathbf{U}_n(\hat{\theta}_n|\mathbf{X}) = \mathbf{0}$ , that is

$$\sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(X_i|\hat{\theta}_n) = \mathbf{0}.$$

This system is called the *likelihood equations*.

The solution to the likelihood equations need not exist. Sometimes there may be multiple solutions, at most one of which is the MLE. If  $I_n(\hat{\theta}_n|\mathbf{X}) > 0$  (the observed information is positive definite at  $\hat{\theta}_n$ ), we know that  $\hat{\theta}_n$  is at least a local maximum. If  $I_n(\theta|\mathbf{X}) > 0$  for every  $\theta \in \Theta$ , the log-likelihood function is concave and the solution to the likelihood equations must be the global maximum and hence the MLE.

In most cases no explicit solution can be found and the MLE must be calculated by numerical methods. There are two commonly used numerical methods for solving the likelihood equations. Let  $\hat{\theta}^{(r)}$  be the  $r$ -th iteration to the solution.

- **The Newton-Raphson method:**  $\hat{\theta}^{(r+1)} = \hat{\theta}^{(r)} + [nI_n(\hat{\theta}^{(r)}|\mathbf{X})]^{-1}\mathbf{U}_n(\hat{\theta}^{(r)}|\mathbf{X})$
- **The Fisher Scoring method:**  $\hat{\theta}^{(r+1)} = \hat{\theta}^{(r)} + [nI(\hat{\theta}^{(r)})]^{-1}\mathbf{U}_n(\hat{\theta}^{(r)}|\mathbf{X})$

They are iterated until the change in  $\hat{\theta}$  from one iteration to the next is sufficiently small or until  $\mathbf{U}_n(\hat{\theta})$  is sufficiently close to  $\mathbf{0}$ . The only difference between the two methods is in the information matrix: N-R uses the observed information, FS uses the expected information.

Both require setting  $\hat{\theta}^{(1)}$ , the starting value for numerical approximation, and are sensitive to its choice.

## Properties of the maximum likelihood estimator

Maximum likelihood estimators are consistent and asymptotically normal as long as so called *regularity conditions* are satisfied.

**Conditions** (Regularity conditions for maximum likelihood estimators).

- R1.** The number of parameters  $d$  in the model  $\mathcal{F}$  is constant.
- R2.** The support set  $S = \{x \in \mathbb{R} : f(x|\theta) > 0\}$  does not depend on the parameter  $\theta$ .
- R3.** The parameter space  $\Theta$  is an open set.
- R4.** The density  $f(x|\theta)$  is sufficiently smooth function of  $\theta$  (at least twice continuously differentiable).
- R5.** The Fisher information matrix  $I(\theta)$  is finite, regular, and positive definite in a neighborhood of  $\theta_X$ .
- R6.** The order of differentiation and integration can be interchanged in expressions such as

$$\frac{\partial}{\partial \theta} \int h(x, \theta) d\mu(x) = \int \frac{\partial}{\partial \theta} h(x, \theta) d\mu(x),$$

where  $h(x, \theta)$  is either  $f(x|\theta)$  or  $\partial f(x|\theta)/\partial \theta$ .

**Note.** Take the identity

$$\int_{-\infty}^{\infty} f(x|\boldsymbol{\theta}) d\mu(x) = 1$$

and differentiate both sides of the equation twice with respect to  $\boldsymbol{\theta}$ . Regularity condition R6 implies

$$\int_{-\infty}^{\infty} \frac{\partial}{\partial \boldsymbol{\theta}} f(x|\boldsymbol{\theta}) d\mu(x) = \int_{-\infty}^{\infty} \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} f(x|\boldsymbol{\theta}) d\mu(x) = \mathbf{0}. \quad (1)$$

**Theorem 2** (consistency of the MLE). Let conditions R1–R6 hold. Then there exists  $n_0$  and a sequence  $\hat{\boldsymbol{\theta}}_n$  ( $n \geq n_0$ ) of solutions to the likelihood equations  $\mathbf{U}_n(\hat{\boldsymbol{\theta}}_n|\mathbf{X}) = \mathbf{0}$  such that  $\hat{\boldsymbol{\theta}}_n \xrightarrow{P} \boldsymbol{\theta}_X$ .

