

Inférence Statistique et Vraisemblance

LSTAT2040

Anouar El Gouch

Contents

Notations and Abbreviations	5
1 Parametric models and exponential families	9
1.1 Motivation and formalization	9
1.1.1 General formalization	10
1.1.2 Parametric models	10
1.1.3 Identifiability	12
1.1.4 Purpose of inferential statistics	13
1.2 Exponential family	14
1.2.1 One-parameter exponential family	14
1.2.2 Properties of the exponential family	16
1.2.3 Multiparameter exponential family	17
1.3 Some useful tools	18
2 Basic concepts of point estimation	23
2.1 Statistic, estimator and estimation	23
2.2 Risk and loss function	24
2.3 MSE, bias, variance, and relative efficiency	25
2.4 Simulation-Based Evaluation of Estimator Performance	28
2.5 The best unbiased estimator (MVUE)	29
3 Fisher information and Cramer-Rao bound	33
3.1 Score and Fisher information	33
3.2 FI contained in a statistic	36
3.3 Sufficient statistic	37
3.4 FI and re-parametrization	40
3.5 Information Inequality: The Cramer-Rao Lower bound (CRLB)	42
3.6 Efficiency in exponential families	44
3.7 CRLB Attainment	45
3.8 Multiparameter case	47
4 Asymptotic evaluations	55
4.1 Convergence in probability and consistency	55

4.2	Convergence in law and asymptotic distribution	59
4.3	Tools for proving asymptotic results	63
4.3.1	Continuous mapping theorem	63
4.3.2	Slutsky's theorem	63
4.3.3	Delta method	64
4.4	Asymptotic efficiency	67
5	Estimation methods	71
5.1	The method of moments (MoM)	71
5.2	The method of maximum likelihood (ML)	72
5.2.1	Likelihood: definition and meaning	72
5.2.2	MLE implementation	74
6	Finite and large sample properties of MLE	81
6.1	Finite sample properties	81
6.1.1	Efficiency	81
6.1.2	Invariance to re-paramerization	81
6.2	Large sample properties	82
6.2.1	Consistency	83
6.2.2	Asymptotic normality and asymptotic efficiency	84
6.2.3	Observed Fisher information	86
7	More about likelihood	89
7.1	Numerical maximization of the likelihood	89
7.1.1	The Newton-Raphson (NR) method	89
7.1.2	Maximum Likelihood in R	93
7.2	Profile likelihood	97
7.3	Likelihood for regression models	99
8	Hypothesis testing	103
8.1	Basic concepts	103
8.2	Evaluating a test	105
8.3	Testing strategy: the Neyman-Pearson approach	107
8.4	The p -Value	110
8.5	Likelihood Ratio Test (LRT)	111
8.6	Asymptotic Tests	115
8.6.1	Simple null hypothesis	115
8.6.2	Composite null hypothesis	120

Notations and Abbreviations

We introduce here the notations and abbreviations that will be used consistently throughout the course. Additional notations will be defined later as they are needed and first appear.

- Unless otherwise specified, vectors and matrices are emphasized using a **bold** font. The null vector is denoted by $\mathbf{0}$.
- The notation $\text{diag}(a_1, \dots, a_n)$ is used for a matrix with diagonal elements a_1, \dots, a_n and all off-diagonal elements zero. The symbol $\mathbb{1}$ is used to denote an identity matrix, i.e. the matrix $\text{diag}(1, \dots, 1)$. The determinant of a matrix A is denoted by $\det(A)$.
- Unless otherwise specified, *any vector that appears in a mathematical operation—such as addition, subtraction, or multiplication—is implicitly interpreted as a **column vector***. For readability, we may write a vector in inline form, for example $\mathbf{x} = (x_1, x_2, x_3)$, but whenever \mathbf{x} enters a vector or matrix expression, it is understood to mean the column vector $\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$. The index t denotes vector/matrix **transpose**. For example,
$$\begin{pmatrix} x_1 & x_2 & x_3 \end{pmatrix}^t = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}.$$
- Random variables are assigned to letters from the last part of the alphabet (X, Y, Z, U, V, \dots), while corresponding observations are assigned to lowercase letters (x, y, z, u, v, \dots).
- Constants are assigned to letters from the first part of the alphabet (a, b, c, \dots).
- The symbol $:=$ indicates an **assignment**. For example, if $x = 2$ and we write $x = y := 2$, then it means that x and y are equal, as y is *defined* to be 2. The symbol \equiv denotes **equivalence**, i.e. things that have exactly the same meaning. These symbols will only be used when they are really relevant to the understanding of a particular concept or situation, otherwise “=” will be used.
- The *probability density function* (pdf) of a continuous random variable/vector, or the *probability mass function* (pmf) of a discrete random variable/vector, will be simply referred to as a *probability distribution* (**pd**), and will be typically denoted f .
- rv/rve: Random variable(s)/Random vector(s).
- cdf: Cumulative distribution function; typically denoted by F .
- iid: Independent and identically distributed rvs (or rves).
- $X \perp\!\!\!\perp Y$: X and Y are independent of each other.
- \log and \ln will be used interchangeably to refer to the **natural logarithm**.

- Any integral sign \int , without limits, should be understood as $\int_{-\infty}^{\infty}$ in the univariate case and as $\int_{\mathbb{R}^d}$ in the multivariate case.
- The **indicator function** of a given set A is denoted by $I(x \in A)$. It takes the value 1 if $x \in A$, and 0 otherwise.
- $a = \arg \min_{x \in I} f(x)$ means that a is a value for which $x \mapsto f(x)$ reaches its minimum on I . The same applies to $\arg \max$.
- For a d -dimensional rve $\mathbf{X} = (X_1, \dots, X_d)^t$. $\boldsymbol{\mu} = E(\mathbf{X})$ denotes the vector of first moments (means), i.e. $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)^t$, with $\mu_j = E(X_j)$. The variance-covariance, or simply the variance, of \mathbf{X} is the $d \times d$ symmetric matrix denoted by $\text{Var}(\mathbf{X}) := E((\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^t) = E(\mathbf{X}\mathbf{X}^t) - \boldsymbol{\mu}\boldsymbol{\mu}^t$. The (j, k) element of this matrix is nothing but $\text{Cov}(X_j, X_k)$. Note that, if \mathbf{A} and \mathbf{B} are two constant objects, then $E(\mathbf{A} + \mathbf{B}\mathbf{X}) = \mathbf{A} + \mathbf{B}E(\mathbf{X})$, and $\text{Var}(\mathbf{A} + \mathbf{B}\mathbf{X}) = \mathbf{B}\text{Var}(\mathbf{X})\mathbf{B}^t$.
- We will use the notation N_d to designate a multivariate normal distribution of dimension d , $d \geq 2$. For a univariate normal ($d = 1$), we simply write N (without any subscript).
- For a scalar function $f(x) : \mathbb{R} \rightarrow \mathbb{R}$, the first **derivative** evaluated at $x = a$ is denoted $f'(a) = \left. \frac{df(x)}{dx} \right|_{x=a}$, and the second derivative is denoted by $f''(a)$. For $k \geq 3$, the k -th order derivative is denoted by $f^{(k)}(a)$.
- For a multivariate function $f(\mathbf{x}) = f(x_1, \dots, x_d) : \mathbb{R}^d \rightarrow \mathbb{R}$:
 - The **partial derivative**, at $\mathbf{x} = \mathbf{a}$, with respect to x_i is denoted by $\partial_i f(\mathbf{a}) = \partial_{x_i} f(\mathbf{a}) = \left. \frac{\partial f(\mathbf{x})}{\partial x_i} \right|_{\mathbf{x}=\mathbf{a}}$. Similarly, we denote the second-order partial derivative by $\partial_{ij} f(\mathbf{a}) = \partial_{x_i x_j} f(\mathbf{a}) = \left. \frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} \right|_{\mathbf{x}=\mathbf{a}}$, and $\partial_i^2 f(\mathbf{a}) = \partial_{x_i}^2 f(\mathbf{a}) = \left. \frac{\partial^2 f(\mathbf{x})}{\partial x_i^2} \right|_{\mathbf{x}=\mathbf{a}}$.
 - The **gradient**, at a point \mathbf{a} , is the column vector given by $\nabla f(\mathbf{a}) = (\partial_1 f(\mathbf{a}), \dots, \partial_d f(\mathbf{a}))^t$.
 - The **Hessian**, at a point \mathbf{a} , is the symmetric matrix given by $\mathbf{H}_f(\mathbf{a}) = \nabla^2 f(\mathbf{a}) = [\partial_{ij} f(\mathbf{a})]_{i,j=1,\dots,d}$.

For example, for $f(x_1, x_2) = x_1^2 + 2x_2^3 + x_1 x_2$, $\nabla f(x_1, x_2) = \begin{pmatrix} 2x_1 + x_2 \\ x_1 + 6x_2^2 \end{pmatrix}$ and

$$\nabla^2 f(x_1, x_2) = \begin{pmatrix} 2 & 1 \\ 1 & 12x_2 \end{pmatrix}.$$
- For a multivariate vector-valued function $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_p(\mathbf{x})) : \mathbb{R}^d \rightarrow \mathbb{R}^p$, where $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$, $i = 1, \dots, p$, the **Jacobian** is the $p \times d$ matrix $\mathbf{J}_f(\mathbf{a}) = \dot{\mathbf{f}}(\mathbf{a}) = [\partial_j f_i(\mathbf{a})]_{i=1,\dots,p; j=1,\dots,d}$.

For example, for $f(x_1, x_2) = (x_1^3, x_1^2 + 2x_2^3 + x_1x_2, x_2^2)$, $\dot{f}(x_1, x_2) = \begin{pmatrix} 3x_1^2 & 0 \\ 2x_1 + x_2 & x_1 + 6x_2^2 \\ 0 & 2x_2 \end{pmatrix}$.

Note that Hessian is the Jacobian matrix of the gradient, thus $\nabla^2 f(\mathbf{a}) = J_{\nabla f}(\mathbf{a})$.

Chapter 1

Parametric models and exponential families

1.1 Motivation and formalization

In order to obtain an estimate of an unknown quantity, say μ_0 , say, for example, a speed, it is common to take n measurements x_1, \dots, x_n and calculate their mean:

$$\bar{x}_n = \frac{x_1 + x_2 + \dots + x_n}{n}.$$

But why should the observations be combined in this way ?

→ the empirical mean is a relevant measure of the center of the observations, since

$$\bar{x}_n = \arg \min_a \sum_{i=1}^n (x_i - a)^2.$$

But (at this level) we can't justify why \bar{x}_n is a good estimate (approximation) of the true value μ_0 since no explicit assumption has been made to connect the data (x_1, \dots, x_n) and μ_0 .

To establish such a connection, we can, for example, presume that:

- (i) each x_i is an observed value of a rv X_i , and
- (ii) $X_i, i = 1, \dots, n$, have a common mean μ_0 .

Even more specifically, we can, for example, assume the following *additive error model*

$$X_i = \mu_0 + \epsilon_i, \epsilon_i \sim N(0, \sigma_0^2),$$

or equivalently $X_i \sim N(\mu_0, \sigma_0^2)$. In this way, the problem we face is the estimation of $\mu_0 = E(X_i)$ from the sample (X_1, \dots, X_n) .

1.1.1 General formalization

The data $\mathbf{x} = (x_1, \dots, x_n)$ that we observe are believed to be generated by a rve $\mathbf{X} = (X_1, \dots, X_n)$, which represents our random sample. We assume that \mathbf{X} follows some joint distribution which is (partly) unknown. The set of assumptions made about this underlying joint distribution is what we call a *statistical model*.

X_1, \dots, X_n are typically assumed to be *iid* copies of some population rv, which we denote hereafter by X . In this case the statistical model reduces to the set of assumptions about the distribution of X . To describe this latter, we usually use a pd f or a cdf F . Under the iid assumption, the joint pd and the joint cdf of \mathbf{X} , denoted by f_n and F_n , respectively, are given by

$$f_n(\mathbf{x}) = f_n(x_1, \dots, x_n) = \prod_{i=1}^n f(x_i) \text{ and } F_n(\mathbf{x}) = P(X_1 \leq x_1, \dots, X_n \leq x_n) = \prod_{i=1}^n F(x_i).$$

The IID hypothesis plays a crucial role

Without the iid assumption the statistical analysis of the data becomes much more complicated. For example, as a statistical model, we could assume that

$$\mathbf{X} \sim N_n(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

with an unknown $\boldsymbol{\theta} = (\mu_1, \dots, \mu_n, \sigma_{11}, \dots, \sigma_{nn}) \in \mathbb{R}^{n + \frac{n \times (n+1)}{2}}$, where $\mu_i = E(X_i)$ and $\sigma_{ij} = \text{Cov}(X_i, X_j)$.

Now, under the iid assumption, $\boldsymbol{\theta}$ reduces to $(\mu_1, \sigma_{11}) \in \mathbb{R}^2$. \square

1.1.2 Parametric models

A parametric model (or parametric family) is a set of distributions indexed by a *finite* dimensional parameter $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d)$, $d \geq 1$. That is to say that the pd of X – the random variable that generated the observed data – is known up to the unknown parameter $\boldsymbol{\theta}$. In which case, we denote the pd of X by $f(x; \boldsymbol{\theta})$ and its cdf by $F(x; \boldsymbol{\theta})$.

We may write $X \sim f(x; \boldsymbol{\theta})$ or, less commonly, $X \sim F(x; \boldsymbol{\theta})$. When the distribution has a well-known name (e.g., Normal, Poisson, Binomial), the notation \sim is usually followed by the distribution name together with its parameters, as in $X \sim N(\mu, \sigma^2)$ or $X \sim \text{Poisson}(\lambda)$.

The set of possible values for the parameter $\boldsymbol{\theta}$, that we denote by $\boldsymbol{\Theta}$, is called *the parameter space*. Within $\boldsymbol{\Theta}$, the particular value $\boldsymbol{\theta}_0$ that actually generated the observed data is referred to as the true parameter value (or simply the *true value*).

Example 1.1.

- The Bernoulli model: $X \sim f(x; \theta)$, with

$$f(x; \theta) \equiv P_\theta(X = x) = \theta^x (1 - \theta)^{(1-x)} I(x \in \{0, 1\}), \theta \in (0, 1).$$

$\iff X \sim \text{Ber}(\theta_0)$, for some specific but unknown $\theta_0 \in (0, 1) = \Theta$.

- The Exponential model: $X \sim f(x; \theta)$, with

$$f(x; \theta) = \theta^{-1} e^{-x/\theta} I(x \geq 0), \theta > 0.$$

$\iff X \sim \text{Exp}(\theta_0)$, for some unknown $\theta_0 \in (0, \infty) = \Theta$.

- The Normal model: $X \sim f(x; \theta)$, with

$$f(x; \theta) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \theta^t = (\mu, \sigma^2) \in \mathbb{R} \times (0, \infty).$$

$\iff X \sim N(\mu_0, \sigma_0^2)$, for some unknown $(\mu_0, \sigma_0^2) \in \mathbb{R} \times (0, \infty) = \Theta$. \square

Attention. Hereafter, for simplicity and when no ambiguity arises, we drop the subscript 0 in θ_0 and use θ to denote both the generic parameter and the unknown true parameter value.

Misspecified model

If there is no $\theta \in \Theta$ such that $X \sim f(x; \theta)$, then the model $\{f(x; \theta), \theta \in \Theta\}$ is said to be *misspecified*. An example of misspecification is when a normal distribution is used for exponential data. A conclusion drawn from a statistical model is *valid only if the chosen model is correctly specified*. In reality, no model is 100% correct, but some models are more useful than others in approximating the true underlying distribution of the data.

In what follows, unless stated otherwise, we assume that the statistical models under consideration are correctly specified. \square

Parametric vs nonparametric models ?

If the distribution of X is not completely determined by a finite number of parameters, then the model is *nonparametric* or *semiparametric* (i.e., a mix of finite- and infinite-dimensional parameters).

Example 1.2.

- X has a pd f , with $\int f''(x)dx < \infty$ and/or $\int x^2 f(x)dx < \infty$.
- X has a symmetric distribution about 0, i.e. it has a pd satisfying $f(-x) = f(x)$, $\forall x$.
- X has a pd f satisfying $f(x; \mu, \sigma) = \frac{1}{\sigma} f_0\left(\frac{x - \mu}{\sigma}\right)$, where μ and $\sigma > 0$ are unknown parameters and f_0 an unknown pd symmetric about 0.

These models cannot be indexed by a finite dimensional parameter. The first two cases are examples of (fully) nonparametric models, while the last case is an example of a semiparametric model. \square

The question of whether to use a parametric or a nonparametric model depends primarily on prior knowledge about the data-generating process and on the associated risk-benefit

trade-off: efficiency/interpretability of parametric models against the flexibility/robustness of nonparametric ones.

For example, suppose we want to estimate $\theta := F(s) = P(X \leq s) \in [0, 1]$ for a given s .

- If we assume that $X \sim N(\mu, 1)$, then it is reasonable to estimate θ by $\hat{\theta}_{\text{para}} = \Phi(s - \hat{\mu})$, where Φ is the cdf of a $N(0, 1)$ and $\hat{\mu}$ is any given estimator of μ , e.g., the sample mean \bar{X}_n .
 - Advantage: Efficient (makes optimal use of available information \rightarrow small variance), if normality holds.
 - Risk: If the distribution is not normal, $\hat{\theta}_{\text{para}}$ can be severely biased.
 - Without making any assumption about the distribution of X , a reasonable estimator is the empirical cdf $\hat{F}(s) = n^{-1} \sum_{i=1}^n I(X_i \leq s)$.
 - Advantage: Consistent for any distribution.
 - Risk: Larger variance than the parametric estimator when normality is true.
- \rightarrow Incorrect assumptions about the underlying distribution of X produce biased conclusions, whereas correct assumptions yield more efficient estimation. \square

1.1.3 Identifiability

For a given parametric model, each parameter value θ determines exactly the distribution of X . However, this does not exclude the possibility that two distinct parameter values $\theta_1 \neq \theta_2$ may generate exactly the same distribution, that is, $f(x; \theta_1) = f(x; \theta_2)$, $\forall x$.

In such a situation, the two parameter values are indistinguishable from the data, even with an infinite sample, since both produce identical distributions. This phenomenon is known as an *identifiability problem*.

Identifiability is an important property of a statistical model, which determines whether the parameter of interest can be recovered (estimated) from the observed data, which is only possible if different values of θ lead to different distributed samples.

Mathematically, this can be formulated by saying that, in a given model $\{f(x; \theta), \theta \in \Theta\}$, the parameter θ (or the model) is identifiable, if, for any θ_1 and θ_2 in Θ ,

$$f(x; \theta_1) = f(x; \theta_2), \forall x \Rightarrow \theta_1 = \theta_2.$$

Example 1.3.

- The Bernoulli, the Exponential, and the Normal models, as defined above, are identifiable. Why?
- Let $X = \mu_1 + \epsilon$, where $\epsilon \sim N(\mu_2, 1)$ and μ_1 and μ_2 are unknown. Suppose that we observe X (and not ϵ), then $\theta = \mu_1 + \mu_2$ is identifiable but $\theta = (\mu_1, \mu_2)$ is not.
- Let $X = |Y|$, where $Y \sim N(\mu, 1)$ and μ is unknown. Suppose that we observe X , then μ is not identifiable.

Let's verify the identifiability of the Normal model. For that, observe that $f(x, \theta_1) = f(x, \theta_2), \forall x$, is equivalent to

$$\frac{(x - \mu_1)^2}{\sigma_1^2} - \frac{(x - \mu_2)^2}{\sigma_2^2} = 2 \log \frac{\sigma_2}{\sigma_1}, \forall x.$$

Since a parabolic function $ax^2 + bx + c$ vanishes (i.e., $ax^2 + bx + c = 0, \forall x$) if and only if $a = b = c = 0$, and since in our case $a = \frac{1}{\sigma_1^2} - \frac{1}{\sigma_2^2}$, we have that $\sigma_1 = \sigma_2$. As consequence, $f(x; \theta_1) = f(x; \theta_2), \forall x$, is equivalent to

$$(x - \mu_1)^2 - (x - \mu_2)^2 = 0, \forall x \iff \mu_1 = \mu_2.$$

To verify the last example above (with $X = |Y|$), observe that $P_\mu(X \leq x) = \Phi(x - \mu) + \Phi(x + \mu) - 1$, where Φ is the cdf of $N(0, 1)$. It follows that, $P_1(X \leq x) = P_{-1}(X \leq x), \forall x$, i.e., $\mu = 1$ and $\mu = -1$ lead to the same distribution for X . This demonstrates that μ is not identifiable. \square

If a model is not identifiable, it is common to introduce additional constraints/assumptions on it in order to make it identifiable. In that case, the set of these requirements is called the **identifiability conditions**. For instance, in our example above with $X = |Y|$, if we assume that $\mu > 0$, i.e., $\Theta = (0, \infty)$, then μ becomes identifiable; can you prove this ?

1.1.4 Purpose of inferential statistics

Statistical inference is the process of learning about a given probability model using observed data. To be more precise, suppose we are given a data set $\mathbf{x} = (x_1, \dots, x_n)$ which we assume to be generated from the model $\{f(x; \theta), \theta \in \Theta\}$. The aim of parametric statistical inference is to gain knowledge about the unknown parameter θ from \mathbf{x} .

There are three major parametric statistical inference procedures:

1. **Point estimation:** A single value is computed from the data \mathbf{x} and used as an estimate (approximation) of the true parameter value θ .
2. **Hypothesis testing:** Sets up some specific hypotheses regarding θ and evaluate the degree to which the data \mathbf{x} support these hypotheses.
3. **Confidence set estimation:** Use the observed data \mathbf{x} to construct a set of possible values for θ . The resulting set must have a high (predetermined) probability of including the true value.

Other well-known topics in statistical inference include *model selection*, *model validation* and *prediction*.

1.2 Exponential family

One important class of statistical models is the exponential family models. These models are widely used in statistics and machine learning. They are characterized by a simple and elegant mathematical structure, which makes them analytically tractable and computationally efficient. Exponential family contains most of the standard discrete and continuous distributions that are used for modeling, such as (multivariate) normal, poisson, binomial, multinomial, exponential, and gamma.

The reason for the special status of the exponential family is that a number of important and useful results in inference can be unified within it. This family also forms the basis for an important class of regression models, known as generalized linear models.

1.2.1 One-parameter exponential family

A family of probability distributions that depend on a single (scalar) parameter θ is a one-parameter exponential family if it can be expressed as

$$f(x; \theta) = h(x) \exp(g(\theta)T(x) - B(\theta)), \quad \forall x.$$

Here x can be a scalar or vector, $h(x) \geq 0$ and $T(x)$ are functions of x only (*cannot depend on θ*), and $g(\theta)$ and $B(\theta)$ are functions of θ only (*cannot depend on x*). $B(\theta)$ is a normalizing constant, ensuring that $f(x; \theta)$ sums or integrates to 1.

The set $\Theta = \{\theta : \int h(x) \exp(g(\theta)T(x))dx < \infty\}$, in the continuous case, and $\Theta = \{\theta : \sum_x h(x) \exp(g(\theta)T(x)) < \infty\}$, in the discrete case, is the parameter space of the family. $T(X)$ is referred to as **natural sufficient** statistic or simply the sufficient statistic. In many cases $T(x) = x$.

The exponential-family representation above is not unique. In particular, the functions $g(\theta)$ and $T(x)$ are only defined up to nonzero linear rescaling: one may multiply g by a nonzero constant a and divide T by the same constant without changing the resulting distribution. This ambiguity is purely algebraic and has no impact on the underlying probability distributions or any statistical conclusions drawn from them.

An exponential family can be reparameterized as

$$h(x) \exp(\eta T(x) - A(\eta)).$$

This expression is called the **canonical (or natural) representation**, and $\eta = g(\theta)$ is the **canonical parameter**. Here $h(x)$ and $T(x)$ are the same as in the original parameterization, while the normalizing function becomes $A(\eta) = B(g^{-1}(\eta))$, provided that g is invertible. The set $\Lambda = \{\eta : \int h(x) \exp(\eta T(x))dx < \infty\}$ is called the natural parameter space, with the convention that the integral is replaced by a sum in the discrete case.

It's analytically convenient and easier to work with an exponential family in its canonical

form. Once a result has been derived for the canonical form, we can rewrite it in terms of the original parameter θ if desired.

To verify that a given pd is a member of the exponential family, we must identify all the functions h , T , g , and B (or equivalently A). The key step is to rewrite the density so that: (1) all terms depending only on x are absorbed into $h(x)$; (2) all terms depending only on θ are collected into $B(\theta)$; and (3) the remaining mixed term factors as $g(\theta) T(x)$ inside the exponential. The next example illustrates this procedure.

Example 1.4.