**Note.** If the log-likelihood is strictly concave, the likelihood equations have a unique solution, which is the MLE. It converges in probability to the true parameter. If the log-likelihood is not strictly concave, the likelihood equations may have multiple solutions representing local maxima and minima of the log-likelihood. There is one solution among them (the closest to  $\boldsymbol{\theta}_X$ ), which provides a consistence sequence of estimators. Other solutions may not be close to  $\boldsymbol{\theta}_X$  and may not converge to it.

**Note.** If there exists a sequence  $\tilde{\boldsymbol{\theta}}_n$  of other estimators that are guaranteed to be consistent (for example, moment estimators of  $\boldsymbol{\theta}_X$ ), a consistent MLE can be obtained by taking the root of the likelihood equations, which is closest to  $\tilde{\boldsymbol{\theta}}_n$ . Alternatively, one can perform one step of the Newton-Raphson algorithm with  $\tilde{\boldsymbol{\theta}}_n$  as the starting value.

**Theorem 3** (Score function properties). Let conditions R1–R6 hold. Then

(i)  $E \mathbf{U}(\boldsymbol{\theta}_X|X_i) = 0$ ,  $\text{var} \mathbf{U}(\boldsymbol{\theta}_X|X_i) = I(\boldsymbol{\theta}_X)$ .

(ii)  $\frac{1}{\sqrt{n}} \mathbf{U}_n(\boldsymbol{\theta}_X|\mathbf{X}) \xrightarrow{D} N_d(\mathbf{0}, I(\boldsymbol{\theta}_X))$ .

**Note.** The Fisher information matrix at  $\boldsymbol{\theta}_X$  can be calculated in two different ways: from Definition 4 (the expectation of minus the second derivative of the log density) or from Theorem 3 (the score function variance).

**Theorem 4** (asymptotic normality of the MLE). Suppose conditions R1–R6 hold. Let  $\hat{\boldsymbol{\theta}}_n$  be a consistent sequence of solutions to the likelihood equations. Then

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_X) \xrightarrow{D} N_d(\mathbf{0}, I^{-1}(\boldsymbol{\theta}_X)).$$

**Note.**

- The asymptotic variance of the MLE is equal to the inverse of the Fisher information. More information means better precision for estimation.
- The asymptotic variance of the MLE is in a certain sense optimal. Other estimators (e.g., moment estimators) cannot have a smaller asymptotic variance.

**Theorem 5** (asymptotic distribution of the likelihood ratio). Suppose conditions R1–R6 hold. Let  $\hat{\boldsymbol{\theta}}_n$  be a consistent sequence of solutions to the likelihood equations. Then

$$2 \log \frac{L_n(\hat{\boldsymbol{\theta}}_n)}{L_n(\boldsymbol{\theta}_X)} = 2(\ell_n(\hat{\boldsymbol{\theta}}_n) - \ell_n(\boldsymbol{\theta}_X)) \xrightarrow{D} \chi_d^2.$$

**Theorem 6** (the  $\Delta$  method for the MLE). Suppose conditions R1–R6 hold. Let  $\hat{\boldsymbol{\theta}}_n$  be a consistent sequence of solutions to the likelihood equations. Take  $q : \Theta \rightarrow \mathbb{R}^k$  a continuously differentiable function. Denote  $\boldsymbol{\nu}_X = q(\boldsymbol{\theta}_X)$  a  $D(\boldsymbol{\theta}) = \partial q(\boldsymbol{\theta}) / \partial \boldsymbol{\theta}$ . Then  $\hat{\boldsymbol{\nu}}_n = q(\hat{\boldsymbol{\theta}}_n)$  is the MLE of the parameter  $\boldsymbol{\nu}_X$  and

$$\sqrt{n}(\hat{\boldsymbol{\nu}}_n - \boldsymbol{\nu}_X) \xrightarrow{D} N_k(\mathbf{0}, D(\boldsymbol{\theta}_X)I^{-1}(\boldsymbol{\theta}_X)D(\boldsymbol{\theta}_X)^T).$$