- Poisson:

$$\begin{aligned}\frac{\theta^x e^{-\theta}}{x!} &= \frac{1}{x!} \exp(x \log(\theta) - \theta), \quad x = 0, 1, \dots, \text{ and } \theta > 0 \\ &\equiv \frac{1}{x!} \exp(x\eta - e^\eta), \quad \eta \in (-\infty, \infty).\end{aligned}$$

- Binomial:

$$\begin{aligned}C_n^x \theta^x (1 - \theta)^{n-x} &= C_n^x \exp\left(x \log\left(\frac{\theta}{1 - \theta}\right) + n \log(1 - \theta)\right), \quad x = 0, 1, \dots, n \text{ and } \theta \in (0, 1) \\ &\equiv C_n^x \exp(x\eta - n \log(1 + e^\eta)), \quad \eta \in (-\infty, \infty).\end{aligned}$$

- Normal with a *known* $\sigma > 0$ ($\theta = \mu$):

$$\begin{aligned}\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-x^2}{2\sigma^2}\right) \exp\left(\frac{\mu}{\sigma^2}x - \frac{\mu^2}{2\sigma^2}\right), \quad x \in (-\infty, \infty) \text{ and } \mu \in (-\infty, \infty) \\ &\equiv \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-x^2}{2\sigma^2}\right) \exp\left(\eta x - \frac{\eta^2 \sigma^2}{2}\right), \quad \eta \in (-\infty, \infty). \quad \square\end{aligned}$$

Example 1.5 (Counter-example). The following density

$$f(x; \theta) = \exp(-(x - \theta)) I(x \geq \theta) = \begin{cases} \exp(-(x - \theta)), & x \geq \theta, \\ 0, & x < \theta. \end{cases}$$

is not an exponential family. The obstruction comes from the indicator term $I(x \geq \theta)$ which cannot be factored into a function of x times a function of θ , as required for the exponential-family form. \square

In general, for a pd $f(x; \theta)$ to belong to an exponential family, its *support* $S = \{x : f(x; \theta) > 0\}$ must be independent of θ .

1.2.2 Properties of the exponential family

A useful fact about exponential families is that the integral $\int h(x) \exp(g(\theta)T(x) - B(\theta))dx$, in the continuous case, or the sum $\sum_x h(x) \exp(g(\theta)T(x) - B(\theta))$, in the discrete case, can be differentiated, with respect to θ (or with respect η for the canonical representation), any number of times *by moving the derivative inside the integral or sum*. This property is the source of many important results about exponential families. One such result is given next.

Proposition 1.1. *With the canonical parameterization, the mean and variance of $T \equiv T(X)$ are given by*

$$E_\eta(T) = A'(\eta) \text{ and } \text{Var}_\eta(T) = A''(\eta) = \partial_\eta E_\eta(T). \quad (1.1)$$

To show the first equality in (1.1), we can differentiate $\eta \mapsto \int h(x) \exp(\eta T(x) - A(\eta))dx = 1$, with respect to η , and then use the fact that the derivative can be moved inside the integral, which gives

$$\begin{aligned} 0 &= \int \partial_\eta h(x) \exp(\eta T(x) - A(\eta))dx \\ &= \int h(x)(T(x) - A'(\eta)) \exp(\eta T(x) - A(\eta))dx \\ &= E_\eta(T) - A'(\eta). \end{aligned}$$

As for the second equality, we can differentiate two times to get

$$\begin{aligned} 0 &= \int \partial_\eta h(x)(T(x) - A'(\eta)) \exp(\eta T(x) - A(\eta))dx \\ &= \int h(x)(-A''(\eta)) \exp(\eta T(x) - A(\eta))dx + \int h(x)(T(x) - A'(\eta))^2 \exp(\eta T(x) - A(\eta))dx \\ &= -A''(\eta) + E_\eta(T - A'(\eta))^2. \end{aligned}$$

In terms of the original parameterization, with θ , we can write

$$E_\theta(T) = \frac{B'(\theta)}{g'(\theta)} \text{ and } \text{Var}_\theta(T) = \frac{\partial_\theta E_\theta(T)}{g'(\theta)}.$$

These can be proven directly from the definition of $f(x; \theta)$ by following the same derivation as we did above for the canonical parameterization. Another way to obtain these results is to use (1.1) and apply the chain rule to $\eta \mapsto A(\eta) = B(g^{-1}(\eta))$, which yields

$$\begin{aligned} A'(\eta) &= \frac{B'}{g'}(g^{-1}(\eta)) \equiv \frac{B'}{g'}(\theta) \\ A''(\eta) &= \frac{1}{g'(g^{-1}(\eta))} \left(\frac{B'}{g'} \right)'(g^{-1}(\eta)) \equiv \frac{1}{g'(\theta)} \left(\frac{B'}{g'} \right)'(\theta), \end{aligned}$$

with $\theta = g^{-1}(\eta)$.

Attention. In the formulas of expectations and variances above, the subscripts θ and η indicate the parameterization with respect to which the calculations are carried out. We will use this notation whenever necessary; otherwise, the subscripts will be omitted to lighten the notation.

Example 1.6. For $N(\mu, \sigma^2)$, see Example 1.4, we have that

- From the Canonical form :

$$E(X) = \partial_\eta \left(\frac{\eta^2 \sigma^2}{2} \right) = \eta \sigma^2 = \mu \quad \text{and} \quad \text{Var}(X) = \partial_\eta (\eta \sigma^2) = \sigma^2.$$

- From the original form :

$$E(X) = \frac{\partial_\mu \frac{\mu^2}{2\sigma^2}}{\partial_\mu \frac{\mu}{\sigma^2}} = \frac{\frac{\mu}{\sigma^2}}{\frac{1}{\sigma^2}} = \mu \quad \text{and} \quad \text{Var}(X) = \frac{\partial_\mu \mu}{1/\sigma^2} = \sigma^2.$$

Another interesting fact about the exponential family is that its structure is preserved under iid sampling. This is better explained in the following.

Proposition 1.2. If X_1, \dots, X_n are iid rv from the exponential family, as defined above, with sufficient statistic T then the joint distribution of $\mathbf{X} = (X_1, \dots, X_n)$:

$$f_n(\mathbf{x}; \eta) = \left[\prod_{i=1}^n h(x_i) \right] \exp \left(\eta \sum_{i=1}^n T(x_i) - nA(\eta) \right)$$

is also an exponential family with sufficient statistic $\sum_{i=1}^n T(X_i)$.

1.2.3 Multiparameter exponential family

A distribution is said to belong to the J -parameter exponential family ($J \geq 1$) if its density or probability mass function can be represented in the form

$$f(x; \boldsymbol{\theta}) = h(x) \exp(\mathbf{g}^t(\boldsymbol{\theta}) \mathbf{T}(x) - B(\boldsymbol{\theta})) = h(x) \exp \left(\sum_{j=1}^J g_j(\boldsymbol{\theta}) T_j(x) - B(\boldsymbol{\theta}) \right),$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_J)$ is the parameter vector, $\mathbf{g}^t(\boldsymbol{\theta}) = (g_1(\boldsymbol{\theta}), \dots, g_J(\boldsymbol{\theta}))$, with each $g_j : \mathbb{R}^J \mapsto \mathbb{R}$, $B : \mathbb{R}^J \rightarrow \mathbb{R}$, and the vector of statistics $\mathbf{T}^t(x) = (T_1(x), \dots, T_J(x))$ is called the sufficient statistic vector.

If we reparameterize by setting $\eta_j := g_j(\boldsymbol{\theta})$, $j = 1, \dots, J$, the family is called the J -parameter canonical exponential family, and we obtain

$$f^*(x; \boldsymbol{\eta}) = h(x) \exp(\boldsymbol{\eta}^t \mathbf{T}(x) - A(\boldsymbol{\eta})),$$

where $\boldsymbol{\eta} = (\eta_1, \dots, \eta_J)$ is called the canonical (or natural) parameter vector, and $A : \mathbb{R}^J \rightarrow \mathbb{R}$.

The properties established in the one-parameter case extend directly to the multi-parameter setting. Using the shorthand $T_j \equiv T_j(X)$, one obtains

$$E(T_j) = \partial_{\eta_j} A(\boldsymbol{\eta}), \text{ and } Cov(T_j, T_k) = \partial_{\eta_j \eta_k} A(\boldsymbol{\eta}), j, k = 1, \dots, J.$$

In matrix notation, these identities become $E(\mathbf{T}) = \nabla A(\boldsymbol{\eta})$ and $Var(\mathbf{T}) = \nabla^2 A(\boldsymbol{\eta})$.

Example 1.7. Normal distribution with an unknown μ and σ ($\boldsymbol{\theta} = (\mu, \sigma^2)$):

$$\begin{aligned} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) &= \frac{1}{\sqrt{2\pi}} \exp\left(\frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2 - \left(\frac{\mu^2}{2\sigma^2} + \log(\sigma)\right)\right) \\ &\equiv \frac{1}{\sqrt{2\pi}} \exp(\eta_1 x + \eta_2 x^2 - A(\boldsymbol{\eta})), \end{aligned}$$

where $(\eta_1, \eta_2) = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right) \iff (\mu, \sigma^2) = \left(-\frac{\eta_1}{2\eta_2}, -\frac{1}{2\eta_2}\right)$, and $A(\eta_1, \eta_2) = -\frac{\eta_1^2}{4\eta_2} - \frac{1}{2}\log(-2\eta_2)$. Here $T_1(x) = x$ and $T_2(x) = x^2$. And we have that,

$$\begin{aligned} E(X) &= \partial_1 A = -\frac{\eta_1}{2\eta_2} = \mu, \quad E(X^2) = \partial_2 A = \frac{\eta_1^2 - 2\eta_2}{4\eta_2^2} = \mu^2 + \sigma^2 \\ Var(X) &= \partial_1^2 A = -\frac{1}{2\eta_2} = \sigma^2, \quad Var(X^2) = \partial_2^2 A = \frac{\eta_2 - \eta_1^2}{2\eta_2^3} = 2\sigma^2(\sigma^2 + 2\mu^2), \end{aligned}$$

$$\text{and } Cov(X, X^2) = \partial_{1,2} A(\boldsymbol{\eta}) = \frac{\eta_1}{2\eta_2^2} = 2\mu\sigma^2. \quad \square$$

1.3 Some useful tools

In this section, we will go through some important mathematical and statistical properties that will be used later in the course.

Law of total expectation/variance

For any two rv X and Y ,

$$\begin{aligned} E(Y) &= E(E[Y|X]), \\ Var(Y) &= E(Var[Y|X]) + Var(E[Y|X]), \end{aligned}$$

where $Var[Y|X] := E[(Y - E[Y|X])^2|X] = E[Y^2|X] - (E[Y|X])^2$, and $E(Y|X = x) = \int y f_{Y|X}(y|x) dy$.

Expected value of a non-negative rv

If X is a non-negative rv, then $E(X) \geq 0$. Moreover, $E(X) = 0$ if and only if $X = 0$ almost surely (i.e., $P(X = 0) = 1$).

Markov-Chebyshev's inequality

If X is a non-negative rv, then $E(X) \geq kP(X \geq k)$, $\forall k \in \mathbb{R}$. As a consequence, for any rv X and any constant $k > 0$, $P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$, where $\mu = E(X)$ and $\sigma = \sqrt{\text{Var}(X)}$.

To see the first inequality, observe that, $\forall X, \forall k, X = XI(X \geq k) + XI(X < k)$. Which, by the fact that $X \geq 0$, implies that, $X \geq kI(X \geq k)$, and hence $E(X) \geq kP(X \geq k)$. Applying this last equality to $(X - \mu)^2$, instead of X , we obtain the second inequality.

Jensen's inequality

If f is a **convex function** (Reminder: $f''(x) \geq 0, \forall x \in I \Rightarrow f$ is convex in I), then

$$f(E(X)) \leq E(f(X)).$$

Moreover, if f is strictly convex ($f'' > 0$), then this inequality is strict unless X is almost surely constant (i.e., X takes the same value all the time; thus $\exists c$ such that $P(X = c) = 1$).

The opposite holds for concave functions (Reminder: f is concave if and only if $-f$ is convex).

Example. Let X be a non-constant rv. Since $x \mapsto |x|^a$ is convex for $a \geq 1$, and strictly convex for $a > 1$, we have that, for example, $|E(X)| \leq E|X|$ and $(E(X))^2 < E(X^2)$. And since $x \mapsto \sqrt{x}$ and $x \mapsto \log(x)$ are strictly concave in $(0, \infty)$, we have that $\sqrt{E(X)} > E(\sqrt{X})$ and $\log(EX) > E(\log X)$, provided that $X > 0$ almost surely. \square

Cauchy-Schwarz's inequality

For any two rv X and Y ,

$$(E(XY))^2 \leq E(X^2)E(Y^2).$$

As a consequence, we get the inequality

$$(\text{Cov}(X, Y))^2 \leq \text{Var}(X)\text{Var}(Y).$$

Moreover, if X and Y are not constant random variables, the last inequality becomes strict unless X and Y are linearly dependent; that is, equality holds if and only if there exist constants $a \neq 0$ and b such that $Y = aX + b$ almost surely.

The composite function rule (chain rule)

- Basic version : If $h(x) = f(g(x))$, then

$$h'(x) = f'(g(x))g'(x).$$

By putting $z = f(y)$ and $y = g(x)$, the above formula can be expressed as $\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$.

This formula can be generalized in several ways.

- Case of multiple compositions: If $z = f(y)$, $y = g(x)$, and $x = h(t)$, then

$$\frac{dz}{dt} = \frac{dz}{dy} \cdot \frac{dy}{dx} \cdot \frac{dx}{dt}.$$

- Case of a function of two variables: If $z = f(x, y)$, $x = g(t)$, and $y = h(t)$, then

$$\frac{dz}{dt} = \frac{\partial z}{\partial x} \cdot \frac{dx}{dt} + \frac{\partial z}{\partial y} \cdot \frac{dy}{dt}.$$

Taylor's theorem

Suppose f is a function such that $f^{(n+1)}$ ($n \geq 0$) is continuous on some interval I . Then, for any $x, a \in I$, there exists a $\theta \in [0, 1]$ such that

$$f(x) = \sum_{i=0}^n \frac{(x-a)^i}{i!} f^{(i)}(a) + \frac{(x-a)^{(n+1)}}{(n+1)!} f^{(n+1)}(a + \theta(x-a)).$$

We can use this result to approximate the function f , and write that, in a sufficiently small neighbourhood of a ,

$$f(x) \approx \sum_{i=0}^n \frac{(x-a)^i}{i!} f^{(i)}(a) = f(a) + (x-a)f'(a) + \frac{(x-a)^2}{2} f''(a) + \dots + \frac{(x-a)^n}{n!} f^{(n)}(a).$$

This is called the n th order Taylor polynomial approximation of f around a .

A similar result holds for functions of several variables. For example, the second-order Taylor polynomial approximation of $f : \mathbb{R}^d \mapsto \mathbb{R}$ around a point a is

$$f(x) \approx f(a) + \nabla^t f(a)(x-a) + \frac{1}{2}(x-a)^t \nabla^2 f(a)(x-a).$$

Definite matrix

Definite matrices play a very important role in statistics and optimization

Let A be a $d \times d$ symmetric matrix ($A^t = A$). A is said to be

- *positive definite*, if $x^t A x > 0, \forall x \in \mathbb{R}^d \neq \mathbf{0}$.
- *positive semidefinite*, if $x^t A x \geq 0, \forall x \in \mathbb{R}^d$.

If the inequalities are reversed, then A is *negative definite* or *negative semidefinite*, respectively.

Here are some examples. The matrix $A = \begin{pmatrix} 3 & 0 \\ 0 & 2 \end{pmatrix}$ is positive definite, since $x^t A x = 3x_1^2 +$

$2x_2^2 > 0, \forall x \neq \mathbf{0}$. The matrix $A = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$ is positive semidefinite since $x^t A x = (x_1 + x_2)^2 \geq$

$0, \forall x$. The matrix $A = \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}$ is positive definite since $x^t A x = 2(x_1^2 - x_1 x_2 + x_2^2) >$

0, $\forall \mathbf{x} \neq \mathbf{0}$. The matrix $A = \begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix}$ is *indefinite* since $\mathbf{x}^t A \mathbf{x} = x_1^2 + 4x_1x_2 + x_2^2$ can be positive or negative.

Different methods exist to check if a matrix is positive definite, such as the Cholesky decomposition, the eigenvalues, or the principal minors. In this course we will not go into the details of these methods. It is sufficient to know that a 2×2 matrix $A = \begin{pmatrix} a & b \\ b & c \end{pmatrix}$ is positive definite if and only if $a > 0$ and $\det(A) := ac - b^2 > 0$.

There also are many interesting properties of positive (semi)definite matrices, such as the fact that :

- A matrix is positive definite if and only if it is positive semidefinite and invertible; its inverse is then also positive definite.
- For every positive (semi)definite matrix A , there exists a unique positive (semi)definite matrix B such that $B^2 := B \times B = A$. This B is called the (natural or principal) **square-root** of A and is denoted by $A^{1/2}$.

Some interesting properties of the Variance-Covariance matrix

Let X be a rve in \mathbb{R}^d and $\Sigma = \text{Var}(X) = (\text{Cov}(X_j, X_k))_{1 \leq j, k \leq d}$. Σ is **symmetric** ($\Sigma^t = \Sigma$) and **positive semidefinite**. This last property follows directly from the fact that $\mathbf{a}^t \Sigma \mathbf{a} = \text{Var}(\mathbf{a}^t X)$.

Σ is **positive definite** if and only if the components of X are **linearly independent** (almost surely); i.e. $\nexists \mathbf{a} \neq \mathbf{0}$ such that $\mathbf{a}^t X = \text{constant}$.

Some properties of the multivariate normal

- If B is a $p \times d$ matrix, \mathbf{a} is a p -dimensional vector, and \mathbf{b} a d -dimensional vector, then

$$\mathbf{b}^t \times N_d(\boldsymbol{\mu}, \Sigma) = N(\mathbf{b}^t \boldsymbol{\mu}, \mathbf{b}^t \Sigma \mathbf{b}), \text{ and } \mathbf{a} + B \times N_d(\boldsymbol{\mu}, \Sigma) = N_p(\mathbf{a} + B \boldsymbol{\mu}, B \Sigma B^t).$$

- $X \sim N_d(\boldsymbol{\mu}, \Sigma) \Rightarrow \Sigma^{-1/2}(X - \boldsymbol{\mu}) \sim N_d(\mathbf{0}, \mathbb{I})$, where $\Sigma^{-1/2}$ is the square-root of Σ^{-1} .
- $X \sim N_d(\mathbf{0}, \mathbb{I}) \Rightarrow X^t X \sim \chi_d^2$, where χ_d^2 the chi-squared distribution with d degrees of freedom.

Chapter 2

Basic concepts of point estimation

2.1 Statistic, estimator and estimation

Definition 2.1. Let $X = (X_1, \dots, X_n)$ be a sample. Any (measurable) function $T(X)$ of X , i.e. any quantity that can be calculated solely from the observed data, is called a *statistic*.

An estimator is any statistic used to estimate a given parameter. Typically, we use the notation $\hat{\theta}_n(X) \equiv \hat{\theta}_n \equiv \hat{\theta}$ to denote an *estimator* of θ .

Any realization $\hat{\theta}_n(x)$ of $\hat{\theta}_n(X)$ is an *estimation* (a guess) of θ .

Example 2.1 (Example of statistics).

$$X_1, (X_1, \dots, X_n), \sum_i X_i, \bar{X}_n = n^{-1} \sum_i X_i, n^{-1} \sum_i X_i^2, n^{-1} \sum_i I(X_i \geq 0),$$

$$X_{(1)} = \min_i X_i, X_{(n)} = \max_i X_i, X_{(k)} \text{ the } k\text{th order statistic, i.e. } k\text{th-smallest value,}$$

$$\tilde{\sigma}_n^2 = n^{-1} \sum_i (X_i - \mu)^2 \text{ (assuming that } \mu \text{ is known), } S_n^2 = (n-1)^{-1} \sum_i (X_i - \bar{X}_n)^2, (\bar{X}_n, S_n^2),$$

$$\arg \min_a \sum_i |X_i - a|. \quad \square$$

To evaluate the performance of an estimation procedure, one examines the theoretical properties of the estimator on which it is based. Since an estimator is a random variable, it possesses a probability distribution, referred to as the *sampling distribution*. The analytical characteristics of this distribution are fundamental for determining the adequacy of an estimator in a given inferential context and for identifying those that should be avoided due to undesirable statistical properties.

Example 2.2 (Uniform model).

Consider the uniform density $f(x; \theta) = \frac{1}{\theta} I(0 \leq x \leq \theta); \theta > 0$.

As an estimator of θ , based on a sample X_1, \dots, X_n , one can, for example, consider one of

the following :

$$\begin{aligned}\hat{\theta}_1 &= X_{(n)}, & \hat{\theta}_2 &= \frac{n+1}{n}X_{(n)} \\ \hat{\theta}_3 &= X_{(1)} + X_{(n)}, & \hat{\theta}_4 &= 2\bar{X}_n \\ \hat{\theta}_5 &= 2\hat{q}_{0.5},\end{aligned}$$

where $\hat{q}_{0.5}$ is the sample median, i.e.

$$\hat{q}_{0.5} = \begin{cases} X_{(k+1)} & \text{if } n = 2k + 1 \text{ is odd,} \\ \frac{X_{(k)} + X_{(k+1)}}{2} & \text{if } n = 2k \text{ is even. } \square \end{cases}$$

This example suggests questions like:

- If many estimators are available—as is always the case—how can we compare them ?
- Are there any general and reliable methods for constructing estimators ?
- Given a statistical model, does a “best estimator” exist, and if so, how can it be identified ?

These questions (and many others of the same nature) will be the subject of our next readings.

2.2 Risk and loss function

It seems reasonable that we want an estimate $\hat{\theta}$ which generally comes quite close to the true value of θ , and dislike an estimate $\hat{\theta}$ which generally misses the true value of θ by a large amount. The question is how to make this precise and quantifiable?

We quantify the idea of $\hat{\theta}$ being close to θ , by measuring the **risk**, that is the *average distance*, between these two quantities. The distance is measured using what is called a **loss function**.

Examples of loss functions include:

- squared error loss: $L(\hat{\theta}, \theta) = (\hat{\theta} - \theta)^2$
- absolute error loss: $L(\hat{\theta}, \theta) = |\hat{\theta} - \theta|$
- absolute relative loss: $L(\hat{\theta}, \theta) = |\hat{\theta}/\theta - 1|$

Once the loss function is chosen, we calculate the risk as follows

$$E_{\theta}(L(\hat{\theta}, \theta)) = \int L(\hat{\theta}(\mathbf{x}), \theta) f_n(\mathbf{x}, \theta) d\mathbf{x},$$

where E_{θ} means that *the expectation is taken under the assumption that θ is the true parameter*; that is, the pd of X is $f_n(\mathbf{x}, \theta)$.

2.3 MSE, bias, variance, and relative efficiency

When the squared-error loss is used, the risk associated with an estimator is precisely its *mean squared error*:

$$\begin{aligned} \text{MSE}_\theta(\hat{\theta}) &:= E_\theta[(\hat{\theta} - \theta)^2] \\ &= E_\theta[(\hat{\theta} - E_\theta(\hat{\theta})) + (E_\theta(\hat{\theta}) - \theta)]^2 \\ &= \text{Bias}_\theta^2(\hat{\theta}) + \text{Var}_\theta(\hat{\theta}), \end{aligned}$$

where $\text{Bias}_\theta(\hat{\theta}) = E_\theta(\hat{\theta}) - \theta$ is the bias of the estimator $\hat{\theta}$.

A large bias indicates low accuracy ($\hat{\theta}$ lies far from θ , i.e. some systematic error), while a large variance indicates low precision (too much fluctuation). If the bias of $\hat{\theta}$ is always zero, i.e. $\text{Bias}_\theta(\hat{\theta}) = 0 \forall \theta \in \Theta$, then $\hat{\theta}$ is called *unbiased*. This means that *on average* the estimator will yield the true value of the unknown parameter (whatever the true value is). *In this case, the MSE reduces to the variance.*

Attention. In the following, in order to ease the notation, if no confusion is possible, we drop the index θ and write E , Var and MSE instead of E_θ , Var_θ and MSE_θ , respectively.

The choice of an estimator is very often restricted to the class of unbiased estimators. But there are cases where a small bias is accepted, in particular if the bias converges to zero when the sample size tends to infinity. Moreover, there are cases where no unbiased estimator exists.

Example: $X_1, \dots, X_n \sim \text{Bernoulli}(p)$, and $\theta = \frac{1}{p}$.

Example 2.3 (Unbiased does not necessarily mean a good estimator).

Let X_i , $i = 1, \dots, n$, be iid rv with $E(X_i) = \mu$ and $\text{Var}(X_i) = \sigma^2$. Then X_1 , \bar{X}_n and $\frac{X_1 + \bar{X}_n}{2}$ are unbiased estimators of μ . Which one should we use?