## Tests based on maximum likelihood theory

The theory of the MLE can be used to derive tests of simple and composite hypotheses about the parameter  $\theta_X$ .

### Testing of simple hypotheses

We want to test the null hypothesis  $H_0 : \theta_X = \theta_0$  against the alternative  $H_1 : \theta_X \neq \theta_0$ , where  $\theta_0 \in \Theta$ . It is a simple hypothesis because there is just a single distribution in the model  $\mathcal{F}$  with the density  $f(x|\theta_0)$ .

We will introduce three different test statistics for testing  $H_0$ .

#### Definition 5.

(i) The statistic

$$\lambda_n = \frac{L_n(\hat{\theta}_n)}{L_n(\theta_0)}$$

is called *the likelihood ratio*.

(ii) The statistic

$$W_n = n(\hat{\theta}_n - \theta_0)^\top \hat{I}_n(\hat{\theta}_n)(\hat{\theta}_n - \theta_0)$$

is called *the Wald statistic*.

(iii) The statistic

$$R_n = \frac{1}{n} \mathbf{U}_n(\theta_0|\mathbf{X})^\top \hat{I}_n^{-1}(\theta_0) \mathbf{U}_n(\theta_0|\mathbf{X})$$

is called *the Rao (score) statistic*.

**Note.** The symbol  $\hat{I}_n$  denotes any consistent estimator of the Fisher information matrix. Three different estimators can be used in Wald and Rao statistics:

1.  $\hat{I}_n(\theta) = I_n(\theta|\mathbf{X}) = -\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta \partial \theta^\top} \log f(\theta|X_i)$  (the observed information matrix)
2.  $\hat{I}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \mathbf{U}(\theta|X_i)^{\otimes 2}$  (the empirical variance of the score function)
3.  $\hat{I}_n(\theta) = I(\theta)$  (the Fisher information matrix)

The most common choice for the Wald statistic is  $\hat{I}_n(\hat{\theta}_n) = I_n(\hat{\theta}_n|\mathbf{X})$ . The most common choice for the Rao statistic is  $\hat{I}_n(\theta_0) = \frac{1}{n} \sum_{i=1}^n \mathbf{U}(\theta_0|X_i)^{\otimes 2}$ .

#### Note.

- The likelihood ratio requires the calculation of  $\hat{\theta}_n$  and  $L_n$  or  $\ell_n$ . It does not require the calculation of  $\mathbf{U}_n$  and  $\hat{I}_n$ .
- The Wald statistic requires the calculation of  $\hat{\theta}_n$  and  $\hat{I}_n$ . It does not require the calculation of  $L_n$  and  $\mathbf{U}_n$ .
- Rao statistic requires the calculation of  $\mathbf{U}_n$  and  $\hat{I}_n$ . It does not require the calculation of  $\hat{\theta}_n$  and  $L_n$ .

**Note.** If  $d = 1$  (one parameter) and  $\theta_0 = 0$ , then the Wald statistic can be written as

$$W_n = \left[ \frac{\hat{\theta}_n}{\sqrt{n^{-1}\hat{I}_n^{-1}(\hat{\theta}_n)}} \right]^2,$$

where  $n^{-1}\hat{I}_n^{-1}(\hat{\theta}_n)$  is the estimator of the asymptotic variance of  $\hat{\theta}_n$ .

**Theorem 7.** Suppose conditions R1–R6 are satisfied. Let the hypothesis  $H_0 : \theta_X = \theta_0$  hold. Then:

(i)

$$2 \log \lambda_n = 2(\ell_n(\hat{\theta}_n) - \ell_n(\theta_0)) \xrightarrow{D} \chi_d^2$$

(ii)

$$W_n \xrightarrow{D} \chi_d^2$$

(iii)

$$R_n \xrightarrow{D} \chi_d^2$$

**Note.** If  $H_0$  holds,  $\hat{\theta}_n$  should be close to  $\theta_0$ ,  $L_n(\hat{\theta}_n)$  should be close to  $L_n(\theta_0)$ , and  $\mathbf{U}_n(\theta_0|\mathbf{X})$  should be close to  $\mathbf{0}$ . Under  $H_0$ , all three test statistics have values close to 0. Their large values testify against  $H_0$ .

**Corollary.** Denote by  $\chi_d^2(1 - \alpha)$  the  $(1 - \alpha)$ -quantile of  $\chi_d^2$  distribution. Consider tests of  $H_0 : \theta_X = \theta_0$  against  $H_1 : \theta_X \neq \theta_0$  defined by the rule: reject  $H_0$  in favor of  $H_1$ , if

(i)  $2 \log \lambda_n \geq \chi_d^2(1 - \alpha)$  (*likelihood ratio test*)

(ii)  $W_n \geq \chi_d^2(1 - \alpha)$  (*Wald test*)

(iii)  $R_n \geq \chi_d^2(1 - \alpha)$  (*score test*)

Each of these tests has asymptotically (for  $n \rightarrow \infty$ ) the level  $\alpha$ .

**Note.** It can be shown that these three tests are asymptotically equivalent. For large sample sizes, their results are almost identical. With smaller sample sizes, their results can differ. Investigations of small sample behavior of these test statistics revealed that the likelihood ratio test has the best properties, the Wald test is the worst of the three.

Thus, in practical applications, the likelihood ratio test should be preferred.

**Note.** Under normality, the three test statistics are identical.

### Estimation in the presence of nuisance parameters and testing of composite hypotheses

It is frequently desirable to estimate and test just a small number of parameters in a model that contains a much larger number of parameters. We divide the parameter vector into two subsets: the parameters of interest and the other parameters – *nuisance parameters*.

Let  $\theta$  be divided into  $\theta_A$  containing the first  $m$  components of  $\theta$ , and  $\theta_B$  containing the remaining  $d - m$  components of  $\theta$ . We have

$$\theta = (\theta_A, \theta_B)^T = (\theta_1, \dots, \theta_m, \theta_{m+1}, \dots, \theta_d)^T$$

We want to test the hypothesis  $H_0^* : \boldsymbol{\theta}_X \in \Theta_0$  against  $H_1^* : \boldsymbol{\theta}_X \notin \Theta_0$ , where  $\Theta_0 = \{\boldsymbol{\theta} : \boldsymbol{\theta}_A = \boldsymbol{\theta}_{A0}\} \subset \Theta$ . We want to know whether the first  $m$  components of  $\boldsymbol{\theta}_X$  are equal to the vector of constants  $\boldsymbol{\theta}_{A0}$  regardless of the other  $d - m$  components of  $\boldsymbol{\theta}_X$ .

This is not a simple null hypothesis because there are many distributions in the model  $\mathcal{F}$  that satisfy  $H_0^*$ .

All the vectors and matrices appearing in the notation of maximum likelihood estimation theory are decomposed into the first  $m$  components (part  $A$ ) and the remaining  $d - m$  components (part  $B$ ). For example,

$$\hat{\boldsymbol{\theta}}_n = \begin{pmatrix} \hat{\boldsymbol{\theta}}_{An} \\ \hat{\boldsymbol{\theta}}_{Bn} \end{pmatrix}, \quad \mathbf{u}_n(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{u}_{An}(\boldsymbol{\theta}) \\ \mathbf{u}_{Bn}(\boldsymbol{\theta}) \end{pmatrix}, \quad I(\boldsymbol{\theta}) = \begin{pmatrix} I_{AA}(\boldsymbol{\theta}) & I_{AB}(\boldsymbol{\theta}) \\ I_{BA}(\boldsymbol{\theta}) & I_{BB}(\boldsymbol{\theta}) \end{pmatrix}, \quad \text{etc.}$$

The following lemma is useful for inverting the decomposed information matrix.