It is clear that all three are unbiased. So to compare these estimators, we have to compare their variances (i.e. their MSE). It is easy to see that

$$\text{Var}(X_1) = \sigma^2, \quad \text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}, \quad \text{and} \quad \text{Var}\left(\frac{X_1 + \bar{X}_n}{2}\right) = \frac{1}{4} \left(1 + \frac{3}{n}\right) \sigma^2.$$

To check the last equality, observe that

$$\text{Var}(X_1 + \bar{X}_n) = \text{Var}\left(\frac{n+1}{n}X_1 + \frac{1}{n}\sum_{i=2}^n X_i\right) = \frac{(n+1)^2}{n^2}\sigma^2 + \frac{n-1}{n^2}\sigma^2 = \left(1 + \frac{3}{n}\right)\sigma^2.$$

$\rightarrow \bar{X}_n$ is better than the other two. \square

Bias and transformation

If $\hat{\theta}$ is an unbiased estimator of θ , then for any constants a and b , the transformed estimator $a + b\hat{\theta}$ remains unbiased for $a + b\theta$; in other words, *unbiasedness is preserved under linear*

transformations. This property allows us to correct certain biased estimators: if $E(\hat{\theta}) = a\theta + b$, then the adjusted estimator $(\hat{\theta} - b)/a$ is unbiased for θ .

In contrast, nonlinear transformations do not generally preserve unbiasedness: even when $\hat{\theta}$ is unbiased for θ , a nonlinear function $g(\hat{\theta})$ will typically fail to be unbiased for $g(\theta)$. For instance, although the sample mean \bar{X}_n is an unbiased estimator of μ , Jensen's inequality implies that, provided \bar{X}_n is not almost surely constant, $(E \bar{X}_n)^2 = \mu^2 < E(\bar{X}_n^2)$; consequently, \bar{X}_n^2 is a biased estimator of μ^2 . \square

Example 2.4. Let $X_i, i = 1, \dots, n$, be iid rv with $E(X_i) = \mu$ and $Var(X_i) = \sigma^2$. Let's find an unbiased estimator for μ^2 . To do so, observe that

$$E(\bar{X}_n^2) = (E(\bar{X}_n))^2 + Var(\bar{X}_n) = \mu^2 + \frac{\sigma^2}{n}.$$

This means that $Bias(\bar{X}_n^2) = \sigma^2/n$. This also implies that an unbiased estimator of μ^2 is given by

$$\bar{X}_n^2 - \frac{S_n^2}{n},$$

provided that S_n^2 is an unbiased estimator of σ^2 (see the next example). \square

Example 2.5 (The sample variance). Let $X_i, i = 1, \dots, n$, be iid rv with $E(X_i) = \mu$ and $Var(X_i) = \sigma^2$. A natural estimator of this latter is given by $\hat{\sigma}_n^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$. We have that

$$E(\hat{\sigma}_n^2) = Var(X_1 - \bar{X}_n) = Var\left(\frac{n-1}{n}X_1 - \frac{1}{n}\sum_{i=2}^n X_i\right) = \frac{n-1}{n}\sigma^2$$

With the correction factor $\frac{n}{n-1}$ we obtain an unbiased estimator, namely the well-known empirical variance (or sample variance):

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2. \square$$

The MSE is a global measure of estimation quality that combines variance and squared bias, two components of different nature. *Comparing the MSE of two estimators is meaningful only when their biases are comparable* (i.e., both are unbiased or have roughly the same bias); otherwise, differences in MSE become difficult to interpret. Moreover, comparing MSEs is valid only when the two estimators address the same estimation problem—estimating the same parameter from the same data, under the same model and with the same sample size.

A standard tool for such comparison is the *relative efficiency* :

$$RE(\hat{\theta}_1, \hat{\theta}_2) = \frac{MSE(\hat{\theta}_2)}{MSE(\hat{\theta}_1)}.$$

This quantity expresses the performance of $\hat{\theta}_2$ relative to $\hat{\theta}_1$. When both estimators are unbiased, the relative efficiency reduces to

$$RE(\hat{\theta}_1, \hat{\theta}_2) = \frac{Var(\hat{\theta}_2)}{Var(\hat{\theta}_1)}.$$

If this quantity is less than one $\forall \theta \in \Theta$, then $\hat{\theta}_1$ has uniformly a larger variance than $\hat{\theta}_2$, and the latter is said to be *more efficient* than the former.

Example 2.6.

Let X_1, \dots, X_n , be an iid sample from $Unif[0, \theta]$. Let $\hat{\theta}_1 = 2\bar{X}_n$ and $\hat{\theta}_2 = \frac{n+1}{n}X_{(n)}$, two estimators of θ .

Recall that an uniform distribution in $[a, b]$ is characterized by its cdf

$$F(x) = \begin{cases} 0, & \text{if } x < a \\ \frac{x-a}{b-a}, & \text{if } a \leq x \leq b \\ 1, & \text{if } x > b. \end{cases}$$

Its mean and variance are given by

$$E(X) = \frac{a+b}{2} \text{ and } Var(X) = \frac{(b-a)^2}{12}.$$

From this, it follows that $\hat{\theta}_1$ is unbiased and $Var(\hat{\theta}_1) = \frac{\theta^2}{3n}$.

To find the expectation and variance of $\hat{\theta}_2$ first observe that

$$P(X_{(n)} \leq x) = (P(X_1 \leq x))^n.$$

So the cdf of $X_{(n)}$ is given by

$$F_{X_{(n)}}(x) = \begin{cases} 0, & \text{if } x < 0 \\ (x/\theta)^n, & \text{if } 0 \leq x \leq \theta \\ 1, & \text{if } x > \theta. \end{cases}$$

Hence, the pd of $X_{(n)}$ is given by $f_{X_{(n)}}(x) = n \frac{x^{n-1}}{\theta^n} I(0 \leq x \leq \theta)$. It follows that $E(X_{(n)}) = \frac{n}{n+1}\theta$ and $E(X_{(n)}^2) = \frac{n}{n+2}\theta^2$. Therefore, $\hat{\theta}_2$ is unbiased and its variance is $Var(\hat{\theta}_2) = \frac{\theta^2}{n(n+2)}$. Finally, the relative efficiency is

$$\frac{Var(\hat{\theta}_2)}{Var(\hat{\theta}_1)} = \frac{3}{n+2}.$$

Indicating that for $n > 1$, $\hat{\theta}_2$ is more efficient than $\hat{\theta}_1$. \square

Note that, in the example above, the relative efficiency does not depend on the parameter θ . This feature is specific to that setting and does not hold in general; a point to which we will return in the subsequent section (after the next).

Remark. To compare estimators that are not defined on the same scale, the Mean Absolute Relative Error (MARE)

$$MARE(\hat{\theta}) = E\left(\left|\frac{\hat{\theta} - \theta}{\theta}\right|\right)$$

is particularly useful because it expresses the average absolute relative deviations making the metric **unit-free** and comparable across heterogeneous settings. Its interpretation is straightforward: for example, a MARE of 0.10 indicates that, *on average*, the estimator deviates from the target by 10%, with lower values reflecting better performance and higher values indicating proportionally larger errors. \square

2.4 Simulation-Based Evaluation of Estimator Performance

A practical way to assess the performance of an estimator is to approximate its bias, variance, MSE, etc., through simulation. The idea is to repeatedly generate data from a model with a known parameter value θ , compute the estimator $\hat{\theta}$ for each simulated dataset, and then summarize the resulting collection of estimates. If $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(N)}$ denote the estimates obtained from N independent simulations, the empirical bias, empirical variance, empirical MSE, and empirical MARE are given, respectively, by

$$\begin{aligned}\widehat{\text{Bias}} &= \bar{\hat{\theta}} - \theta \\ \widehat{\text{Var}} &= \frac{1}{N} \sum_{s=1}^N \left(\hat{\theta}^{(s)} - \bar{\hat{\theta}} \right)^2 \\ \widehat{\text{MSE}} &:= \frac{1}{N} \sum_{s=1}^N \left(\hat{\theta}^{(s)} - \theta \right)^2 = \widehat{\text{Bias}}^2 + \widehat{\text{Var}} \\ \widehat{\text{MARE}} &:= \frac{1}{N} \sum_{s=1}^N \left| \frac{\hat{\theta}^{(s)}}{\theta} - 1 \right|\end{aligned}$$

where $\bar{\hat{\theta}} = \frac{1}{N} \sum_{s=1}^N \hat{\theta}^{(s)}$ is the average of the simulated estimates.

The number of simulations N must be chosen sufficiently large to ensure that these empirical estimates are stable and accurately reflect their theoretical counterparts.

To obtain a reliable picture of the estimator's performance, and since these measures typically depend on the value of θ , one must repeat the entire simulation procedure across a broad range of parameter values. This can become computationally demanding, especially when the estimator itself is costly to compute.

The R code below demonstrates the simulation procedure outlined above, applied to the setting of Example 2.6.

```
mse <- function(estimates, theta) {
  mean <- mean(estimates)
  bias <- mean - theta          # empirical bias
  var <- mean((estimates - mean)^2) # empirical variance
  mse <- bias^2 + var           # empirical mean squared error
  mare <- mean(abs(estimates / theta - 1)) # empirical mean absolute relative error
  c(bias = bias, var = var, mse = mse, mare = mare)
}

estSim <- function(dataFun, estFun, N = 10000) {
  estimates <- replicate(N, dataFun() |> estFun())
  return(estimates)
}
```

- Performances of $\hat{\theta}_1$:

```
n <- 10; theta <- 3
estSim(dataFun = \() runif(n, min = 0, max = theta),
       estFun = \ (data) 2 * mean(data)) |> mse(theta)
```

bias	var	mse	mare
0.000695	0.301129	0.301130	0.146706

- Performances of $\hat{\theta}_2$:

```
n <- 10; theta <- 3
estSim(dataFun = \() runif(n, min = 0, max = theta),
       estFun = \ (data) ((n + 1) / n) * max(data)) |> mse(theta)
```

bias	var	mse	mare
-0.0049	0.0772	0.0773	0.0706

These results closely match the theory which tells us that, for $n = 10$, the MSE of $\hat{\theta}_1$ is $3^2/(3 \times 10) = 0.3$ and that the MSE of $\hat{\theta}_2$ is $3^2/(10 \times (10 + 2)) = 0.075$.

2.5 The best unbiased estimator (MVUE)

It is natural to prefer estimators that have a small MSE.

If an estimator $\hat{\delta}$ has a larger MSE than another estimator $\hat{\theta}$ for every possible value of the parameter – that is, if

$$MSE_{\theta}(\hat{\theta}) \leq MSE_{\theta}(\hat{\delta}), \forall \theta \in \Theta,$$

then $\hat{\delta}$ is called *inadmissible*.

One might hope to find an estimator that has the smallest possible MSE for every value of θ . However, this turns out to be impossible. For an estimator $\hat{\theta}$ to be the “best” for every θ , it would need to satisfy

$$E_{\theta}[(\hat{\theta} - \theta)^2] = 0 \quad \forall \theta \in \Theta.$$

But the only way for this expectation to be zero is if $\hat{\theta} = \theta$ with probability 1. This would mean the estimator always returns the true parameter exactly, which is impossible in any non-trivial statistical model.

Typically, two estimators $\hat{\theta}$ and $\hat{\delta}$ are not uniformly comparable. In general, one can find parameter values $\theta_1, \theta_2 \in \Theta$ such that

$$\text{MSE}_{\theta_1}(\hat{\theta}) < \text{MSE}_{\theta_1}(\hat{\delta}) \quad \text{and} \quad \text{MSE}_{\theta_2}(\hat{\theta}) > \text{MSE}_{\theta_2}(\hat{\delta}).$$

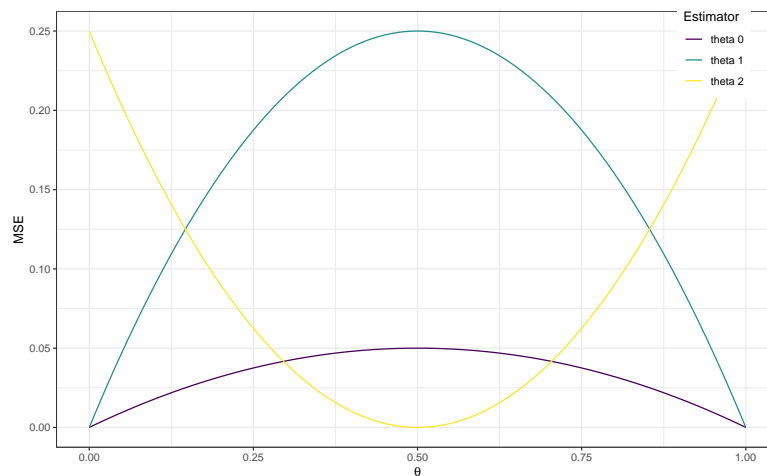
This means that each estimator performs better for some values of θ and worse for others, so neither dominates the other uniformly.

Example 2.7.

Suppose an iid sample X_1, \dots, X_n from a Bernoulli distribution with an unknown parameter θ , $0 \leq \theta \leq 1$. Let $\hat{\theta}_0 = \bar{X}_n$, $\hat{\theta}_1 = X_1$ and $\hat{\theta}_2 = 1/2$. It is easy to see that

$$\text{MSE}(\hat{\theta}_0) = \frac{\theta(1-\theta)}{n}, \quad \text{MSE}(\hat{\theta}_1) = \theta(1-\theta), \quad \text{MSE}(\hat{\theta}_2) = (\theta - 1/2)^2.$$

These three MSE (quadratic risk) functions, of θ , are plotted below (for $n = 5$).



$\hat{\theta}_1$ is inadmissible as it is uniformly less efficient than $\hat{\theta}_0$. The estimators $\hat{\theta}_0$ and $\hat{\theta}_2$ are not uniformly comparable. Near $\theta = 1/2$, $\hat{\theta}_2$ is the best, and away from $\theta = 1/2$, $\hat{\theta}_0$ is the best. \square

In view of the fact that there is no estimator that has the smallest possible MSE for every value of the parameter, statisticians adopt other strategies to choose good estimators. One such strategy is to restrict attention to the *class of unbiased estimators* and then search for the best estimator within this restricted group. By focusing on unbiased estimators, it becomes somewhat easier to compare their performance and identify the one that performs best.

Definition 2.2 (MVUE). An unbiased estimator $\hat{\theta}$ of $\theta \in \Theta$ is the uniformly *Minimum Variance Unbiased Estimator* (MVUE) if for any other unbiased estimator $\hat{\delta}$

$$\text{Var}_{\theta}(\hat{\theta}) \leq \text{Var}_{\theta}(\hat{\delta}), \forall \theta \in \Theta.$$

In other words, the MVUE is the best (*most efficient*) *unbiased* estimator that can be found.

Facts to know

- The MVUE may not exist (in some problems, even an unbiased estimator does not exist). However, *when an MVUE does exist, it is unique*.
- In terms of MSE, the MVUE is not necessarily the best estimator. There may be *biased estimators* whose MSE is smaller than that of the MVUE. In fact, a small increase in bias can sometimes lead to a large reduction in variance, resulting in a lower overall MSE.

□

The question now is how to find the MVUE (when it exists). To address this, several techniques have been developed in the statistical literature. One important approach relies on a variance inequality known as the *Cramér–Rao bound*. Before introducing this method, we first need to define the concept of *Fisher information*, which plays a central role in statistical inference.

Chapter 3

Fisher information and Cramer-Rao bound

3.1 Score and Fisher information

Let X be a random variable (or a random vector) with pdf $f(x; \theta)$ indexed by an unknown parameter $\theta \in \Theta \subset \mathbb{R}$. The question of interest here is: how much information about θ can be obtained from observing X ?

To explore this, assume that f is differentiable with respect to θ , and define the **score function** associated with f as

$$S(\theta, x) := \partial_{\theta} \log f(x; \theta) = \frac{\partial_{\theta} f(x; \theta)}{f(x; \theta)}.$$

Observe that, by definition, for any fixed $\theta_0 \in \Theta$,

$$S(\theta_0, x) = \lim_{\epsilon \rightarrow 0} \frac{\frac{1}{\epsilon} [f(x; \theta_0 + \epsilon) - f(x; \theta_0)]}{f(x; \theta_0)}.$$

Thus, the score $S(\theta_0, x)$ can be interpreted as the *relative instantaneous rate of change* of the function $\theta \mapsto f(x; \theta)$ at the point θ_0 . In particular, if f changes rapidly in a neighborhood of θ_0 , the score will have a large absolute value; conversely, if f is relatively flat, the score will be small. In other words, a large value of $|S(\theta_0, x)|$ indicates that the observation x is highly informative for distinguishing θ_0 from nearby parameter values.

A remarkable property of the score function is given in the following proposition.

Proposition 3.1. *Suppose that*

- (I) *the support of f , i.e. the set $\{x : f(x; \theta) > 0\}$, does not depend on θ , and*
- (II) *the operations of integration (or summation) and differentiation by θ can be interchanged in $\int f(x; \theta) dx$. Thus, $\partial_{\theta} \int f(x; \theta) dx = \int \partial_{\theta} f(x; \theta) dx$ ¹.*

¹For more details on this condition, see the discussion of the Leibniz integral rule [here](#)

Then

$$E_{\theta} S(\theta, X) = 0, \forall \theta \in \Theta.$$

The expected score is zero because, when the parameter is set to its correct value, a tiny change in that parameter makes the density increase for some values of x and decreases for others. These opposing changes must perfectly balance out, because the total probability — the area under the density curve — must always sum to 1. Consequently, the average effect of that small parameter shift, weighted by the data probabilities, is zero, which is precisely what the expected score expresses.

Attention. From now on, unless explicitly stated otherwise or unless it is evident that they fail to hold, we will assume that conditions (I) and (II) are satisfied. Notably, these conditions are satisfied by the (regular) exponential family.

Taking the square (of S) and averaging we obtain $I_X(\theta)$:

$$I_X(\theta) := E_{\theta} [S^2(\theta, X)] = \text{Var}_{\theta} [S(\theta, X)]$$

which is known as the (expected) **Fisher information (FI)** that X contains about θ , or the FI for θ based on X .

FI attempts to *quantify the average sensitivity of the random variable X to the value of the parameter θ* . If small changes in θ result in large changes in the values of X , then observing the latter can tell us a lot about θ . In this case the FI would be quite large. In other words, FI attempts to quantify how easy one can guess the θ that produced the observed X .

Attention. Note that the subscript X in $I_X(\theta)$ serves solely as a label indicating that the FI is computed with respect to the random variable X ; it does **not** imply that the Fisher information itself is random. In contrast, the score $S(\theta, X)$ is a random variable, since it is defined as a function of X . In what follows, whenever no ambiguity about the underlying random variable arises, we will simply write $I(\theta)$ instead of $I_X(\theta)$.

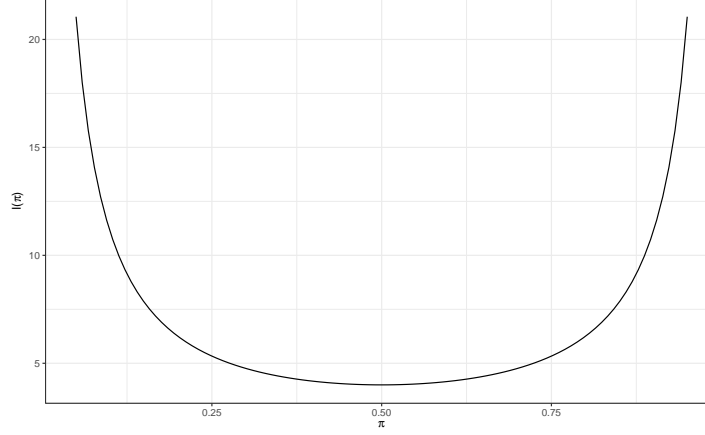
Example 3.1 (Calculating Fisher Information 1).

- Bernoulli distribution: $X \sim \text{Be}(\pi)$ with $\pi \in (0, 1)$. $f(x; \pi) = \pi^x(1 - \pi)^{1-x}$, $x = 0, 1$.

$$\begin{aligned} S(\pi, X) &= \partial_{\pi} \{X \log(\pi) + (1 - X) \log(1 - \pi)\} = \frac{X - \pi}{\pi(1 - \pi)}. \\ \implies I(\pi) &= E_{\pi} [S^2(\pi, X)] = \frac{1}{\pi(1 - \pi)}. \end{aligned}$$

Here, the Fisher information turns out to be the reciprocal of the variance of a Bernoulli. This is not unusual. In fact, as we will see later, the Fisher information is typically inversely proportional to the variance. Intuitively, when the data exhibit greater variability, it becomes harder to infer the true value of the parameter, and the Fisher information decreases. Conversely, when the variance is small, the data are more

concentrated and provide more precise information about the parameter.



- Normal distribution: $X \sim N(\mu, \sigma^2)$ with $\mu \in \mathbb{R}$ and $\sigma^2 > 0$. $f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$.

– For μ :

$$S(\mu, X) = \partial_\mu \left\{ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (X - \mu)^2 \right\} = \frac{X - \mu}{\sigma^2}.$$

$$\implies I(\mu) = E[S^2(\mu, X)] = \frac{1}{\sigma^2}.$$

Note that, in this case, the FI about μ does not depend on μ but only on σ^2 . It decreases with σ^2 .

– For σ^2 :

$$S(\sigma^2, X) = \partial_{\sigma^2} \left\{ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (X - \mu)^2 \right\} = -\frac{1}{2\sigma^2} + \frac{1}{2} \frac{(X - \mu)^2}{\sigma^4}.$$

$$\implies I(\sigma^2) = E[S^2(\sigma^2, X)] = \frac{1}{4} E \left(\frac{(X - \mu)^2}{\sigma^4} - \frac{1}{\sigma^2} \right)^2$$

$$= \frac{1}{4} \left(\frac{1}{\sigma^8} E(X - \mu)^4 + \frac{1}{\sigma^4} - \frac{2}{\sigma^6} E(X - \mu)^2 \right) = \frac{1}{2\sigma^4},$$

where we have used the fact that, for $X \sim N(\mu, \sigma^2)$, $E(X - \mu)^4 = 3\sigma^4$. \square

The following proposition offers a useful alternative for calculating the FI. Rather than requiring the second moment of the score, it expresses $I(\theta)$ as the negative expected derivative of the score function itself.

Proposition 3.2. Let $X \sim f(x; \theta)$. Assume that $f(x; \theta)$ is twice differentiable with respect to θ , and that double differentiation and integration (or summation) can be interchanged, so that $\partial_\theta^2 \int f(x; \theta) dx = \int \partial_\theta^2 f(x; \theta) dx$. Then,

$$I(\theta) = -E_\theta[\partial_\theta S(\theta, X)].$$

Note that the equality above is equivalent to $I(\theta) = -E_\theta[\partial_\theta^2 \log f(X; \theta)]$.

To see why this proposition is true, consider the following calculation where we have simplified the notation by writing S and f instead of $S(\theta, X)$ and $f(X; \theta)$, respectively.

$$\partial_\theta S = \partial_\theta \frac{\partial_\theta f}{f} = \frac{\partial_\theta^2 f \times f - (\partial_\theta f)^2}{f^2} = \frac{\partial_\theta^2 f}{f} - S^2.$$

But $E_\theta \left(\frac{\partial_\theta^2 f(X; \theta)}{f(X; \theta)} \right) = \int \partial_\theta^2 f(x; \theta) dx = \partial_\theta^2 \int f(x; \theta) dx = 0$. Thus, $E_\theta(\partial_\theta S) = -E_\theta S^2$.

Example 3.2 (Calculating Fisher Information 2).

- Bernoulli distribution: $X \sim Be(\pi)$ with $\pi \in (0, 1)$.

$$\begin{aligned} \partial_\pi S(\pi, X) &= \frac{1}{\pi^2(1-\pi)^2} \left(-\pi(1-\pi) - (\pi(1-\pi))'(X-\pi) \right). \\ \implies I(\pi) &= -E[\partial_\pi S(\pi, X)] = \frac{1}{\pi(1-\pi)}. \end{aligned}$$

- Normal distribution: $X \sim N(\mu, \sigma^2)$ with $\mu \in \mathbb{R}$ and $\sigma^2 > 0$.

$$\begin{aligned} \partial_\mu S(\mu, X) &= -\frac{1}{\sigma^2}. & \partial_{\sigma^2} S(\sigma^2, X) &= \frac{1}{2\sigma^4} - \frac{(X-\mu)^2}{\sigma^6}. \\ \implies I(\mu) &= -E[\partial_\mu S(\mu, X)] = \frac{1}{\sigma^2}. & \implies I(\sigma^2) &= -E[\partial_{\sigma^2} S(\sigma^2, X)] = \frac{1}{2\sigma^4}. \end{aligned}$$

3.2 FI contained in a statistic

The above definitions of the score and FI can be directly applied to any statistics. In fact, let $T \equiv T(\mathbf{X}) = T(X_1, \dots, X_n)$ be a statistic whose pd is given by $h_n(t; \theta)$. The score associated with h_n and its corresponding FI are given by

$$\begin{aligned} S(\theta, T) &= \partial_\theta \log h_n(T; \theta) \\ I_T(\theta) &= E_\theta [S^2(\theta, T)] \end{aligned}$$

$I_T(\theta)$ is the (Fisher) information about θ that we can extract from T .

Assuming the interchangeability of integration and differentiation twice, this FI can also be expressed as

$$I_T(\theta) = -E_\theta [\partial_\theta S(\theta, T)].$$

Example 3.3.

- Let $X_i, i = 1, \dots, n$, be an iid sample from $Be(\pi)$. Let us define the statistic $T = \sum_{i=1}^n X_i$. Since $T \sim Bin(n, \pi)$, the pd of T is given by $h_n(t; \pi) = P_\pi(T = t) = C_n^t \pi^t (1-\pi)^{n-t}$. It follows that

$$I_T(\pi) = E[S^2(\pi, T)] = E \left[\frac{T - n\pi}{\pi(1-\pi)} \right]^2 = \frac{n}{\pi(1-\pi)}.$$

- Let $X_i, i = 1, \dots, n$, be an iid sample from a Normal distribution $N(\mu, \sigma^2)$. Since $\bar{X}_n \sim N(\mu, \sigma^2/n)$, it follows that

$$I_{\bar{X}_n}(\mu) = \frac{n}{\sigma^2}. \square$$

An important fact about FI is its **additivity**. Let T_1 and T_2 be two statistics with pd h_1 and h_2 , and with FI $I_{T_1}(\theta)$ and $I_{T_2}(\theta)$, respectively. If T_1 and T_2 are *independent*, i.e. if $h(t_1, t_2; \theta) = h_1(t_1; \theta)h_2(t_2; \theta)$, $\forall t$, with h being the joint pd of (T_1, T_2) , then

$$\begin{aligned} I_{(T_1, T_2)}(\theta) &= E(\partial_\theta \log h(T_1, T_2; \theta))^2 = E(\partial_\theta \log h_1(T_1; \theta) + \partial_\theta \log h_2(T_2; \theta))^2 \\ &= E(\partial_\theta \log h_1(T_1; \theta))^2 + E(\partial_\theta \log h_2(T_2; \theta))^2 + 2E(\partial_\theta \log h_1(T_1; \theta) \partial_\theta \log h_2(T_2; \theta)) \\ &= I_{T_1}(\theta) + I_{T_2}(\theta). \end{aligned}$$

As consequence we have the following result.

Proposition 3.3. Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample from a distribution with density $f(x; \theta)$, and let $f_n(\mathbf{x}; \theta) = \prod_{i=1}^n f(x_i; \theta)$ denote the joint density of \mathbf{X} . Define the individual and joint Score and FI functions by

$$\begin{aligned} S(\theta, X_i) &= \partial_\theta \log f(X_i; \theta) & I_{X_i}(\theta) &= E[S^2(\theta, X_i)] \\ S(\theta, \mathbf{X}) &= \partial_\theta \log f_n(\mathbf{X}; \theta) & I_{\mathbf{X}}(\theta) &= E[S^2(\theta, \mathbf{X})] \end{aligned}$$

Then, $S(\theta, \mathbf{X}) = \sum_{i=1}^n S(\theta, X_i)$, and $I_{\mathbf{X}}(\theta) = \sum_{i=1}^n I_{X_i}(\theta) = n I_{X_1}(\theta)$.

Attention. From now on, if no confusion is possible, we use the notation I_n to denote the joint FI, i.e. $I_n = I_{\mathbf{X}}$, and write I for the individual FI, i.e. $I = I_{X_i} = I_{X_1}$. With this convention, the equality above becomes $I_n(\theta) = nI(\theta)$.

For any statistic $\mathbf{T} = (T_1(\mathbf{X}), T_2(\mathbf{X}), \dots, T_d(\mathbf{X}))$, $d \geq 1$, it can be shown that

$$0 \leq I_{\mathbf{T}}(\theta) \leq I_n(\theta).$$

This inequality expresses a fundamental principle: no statistic computed from the sample can contain more information about θ than the sample itself.

3.3 Sufficient statistic

One may naturally ask under what circumstances a given statistic can attain the maximal possible Fisher information I_n . To answer this question, we must introduce the notion of sufficiency.

The definition below applies to both the single-parameter case and the multi-parameter case, where θ may represent a vector of several parameters.

Definition 3.1 (Sufficiency). Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random sample with joint pd $f_n(\mathbf{x}; \boldsymbol{\theta})$. A statistic T is sufficient for $\boldsymbol{\theta}$ if the conditional distribution of \mathbf{X} given T does not depend on $\boldsymbol{\theta}$.

Put another way, given a sufficient statistic T for $\boldsymbol{\theta}$, the sample \mathbf{X} provides no additional information about $\boldsymbol{\theta}$. Sufficient statistics are especially valuable when their dimension d is much smaller than the sample size n , because they reduce a large dataset to one or a few numbers while preserving all the information about the parameter.

Example 3.4. Let $X_i, i = 1, \dots, n$, be an iid sample from $Be(\pi)$. Define the statistic $T = \sum_{i=1}^n X_i$. We have that

$$\begin{aligned} P(X_1 = x_1, \dots, X_n = x_n | T = t) &= \frac{P(X_1 = x_1, \dots, X_n = x_n)}{P(T = t)} I(\sum_{i=1}^n x_i = t) \\ &= \frac{\prod_i \pi^{x_i} (1 - \pi)^{1-x_i}}{C_n^t \pi^t (1 - \pi)^{n-t}} I(\sum_{i=1}^n x_i = t) = \frac{I(\sum_{i=1}^n x_i = t)}{C_n^t}. \end{aligned}$$

This expression does **not** depend on π . Therefore, by the definition of sufficiency, $\sum_{i=1}^n X_i$ is sufficient for π . \square

Sufficiency and transformation

- If T is sufficient for $\boldsymbol{\theta}$, then T is also sufficient for any transformation $k(\boldsymbol{\theta})$, regardless of the form of k .
- If T is sufficient for $\boldsymbol{\theta}$, then any statistic $U = k(T)$ remains sufficient for $\boldsymbol{\theta}$, provided that the function k is bijective (or at least one-to-one). In such cases, U contains exactly the same information as T , just expressed on a different scale. For example, in the $Ber(\pi)$ model, since $\sum_{i=1}^n X_i$ is sufficient for π , the sample mean $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ is also sufficient for π . \square

The following result, known as the **Factorization Theorem**, makes it very easy to identify sufficient statistics (this theorem applies to both the single-parameter and multi-parameter cases).

Theorem 3.1 (Factorization Theorem). Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random sample with joint pd $f_n(\mathbf{x}; \boldsymbol{\theta})$. A statistic $T(\mathbf{X})$ sufficient for $\boldsymbol{\theta}$ if and only if there exist nonnegative functions φ and h such that f_n can be factorized as

$$f_n(\mathbf{x}; \boldsymbol{\theta}) = \varphi(T(\mathbf{x}); \boldsymbol{\theta}) h(\mathbf{x}), \forall \mathbf{x}, \boldsymbol{\theta}.$$

where h does not depend on $\boldsymbol{\theta}$, and φ depends on \mathbf{x} only through $T(\mathbf{x})$.

Example 3.5.

- Let $X_i, i = 1, \dots, n$, be an iid sample from $Be(\pi)$. The joint pd of (X_1, \dots, X_n) is

$$\prod_{i=1}^n \pi^{x_i} (1 - \pi)^{1-x_i} = \pi^{\sum_{i=1}^n x_i} (1 - \pi)^{n - \sum_{i=1}^n x_i}.$$

It follows that $\sum_{i=1}^n X_i$ is sufficient for π .

- Let $X_i, i = 1, \dots, n$, be an iid sample from $N(\mu, \sigma^2)$. The joint pd of (X_1, \dots, X_n) is

$$\begin{aligned} f(x_1, \dots, x_n; \mu, \sigma^2) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \left(\sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + n\mu^2\right)\right) \end{aligned}$$

It follows that $(\sum_{i=1}^n X_i, \sum_{i=1}^n X_i^2)$ is sufficient for (μ, σ^2) . This implies that (\bar{X}, S^2) is sufficient for θ . \square

It can be shown that, for both the single-parameter and multi-parameter cases,

$$T \text{ is sufficient for } \theta \implies I_T(\theta) = I_n(\theta).$$

Under some regularity conditions, the reverse is also true². In the multi-parameter case, I (in bold) denotes the Fisher information matrix; see Section 3.8 for more details.

Example 3.6. We have seen that for the Bernoulli model, the statistic $T = \sum_{i=1}^n X_i$ is sufficient for π . Consequently, the Fisher information contained in T is equal to the Fisher information in the full sample, that is, $I_{\sum_{i=1}^n X_i}(\pi) = I_n(\pi)$. This equality can be verified directly by the calculations carried out in Example 3.1 and Example 3.3. \square

As a direct consequence of the Factorization theorem, we have the following result that allows one to obtain sufficient statistics when data come from an exponential family.

Proposition 3.4. If $X = (X_1, \dots, X_n)$ is an iid sample from a J -parameter exponential family ($J \geq 1$) with pd

$$h(x) \exp(\eta^t(\theta)T(x) - B(\theta)),$$

then the statistic $\sum_{i=1}^n T(X_i)$ is sufficient for θ .

Example 3.7. Let X_1, \dots, X_n be an iid sample from $f(x; \theta) = \theta x^{\theta-1}$, where $x \in (0, 1)$ and $\theta > 0$. Let's show that $\prod_{i=1}^n X_i$ is sufficient for θ . To see this, we can write $f(x; \theta)$ as

$$f(x; \theta) = I(0 < x < 1) x^{-1} \exp(\theta \log x + \log \theta),$$

which is an exponential family with natural statistic $T = \log X$. It follows that $L = \sum_i \log(X_i)$

²See Pollard, D. (2013). A note on insufficiency and the preservation of Fisher information. DOI: 10.1214/12-IMSCOLL919

is sufficient for θ . And since $\prod_{i=1}^n X_i = \exp(L)$ is a bijective function of L , it is also sufficient for θ . \square

3.4 FI and re-parametrization

We have already seen that statistical models can be parameterized in different ways. It is important to realize that FI depends on the chosen parameterization.

Proposition 3.5 (FI re-parametrization). *Let $\eta : \theta \mapsto \eta(\theta)$ be a differentiable, one-to-one transformation of θ .*

Define the reparameterized pd

$$f^*(x; \eta) = f(x; \theta), \forall x, \text{ where } \eta = \eta(\theta).$$

Let $I(\theta) = E[\partial_\theta \log f(X; \theta)]^2$ be the FI about θ (when the parameterization in θ is used) and $I(\eta) = E[\partial_\eta \log f^(X; \eta)]^2$ be the FI about η (when the parameterization in η is used). Then, $I(\theta)$ and $I(\eta)$ are related by $I(\theta) = I(\eta) (\eta'(\theta))^2$.*

The proof of this proposition is straightforward and is a direct consequence of the chain rule :

$$I(\theta) = E[\partial_\theta \log f(X; \theta)]^2 = E[\partial_\theta \log f^*(X; \eta)]^2 = E[\partial_\theta \eta(\theta) \partial_\eta \log f^*(X; \eta)]^2 = (\eta'(\theta))^2 I(\eta).$$

Attention. In the following, for a given model $\{f(x; \theta) : \theta \in \Theta\}$, the notation $I(\eta(\theta))$ will always refer to the FI computed under the reparameterized model $f^*(x; \eta)$ with $\eta = \eta(\theta)$. It does **not** denote the quantity obtained by simply substituting θ with $\eta(\theta)$ inside $I(\theta)$.

Example 3.8.

- Let $X \sim N(\mu, \sigma^2)$. $f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$. We have seen (see Example 3.1 or 3.2) that $I(\sigma^2) = \frac{1}{2\sigma^4}$.

Let's now compute $I(\sigma)$. Using $I(\sigma) = -E[\partial_\sigma^2 \log f]$, we obtain

$$I(\sigma) = -E\left[\partial_\sigma \left\{ -\frac{1}{\sigma} + \frac{(X-\mu)^2}{\sigma^3} \right\}\right] = -E\left[\frac{1}{\sigma^2} - 3\frac{(X-\mu)^2}{\sigma^4}\right] = \frac{2}{\sigma^2}.$$

This matches the result one obtains by applying Proposition 3.5, in fact :

$$I(\sigma) = I(\sigma^2) \times \left(\partial_t t^2 \Big|_{t=\sigma} \right)^2 = I(\sigma^2) (2\sigma)^2 = \frac{2}{\sigma^2}.$$

Or, equivalently,

$$I(\sigma^2) = I(\sigma) \times \left(\partial_t \sqrt{t} \Big|_{t=\sigma^2} \right)^2 = I(\sigma) \left(\frac{1}{2\sqrt{\sigma^2}} \right)^2 = \frac{1}{2\sigma^4}.$$

- Let $X \sim \text{Pois}(\theta)$. $f(x; \theta) = \frac{\theta^x}{x!} e^{-\theta}$, $x = 0, 1, \dots$ and $\theta > 0$. Direct calculation leads to

$$I(\theta) = -E\{\partial_\theta^2 \log f\} = -E\left\{-\frac{X}{\theta^2}\right\} = \frac{1}{\theta}$$

Now, let's consider the parametrization with $\eta = \log(\theta)$: $f^*(x; \eta) = \frac{e^{x\eta}}{x!} e^{-e^\eta}$, $\eta \in (-\infty, \infty)$.

$$I(\eta) = -E\{\partial_\eta^2 \log f^*\} = -E\{-e^\eta\} = e^\eta$$

$$\iff I(\log(\theta)) = \theta$$

This same result can be obtained by directly applying Proposition 3.5 as follows:

$$I(\theta) = I(\log(\theta)) \times \left(\partial_t \log(t) \Big|_{t=\theta} \right)^2 = \frac{I(\log(\theta))}{\theta^2}.$$

Thus, $I(\log(\theta)) = \theta^2 I(\theta) = \theta$. \square

This last example suggests that, in Poisson model, it is easier to estimate $\eta = \log(\theta)$ than θ when the latter is large. Let's check this out. A natural estimator of θ is $\hat{\theta} := \bar{X}$ and a natural estimator of η is $\hat{\eta} := \log(\hat{\theta}) = \log(\bar{X})$. We have that $MSE(\hat{\theta}) = \frac{\theta}{n}$ and, by first order Taylor polynomial approximation, i.e. $\log(\hat{\theta}) \approx \log(\theta) + (\hat{\theta} - \theta) \frac{1}{\theta}$, we can write $MSE(\hat{\eta}) \approx \frac{1}{n\theta}$. These mean-square errors cannot be compared, as they are of different scales/magnitudes. However, we can see that as θ increases, the MSE performance of $\hat{\theta}$ becomes worse and worse compared to that of $\hat{\eta}$. The following simulation confirms this fact.

```
n <- 100; theta <- 10
hat.theta <- estSim(dataFun = \( ) rpois(n, theta), estFun = mean)
```

- Performances of $\hat{\theta}$:

```
mse(hat.theta, theta)
```

bias	var	mse	mare
-0.00218	0.10112	0.10113	0.02537

- Performances of $\hat{\eta}$:

```
mse(log(hat.theta), log(theta))
```

bias	var	mse	mare
-0.000724	0.001012	0.001013	0.011024

Remark. We have seen that, in general, $I(g(\theta)) \neq I(\theta)$. We may ask the question what happen with FI when the data itself, or a statistic from it, are transformed. The answer

depends on the type of transformation used. For example in the case of strictly monotonic and differentiable transformation of the data, FI does change. More precisely, if $U = k(T)$, where k is a differentiable and strictly monotonic function that does not depend on θ , then $I_T(\theta) = I_U(\theta)$. This is a direct consequence of the *Change of Variable(s) Formula*: $f_T(T; \theta) = f_U(U; \theta) |k'(U)|$. \square

3.5 Information Inequality: The Cramer-Rao Lower bound (CRLB)

We will develop a lower bound for the variance of any given statistic, which can be mainly used (i) as a benchmark for comparing estimator performance, and (ii) to find the MVUE (Minimum Variance Unbiased Estimator). The bound we are interested in is called the Cramer-Rao Lower Bound (**CRLB**), and is given in the following theorem.

Theorem 3.2 (Information Inequality). *Let $\mathbf{X} \equiv \mathbf{X}_n = (X_1, \dots, X_n)$ be a random sample with joint pd $f_n(\mathbf{x}; \theta)$, $\theta \in \Theta$. Assume that assumptions (I) and (II) as given above (see Proposition 3.1) hold for $f_n(\mathbf{x}; \theta)$, the joint dp of \mathbf{X} . Let $T \equiv T(\mathbf{X})$ be a statistic. Assume that (III) $\partial_\theta E_\theta(T)$ exists and can be obtained by differentiating under the integral (or sum) sign. i.e., $\partial_\theta \int T(\mathbf{x}) f_n(\mathbf{x}; \theta) d\mathbf{x} = \int T(\mathbf{x}) \partial_\theta f_n(\mathbf{x}; \theta) d\mathbf{x}$. Then*

$$\text{Var}_\theta(T) \geq \frac{(\partial_\theta E_\theta(T))^2}{I_n(\theta)}, \forall \theta \in \Theta.$$

The inequality above is a direct consequence of the Cauchy-Schwarz inequality. The proof goes as follows. Let $S_n \equiv S(\theta, \mathbf{X}) = \partial_\theta \log f_n(\mathbf{X}; \theta)$ be the Score associated with f_n . By the Cauchy-Schwarz inequality,

$$[\text{Cov}(T, S_n)]^2 = [E(TS_n) - E(T)E(S_n)]^2 = [E(TS_n)]^2 \leq \text{Var}(T)\text{Var}(S_n) = \text{Var}(T)I_n(\theta).$$

The final result is the consequence of the fact that $E(TS_n) = \int T(\mathbf{x}) \partial_\theta \log f_n(\mathbf{X}; \theta) f_n(\mathbf{x}; \theta) d\mathbf{x} = \int T(\mathbf{x}) \partial_\theta f_n(\mathbf{x}; \theta) d\mathbf{x} = \partial_\theta E(T)$.

As an immediate consequence, if $\hat{\theta} \equiv T$ is an *unbiased* estimator of θ , then $\text{Var}_\theta(\hat{\theta}) \geq \frac{1}{I_n(\theta)}$.

More generally, if $\widehat{g(\theta)} \equiv T$ is an unbiased estimator of $g(\theta)$, then, $\forall \theta \in \Theta$,

$$\text{Var}_\theta(\widehat{g(\theta)}) \geq I_n^{-1}(g(\theta)),$$

where, $I_n^{-1}(g(\theta)) = \frac{1}{I_n(g(\theta))}$, and $I_n(g(\theta)) := \frac{I_n(\theta)}{(g'(\theta))^2}$.³

The right-hand side of the inequality above is known as the **Cramér–Rao Lower Bound (CRLB)** for estimating $g(\theta)$, i.e. $\text{CRLB}(g(\theta)) = I_n^{-1}(g(\theta))$. It represents the smallest possible variance that any *unbiased estimator* of $g(\theta)$ can achieve under the stated regularity conditions.

³ g does not need to be one-to-one. If it is one-to-one, then $I_n(g(\theta))$ as defined above coincides with the FI of the reparameterization $\theta \mapsto g(\theta)$.

Note. All (regular) exponential families meet the theorem's assumptions; in particular, Assumption (III) holds for any statistic T , eliminating the need for case-by-case verification. \square

An *unbiased* estimator $\widehat{g(\theta)}$ of $g(\theta)$ is called **efficient** if its variance equals the CRLB for $g(\theta)$, i.e. if $\text{Var}(\widehat{g(\theta)}) = \text{CRLB}(g(\theta)) := I_n^{-1}(g(\theta))$; otherwise its (absolute) efficiency is defined to be

$$\text{Eff}(\widehat{g(\theta)}) := \frac{\text{CRLB}(g(\theta))}{\text{Var}(\widehat{g(\theta)})} = \frac{(g'(\theta))^2}{I_n(\theta) \text{Var}(\widehat{g(\theta)})}.$$

$\text{Eff}(\widehat{g(\theta)}) \in (0, 1]$, and equals 1 if and only if $\widehat{g(\theta)}$ is efficient.

By definition, an efficient estimator is (1) unbiased and (2) its variance is uniformly lower than (or equal to) the variance of any other unbiased estimator. Thus, *an efficient estimator, when it exists, is the uniformly minimum variance unbiased estimator (MVUE).*

Efficiency is a stronger requirement than being the MVUE. Indeed, while every efficient estimator is necessarily an MVUE, the converse is not true. An estimator can be MVUE without attaining the CRLB, either because the bound is not attainable or because regularity conditions fail and the CRLB does not apply. In many models, no unbiased estimator reaches the CRLB, yet a unique MVUE still exists with variance strictly larger than $1/I_n(\theta)$ for some θ . Thus, efficiency implies MVUE, but MVUE does not imply efficiency.

$$\text{Efficient} \implies \text{MVUE}$$

Example 3.9. Let $X = (X_1, \dots, X_n)$ be an iid sample from $N(\mu, \sigma^2)$ (exponential family).

- Suppose that μ is our parameter of interest. We have seen that $I_n(\mu) = n/\sigma^2$. So for any unbiased estimator $\hat{\mu}_n$ of μ ,

$$\text{Var}(\hat{\mu}_n) \geq \frac{\sigma^2}{n}.$$

Now, since \bar{X}_n is an unbiased estimator of μ and $\text{Var}(\bar{X}_n) = \frac{\sigma^2}{n}$, we conclude that \bar{X}_n is efficient for μ . And so, \bar{X}_n is the MVUE of μ .

- Suppose that σ^2 is our parameter of interest and μ is known. We know that $I_n(\sigma^2) = \frac{n}{2\sigma^4}$. So for any unbiased estimator $\hat{\sigma}_n^2$ of σ^2 ,

$$\text{Var}(\hat{\sigma}_n^2) \geq 2\sigma^4/n.$$

On the other hand, we know that $\tilde{\sigma}_n^2 = n^{-1} \sum_i (X_i - \mu)^2$ is an unbiased estimator of σ^2 . And, using the fact that $E(X - \mu)^4 = 3\sigma^4$, we have that $\text{Var}(\tilde{\sigma}_n^2) = n^{-1} \text{Var}(X - \mu)^2 = n^{-1} (E(X - \mu)^4 - \sigma^4) = 2\sigma^4/n$. So, we conclude that $\tilde{\sigma}_n^2$ is efficient for σ^2 . And so it is

the MVUE of σ^2 .

- Suppose that μ^2 is our parameter of interest. By the information inequality, with $g : \mu \mapsto \mu^2$, we have that

$$\text{Var}(\hat{\delta}) \geq \frac{1}{I_n(\mu^2)} = \frac{(2\mu)^2}{I_n(\mu)} = \frac{(2\mu)^2}{n/\sigma^2} = \frac{4\mu^2\sigma^2}{n}.$$

for any unbiased estimator $\hat{\delta}$ of μ^2 . Thus, $\text{CRLB}(\mu^2) = 4n^{-1}\mu^2\sigma^2$. But, for now, we cannot say if this limit is attainable or not and thus if there is an efficient estimator for μ^2 or not.

Example 3.10 (Importance of assumptions). Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample from $\text{Unif}(0, \theta)$, $\theta > 0$. Thus $f(x; \theta) = \frac{1}{\theta}$, $0 < x < \theta$. Since $\partial_\theta \log f(x; \theta) = -1/\theta$, if we apply Theorem 3.2, we could conclude that, for any unbiased estimator $\hat{\theta}$ of θ ,

$$\text{Var}(\hat{\theta}) \geq \frac{\theta^2}{n}.$$

But, we learned earlier that the estimator $\hat{\theta}_2 = \frac{n+1}{n}X_{(n)}$ is unbiased and its variance is $\frac{\theta^2}{n(n+2)}$, which is uniformly smaller than θ^2/n !

The apparent contradiction occurs because Assumption (I) (support independent of θ) is violated. Consequently, the FI computed above is incorrect and the CRLB does not apply to this model, so no conflict exists !

Note that it can be shown that $\hat{\theta}_2$ is actually the MVUE of θ . \square

Remark. Efficiency is not generally preserved under transformations: $\hat{\theta}$ is efficient for $\theta \not\Rightarrow g(\hat{\theta})$ is efficient for $g(\theta)$. The only transformations that preserve efficiency are affine maps $g(\theta) = a\theta + b$; ($a \neq 0$). Specifically,

$$\hat{\theta} \text{ is efficient for } \theta \implies a\hat{\theta} + b \text{ is efficient for } a\theta + b. \square$$

3.6 Efficiency in exponential families

Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample from a one-parameter exponential family, written under two parametrizations:⁴

$$f(x; \theta) = h(x) \exp(\eta(\theta)T(x) - B(\theta)) \quad (\text{original parameterization})$$

$$f^*(x; \eta) = h(x) \exp(\eta T(x) - A(\eta)) \quad (\text{canonical parameterization})$$

⁴For convenience, the symbol η denotes hereafter both the mapping $g : \theta \mapsto \eta(\theta)$ and the natural parameter in the canonical form. The intended meaning will always be clear from the context: when η stands alone it refers to the natural (canonical) parameter, whereas $\eta(\theta)$ explicitly indicates the original parametrisation.