**Lemma 8** (Block matrix inversion). Let the matrix

$$I = \begin{pmatrix} I_{AA} & I_{AB} \\ I_{BA} & I_{BB} \end{pmatrix}$$

be of full rank. Then there exists an inverse matrix to  $I$  and it can be expressed as

$$I^{-1} = \begin{pmatrix} I^{AA} & I^{AB} \\ I^{BA} & I^{BB} \end{pmatrix},$$

where

$$\begin{aligned} I^{AA} &= I_{AA.B}^{-1} \\ I^{AB} &= -I_{AA.B}^{-1} I_{AB} I_{BB}^{-1} \\ I^{BA} &= -I_{BB.A}^{-1} I_{BA} I_{AA}^{-1} \\ I^{BB} &= I_{BB.A}^{-1} \\ I_{AA.B} &= I_{AA} - I_{AB} I_{BB}^{-1} I_{BA} \\ I_{BB.A} &= I_{BB} - I_{BA} I_{AA}^{-1} I_{AB}. \end{aligned}$$

If the null hypothesis  $H_0^* : \boldsymbol{\theta}_X \in \Theta_0$  holds we know that  $\boldsymbol{\theta}_{AX} = \boldsymbol{\theta}_{A0}$ , but we do not know the value of  $\boldsymbol{\theta}_{BX}$ . We can estimate  $\boldsymbol{\theta}_{BX}$  by the maximum likelihood method applied to the nested submodel

$$\mathcal{F}_0 = \{\text{distributions with density } f(x | (\boldsymbol{\theta}_A, \boldsymbol{\theta}_B)), \boldsymbol{\theta}_A = \boldsymbol{\theta}_{A0}, \boldsymbol{\theta}_B \in \Theta_B \subseteq \mathbb{R}^{d-m}\},$$

with  $d - m$  unknown parameters.

Denote the maximum likelihood estimator of  $\boldsymbol{\theta}_X$  in the submodel  $\mathcal{F}_0$  by  $\tilde{\boldsymbol{\theta}}_n = \begin{pmatrix} \tilde{\boldsymbol{\theta}}_{An} \\ \tilde{\boldsymbol{\theta}}_{Bn} \end{pmatrix}$ , where  $\tilde{\boldsymbol{\theta}}_{An} = \boldsymbol{\theta}_{A0}$  and  $\tilde{\boldsymbol{\theta}}_{Bn}$  solves the system of likelihood equations

$$\mathbf{u}_{Bn}(\boldsymbol{\theta}_{A0}, \tilde{\boldsymbol{\theta}}_{Bn}) = \mathbf{0}.$$

The Fisher information matrix for  $\boldsymbol{\theta}_B$  in this model is  $I_{BB}(\boldsymbol{\theta}_X)$ .

By Theorems 3 and 4 applied to the submodel  $\mathcal{F}_0$ , we get

$$\frac{1}{\sqrt{n}} \mathbf{u}_{Bn}(\boldsymbol{\theta}_X) \xrightarrow{D} N_{d-m}(\mathbf{0}, I_{BB}(\boldsymbol{\theta}_X))$$

and

$$\sqrt{n}(\tilde{\boldsymbol{\theta}}_{Bn} - \boldsymbol{\theta}_{BX}) \xrightarrow{D} N_{d-m}(\mathbf{0}, I_{BB}^{-1}(\boldsymbol{\theta}_X)).$$