For the statistic $T_i = T(X_i)$, $i = 1, \dots, n$, we have seen that (for the canonical representation)

$$E(T_i) = A'(\eta), \text{ and } \text{Var}(T_i) = A''(\eta).$$

The FI about η is

$$I(\eta) = E(\partial_\eta \log f^*)^2 = E\left(T_1 - A'(\eta)\right)^2 = A''(\eta).$$

Define the sample mean of the sufficient statistic : $\bar{T} = n^{-1} \sum_i T_i$. Then,

$$E(\bar{T}) = A'(\eta), \text{ and } \text{Var}(\bar{T}) = \frac{A''(\eta)}{n}.$$

When we regard $A'(\eta)$ as the parameter of interest, the CRLB for estimating this quantity is

$$I^{-1}(A'(\eta)) = \frac{(A''(\eta))^2}{I_n(\eta)} = \frac{A''(\eta)}{n}.$$

$\implies \bar{T}$ is efficient for $A'(\eta)$.

Under the original parametrization, this result is equivalent to say that \bar{T} is efficient for $\frac{B'(\theta)}{\eta'(\theta)}$.

Example 3.11. Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample from the exponential distribution. Thus, for some $\theta > 0$, the pd of X_i is given by

$$\begin{aligned} f(x; \theta) &= \frac{1}{\theta} e^{-\frac{x}{\theta}} I(x > 0) \\ &= I(x > 0) \exp\left(-\frac{1}{\theta}x - \log(\theta)\right) \end{aligned} \quad (\text{exponential family})$$

from the result above, we can directly conclude that $\bar{T} := n^{-1} \sum_{i=1}^n X_i$ is efficient for

$$\frac{(\log(\theta))'}{(-1/\theta)'} = \frac{1/\theta}{1/\theta^2} = \theta = E(X_1). \quad \square$$

3.7 CRLB Attainment

A natural question to ask is under what conditions a given unbiased estimator, say $T(\mathbf{X})$, of $g(\theta)$ can attain the CRLB? It turns out that the CRLB is achieved only when the definition of the estimator $T(\mathbf{X})$ has the special form given in the following theorem.

Theorem 3.3 (CRLB Attainment). *Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random sample with a joint pd $f_n(\mathbf{x}; \theta)$, $\theta \in \Theta$. Suppose the regularity conditions of Theorem 3.2 are satisfied. Then a statistic $T(\mathbf{X})$*

is efficient for $g(\theta)$ if and only if, $\forall \theta \in \Theta$, there exists a function $a_n(\theta) \neq 0$, such that

$$\partial_\theta \log f_n(\mathbf{x}; \theta) = a_n(\theta)[T(\mathbf{x}) - g(\theta)]. \quad (3.1)$$

Moreover, when (3.1) holds: (i) $a_n(\theta) = I_n(\theta)/g'(\theta)$, and (ii) $f_n(\mathbf{x}; \theta)$ belongs to a one-parameter exponential family

Let's prove this result. Let $S_n = \partial_\theta \log f_n(\mathbf{X}; \theta)$ be the Score associated with f_n . Remember that, under the stated assumptions, $E(S_n) = 0$, $\text{Var}(S_n) = I_n(\theta)$ and $\text{Cov}(S_n, T) = \partial_\theta E(T)$. Suppose that $T(\mathbf{X})$ is an efficient estimator of $g(\theta)$. Then, $\forall \theta \in \Theta$, $E(T) = g(\theta)$, and $\text{Var}(T) = (g'(\theta))^2 / I_n(\theta)$. So, $\text{Cov}^2(S_n, T) = (g'(\theta))^2 = \text{Var}(T)\text{Var}(S_n)$. Since Cauchy-Schwarz inequality become an equality only in the case of linear dependence, we conclude that $\exists a \equiv a_n(\theta) \neq 0$ and $b \equiv b_n(\theta)$, such that $S_n = aT + b$. But since $E(S_n) = 0$, we have that $b = -ag(\theta)$. Thus, $S_n = a(T - g(\theta))$.

Conversely, suppose that, $\forall \theta \in \Theta$, $\exists a \equiv a_n(\theta) \neq 0$ such that $S_n = a(T - g(\theta))$. Observe that $E(S_n) = a(E(T) - g(\theta)) \Rightarrow E(T) = g(\theta)$, $\text{Var}(S_n) = a^2 \text{Var}(T) \Rightarrow \text{Var}(T) = I_n(\theta)/a^2$, and $\text{Cov}(S_n, T) = \text{Cov}(a(T - g(\theta)), T) = a \text{Var}(T) \Rightarrow g'(\theta) = a \text{Var}(T)$. This implies that $a = I_n(\theta)/g'(\theta)$, which in turn implies that $\text{Var}(T) = (g'(\theta))^2 / I_n(\theta)$. This concludes the proof.

The CRLB attainment theorem above leads to an explicit constructive procedure for deriving the (efficient) MVUE of $g(\theta)$ when it exists. Namely, put

$$T = g(\theta) + \frac{g'(\theta)}{I_n(\theta)} \partial_\theta \log f_n(\mathbf{X}; \theta).$$

If the expression on the right hand side of the equality above *does not depend on θ* , i.e. T as defined above is a statistic, then T is efficient for $g(\theta)$.

Example 3.12. Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample from the exponential distribution. Thus, for some $\theta > 0$, the pd of X_i is given by $f(x; \theta) = \frac{1}{\theta} e^{-x/\theta}$, for $x > 0$. It is easy to check that

$$\begin{aligned} \partial_\theta \log f_n(\mathbf{X}; \theta) &= -\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n X_i \\ \partial_\theta^2 \log f_n(\mathbf{X}; \theta) &= \frac{n}{\theta^2} - \frac{2}{\theta^3} \sum_{i=1}^n X_i \Rightarrow I_n(\theta) = -E[\partial_\theta^2 \log f_n(\mathbf{X}; \theta)] = \frac{n}{\theta^2} \end{aligned}$$

- Let's try to find the efficient estimator for θ . To do so, put

$$\begin{aligned} T &:= \theta + \partial_\theta \log f_n(\mathbf{X}; \theta) / I_n(\theta) \\ &= \theta + \frac{\theta^2}{n} \left(-\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n X_i \right) = n^{-1} \sum_{i=1}^n X_i \end{aligned}$$

$\rightarrow n^{-1} \sum_i X_i$ is the desired estimator.

- Let's now try to find an efficient estimator for $\delta = \frac{1}{\theta}$. Following the same procedure, let

$$\begin{aligned} T &:= \frac{1}{\theta} + \frac{\left(\frac{1}{\theta}\right)'}{\frac{n}{\theta^2}} \left(-\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n X_i \right) \\ &= \frac{2}{\theta} - \frac{1}{\theta^2} n^{-1} \sum_{i=1}^n X_i. \end{aligned}$$

The latter is not a statistic, so there is no efficient estimator for $1/\theta$. Note, however, that this does not mean that there is no MVUE. \square

3.8 Multiparameter case

The above theory extends naturally to parametric models whose pdf $f(x; \theta)$ depend on several parameters $\theta = (\theta_1, \dots, \theta_d) \in \Theta \subset \mathbb{R}^d$. In what follows, we present the main results without working through the full details of the calculations.

In the sequel, we assume that the following two regularity conditions are satisfied:

- (I) The set $\{x : f(x; \theta) > 0\}$ does not depend on θ .
- (II) The operations of integration (or summation) and differentiation by θ_j can be interchanged in $\int f(x; \theta) dx$. i.e., $\partial_{\theta_j} \int f(x; \theta) dx = \int \partial_{\theta_j} f(x; \theta) dx$, $\forall j = 1, \dots, d$.

The Score vector

The Score vector of f is defined as the gradient of $\theta \mapsto \log f(X; \theta)$, i.e.

$$S := \nabla_{\theta} \log f(X; \theta) = (S_1, \dots, S_d)^t,$$

where $S_j = \partial_{\theta_j} \log f(X; \theta)$ is the score for θ_j . It is easy to see that $E(S_j) = 0$, $\forall j$. Thus, $E(S) = \mathbf{0}$.

The FI matrix

The **FI matrix** contained in X about θ is defined as

$$I(\theta) := E(SS^t) = \text{Var}(S).$$

Thus, $I(\theta) = [I_{jk}(\theta)]_{1 \leq j, k \leq d}$, where $I_{jk}(\theta) := E(S_j \times S_k) = \text{Cov}(S_j, S_k)$.

Any FI matrix is symmetric and positive semidefinite by construction (it is the covariance matrix of the score vector). It becomes positive definite when the statistical model is regular and identifiable.

Attention. Here, we restrict our attention to regular models, for which the FI matrix is positive definite—and therefore invertible.

FI matrix–Hessian form

If $f(x; \theta)$ is twice differentiable and double integration and differentiation under the integral sign can be interchanged, i.e., $\partial_{\theta_k \theta_j} \int f(x; \theta) dx = \int \partial_{\theta_k \theta_j} f(x; \theta) dx$, $\forall j, k = 1, \dots, d$, then $I_{jk}(\theta) = -E(\partial_{\theta_k} S_j)$. Thus,

$$I(\theta) = -E(\nabla_{\theta}^2 \log f(X; \theta)),$$

where $\nabla_{\theta}^2 \log f(X; \theta)$ denotes the Hessian of $\theta \mapsto \log f(X; \theta)$.

To be more explicit about the formulas given above. Let's consider the special case of a two-parameter model with $\theta = (\theta_1, \theta_2)$, $S_1 = \partial_{\theta_1} \log f(X; \theta)$, and $S_2 = \partial_{\theta_2} \log f(X; \theta)$. Under the regularity assumptions stated above, we can write the FI matrix in any of the following equivalent expressions:

$$I(\theta_1, \theta_2) = E \begin{pmatrix} S_1^2 & S_1 S_2 \\ S_1 S_2 & S_2^2 \end{pmatrix} = \begin{pmatrix} \text{Var}(S_1) & \text{Cov}(S_1, S_2) \\ \text{Cov}(S_1, S_2) & \text{Var}(S_2) \end{pmatrix} = -E \begin{pmatrix} \partial_{\theta_1} S_1 & \partial_{\theta_2} S_1 \\ \partial_{\theta_1} S_2 & \partial_{\theta_2} S_2 \end{pmatrix}.$$

FI matrix in an iid Sample

Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample with joint pd $f_n(x; \theta)$. The joint Score and associated FI are defined as

$$S_n = \nabla_{\theta} \log f_n(\mathbf{X}; \theta), \quad I_n(\theta) = E(S_n S_n^t).$$

Because of the iid assumption, the FI is additive, yielding

$$I_n(\theta) = nI(\theta).$$

Hereafter, we'll use $I_n^{-1}(\theta)$ to denote $[I_n(\theta)]^{-1} = n^{-1}[I(\theta)]^{-1}$, the inverse matrix of $I_n(\theta)$.

Example 3.13. Normal distribution $N(\mu, \sigma^2)$ with $\log f(x; \mu, \sigma^2) = -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2$.

$$\begin{aligned} S_1 &= \partial_{\mu} \log f = \frac{X - \mu}{\sigma^2} & S_2 &= \partial_{\sigma^2} \log f = -\frac{1}{2\sigma^2} + \frac{(X - \mu)^2}{2\sigma^4} \\ \partial_{\mu} S_1 &= -\frac{1}{\sigma^2} & \partial_{\mu} S_2 &= -\frac{X - \mu}{\sigma^4} \\ \partial_{\sigma^2} S_1 &= -\frac{X - \mu}{\sigma^4} & \partial_{\sigma^2} S_2 &= \frac{1}{2\sigma^4} - \frac{(X - \mu)^2}{\sigma^6} \\ I_{11} &= -E(\partial_{\mu} S_1) = \frac{1}{\sigma^2} & I_{22} &= -E(\partial_{\sigma^2} S_2) = \frac{1}{2\sigma^4} \\ I_{21} &= -E(\partial_{\mu} S_2) = 0 & I_{12} &= -E(\partial_{\sigma^2} S_1) = 0 \end{aligned}$$

$$\Rightarrow I(\mu, \sigma^2) = \begin{pmatrix} I_{11} & I_{12} \\ I_{21} & I_{22} \end{pmatrix} = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix} \text{ and } I_n(\mu, \sigma^2) = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix}. \square$$

FI under reparametrization – Multiparameter Case

Assume the model is reparametrized by a one-to-one differentiable mapping $\boldsymbol{\eta} : \boldsymbol{\theta} \mapsto \boldsymbol{\eta}(\boldsymbol{\theta}) = (\eta_1(\boldsymbol{\theta}), \dots, \eta_d(\boldsymbol{\theta}))$.

Let $\mathbf{I}(\boldsymbol{\theta})$ and $\mathbf{I}(\boldsymbol{\eta}) \equiv \mathbf{I}(\boldsymbol{\eta}(\boldsymbol{\theta}))$ denote the FI matrices under the $\boldsymbol{\theta}$ - and $\boldsymbol{\eta}$ -parametrizations, respectively.

Then the two matrices are related by

$$\mathbf{I}(\boldsymbol{\theta}) = \dot{\boldsymbol{\eta}}^t(\boldsymbol{\theta}) \mathbf{I}(\boldsymbol{\eta}) \dot{\boldsymbol{\eta}}(\boldsymbol{\theta}),$$

where $\dot{\boldsymbol{\eta}}(\boldsymbol{\theta}) = \partial_{\boldsymbol{\theta}} \boldsymbol{\eta}(\boldsymbol{\theta})$ is the $d \times d$ Jacobian matrix of $\boldsymbol{\theta} \mapsto \boldsymbol{\eta}(\boldsymbol{\theta})$, whose (j, k) -entry is $\partial_{\theta_k} \eta_j(\boldsymbol{\theta})$, and $\dot{\boldsymbol{\eta}}^t(\boldsymbol{\theta}) = [\dot{\boldsymbol{\eta}}(\boldsymbol{\theta})]^t$ denotes its transpose.

Example 3.14. Normal distribution $N(\mu, \sigma^2)$.

Let $\boldsymbol{\eta} : (\mu, \sigma) \mapsto (\eta_1, \eta_2)$, with $\eta_1(\mu, \sigma) = \mu$, and $\eta_2(\mu, \sigma) = \sigma^2$. We have that

$$\begin{aligned} \dot{\boldsymbol{\eta}} \equiv \dot{\boldsymbol{\eta}}(\mu, \sigma) &:= \begin{pmatrix} \partial_{\mu} \eta_1 & \partial_{\sigma} \eta_1 \\ \partial_{\mu} \eta_2 & \partial_{\sigma} \eta_2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 2\sigma \end{pmatrix}. \\ \Rightarrow \mathbf{I}(\mu, \sigma) = \dot{\boldsymbol{\eta}}^t \mathbf{I}(\mu, \sigma^2) \dot{\boldsymbol{\eta}} &= \begin{pmatrix} 1 & 0 \\ 0 & 2\sigma \end{pmatrix} \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{1}{2\sigma^4} \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 2\sigma \end{pmatrix} = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{2}{\sigma^2} \end{pmatrix}. \square \end{aligned}$$

Information Inequality Theorem–Multiparameter CRLB

Let $\mathbf{X} = (X_1, \dots, X_n)$ be a random sample with joint pd $f_n(\mathbf{x}; \boldsymbol{\theta})$, $\boldsymbol{\theta} = (\theta_1, \dots, \theta_d) \in \boldsymbol{\Theta} \subseteq \mathbb{R}^d$. Let $\mathbf{T} \equiv (T_1(\mathbf{X}), \dots, T_p(\mathbf{X}))$ be an estimator of $\mathbf{g}(\boldsymbol{\theta})$, where $\mathbf{g}(\cdot) = (g_1(\cdot), \dots, g_p(\cdot))$ is a differentiable mapping from \mathbb{R}^d to \mathbb{R}^p . Assume that

- (I) and (II), as given above (see the beginning of the current section), hold for $f_n(\mathbf{x}; \boldsymbol{\theta})$.
 - (III) $\partial_{\theta_k} E_{\boldsymbol{\theta}}(T_j)$ exists and can be obtained by differentiating under the integral sign, $\forall j, k, \boldsymbol{\theta}$.
- If \mathbf{T} is an unbiased estimator of $\mathbf{g}(\boldsymbol{\theta})$, i.e. if $E_{\boldsymbol{\theta}}(T_j) = g_j(\boldsymbol{\theta})$, $j = 1, \dots, p$, then, $\forall \boldsymbol{\theta} \in \boldsymbol{\Theta}$,

$$\text{Var}_{\boldsymbol{\theta}}(\mathbf{T}) \succeq \mathbf{I}_n^{-1}(\mathbf{g}(\boldsymbol{\theta})),$$

where $\text{Var}_{\boldsymbol{\theta}}(\mathbf{T})$ is the variance-covariance matrix of \mathbf{T} , $\mathbf{I}_n^{-1}(\mathbf{g}(\boldsymbol{\theta})) := \dot{\mathbf{g}}(\boldsymbol{\theta}) \mathbf{I}_n^{-1}(\boldsymbol{\theta}) \dot{\mathbf{g}}^t(\boldsymbol{\theta})$ ⁵, $\dot{\mathbf{g}}$ is the $p \times d$ Jacobian matrix of \mathbf{g} , whose (j, k) -th element is $\partial_{\theta_k} g_j(\boldsymbol{\theta})$, and $\dot{\mathbf{g}}^t(\boldsymbol{\theta}) = [\dot{\mathbf{g}}(\boldsymbol{\theta})]^t$.

In particular, if $d = p$ and $\mathbf{g}(\boldsymbol{\theta}) = \boldsymbol{\theta}$, i.e. if \mathbf{T} is unbiased for $\boldsymbol{\theta}$, then $\text{Var}_{\boldsymbol{\theta}}(\mathbf{T}) \succeq \mathbf{I}_n^{-1}(\boldsymbol{\theta})$.

Above, the notation $\mathbf{A} \succeq \mathbf{B}$ means that $\mathbf{A} - \mathbf{B}$ is positive semi-definite matrix. Consequently, writing $\text{Var}(\mathbf{T}) \succeq \mathbf{I}_n^{-1}$ is equivalent to say that $\text{Var}(\mathbf{a}^t \mathbf{T}) \geq \mathbf{a}^t \mathbf{I}_n^{-1} \mathbf{a}$, $\forall \mathbf{a}$.

This equivalence reveals a key interpretation: each choice of \mathbf{a} selects a direction in the parameter space, corresponding to a particular linear combination of parameters. The CRLB states that, in every such direction, the variance of the corresponding estimator (i.e. $\text{Var}(\mathbf{a}^t \mathbf{T})$)

⁵ $\dot{\mathbf{g}}$ does not need to be one-to-one, but if it is, then $\mathbf{I}_n^{-1}(\mathbf{g}(\boldsymbol{\theta}))$, as defined above, coincides with $[\mathbf{I}_n(\mathbf{g}(\boldsymbol{\theta}))]^{-1}$, the inverse of the FI about $\mathbf{g}(\boldsymbol{\theta})$.

cannot fall below the threshold $\mathbf{a}^t \mathbf{I}_n^{-1} \mathbf{a}$ determined by the FI matrix. In short, the CRLB provides a universal lower bound on how precisely we can estimate all parameters (or several functions thereof) simultaneously.

In the case of $p = 1$, i.e. when $g : \mathbb{R}^d \mapsto \mathbb{R}$ is a scalar-valued function, the CRLB simplifies to the following form: for any unbiased estimator T of $g(\boldsymbol{\theta})$,

$$\text{Var}(T) \geq \nabla^t g(\boldsymbol{\theta}) \mathbf{I}_n^{-1}(\boldsymbol{\theta}) \nabla g(\boldsymbol{\theta}),$$

where $\nabla g(\boldsymbol{\theta}) = (\partial_1 g(\boldsymbol{\theta}), \dots, \partial_d g(\boldsymbol{\theta}))^t$ is the gradient of g .

Example 3.15. Normal distribution $N(\mu, \sigma^2) \equiv N(\theta_1, \theta_2)$. For any unbiased estimator $\hat{\eta}$ of $\eta = g(\mu, \sigma^2) \in \mathbb{R}$, we have

$$\begin{aligned} \text{Var}(\hat{\eta}) &\geq \begin{pmatrix} \partial_\mu g(\mu, \sigma^2) & \partial_{\sigma^2} g(\mu, \sigma^2) \end{pmatrix} \mathbf{I}_n^{-1}(\mu, \sigma^2) \begin{pmatrix} \partial_\mu g(\mu, \sigma^2) \\ \partial_{\sigma^2} g(\mu, \sigma^2) \end{pmatrix} \\ &= \frac{\sigma^2}{n} (\partial_\mu g(\mu, \sigma^2))^2 + \frac{2\sigma^4}{n} (\partial_{\sigma^2} g(\mu, \sigma^2))^2. \end{aligned}$$

For example, if we are interested in estimating the coefficient of variation (CV), taking $g(\mu, \sigma^2) = \sigma/\mu$ and applying the CRLB gives

$$\text{Var}(\hat{\eta}) \geq \frac{\sigma^2}{n} \left(-\frac{\sigma}{\mu^2} \right)^2 + \frac{2\sigma^4}{n} \left(\frac{1}{2\mu\sigma} \right)^2 = \frac{\sigma^2}{n\mu^2} \left(\frac{\sigma^2}{\mu^2} + \frac{1}{2} \right).$$

— \rightarrow no unbiased estimator (if any) of the CV can have variance smaller than $\frac{\sigma^2}{n\mu^2} \left(\frac{\sigma^2}{\mu^2} + \frac{1}{2} \right)$
— this is the best possible precision for any unbiased CV estimator under normality. \square

Effect of nuisance parameters

The multiparameter CRLB provides important insight into *the effect of nuisance parameters on estimation precision*. To see this, consider again the two-parameter case ($d = 2$), with $\boldsymbol{\theta} = (\theta_1, \theta_2)$. Suppose θ_1 is our *primary parameter of interest* and θ_2 is a *nuisance parameter*. In the normal model, for example, we might be interested only in the mean μ , while σ^2 is a nuisance parameter that is required for model correctness but is not of direct interest.

In such a context, the natural question is: can we estimate θ_1 with the same precision we would achieve if θ_2 were known?

To investigate this, we start by writing the FI matrix for the two-parameter model, along with its inverse. Using the same notations for the score components as above (S_1 for θ_1 and S_2 for θ_2), we write

$$\mathbf{I}_n(\boldsymbol{\theta}) = n \begin{pmatrix} \text{Var}(S_1) & \text{Cov}(S_1, S_2) \\ \text{Cov}(S_2, S_1) & \text{Var}(S_2) \end{pmatrix} \equiv n \begin{pmatrix} I_{11} & I_{12} \\ I_{12} & I_{22} \end{pmatrix}.$$

$$\mathbf{I}_n^{-1}(\boldsymbol{\theta}) = \frac{1}{n(I_{11}I_{22} - I_{12}^2)} \begin{pmatrix} I_{22} & -I_{12} \\ -I_{12} & I_{11} \end{pmatrix}.$$

We can think of two situations :

- θ_2 is known: In this case, the model reduces to a single parameter to $\theta = \theta_1$, and we return to the univariate setting. Applying the univariate CRLB theorem (Theorem 3.2), we have, for any unbiased estimator T_1 of θ_1 ,

$$\text{Var}(T_1) \geq \frac{1}{nI_{11}}.$$

- θ_2 is unknown: In this case, applying the multiparameter version of the CRLB theorem with $g : (\theta_1, \theta_2) \mapsto \theta_1$, we have, for any unbiased estimator T_1 of θ_1 ,

$$\text{Var}(T_1) \geq (1, 0) \mathbf{I}_n^{-1}(\boldsymbol{\theta}) (1, 0)^t = \frac{I_{22}}{n(I_{11}I_{22} - I_{12}^2)} = \frac{1}{nI_{11}^*},$$

$$\text{where } I_{11}^* := \frac{I_{11}I_{22} - I_{12}^2}{I_{22}} = I_{11} \left(1 - \frac{I_{12}^2}{I_{11}I_{22}} \right) = I_{11} (1 - \text{Cor}^2(S_1, S_2)).$$

Observe that $I_{11}^* \leq I_{11}$, with equality if and only if S_1 and S_2 are uncorrelated—that is, when the Fisher Information matrix is diagonal.

This shows that the presence of nuisance parameters can reduce the effective information available for estimating the parameters of interest. In other words, nuisance parameters not only complicate the estimation process but also diminish the attainable precision by lowering the corresponding FI.

Example 3.16. Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample from $N(\mu, \sigma^2)$, where $\boldsymbol{\theta} = (\mu, \sigma^2)$ is unknown. We have seen that the FI matrix for this model is

$$\mathbf{I}_n(\mu, \sigma^2) = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix},$$

which is *diagonal* \longrightarrow no loss from nuisance parameters.