On the other hand, Theorems 3 and 4 and Lemma 8 applied to the larger model imply

$$\frac{1}{\sqrt{n}}\mathbf{U}_{Bn}(\boldsymbol{\theta}_X) \xrightarrow{D} N_{d-m}(\mathbf{0}, I_{BB}(\boldsymbol{\theta}_X))$$

and

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_{Bn} - \boldsymbol{\theta}_{BX}) \xrightarrow{D} N_{d-m}(\mathbf{0}, I_{BB.A}^{-1}(\boldsymbol{\theta}_X)),$$

where (dropping the arguments  $\boldsymbol{\theta}_X$ )

$$I_{BB.A}^{-1} = (I_{BB} - I_{BA}I_{AA}^{-1}I_{AB})^{-1} \geq I_{BB}^{-1}.$$

Thus, the asymptotic variance of the MLE of the parameter  $\boldsymbol{\theta}_{BX}$  depends on whether or not  $\boldsymbol{\theta}_{AX}$  is known. If  $\boldsymbol{\theta}_{AX}$  is known (which is true if  $H_0^*$  holds), the asymptotic variance of the MLE  $\tilde{\boldsymbol{\theta}}_{Bn}$  is generally larger than the asymptotic variance of the MLE  $\hat{\boldsymbol{\theta}}_{Bn}$  that does not assume a known  $\boldsymbol{\theta}_{AX}$ .

However, when  $I_{BA} = 0$  (the estimators of  $\boldsymbol{\theta}_{AX}$  and  $\boldsymbol{\theta}_{BX}$  are asymptotically independent), then the asymptotic variances of  $\tilde{\boldsymbol{\theta}}_{Bn}$  and  $\hat{\boldsymbol{\theta}}_{Bn}$  are the same. Then it does not matter whether or not  $\boldsymbol{\theta}_{AX}$  is known.

Let us generalize the three test statistics introduced in Definition 5 of the previous section to testing the composite hypothesis  $H_0^* : \boldsymbol{\theta}_X \in \Theta_0$  against  $H_1^* : \boldsymbol{\theta}_X \notin \Theta_0$ , where  $\Theta_0 = \{\boldsymbol{\theta} : \boldsymbol{\theta}_A = \boldsymbol{\theta}_{A0}\} \subset \Theta$ .

**Definition 6.**

(i) The statistic

$$\lambda_n^* = \frac{L_n(\hat{\boldsymbol{\theta}}_n)}{L_n(\tilde{\boldsymbol{\theta}}_n)}$$

is called *the likelihood ratio*.

(ii) The statistic

$$W_n^* = n(\hat{\boldsymbol{\theta}}_{An} - \boldsymbol{\theta}_{A0})^\top \hat{I}_{AA.B}(\hat{\boldsymbol{\theta}}_n)(\hat{\boldsymbol{\theta}}_{An} - \boldsymbol{\theta}_{A0})$$

is called *the Wald statistic*.

(iii) The statistic

$$R_n^* = \frac{1}{n}\mathbf{U}_n(\tilde{\boldsymbol{\theta}}_n)^\top \hat{I}_n^{-1}(\tilde{\boldsymbol{\theta}}_n)\mathbf{U}_n(\tilde{\boldsymbol{\theta}}_n)$$

is called *the Rao (score) statistic*.

**Note.**

- Obviously,  $\lambda_n^* \geq 1$ .
- The expression  $\hat{I}_{AA.B}$  in the Wald statistic means the inverse of the upper left block of the the matrix  $\hat{I}_n^{-1}$ .
- Since  $\mathbf{U}_{Bn}(\tilde{\boldsymbol{\theta}}_n) = \mathbf{0}$ , the Rao statistic can be written as

$$R_n^* = \frac{1}{n}\mathbf{U}_{An}(\tilde{\boldsymbol{\theta}}_n)^\top \hat{I}_{AA.B}^{-1}(\tilde{\boldsymbol{\theta}}_n)\mathbf{U}_{An}(\tilde{\boldsymbol{\theta}}_n).$$

- Theh Rao statistic does not require the calculation of the MLE  $\hat{\boldsymbol{\theta}}_n$  in the larger model, it only needs the MLE  $\tilde{\boldsymbol{\theta}}_n$  in the submodel. This is often much easier to get.