- Suppose that μ is our parameter of interest. The information inequality states that $\text{Var}(\hat{\mu}) \geq \frac{\sigma^2}{n}$, for any unbiased estimator $\hat{\mu}$ of μ , whether σ^2 is known or not. On the other hand, since $E(\bar{X}_n) = \mu$ and $\text{Var}(\bar{X}_n) = \sigma^2/n$, \bar{X}_n is efficient for μ , regardless of whether σ^2 is known.
- Suppose now that σ^2 is our parameter of interest. The information inequality tells us that $\text{Var}(\hat{\sigma}^2) \geq \frac{2\sigma^4}{n}$, for any unbiased estimator $\hat{\sigma}^2$ of σ^2 , whether μ is known or not.
 - If μ is known, a natural estimator of σ^2 is $\tilde{\sigma}_n^2 = n^{-1} \sum_i (X_i - \mu)^2$. We have seen that $E(\tilde{\sigma}_n^2) = \sigma^2$ and $\text{Var}(\tilde{\sigma}_n^2) = 2\sigma^4/n$. Hence, in this setting, $\tilde{\sigma}_n^2$ attains the CRLB and is therefore an **efficient estimator** of σ^2 .
 - If μ is unknown, $\tilde{\sigma}_n^2$ is no longer usable because it depends on the (unknown)

true value of μ . In this case, an unbiased estimate of σ^2 is $S_n^2 = \frac{1}{n-1} \sum_i (X_i - \bar{X}_n)^2$. However, $\text{Var}(S_n^2) = \frac{2\sigma^4}{n-1} > \frac{2\sigma^4}{n}$. So, S_n^2 is unbiased but **slightly inefficient**. Nevertheless, it can be shown that S_n^2 is the MVUE of σ^2 . \square

Example 3.17. Let $\mathbf{X} = (X_1, X_2, X_3)$ be a random vector drawn from the [trinomial distribution](#) $\text{Mult}(n, (\pi_1, \pi_2, \pi_3))$. Put $\mathbf{x} = (x_1, x_2, x_3)$, and $\boldsymbol{\pi} = (\pi_1, \pi_2, \pi_3)$. The joint pd of \mathbf{X} is

$$f_n(\mathbf{x}; \boldsymbol{\pi}) := P(X_1 = x_1, X_2 = x_2, X_3 = x_3) = \frac{n!}{x_1!x_2!x_3!} \pi_1^{x_1} \pi_2^{x_2} \pi_3^{x_3},$$

with $x_1, x_2, x_3 \in \{0, \dots, n\}$, $x_3 = n - x_1 - x_2$, and $\pi_3 = 1 - \pi_1 - \pi_2$. Because of these constraints, $f(\mathbf{x}; \boldsymbol{\pi})$ depends effectively only on (x_1, x_2) and (π_1, π_2) . The parameter space is $\Theta = \{(\pi_1, \pi_2) \in [0, 1]^2 : \pi_1 \geq 0, \pi_2 \geq 0, \pi_1 + \pi_2 \leq 1\}$. It is known that $E(\mathbf{X}) = n\boldsymbol{\pi}$ and $\text{Var}(\mathbf{X}) = n(\text{diag}(\boldsymbol{\pi}) - \boldsymbol{\pi}\boldsymbol{\pi}^t)$.

We have that

$$\begin{aligned} S_1 &= \partial_{\pi_1} \log f_n(\mathbf{X}; \boldsymbol{\pi}) = X_1 \partial_{\pi_1} \log(\pi_1) + X_2 \partial_{\pi_1} \log(\pi_2) + X_3 \partial_{\pi_1} \log(\pi_3) \\ &= \frac{X_1}{\pi_1} + X_3 \partial_{\pi_1}(\pi_3) \partial_{\pi_3}(\log(\pi_3)) = \frac{X_1}{\pi_1} - \frac{X_3}{\pi_3}. \\ \partial_{\pi_1} S_1 &= X_1 \partial_{\pi_1} \pi_1^{-1} - X_3 \partial_{\pi_1} \pi_3^{-1} = -\frac{X_1}{\pi_1^2} - X_3 \partial_{\pi_1}(\pi_3) \partial_{\pi_3}(\pi_3^{-1}) = -\frac{X_1}{\pi_1^2} - \frac{X_3}{\pi_3^2}. \\ \partial_{\pi_2} S_1 &= -X_3 \partial_{\pi_2} \pi_3^{-1} = -\frac{X_3}{\pi_3^2}. \end{aligned}$$

Similarly, by simetry, $S_2 := \partial_{\pi_2} \log f_n(\mathbf{X}; \boldsymbol{\pi}) = \frac{X_2}{\pi_2} - \frac{X_3}{\pi_3}$, $\partial_{\pi_1} S_2 = -\frac{X_3}{\pi_3^2}$, and $\partial_{\pi_2} S_2 = -\frac{X_2}{\pi_2^2} - \frac{X_3}{\pi_3^2}$.

Using the fact that $E(X_j) = n\pi_j$, $j = 1, 2, 3$, we get the FI matrix

$$\mathbf{I}_n(\pi_1, \pi_2) = -E \begin{pmatrix} \partial_{\pi_1} S_1 & \partial_{\pi_2} S_1 \\ \partial_{\pi_1} S_2 & \partial_{\pi_2} S_2 \end{pmatrix} = n \begin{pmatrix} \frac{1}{\pi_1} + \frac{1}{\pi_3} & \frac{1}{\pi_3} \\ \frac{1}{\pi_3} & \frac{1}{\pi_2} + \frac{1}{\pi_3} \end{pmatrix}.$$

The inverse of thi matrix is

$$\mathbf{I}_n^{-1}(\pi_1, \pi_2) = n^{-1}(\pi_1 \pi_2 \pi_3) \begin{pmatrix} \frac{1}{\pi_2} + \frac{1}{\pi_3} & -\frac{1}{\pi_3} \\ -\frac{1}{\pi_3} & \frac{1}{\pi_1} + \frac{1}{\pi_3} \end{pmatrix} = n^{-1} \begin{pmatrix} \pi_1(1 - \pi_1) & -\pi_1 \pi_2 \\ -\pi_1 \pi_2 & \pi_2(1 - \pi_2) \end{pmatrix}.$$

Suppose that π_1 is our parameter of interest and π_2 is a nuisance parameter. Given the FI matrix above and its inverse, according and the information inequality, we can say that, for any unbiased estimator $\hat{\pi}_1$ of π_1 ,

$$\text{Var}(\hat{\pi}_1) \geq \frac{1}{n \left(\frac{1}{\pi_1} + \frac{1}{\pi_3} \right)} = \frac{\pi_1(1 - \pi_1 - \pi_2)}{n(1 - \pi_2)}, \text{ if } \pi_2 \text{ is known.}$$

$$\text{Var}(\hat{\pi}_1) \geq \frac{\pi_1(1 - \pi_1)}{n}, \text{ if } \pi_2 \text{ is unknown.}$$

And because $\frac{\pi_1(1-\pi_1-\pi_2)}{n(1-\pi_2)} < \frac{\pi_1(1-\pi_1)}{n}$, knowing π_2 gives a smaller (tighter) lower bound, meaning that knowing π_2 allows potentially more precise estimation of π_1 . \square

Efficiency in the Multiparameter Case

An unbiased estimator $T = (T_1, \dots, T_k)^\top$ of a vector-valued parameter $g(\theta)$ is efficient if its covariance matrix attains the CRLB for all θ in the parameter space, i.e.

$$\text{Var}_\theta(T) = \dot{g}(\theta) I_n^{-1}(\theta) \dot{g}^t(\theta)$$

The multiparameter version of the **CRLB attainment theorem** provides a practical criterion for establishing the existence of an efficient estimator and, when one exists, determining its explicit form. Under regularity conditions (I)–(III), the theorem states that if the quantity

$$T := g(\theta) + \dot{g}(\theta) I_n^{-1}(\theta) \nabla_\theta \log f_n(X; \theta),$$

is free of θ , then this T constitutes an efficient estimator of $g(\theta)$.

Example 3.18. Consider again the trinomial model as in our previous example. We seek to determine whether an efficient estimator exists for $\pi = (\pi_1, \pi_2)$. For that, we need to compute

$$\pi + I_n^{-1}(\pi) \nabla_\pi \log f_n(X; \pi) = \begin{pmatrix} \pi_1 \\ \pi_2 \end{pmatrix} + n^{-1} \begin{pmatrix} \pi_1(1-\pi_1) & -\pi_1\pi_2 \\ -\pi_1\pi_2 & \pi_2(1-\pi_2) \end{pmatrix} \begin{pmatrix} \frac{X_1}{\pi_1} - \frac{X_3}{\pi_3} \\ \frac{X_2}{\pi_2} - \frac{X_3}{\pi_3} \end{pmatrix}$$

The first component of the matrix by the vector product above is

$$\begin{aligned} \pi_1(1-\pi_1) \left(\frac{X_1}{\pi_1} - \frac{X_3}{\pi_3} \right) - \pi_1\pi_2 \left(\frac{X_2}{\pi_2} - \frac{X_3}{\pi_3} \right) &= (1-\pi_1)X_1 - \pi_1X_2 + (\pi_1(\pi_1+\pi_2) - \pi_1) \frac{X_3}{\pi_3} \\ &= (1-\pi_1)X_1 - \pi_1X_2 + (\pi_1(1-\pi_3) - \pi_1) \frac{X_3}{\pi_3} \\ &= X_1 - \pi_1(X_1 + X_2 + X_3) = X_1 - n\pi_1. \end{aligned}$$

As for the second component, exactly the same algebra yields to

$$-\pi_1\pi_2 \left(\frac{X_1}{\pi_1} - \frac{X_3}{\pi_3} \right) + \pi_2(1-\pi_2) \left(\frac{X_2}{\pi_2} - \frac{X_3}{\pi_3} \right) = X_2 - n\pi_2$$

Thus,

$$\pi + I_n^{-1}(\pi) \nabla_\pi \log f_n(X; \pi) = \begin{pmatrix} \pi_1 \\ \pi_2 \end{pmatrix} + n^{-1} \begin{pmatrix} X_1 - n\pi_1 \\ X_2 - n\pi_2 \end{pmatrix} = \begin{pmatrix} \frac{X_1}{n} \\ \frac{X_2}{n} \end{pmatrix}.$$

$$\implies \begin{pmatrix} \frac{X_1}{n} \\ \frac{X_2}{n} \end{pmatrix} \text{ is efficient for } (\pi_1, \pi_2). \quad \square$$

Chapter 4

Asymptotic evaluations

All the criteria we have examined so far are finite sample size criteria. In contrast, we might consider asymptotic properties, properties describing the behavior of a procedure *as the sample size becomes infinite* or very large. Describing the distribution of an estimator for a fixed (finite) sample size n is usually very difficult (if not impossible). But in the limit (or asymptotic) regime, things become more structured and generally lead to practical and powerful solutions.

4.1 Convergence in probability and consistency

Definition 4.1 (Convergence in probability). A sequence of random variables $X_n, n = 1, 2, \dots$, is said to converge to X in probability, in symbols $X_n \xrightarrow{p} X$, if for every $\epsilon > 0$,

$$P(|X_n - X| \geq \epsilon) \rightarrow 0 \quad \text{when } n \rightarrow \infty.$$

Roughly speaking, $X_n \xrightarrow{p} X$ means that, as n increases, X_n and X get closer to each other (*with a probability approaching 1*), i.e. $P(X_n \approx X) \approx 1$, when n becomes large. In most practical situations, X is a constant (non-random).

Example 4.1.

- Let X_n be a sequence of random variables such that $X_n \sim \text{Ber}(1/n)$, $n = 1, 2, \dots$

$$P(|X_n - 0| \geq \epsilon) = P(X_n \geq \epsilon) = P(X_n = 1) = 1/n \rightarrow 0.$$

So $X_n \xrightarrow{p} 0$.

- Let X_n be a sequence of random variables such that $X_n = (1 + 1/n)X$, $n = 1, 2, \dots$, with $P(|X| < 10) = 1$.

$$P(|X_n - X| \geq \epsilon) = P(|X| \geq n\epsilon) \leq P(|X| \geq 10) = 0, \forall n \geq 10\epsilon^{-1}.$$

So $X_n \xrightarrow{p} X$. \square

Note also that the definition above says nothing about the rate of convergence (how fast X_n converges to X). The figure below illustrates the concept of convergence of probability and convergence speed.

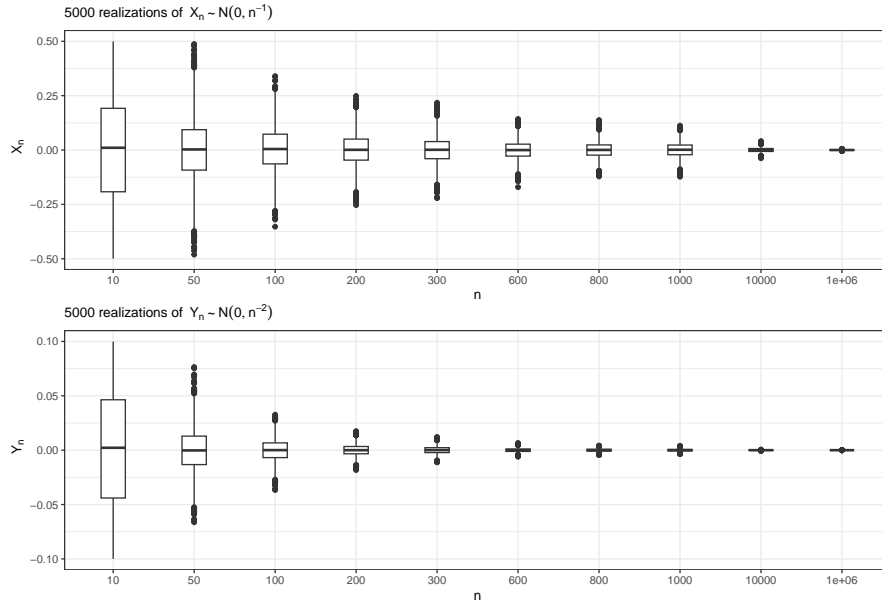


Figure 4.1: Both X_n and $Y_n \xrightarrow{p} 0$ but Y_n goes faster to 0.

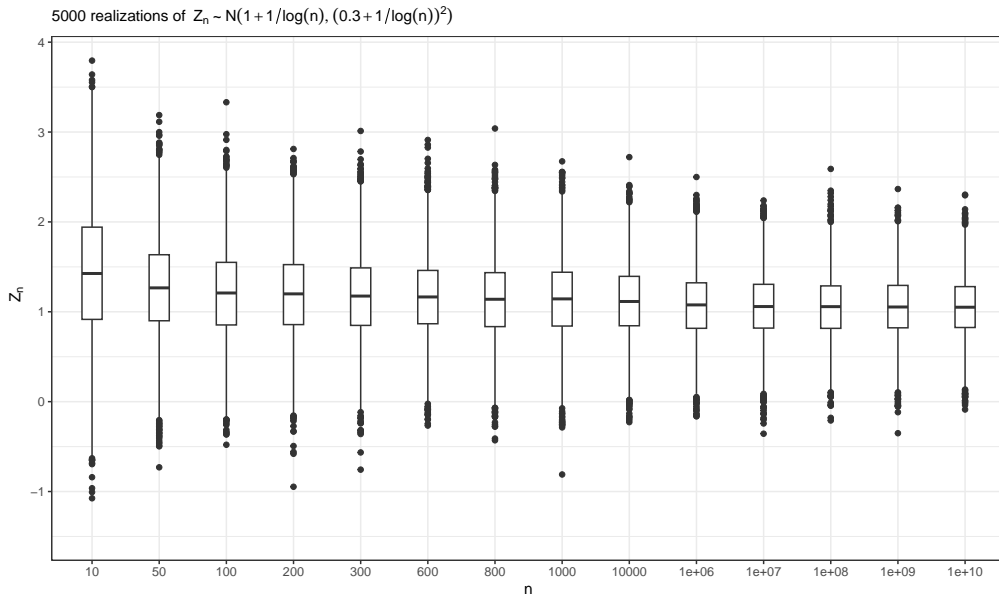


Figure 4.2: $Z_n \xrightarrow{p} N(1, 0.3^2)$.

Facts to know

Convergence in probability is closed under all arithmetic operations. Thus, if $X_n \xrightarrow{p} X$ and $Y_n \xrightarrow{p} Y$, then

- $a_n X_n + b_n \xrightarrow{p} aX + b$, if $a_n \rightarrow a$ and $b_n \rightarrow b$.
- $X_n \pm Y_n \xrightarrow{p} X \pm Y$.
- $X_n Y_n \xrightarrow{p} XY$.
- $1/X_n \xrightarrow{p} 1/X$, if $X \neq 0$ (with probability one).

As we will see later, these results are mainly a direct consequence of a more general result known as *continuous mapping Theorem* (see Theorem 4.4).

For now, as an example, let's take X_n to be a sequence of random variables such that $X_n \sim N(\mu_n, \sigma_n^2)$, $n = 1, 2, \dots$, with $\sigma_n^2 > 0$ and $\mu_n \rightarrow \mu$. Let $Z \sim N(0, 1)$. Applying the first of the above properties, we can immediately conclude that if $\sigma_n^2 \rightarrow \sigma^2 \geq 0$, then $X_n \xrightarrow{p} X := \mu + \sigma Z$, with $X \sim N(\mu, \sigma^2)$. In short, we write $X_n \xrightarrow{p} N(\mu, \sigma^2)$. If $\sigma = 0$, then $X_n \xrightarrow{p} \mu$.

The concept of convergence in probability is easily generalized to the multivariate case as follows

$$(X_n, Y_n) \xrightarrow{p} (X, Y) \iff X_n \xrightarrow{p} X \text{ and } Y_n \xrightarrow{p} Y.$$

Definition 4.2 (Consistency). Let $\hat{\theta}_n \equiv \hat{\theta}_n(X_1, \dots, X_n)$ be an estimator of a parameter θ , $\theta \in \Theta \subset \mathbb{R}$. $\hat{\theta}_n$ is said to be consistent for θ if $\hat{\theta}_n \xrightarrow{p} \theta$, $\forall \theta \in \Theta$.

In the case of *multiple parameters*, the consistency of a parameter vector estimator is defined as the consistency of each of its components. Thus, for example, we say $\hat{\theta} = (\hat{\theta}_{1n}, \hat{\theta}_{2n})$ is consistent for $\theta = (\theta_1, \theta_2)$ if $\hat{\theta}_{1n}$ is consistent for θ_1 , and $\hat{\theta}_{2n}$ is consistent for θ_2 .

The following theorem states one of the most useful result for verifying the consistency of an estimator provided its MSE exists.

Theorem 4.1 (Consistent in mean square error). $\hat{\theta}_n$ is a consistent estimator for θ if

$$MSE_{\theta}(\hat{\theta}_n) := E_{\theta}(\hat{\theta}_n - \theta)^2 \xrightarrow{n \rightarrow \infty} 0, \forall \theta \in \Theta.$$

Note that, by definition, $MSE(\hat{\theta}_n) = Bias^2(\hat{\theta}_n) + Var(\hat{\theta}_n)$. Thus, this theorem can be equivalently restated as $E(\hat{\theta}_n) \rightarrow \theta$ and $Var(\hat{\theta}_n) \rightarrow 0$ or as $E(\hat{\theta}_n) \rightarrow \theta$ and $E(\hat{\theta}_n^2) \rightarrow \theta^2$. An estimator such that $MSE_{\theta}(\hat{\theta}_n) \rightarrow 0$ is referred to as *consistent in mean square* or **MSE-consistent**. And so the above theorem states that the MSE-consistency implies consistency.

Example 4.2.

Let X_1, \dots, X_n , be an iid sample from $Unif(0, \theta)$, $\theta > 0$. We have seen previously that the cdf of $X_{(n)}$ is given by

$$F_{X_{(n)}}(x) = \begin{cases} 0, & \text{if } x < 0 \\ (x/\theta)^n, & \text{if } 0 \leq x \leq \theta \\ 1, & \text{if } x > \theta. \end{cases}$$

So, for any $\epsilon > 0$,

$$P(|X_{(n)} - \theta| \geq \epsilon) = P(X_{(n)} \leq \theta - \epsilon) = \begin{cases} 0 & \text{if } \epsilon > \theta \\ (1 - \frac{\epsilon}{\theta})^n \rightarrow 0 & \text{if } \epsilon \leq \theta. \end{cases}$$

Thus, we conclude that $X_{(n)}$ is a consistent estimator for θ .

Instead of using the definition, we can prove the consistency of $X_{(n)}$ by using the above Theorem 4.1. In fact, we have seen that $E(X_{(n)}) = \frac{n}{n+1}\theta$ and $E(X_{(n)}^2) = \frac{n}{n+2}\theta^2$. And since $E(X_{(n)}) \rightarrow \theta$ and $E(X_{(n)}^2) \rightarrow \theta^2$, $X_{(n)}$ is a MSE-consistent estimator for θ . Another way to get to this result is to observe that

$$MSE(X_{(n)}) = \frac{2\theta^2}{(n+1)(n+2)} \rightarrow 0.$$

Another estimator for θ that we previously studied is $\hat{\theta}_1 = 2\bar{X}_n$. We have seen that

$$MSE(\hat{\theta}_1) = Var(\hat{\theta}_1) = \frac{\theta^2}{3n} \rightarrow 0.$$

So, $\hat{\theta}_1$ is also a consistent for θ .

But clearly $X_{(n)}$ is better as its MSE converges to 0 faster than the one of $\hat{\theta}_1$ (rate of $1/n^2$ vs $1/n$). \square

Statisticians are often interested in the limiting/asymptotic properties of estimators that can be expressed as *arithmetic means or functions of means*. The most basic and most famous result of this type is the following.

Theorem 4.2 (Weak Law of Large Numbers (WLLN)). *Suppose $X_i, i = 1, \dots, n$, is a sequence of iid random variables with finite mean μ , then $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i \xrightarrow{p} \mu$.*

It is important to note that for the above result to hold, it is necessary that $\mu = E(X_i)$ exists and is finite.

Example 4.3 (Consistency of the empirical distribution function). Let $X_i, i = 1, \dots, n$, be an iid sample from a cdf $F(x)$. The empirical distribution function is given by $F_n(x) = n^{-1} \sum_{i=1}^n I(X_i \leq x)$. By the WLLN, we have that

$$F_n(x) \xrightarrow{p} E(I(X_1 \leq x)) = F(x). \square$$

Example 4.4 (Consistency of the sample variance). Let $X_i, i = 1, \dots, n$, be an iid sample with $\mu = E(X_1)$ and $\sigma^2 = Var(X_1)$. The empirical variance is given by

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n}{n-1} \left(\frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - (\bar{X}_n - \mu)^2 \right),$$

By the WLLN, $\frac{1}{n} \sum_i (X_i - \mu)^2 \xrightarrow{p} \sigma^2$, and $\bar{X}_n \xrightarrow{p} \mu$, the fact that the convergence in probability is closed under the arithmetic operations, and $\frac{n}{n-1} \rightarrow 1$, we get that $S_n^2 \xrightarrow{p} \sigma^2$. \square

4.2 Convergence in law and asymptotic distribution

Convergence in law, also called *convergence in distribution*, is another widely used concept in statistical inference. The definition is given here for the bivariate case, the same applies to any dimension (univariate or multivariate).

Definition 4.3. Let (X_n, Y_n) , $n = 1, 2, \dots$, be a sequence of rv's with cdf $F_n(x, y) = P(X_n \leq x, Y_n \leq y)$. If there exists a cdf $F(x, y)$ such that

$$F_n(x, y) \xrightarrow{n \rightarrow \infty} F(x, y), \quad \forall (x, y) \text{ at which } F \text{ is continuous}$$

then F is called the *limiting cdf* of (X_n, Y_n) .

Letting (X, Y) have the cdf F , i.e. $F(x, y) = P(X \leq x, Y \leq y)$, then we say that (X_n, Y_n) converges in distribution (or converges in law) to (X, Y) , and we denote this by $(X_n, Y_n) \xrightarrow{d} (X, Y)$ or $(X_n, Y_n) \xrightarrow{L} (X, Y)$.

The convergence in distribution is a property of the distribution of X_n rather than X_n itself. $X_n \xrightarrow{d} X$ means that $P(X_n \leq x) \approx P(X \leq x)$, for large n , but *this does not imply convergence in probability*, i.e. $X_n \xrightarrow{d} X \not\Rightarrow X_n \xrightarrow{p} X$. In other words, $X_n \xrightarrow{d} X$ does not imply that X_n gets closer to X as n goes to infinity.