**Theorem 9.** Let the null hypothesis  $H_0^* : \boldsymbol{\theta}_X \in \Theta_0$ , where  $\Theta_0 = \{\boldsymbol{\theta} : \boldsymbol{\theta}_A = \boldsymbol{\theta}_{A0}\}$ , hold. Then

(i)

$$2 \log \lambda_n^* = 2(\ell_n(\hat{\boldsymbol{\theta}}_n) - \ell_n(\tilde{\boldsymbol{\theta}}_n)) \xrightarrow{D} \chi_m^2;$$

(ii)

$$W_n^* \xrightarrow{D} \chi_m^2;$$

(iii)

$$R_n^* \xrightarrow{D} \chi_m^2.$$

**Note.** Under  $H_0^*$ , we expect  $\hat{\boldsymbol{\theta}}_n$  to be close to  $\tilde{\boldsymbol{\theta}}_n$ ,  $L_n(\hat{\boldsymbol{\theta}}_n)$  to be close to  $L_n(\tilde{\boldsymbol{\theta}}_n)$ , and  $\mathbf{U}_n(\tilde{\boldsymbol{\theta}}_n)$  to be close to  $\mathbf{0}$ . The large values of the three test statistics testify against the null hypothesis.

**Corollary.** Let  $\chi_m^2(1 - \alpha)$  be  $(1 - \alpha)$ -quantile of the  $\chi_m^2$  distribution. Consider tests of  $H_0^* : \boldsymbol{\theta}_X \in \Theta_0$ , where  $\Theta_0 = \{\boldsymbol{\theta} : \boldsymbol{\theta}_A = \boldsymbol{\theta}_{A0}\}$ , against  $H_1^* : \boldsymbol{\theta}_X \notin \Theta_0$  given by the rule: reject  $H_0^*$  in favor of  $H_1^*$  if

- (i)  $2 \log \lambda_n^* \geq \chi_m^2(1 - \alpha)$  (the likelihood ratio test)
- (ii)  $W_n^* \geq \chi_m^2(1 - \alpha)$  (the Wald test)
- (iii)  $R_n^* \geq \chi_m^2(1 - \alpha)$  (the score test)

Then each of these three tests has asymptotically (for  $n \rightarrow \infty$ ) the level  $\alpha$ .

**Note.** The number of degrees of freedom in the reference  $\chi_m^2$  distribution is equal to the number of tested parameters.

**Note.** These three tests are asymptotically equivalent under the null hypothesis as well as under local alternatives. With small or moderate sample sizes, the likelihood ratio test has the best properties and the Wald test is the worst of the three. In practical applications, the likelihood ratio test should be preferred.

**Note.** Let  $m = 1$ ,  $\boldsymbol{\theta}_{AX} = \theta_{Xj}$ , and  $\boldsymbol{\theta}_{A0} = 0$ . Consider the test of the hypothesis  $H_0^* : \theta_{Xj} = 0$  against  $H_1^* : \theta_{Xj} \neq 0$  (zero value of the  $j$ -th parameter in the presence of other parameters that are unspecified by the hypothesis). Then the Wald statistic can be written as

$$W_n = \left[ \frac{\hat{\theta}_{jn}}{\sqrt{n^{-1} \hat{I}_{jj}^{-1}}} \right]^2,$$

where  $n^{-1} \hat{I}_{jj}^{-1}$  is the estimator of the asymptotic variance of  $\hat{\theta}_{jn}$ . This is the square of the test statistic that statistical software typically evaluates to test zero value of a single model parameter.