Example 4.5. Let $X \sim \text{Ber}(1/2)$, i.e. $P(X = 0) = P(X = 1) = 1/2$, and let $X_n = 1 + \frac{1}{n} - X$, $n \geq 1$. We have that

$$F_n(y) = P(X_n \leq x) = P(1 - X \leq x - 1/n) = \begin{cases} 0, & \text{if } x < 1/n \\ 1/2, & \text{if } 1/n \leq x < 1 + 1/n \\ 1, & \text{if } x \geq 1 + 1/n, \end{cases}$$

$$\xrightarrow{n \rightarrow \infty} \begin{cases} 0, & \text{if } x \leq 0 \\ 1/2, & \text{if } 0 < x \leq 1 \\ 1, & \text{if } x > 1. \end{cases}$$

This limit function is continuous everywhere but at $x = 0$ and $x = 1$, and with the exception of these two points, it coincides with the cdf of $\text{Ber}(1/2)$. So, by definition, we can write that $X_n \xrightarrow{d} \text{Ber}(1/2)$ or, equivalently, that $X_n \xrightarrow{d} X$. Observe that $|X_n - X| = |\pm 1 + 1/n| \geq 1/2$, $\forall n \geq 2$. Hence, $X_n \not\xrightarrow{p} X$. \square

Example 4.6. Let X_n be a rv with cdf $F_n(x) = (1 - 1/x^n)I(x \geq 1)$. Put $Y_n = nX_n - n$. We have that

$$G_n(y) := P(Y_n \leq y) = P(X_n \leq 1 + y/n) = \left(1 - \frac{1}{(1 + y/n)^n}\right) I(y \geq 0).$$

L'Hôpital's rule can be used to verify that $(1 + y/n)^n \rightarrow e^y$. So,

$$G_n(y) \xrightarrow{n \rightarrow \infty} \begin{cases} 0, & \text{if } y < 0 \\ 1 - e^{-y}, & \text{if } y \geq 0. \end{cases}$$

This latter is the cdf of the **exponential distribution** with rate 1. We therefore conclude that $Y_n \xrightarrow{d} \text{Expo}(1)$. \square

Example 4.7. Let X_i , $i = 1, \dots, n$, be an iid sample from $\text{Unif}[0, 1]$. Put $Y_n = nX_{(1)}$, we have that

$$P(Y_n \leq y) = 1 - (1 - P(X_1 \leq y/n))^n = \begin{cases} 1 - (1 - 0)^n = 0, & \text{if } y < 0 \\ 1 - (1 - y/n)^n, & \text{if } 0 \leq y < n \\ 1 - (1 - 1)^n = 1, & \text{if } y \geq n. \end{cases}$$

$$\xrightarrow{n \rightarrow \infty} \begin{cases} 0, & \text{if } y < 0 \\ 1 - e^{-y}, & \text{if } y \geq 0. \end{cases}$$

So, $Y_n \xrightarrow{d} \text{Expo}(1)$. \square

Joint versus marginal convergence in distribution

- If $(X_n, Y_n) \xrightarrow{d} (X, Y)$, then $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{d} Y$.
- But, unlike convergence in probability, marginal convergence in distribution does not (generally) imply joint convergence in distribution. In other words,

$$X_n \xrightarrow{d} X \text{ and } Y_n \xrightarrow{d} Y \not\Rightarrow (X_n, Y_n) \xrightarrow{d} (X, Y).$$

- Sometimes, marginal convergence in distribution does imply joint convergence in distribution. This is for example the case if X (or Y) is a constant or $X_n \perp\!\!\!\perp Y_n$, $\forall n$.

Example 4.8. Let $X_n = X + 1/n$ with $X \sim N(0, 1)$, and $Y_n = -X_n$. Clearly, $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{d} X$, but $(X_n, Y_n) \not\xrightarrow{d} (X, X)$. \square

Convergence in probability is stronger than convergence in distribution. Thus,

$$X_n \xrightarrow{p} X \Rightarrow X_n \xrightarrow{d} X.$$

A special case of convergence in distribution occurs when the limiting distribution degenerates, thus when $F(x) = I(x \geq c)$, for some $c \in \mathbb{R}$, i.e. X is constant with $P(X = c) = 1$. In this case, X_n converges in distribution to a constant, and we write $X_n \xrightarrow{d} c$, i.e. $P(X_n \leq x) \rightarrow$

1, $\forall x > c$ and $P(X_n \leq x) \rightarrow 0, \forall x < c$. It can be shown that

$$X_n \xrightarrow{p} c \Leftrightarrow X_n \xrightarrow{d} c.$$

Example 4.9. Let $X_n \sim \text{Unif}(0, 1/n)$. We have that

$$\begin{aligned} P(X_n \leq x) &= \begin{cases} 0, & \text{if } x < 0 \\ nx, & \text{if } 0 \leq x < 1/n \\ 1, & \text{if } x \geq 1/n, \end{cases} \\ &\xrightarrow{n \rightarrow \infty} \begin{cases} 0, & \text{if } x \leq 0 \\ 1, & \text{if } x > 0. \end{cases} \\ &= I(x \geq 0), \forall x \neq 0. \end{aligned}$$

So, $X_n \xrightarrow{d} 0$. Thus, $X_n \xrightarrow{p} 0$. In fact, $P(|X_n - 0| \geq \epsilon) = P(X_n \geq \epsilon) = 1 - P(X_n \leq \epsilon) = 0, \forall n \geq 1/\epsilon$. \square

Convergence in probability or in distribution does not imply convergence in mean. Thus,

$$X_n \xrightarrow{p \text{ or } d} X \not\Rightarrow E(X_n) \rightarrow E(X).$$

Example 4.10. Let X_n be a sequence of rv such that $P(X_n = 0) = 1 - 1/n$ and $P(X_n = n) = 1/n, n = 1, 2, \dots$. For any $\epsilon > 0$, we have that $P(|X_n| \geq \epsilon) = P(X_n = n) = 1/n \rightarrow 0$. So $X_n \xrightarrow{p} 0$. But $E(X_n) = 1, \forall n$. \square

A sufficient condition to ensure that $X_n \xrightarrow{p \text{ or } d} X \Rightarrow E(X_n) \rightarrow E(X)$ is that $\sup_n E(|X_n|^{1+\delta}) < \infty$, for some $\delta > 0$.

The behavior of the sample mean is very important, and especially its limiting distribution. In this respect, one of the most remarkable theorems in statistics is the central limit theorem (CLT).

Theorem 4.3 (CLT).

- (Univariate case). Let $X_i, i = 1, \dots, n$, be an iid rv's with mean μ and with variance $\sigma^2 \in (0, \infty)$. Then

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2).$$

- (Multivariate case). Let $\mathbf{X}_i, i = 1, \dots, n$, be a d -dimensional iid rve's with mean $\boldsymbol{\mu}$ and with a variance-covariance matrix $\text{Var}(\mathbf{X}_1) = \boldsymbol{\Sigma}$, where $\boldsymbol{\Sigma}$ has finite entries. Then

$$\sqrt{n}(\bar{\mathbf{X}}_n - \boldsymbol{\mu}) \xrightarrow{d} N_d(\mathbf{0}, \boldsymbol{\Sigma}).$$

In particular, for the case of bivariate iid sample $(X_i, Y_i), i = 1, \dots, n$, say from sample of

(X, Y) . The CLT states that

$$\sqrt{n} \left(\begin{pmatrix} \bar{X}_n \\ \bar{Y}_n \end{pmatrix} - \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \right) = \begin{pmatrix} \sqrt{n}(\bar{X}_n - \mu_1) \\ \sqrt{n}(\bar{Y}_n - \mu_2) \end{pmatrix} \xrightarrow{d} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} \right),$$

where $\mu_1 = E(X)$, $\mu_2 = E(Y)$, $\sigma_1^2 = \text{Var}(X)$, $\sigma_2^2 = \text{Var}(Y)$, and $\sigma_{12} = \text{Cov}(X, Y)$.

Since joint convergence in distribution implies marginal convergence, the result above implies that $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma_1^2)$, which is the univariate CLT. In other words, Multivariate CLT \Rightarrow Univariate CLT.

Also, it is interesting to note the fact that $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$ implies that $\bar{X}_n \xrightarrow{p} \mu$. So CLT is stronger than WLLN. $\bar{X}_n \xrightarrow{p} \mu$ means that \bar{X}_n gets closer to μ as n increases but it says nothing about how \bar{X}_n varies around μ which is what CLT is about.

Asymptotic distribution.

According to the (univariate) CLT, using the definition of convergence in distribution, with some abuse of notations, we can write that

$$P(\bar{X}_n \leq x) \approx P(N(\mu, \sigma^2/n) \leq x), \text{ for large } n.$$

$N(\mu, \sigma^2/n)$ is called the **asymptotic distribution** of \bar{X}_n , and we write this down as

$$\bar{X}_n \sim_a N(\mu, \sigma^2/n), \text{ or } \bar{X}_n \sim AN(\mu, \sigma^2/n).$$

Example 4.11. Let $X_i, i = 1, \dots, n$, be an iid sample from $\text{Pois}(\lambda)$ with pd $f(x; \lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$, $x = 0, 1, \dots$, and $\lambda > 0$. We have that $E(X_1) = \text{Var}(X_1) = \lambda$, $E(I(X_1 = 0)) = e^{-\lambda}$, $\text{Var}(I(X_1 = 0)) = e^{-\lambda}(1 - e^{-\lambda})$, and $\text{Cov}(X_1, I(X_1 = 0)) = -\lambda e^{-\lambda}$. Let's find the asymptotic distribution of \bar{X}_n , the asymptotic distribution of \bar{Z}_n , where $Z_i = I(X_i = 0)$, and the joint asymptotic distribution of (\bar{X}_n, \bar{Z}_n) .

By the (univariate) CLT, $\sqrt{n}(\bar{X}_n - \lambda) \xrightarrow{d} N(0, \lambda)$ and $\sqrt{n}(\bar{Z}_n - e^{-\lambda}) \xrightarrow{d} N(0, e^{-\lambda}(1 - e^{-\lambda}))$. Thus, $\bar{X}_n \sim_a N(\lambda, \lambda/n)$ and $\bar{Z}_n \sim_a N(e^{-\lambda}, e^{-\lambda}(1 - e^{-\lambda})/n)$.

By the (multivariate) CLT

$$\begin{pmatrix} \sqrt{n}(\bar{X}_n - \lambda) \\ \sqrt{n}(\bar{Z}_n - e^{-\lambda}) \end{pmatrix} \xrightarrow{d} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \lambda & -\lambda e^{-\lambda} \\ -\lambda e^{-\lambda} & e^{-\lambda}(1 - e^{-\lambda}) \end{pmatrix} \right). \quad (4.1)$$

Thus,

$$\begin{pmatrix} \bar{X}_n \\ \bar{Z}_n \end{pmatrix} \sim_a N_2 \left(\begin{pmatrix} \lambda \\ e^{-\lambda} \end{pmatrix}, n^{-1} \begin{pmatrix} \lambda & -\lambda e^{-\lambda} \\ -\lambda e^{-\lambda} & e^{-\lambda}(1 - e^{-\lambda}) \end{pmatrix} \right).$$

Equation (4.1) can also be written as

$$\begin{pmatrix} \sqrt{n}(\bar{X}_n - \lambda) \\ \sqrt{n}(\bar{Z}_n - e^{-\lambda}) \end{pmatrix} \xrightarrow{d} \begin{pmatrix} U \\ V \end{pmatrix}, \text{ with } \begin{pmatrix} U \\ V \end{pmatrix} \sim N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \lambda & -\lambda e^{-\lambda} \\ -\lambda e^{-\lambda} & e^{-\lambda}(1 - e^{-\lambda}) \end{pmatrix} \right)$$

□

4.3 Tools for proving asymptotic results

There are many techniques available to check the consistency of an estimator and to find its asymptotic distribution. We present the main useful techniques below.

4.3.1 Continuous mapping theorem

Theorem 4.4 (Continuous Mapping Theorem (CMT)). *Let X_n be a sequence of d -dimensional random vectors and X a d -dimensional random vector. Let $g : \mathbb{R}^d \rightarrow \mathbb{R}^p$ be a function **continuous** on \mathbb{R}^d except maybe in a subset I for which $P(X \in I) = 0$. Then,*

$$X_n \rightarrow X \Rightarrow g(X_n) \rightarrow g(X).$$

Above, \rightarrow can be either the convergence in probability (\xrightarrow{p}) or in distribution (\xrightarrow{d}).

The CMT allows the function g to be discontinuous but the probability of X being at a discontinuity point of g must be zero. For example, suppose $X_n \rightarrow X$. Since the function $x \mapsto 1/x$ is discontinuous at 0, we can apply the CMT to conclude that $1/X_n \rightarrow 1/X$, provided that $P(X = 0) = 0$.

CMT implies that the convergence in probability (or in distribution) is closed under all arithmetic operations. So, if $(X_n, Y_n) \rightarrow (X, Y)$, then $aX_n + bY_n \rightarrow aX + bY$, $\forall a, b$, $X_n Y_n \rightarrow XY$, and $X_n/Y_n \rightarrow X/Y$, if $P(Y = 0) = 0$.

Example 4.12.

From the examples studied above, we can conclude that $S_n \xrightarrow{p} \sigma$, $S_n/\bar{X}_n \xrightarrow{p} \sigma/\mu$ (if $\mu \neq 0$), $n \frac{(\bar{X}_n - \mu)^2}{\sigma^2} \xrightarrow{d} \chi_1^2$, and, from Example 4.11, that $\frac{\bar{X}_n - \lambda}{\bar{Z}_n - e^{-\lambda}} \xrightarrow{d} \frac{U}{V}$, where (U, V) is bivariate-normal with mean and variance as defined above. □

4.3.2 Slutsky's theorem

Theorem 4.5 (Slutsky). *If $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{p} C$, where C is a constant, then*

$$\begin{pmatrix} X_n \\ Y_n \end{pmatrix} \xrightarrow{d} \begin{pmatrix} X \\ C \end{pmatrix}$$

This theorem along with the CMT implies that if $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{p} C$, then $X_n \pm Y_n \xrightarrow{d} X \pm C$, $X_n Y_n \xrightarrow{d} XC$, and $Y_n^{-1} X_n \xrightarrow{d} C^{-1}X$ (when $C \neq 0$). People often call such results (CMT + Slutsky) “Slutsky’s theorem”.

Example 4.13.

- By CLT we have that, $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$. And we known that $S_n^2 \xrightarrow{p} \sigma^2$. So, by CMT + Slutsky, we deduce that

$$\sqrt{n} \frac{\bar{X}_n - \mu}{S_n} \xrightarrow{d} N(0, 1).$$

- Let $X_i, i = 1, \dots, n$ be an iid sample from $Ber(\pi)$. Put $\hat{\pi}_n = n^{-1} \sum_{i=1}^n X_i$. The same reasoning as above leads to

$$\sqrt{n} \frac{\hat{\pi}_n - \pi_n}{\sqrt{\hat{\pi}_n(1 - \hat{\pi}_n)}} \xrightarrow{d} N(0, 1).$$

- In Example 4.4, we have seen that $S_n^2 = \frac{n}{n-1} \hat{\sigma}_n^2$, where $\hat{\sigma}_n^2 := n^{-1} \sum_i (X_i - \bar{X}_n)^2 = \tilde{\sigma}_n^2 - (\bar{X}_n - \mu)^2$, with $\tilde{\sigma}_n^2 = n^{-1} \sum_i (X_i - \mu)^2$. By CLT, we known that

$$\begin{aligned} \sqrt{n}(\bar{X}_n - \mu) &\xrightarrow{d} N(0, \sigma^2), \text{ and} \\ \sqrt{n}(\tilde{\sigma}_n^2 - \sigma^2) &\xrightarrow{d} N(0, v^2), \text{ with } v^2 := \text{Var}(X_1 - \mu)^2. \end{aligned}$$

Hence, $\sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) = \sqrt{n}(\tilde{\sigma}_n^2 - \sigma^2) - \sqrt{n}(\bar{X}_n - \mu)^2 \xrightarrow{d} N(0, v^2)$. And since,

$$\sqrt{n}(S_n^2 - \sigma^2) = \frac{n}{n-1} \sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) + \frac{1}{\sqrt{n}} \frac{n}{n-1} \sigma^2,$$

we conclude that $\sqrt{n}(S_n^2 - \sigma^2) \xrightarrow{d} N(0, v^2)$. \square

4.3.3 Delta method

Theorem 4.6 (Delta method). *Let $g : \mathbb{R}^d \mapsto \mathbb{R}^p$ such that \dot{g} is continuous in a neighborhood of θ and let $a_n \xrightarrow[n \rightarrow \infty]{} \infty$. Then,*

$$a_n(\hat{\theta}_n - \theta) \xrightarrow{d} W \Rightarrow a_n(g(\hat{\theta}_n) - g(\theta)) \xrightarrow{d} \dot{g}(\theta)W,$$

where \dot{g} is the $p \times d$ Jacobian matrix of g , whose (i, j) -th element is $\partial_{\theta_j} g_i(\theta)$. In particular, if $W \sim N_d(\mu, \Sigma)$, then $a_n(g(\hat{\theta}_n) - g(\theta)) \xrightarrow{d} N_p(\dot{g}(\theta)\mu, \dot{g}(\theta)\Sigma\dot{g}^t(\theta))$, where $\dot{g}^t(\theta)$ denote $[\dot{g}(\theta)]^t$.

As a consequence, we have the following results:

- **Scalar case with real-valued function:** $\theta \in \mathbb{R}$ and $g(\theta) : \mathbb{R} \mapsto \mathbb{R}$. Suppose that

$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, \sigma^2)$. Then

$$\sqrt{n}(g(\hat{\theta}) - g(\theta)) \xrightarrow{d} g'(\theta) \times N(0, \sigma^2) = N\left(0, (g'(\theta))^2 \sigma^2\right).$$

- **Scalar case with vector-valued function:** $\theta \in \mathbb{R}$ and $g(\theta) := (g_1(\theta), g_2(\theta)) : \mathbb{R} \mapsto \mathbb{R}^2$. Suppose that $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, \sigma^2)$. Then

$$\sqrt{n} \begin{pmatrix} g_1(\hat{\theta}) - g_1(\theta) \\ g_2(\hat{\theta}) - g_2(\theta) \end{pmatrix} \xrightarrow{d} \begin{pmatrix} \partial_{\theta} g_1(\theta) \\ \partial_{\theta} g_2(\theta) \end{pmatrix} \times N(0, \sigma^2) = N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \sigma^2 \mathbf{\Lambda} \right),$$

$$\text{where } \mathbf{\Lambda} = \begin{pmatrix} \partial_{\theta} g_1(\theta) \\ \partial_{\theta} g_2(\theta) \end{pmatrix} \begin{pmatrix} \partial_{\theta} g_1(\theta) & \partial_{\theta} g_2(\theta) \end{pmatrix} = \begin{pmatrix} (\partial_{\theta} g_1(\theta))^2 & \partial_{\theta} g_1(\theta) \partial_{\theta} g_2(\theta) \\ \partial_{\theta} g_1(\theta) \partial_{\theta} g_2(\theta) & (\partial_{\theta} g_2(\theta))^2 \end{pmatrix}.$$

- **Bivariate case with real-valued function:** $\theta := (\theta_1, \theta_2) \in \mathbb{R}^2$ and $g(\theta) = g(\theta_1, \theta_2) : \mathbb{R}^2 \mapsto \mathbb{R}$. Suppose that $\sqrt{n}(\hat{\theta} - \theta) = \sqrt{n} \begin{pmatrix} \hat{\theta}_1 - \theta_1 \\ \hat{\theta}_2 - \theta_2 \end{pmatrix} \xrightarrow{d} N_2(\mathbf{0}, \mathbf{\Sigma})$, where $\mathbf{\Sigma} = [\sigma_{ij}]_{i,j=1,2}$. Then,

$$\sqrt{n}(g(\hat{\theta}) - g(\theta)) \xrightarrow{d} \begin{pmatrix} \partial_{\theta_1} g(\theta) & \partial_{\theta_2} g(\theta) \end{pmatrix} \times N_2(\mathbf{0}, \mathbf{\Sigma}) = N(0, \sigma^2),$$

$$\text{where } \sigma^2 = (\partial_{\theta_1} g(\theta))^2 \sigma_{11} + (\partial_{\theta_2} g(\theta))^2 \sigma_{22} + 2 \partial_{\theta_2} g(\theta) \partial_{\theta_1} g(\theta) \sigma_{12}.$$

- **Bivariate case with vector-valued function:** $\theta := (\theta_1, \theta_2) \in \mathbb{R}^2$ and $g(\theta) := (g_1(\theta), g_2(\theta)) : \mathbb{R}^2 \mapsto \mathbb{R}^2$. Suppose that $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N_2(\mathbf{0}, \mathbf{\Sigma})$. Then,

$$\sqrt{n} \begin{pmatrix} g_1(\hat{\theta}) - g_1(\theta) \\ g_2(\hat{\theta}) - g_2(\theta) \end{pmatrix} \xrightarrow{d} \dot{g}(\theta) \times N_2(\mathbf{0}, \mathbf{\Sigma}) = N_2(\mathbf{0}, \mathbf{\Lambda}),$$

$$\text{where } \mathbf{\Lambda} = \dot{g}(\theta) \mathbf{\Sigma} \dot{g}^t(\theta), \text{ with } \dot{g}(\theta) = \begin{pmatrix} \partial_{\theta_1} g_1(\theta) & \partial_{\theta_2} g_1(\theta) \\ \partial_{\theta_1} g_2(\theta) & \partial_{\theta_2} g_2(\theta) \end{pmatrix}.$$

Example 4.14.

- Let $X_i, i = 1, \dots, n$, be an iid sample from $Pois(\lambda)$ with pd $f(x, \lambda) = \frac{\lambda^x}{x!} e^{-\lambda}, x = 0, 1, \dots$, and $\lambda > 0$. By CLT, $\sqrt{n}(\bar{X}_n - \lambda) \xrightarrow{d} N(0, \lambda)$. So, by Delta method,

$$\sqrt{n}(e^{-\bar{X}_n} - e^{-\lambda}) \xrightarrow{d} N\left(0, \lambda((e^{-\lambda})')^2\right) = N\left(0, \lambda e^{-2\lambda}\right)$$

- Let $X_i, i = 1, \dots, n$ be an iid sample from $Ber(\pi)$. Let $\hat{\pi} \equiv \hat{\pi}_n = n^{-1} \sum_i X_i$. By CLT, we know that $\sqrt{n}(\hat{\pi} - \pi) \xrightarrow{d} N(0, \pi(1 - \pi))$. Put $O = \frac{\pi}{1 - \pi}$ (this is called the odds) and $\hat{O} = \frac{\hat{\pi}}{1 - \hat{\pi}}$. By Delta method,

$$\sqrt{n}(\hat{O} - O) \xrightarrow{d} N\left(0, \pi(1 - \pi) \left(\left(\frac{\pi}{1 - \pi} \right)' \right)^2 \right) = N\left(0, \frac{\pi}{(1 - \pi)^3}\right) = N(0, O(1 + O)^2).$$

- Let $X_i, i = 1, \dots, n$, be an iid sample from $Pois(\lambda)$ with pd $f(x, \lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$, $x = 0, 1, \dots$, and $\lambda > 0$. By CLT, $\sqrt{n}(\bar{X}_n - \lambda) \xrightarrow{d} N(0, \lambda)$. So, by the Delta method, with $g : \lambda \mapsto (\lambda, e^{-\lambda})$,

$$\sqrt{n} \begin{pmatrix} \bar{X}_n - \lambda \\ e^{-\bar{X}_n} - e^{-\lambda} \end{pmatrix} \xrightarrow{d} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \lambda \begin{pmatrix} 1 \\ -e^{-\lambda} \end{pmatrix} \begin{pmatrix} 1 & -e^{-\lambda} \end{pmatrix} \right) = N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \lambda & -\lambda e^{-\lambda} \\ -\lambda e^{-\lambda} & \lambda e^{-2\lambda} \end{pmatrix} \right).$$

- Let $(X_i, Y_i), i = 1, \dots, n$, be an iid sample of (X, Y) . Let $\mu_1 = E(X)$, $\mu_2 = E(Y)$, $\sigma_1^2 = Var(X)$, $\sigma_2^2 = Var(Y)$, and $\sigma_{12} = Cov(X, Y)$. By CLT

$$\sqrt{n} \left(\begin{pmatrix} \bar{X}_n \\ \bar{Y}_n \end{pmatrix} - \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \right) \xrightarrow{d} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} \right)$$

So, by the Delta method, with $g : (\mu_1, \mu_2) \mapsto \mu_1 - \mu_2$,

$$\sqrt{n}((\bar{X}_n - \bar{Y}_n) - (\mu_1 - \mu_2)) \xrightarrow{d} N(0, \sigma_1^2 + \sigma_2^2 - 2\sigma_{12}).$$

Using again the Delta method, but this time with $g : (\mu_1, \mu_2) \mapsto \mu_1/\mu_2$, we get

$$\sqrt{n}(\bar{X}_n/\bar{Y}_n - \mu_1/\mu_2) \xrightarrow{d} N(0, v^2),$$

$$\text{with } v^2 = \frac{1}{\mu_2^2} \left(\sigma_1^2 + \frac{\mu_1^2}{\mu_2^2} \sigma_2^2 - 2\frac{\mu_1}{\mu_2} \sigma_{12} \right).$$

- Let $X_i, i = 1, \dots, n$, be an iid sample from $Ber(\pi_1)$ and $Y_i, i = 1, \dots, n$, be another iid sample from $Ber(\pi_2)$. Suppose that these two samples are independent. Let $\hat{\pi}_1 = n^{-1} \sum_i X_i$, and $\hat{\pi}_2 = n^{-1} \sum_i Y_i$. By CLT,

$$\sqrt{n} \left(\begin{pmatrix} \hat{\pi}_1 \\ \hat{\pi}_2 \end{pmatrix} - \begin{pmatrix} \pi_1 \\ \pi_2 \end{pmatrix} \right) \xrightarrow{d} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \pi_1(1-\pi_1) & 0 \\ 0 & \pi_2(1-\pi_2) \end{pmatrix} \right)$$

So, by the Delta method, with $g : (\pi_1, \pi_2) \mapsto (\pi_1 - \pi_2, \pi_1/\pi_2)$,

$$\sqrt{n} \left(\begin{pmatrix} \hat{\pi}_1 - \hat{\pi}_2 \\ \hat{\pi}_1/\hat{\pi}_2 \end{pmatrix} - \begin{pmatrix} \pi_1 - \pi_2 \\ \pi_1/\pi_2 \end{pmatrix} \right) \xrightarrow{d} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Lambda \right),$$

where

$$\Lambda = \begin{pmatrix} 1 & -1 \\ 1/\pi_2 & -\pi_1/\pi_2^2 \end{pmatrix} \begin{pmatrix} \pi_1(1-\pi_1) & 0 \\ 0 & \pi_2(1-\pi_2) \end{pmatrix} \begin{pmatrix} 1 & 1/\pi_2 \\ -1 & -\pi_1/\pi_2^2 \end{pmatrix} = \begin{pmatrix} \pi_1(1-\pi_1) + \pi_2(1-\pi_2) & \frac{\pi_1(2-\pi_1-\pi_2)}{\pi_2} \\ \frac{\pi_1(2-\pi_1-\pi_2)}{\pi_2} & \frac{\pi_1(\pi_1-\pi_2)}{\pi_2^2} \end{pmatrix}$$

□

4.4 Asymptotic efficiency

In general, for a given estimator $\hat{\theta}_n$ of a parameter θ , if one can show that $\frac{\hat{\theta}_n - \theta_n}{\sigma_n} \xrightarrow{d} N(0, 1)$, then we say that $\hat{\theta}_n$ is *asymptotically distributed as* $N(\theta_n, \sigma_n^2)$, or that $\hat{\theta}_n$ is **asymptotically normal** with **asymptotic mean** $Amean(\hat{\theta}_n) := \theta_n$ and **asymptotic variance** $Avar(\hat{\theta}_n) := \sigma_n^2$, and write

$$\hat{\theta}_n \sim_a N(\theta_n, \sigma_n^2), \text{ or } \hat{\theta}_n \sim AN(\theta_n, \sigma_n^2).$$

The quantity $Abias(\hat{\theta}_n) := \theta_n - \theta$ is called the **asymptotic bias**. If $\theta_n = \theta$, then we say that $\hat{\theta}_n$ is *asymptotically unbiased*.

For two asymptotically unbiased estimators of θ , the *asymptotic relative efficiency* of $\hat{\theta}_1$ to $\hat{\theta}_2$ is the ratio of the asymptotic variance of $\hat{\theta}_2$ to the asymptotic variance of $\hat{\theta}_1$:

$$ARE(\hat{\theta}_1, \hat{\theta}_2) = \frac{Avar(\hat{\theta}_2)}{Avar(\hat{\theta}_1)}.$$

Example 4.15. Let $X_i, i = 1, \dots, n$, be an iid sample from $Pois(\lambda)$ with pd $f(x, \lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$, $x = 0, 1, \dots$, and $\lambda > 0$. Let's say that our parameter of interest is $\delta := P(X = 0) = e^{-\lambda} \in (0, 1)$. Put $\hat{\delta}_1 = n^{-1} \sum_{i=1}^n I(X_i = 0)$, and $\hat{\delta}_2 = e^{-\bar{X}_n}$. By CLT,

$$\sqrt{n} (\hat{\delta}_1 - \delta) \xrightarrow{d} N(0, \delta(1 - \delta)).$$

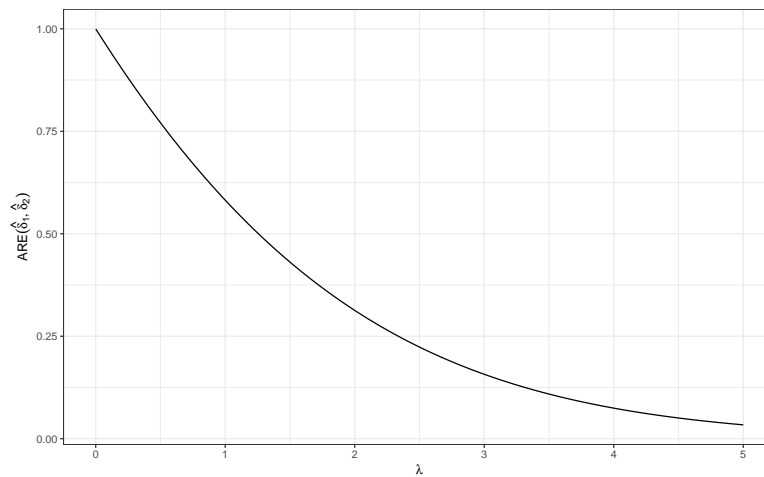
Thus, $\hat{\delta}_1 \sim_a N(\delta, e^{-\lambda}(1 - e^{-\lambda})/n)$. On the other hand, we have that (see Example 4.14)

$$\sqrt{n} (e^{-\bar{X}_n} - e^{-\lambda}) \xrightarrow{d} N(0, \lambda e^{-2\lambda}).$$

Thus, $\hat{\delta}_2 \sim_a N(\delta, \lambda e^{-2\lambda}/n)$. We conclude that

$$ARE(\hat{\delta}_1, \hat{\delta}_2) = \frac{\lambda}{e^\lambda - 1}.$$

This is a strictly decreasing function of λ with a maximum of 1.



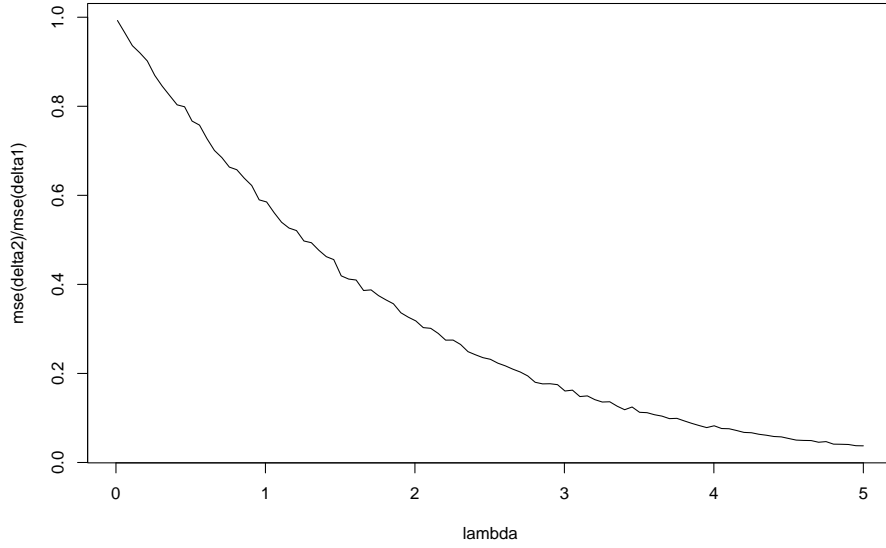
So, $\hat{\delta}_2$ is asymptotically uniformly more efficient than $\hat{\delta}_1$. And, the larger λ is, the better $\hat{\delta}_2$ is.

The following simulation confirms the calculations above (for $n = 100$).

```
REpois <- function(lambda, n, N) {
  delta <- exp(-lambda)
  hat <- replicate(N, {
    samp <- rpois(n, lambda)
    c(delta1 = mean(samp == 0), delta2 = exp(-mean(samp)))
  })
  mse(hat["delta2", ], delta) / mse(hat["delta1", ], delta) # Relative MSE
}

VecREpois <- REpois |> Vectorize(vectorize.args = "lambda")

curve(VecREpois(x, n = 100, N = 10^4),
      from = 0.01, to = 5, xlab = "lambda",
      ylab = "mse(delta2)/mse(delta1)"
)
```



□

We have seen that if $X_i, i = 1, \dots, n$, is an iid sample from a pd $f(x, \theta)$ and $\hat{\theta}$ is any *unbiased estimator* of θ , then

$$\text{Var}(\hat{\theta}) \geq \frac{1}{I_n(\theta)}. \quad (4.2)$$

Under some regularity conditions (see the previous chapter), this property holds *for any finite sample size* n . The bound $I_n^{-1}(\theta)$ is attainable if $f(x, \theta)$ belongs to the exponential family (see the CRLB Attainment theorem). Although this family is very rich, this considerably limits the applicability of such a result.

In asymptotic regime, under some some regularity conditions, it can be shown that for *any asymptotically normal and asymptotically unbiased estimator* $\hat{\theta}$ of θ ,

$$\text{Avar}(\hat{\theta}) \geq \frac{1}{I_n(\theta)}. \quad (4.3)$$

The regularity conditions that guarantee the validity of this asymptotic result are less restrictive than those required for the finite-sample variant (Equation (4.2)). Furthermore, the number of models/families for which the limit in (4.3) is attainable is much larger.

For an asymptotically normal and asymptotically unbiased estimator $\hat{\theta}$, we define the asymptotic efficiency by

$$\text{Aeff}(\hat{\theta}) = \frac{I_n^{-1}(\theta)}{\text{Avar}(\hat{\theta})}$$

If $\text{Aeff}(\hat{\theta}) = 1$, then $\hat{\theta}$ is said to be *asymptotically efficient* for θ . To put it another way, $\hat{\theta}$ is asymptotically efficient for θ if $\hat{\theta} \sim_a N(\theta, I_n^{-1}(\theta))$, i.e. if $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, I^{-1}(\theta))$.

The same concept applies to the case of multiple parameters, resulting in the following

general definition: $\hat{\theta}$ is asymptotically efficient for $\theta \in \mathbb{R}^d$ if

$$\hat{\theta} \sim_a N_d(\theta, I_n^{-1}(\theta)).$$

An interesting property of asymptotic efficiency is as follows

$\hat{\theta}$ is asymptotically efficient for $\theta \implies g(\theta)$ is asymptotically efficient for $g(\theta)$.

$g : \mathbb{R}^d \mapsto \mathbb{R}^p$ is any function with continuous Jacobian. This result is a direct consequence of the Delta method. In fact, if $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N_d(0, I^{-1}(\theta))$, then, by Delta method,

$$\sqrt{n}(g(\hat{\theta}) - g(\theta)) \xrightarrow{d} N_p(0, \dot{g}(\theta) I^{-1}(\theta) \dot{g}^t(\theta)).$$

Thus, $g(\hat{\theta}) \sim_a N_p(g(\theta), I_n^{-1}(g(\theta)))$.

Example 4.16. Let $X_i, i = 1, \dots, n$, be an iid sample from $Pois(\lambda)$. $f(x; \lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$, $x = 0, 1, \dots$, and $\lambda > 0$.

We have seen that $\sqrt{n}(\bar{X} - \lambda) \xrightarrow{d} N(0, \lambda)$ and that $I(\lambda) = 1/\lambda$. So, \bar{X} is asymptotically efficient for λ (actually, \bar{X} is efficient for λ). As consequence, we can, for example, say that $e^{-\bar{X}}$ is asymptotically efficient for $e^{-\lambda}$. In fact, by Delta method, $\sqrt{n}(e^{-\bar{X}} - e^{-\lambda}) \xrightarrow{d} N(0, \lambda)(e^{-\lambda})' = N(0, I^{-1}(e^{-\lambda}))$. Note that it can be shown that $e^{-\bar{X}}$ is not efficient for $e^{-\lambda}$.

The asymptotic efficiency of \bar{X} also implies that of, for example, $(\bar{X}, e^{-\bar{X}})$ as an estimator of $(\lambda, e^{-\lambda}) = (E(X), P(X = 0))$. \square

Chapter 5

Estimation methods

In this section, we examine some general methods that can be used to construct estimators that often have good properties.

5.1 The method of moments (MoM)

Let X_i , $i = 1, \dots, n$, be an iid sample of X . Let $\mu_j = E(X^j)$ denote the j -th moment of X , and $\hat{\mu}_j = n^{-1} \sum_{i=1}^n X_i^j$ the j -th sample moment.

To apply the method of moments to the problem of estimating a parameter θ , we need to be able to express θ as a function of μ_1, μ_2, \dots . Thus, we need to find a (known) function, say g , such that

$$\theta = g(\mu_1, \mu_2, \dots).$$

A simple estimation method consists in replacing the μ_j 's in the equation above by their empirical versions. This leads to the following MoM estimator:

$$\hat{\theta} = g(\hat{\mu}_1, \hat{\mu}_2, \dots).$$

Broadly speaking, MoM estimators do not necessarily offer the best performances, but they are typically easy to obtain and, under reasonable conditions, they are consistent and are asymptotically normal. In fact,

- By the WLLN + CMT, a MoM estimator is consistent provided (i) the population moments exist, and (ii) g is a continuous.
- By the CLT + Delta method, a MoM estimator can, typically, be shown to be asymptotically normal.

Example 5.1.

- Let X_i , $i = 1, \dots, n$, be an iid sample from the exponential distribution with pdf $f(x; \lambda) = \lambda e^{-\lambda x} I(x \geq 0)$, $\lambda > 0$. Since $E(X_1) = 1/\lambda$, the MoM estimator of λ is

$\hat{\lambda} = 1/\bar{X}_n$. This is a consistent estimator. Moreover, since $\text{Var}(X_1) = 1/\lambda^2$, $\hat{\lambda}$ is asymptotically normal with limiting distribution given by $N(\lambda, \lambda^2/n)$.

- Let $X_i, i = 1, \dots, n$, be an iid sample of X . Let $\mu_j = E(X^j)$, $\sigma_{jk} = \text{Cov}(X^j, X^k)$, $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$, and $\bar{X}_n^2 = n^{-1} \sum_{i=1}^n X_i^2$. Let's find the MoM estimator of $\sigma^2 \equiv \sigma_{11} = \text{Var}(X)$ and its asymptotic distribution. Since $\sigma^2 = \mu_2 - \mu_1^2$, the MoM estimator of σ^2 is

$$\hat{\sigma}_n^2 = \hat{\mu}_2 - \hat{\mu}_1^2 = \bar{X}_n^2 - \bar{X}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

By CLT,

$$\sqrt{n} \left(\begin{pmatrix} \bar{X}_n \\ \bar{X}_n^2 \end{pmatrix} - \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \right) \xrightarrow{d} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{pmatrix} \right)$$

So, by the Delta method, with $g(x, y) = y - x^2$,

$$\sqrt{n}(\hat{\sigma}_n^2 - \sigma^2) \xrightarrow{d} N(0, \nu^2),$$

where $\nu^2 := 4\mu_1^2\sigma_{11} + \sigma_{22} - 4\mu_1\sigma_{12} = \text{Var}(X - \mu_1)^2$.

- Let $Y_i, i = 1, \dots, n$, be an iid sample from the **Log-normal distribution** (i.e. $\log(Y) \sim N(\mu, \sigma^2)$) with parameters μ and σ^2 . Y_i has the pd

$$f(y; \mu, \sigma^2) = \frac{1}{y\sigma\sqrt{2\pi}} \exp\left(-\frac{(\log y - \mu)^2}{2\sigma^2}\right) I(y > 0), \mu \in (-\infty, \infty), \text{ and } \sigma > 0.$$

Since

$$\mu_1 := E(Y_1) = \exp(\mu + \sigma^2/2) \text{ and } \mu_2 := E(Y_1^2) = \exp(2\mu + 2\sigma^2),$$

the MoM estimator of μ and σ^2 are

$$\hat{\mu} = 2 \log(\hat{\mu}_1) - \frac{1}{2} \log(\hat{\mu}_2) \text{ and } \hat{\sigma}^2 = \log(\hat{\mu}_2) - 2 \log(\hat{\mu}_1).$$

Can you prove that $(\hat{\mu}, \hat{\sigma}^2)$ is consistent? What about asymptotic normality? \square

5.2 The method of maximum likelihood (ML)

The method of maximum likelihood is, by far, the most popular technique for deriving estimators and performing inference. It has an intuitive motivation and usually has fairly very good properties (at least asymptotically).

5.2.1 Likelihood: definition and meaning

Definition 5.1 (The likelihood function). Let $f_n(\cdot; \theta)$ denote the joint pd of the sample $\mathbf{X} = (X_1, \dots, X_n)$, where $\theta \in \Theta$ is the parameter of interest (θ may be vector valued; we don't bold it here for ease of notation).

Given that $\mathbf{X} = \mathbf{x} := (x_1, \dots, x_n)$ is observed, the function of θ defined by

$$L_n(\theta|\mathbf{x}) = f_n(\mathbf{x};\theta),$$

is called the likelihood function of θ (given the observation $\mathbf{X} = \mathbf{x}$).

The parameter θ is listed first in L_n because, unlike f_n , L_n is viewed as a function of θ , for a given sample point \mathbf{x} of \mathbf{X} .

If X_1, \dots, X_n are **iid** with marginal pd $f(x;\theta)$, the likelihood factorizes into

$$L_n(\theta|\mathbf{x}) = \prod_{i=1}^n f(x_i;\theta)$$

For the discrete case, $L_n(\theta|\mathbf{x}) = P_\theta(\mathbf{X} = \mathbf{x})$, where P_θ means that *the probability is taken under the assumption that θ is the true parameter*. So, if we compare the likelihood function at two parameter points, say θ_1 and θ_2 , and find, for example, that

$$L_n(\theta_1|\mathbf{x}) = P_{\theta_1}(\mathbf{X} = \mathbf{x}) > L_n(\theta_2|\mathbf{x}) = P_{\theta_2}(\mathbf{X} = \mathbf{x}),$$

then the sample \mathbf{x} we actually observe is *more likely* to have occurred if $\theta = \theta_1$ than if $\theta = \theta_2$. This can be interpreted by saying that, given the data \mathbf{x} , θ_1 is a more plausible value for θ than θ_2 .

In general, we can think of $L_n(\theta|\mathbf{x})$ as a measure of how “likely” θ has produced the observed \mathbf{x} . A similar interpretation applies to the continuous case.

Example 5.2. Let X_i , $i = 1, \dots, n$, be an iid sample from $Ber(\pi)$, where π is our parameter of interest. Suppose that we have $n = 6$ observations $x_i = 1, 1, 0, 1, 1, 0$, and $\pi \in \{0.2, 0.3, 0.7, 0.8, 0.9\} = \Theta$.

The likelihood function of the observed data is

$$L_n(\pi|\mathbf{x}) = P(X_1 = 1, X_2 = 1, \dots, X_6 = 0) = \prod_{i=1}^6 P(X_i = x_i) = \pi^4(1 - \pi)^2.$$

π	$L_n(\pi)$
0.2	0.001
0.3	0.004
0.7	0.0216
0.8	0.0164
0.9	0.0066

Given the actual sample, we can say that 0.7 is the most plausible value of π among $\{0.2, 0.3, 0.7, 0.8, 0.9\}$. \square

Definition 5.2 (Maximum Likelihood Estimator). A statistic $\hat{\theta} \equiv \hat{\theta}(X)$ is called a maximum likelihood estimator (MLE) of $\theta \in \Theta$, if (i) $\hat{\theta} \in \Theta$, and (ii) for each sample point x ,

$$L_n(\hat{\theta}(x)|x) \geq L_n(\theta|x), \quad \forall \theta \in \Theta.$$

In other words, a MLE $\hat{\theta}(x)$ is a parameter point at which the likelihood $L_n(\theta|x)$, as a function of θ , attains its maximum: $L_n(\hat{\theta}(x)|x) = \max_{\theta \in \Theta} L_n(\theta|x)$. In mathematical notation, we write

$$\hat{\theta}_n(x) = \arg \max_{\theta \in \Theta} L_n(\theta|x),$$

where $\arg \max$ is a shortcut for “the arguments of the maxima”, i.e. any point at which $\theta \mapsto L_n(\theta|x)$ is maximized over Θ .

Facts to know

- MLE *may not exist* and *may not be unique*.
- By construction, the range of MLE coincides with Θ , the range of the parameter.
- MLE is a judicious choice in that it represents the parameter point that is most likely to have generated the data. In general, this is a good point estimator in that it possesses some *optimal properties* that will be examined later.
- However, this method has the inherent drawback of having to find the maximum of the likelihood function, which is often a difficult problem. In fact, it can be challenging to find (analytically or numerically) a *global maximizer* and to ensure that it is indeed a global (and not local) maximizer, especially in the multi-parameter case; more on this later.

5.2.2 MLE implementation

For the sake of generality, we consider here the multi-parameter case, where the parameter of interest θ is a vector of dimension d .

Maximum likelihood estimators are often found by maximizing the log-likelihood function

$$\ell_n(\theta|x) = \log(L_n(\theta|x)) = \sum_{i=1}^n \log f(x_i; \theta).$$

Since the logarithmic transformation is strictly monotonically increasing, it does not make any difference to maximize ℓ_n or L_n .

If the likelihood function is differentiable, possible candidates for the MLE are the θ 's that solve the *likelihood equation*:

$$\partial_{\theta_k} \ell_n(\theta|x) = 0, \quad \forall k = 1, \dots, d.$$

This is equivalent to solve the *Score equation* $S_n(\theta, x) = \mathbf{0}$, where $S_n(\theta, x)$ is the *score function* associated with $f_n(x; \theta)$, i.e. $S_n(\theta, x) = (\partial_{\theta_1} \ell_n(\theta|x), \dots, \partial_{\theta_d} \ell_n(\theta|x))^t$.

Let's first focus on the one-parameter case. When looking for a maximum using derivatives, remember that :

- **Stationary points** (i.e. any θ for which $\ell'_n(\theta|\mathbf{x}) = 0$) may be local or global minimizer of $\theta \mapsto \ell_n(\theta|\mathbf{x})$, local or global maximizer, or inflection points (i.e. points at which the concavity changes). Our goal is to find a **global maximizer**.
- The **second derivative test** can be used to determine the nature of a stationary point (min/max/inflection). In fact, *if the second derivative at a stationary point is negative, then this point is a local maximizer*. Furthermore, if $\ell''_n(\theta|\mathbf{x}) < 0, \forall \theta \in \Theta$, then ℓ_n is **strictly concave**, and hence any stationary point that can be found will be the unique global maximizer.
- Maximum (or minimum) can occur where the derivative does not exist. \rightarrow *Check all the points where the derivative dose not exists.*
- The zeros of the first derivative locate only extremum points *in the interior of the domain* of a function (here Θ). If the extremum occurs on the boundary, the first derivative may not be 0. \rightarrow *Check the endpoints of Θ .*
- When all maxima candidates (if any) have been identified, the one(s) with the highest likelihood is/are the MLE.

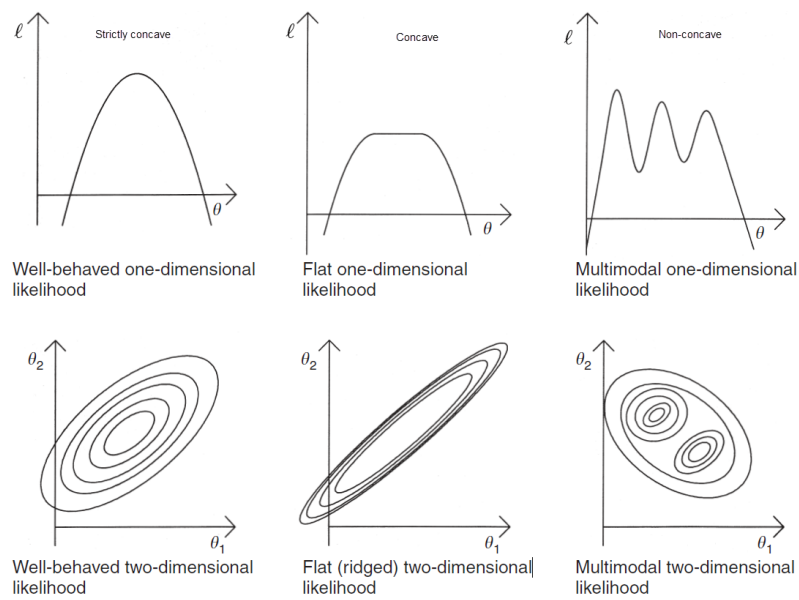
For most of the models considered here after, only one maximum exists, and it corresponds to the solution of the score equation.

Basically, the same approach can be applied for the multi-parameter case but the computational complexity increases with the number of parameters to be estimated. For example, in the case of a model with two parameters, i.e. $\theta = (\theta_1, \theta_2)$, a stationary point, say θ^* , is a **local maxima** if the *Hessian matrix*

$$\nabla^2 \ell_n(\theta^*|\mathbf{x}) = \begin{pmatrix} \partial_{\theta_1}^2 \ell_n(\theta^*|\mathbf{x}) & \partial_{\theta_1 \theta_2} \ell_n(\theta^*|\mathbf{x}) \\ \partial_{\theta_1 \theta_2} \ell_n(\theta^*|\mathbf{x}) & \partial_{\theta_2}^2 \ell_n(\theta^*|\mathbf{x}) \end{pmatrix}$$

is *negative definite*. This is the case if and only if (i) $\partial_{\theta_1}^2 \ell_n(\theta^*|\mathbf{x}) < 0$, and (ii) $\det(\nabla^2 \ell_n(\theta^*|\mathbf{x})) := \partial_{\theta_1}^2 \ell_n(\theta^*|\mathbf{x}) \partial_{\theta_2}^2 \ell_n(\theta^*|\mathbf{x}) - (\partial_{\theta_1 \theta_2} \ell_n(\theta^*|\mathbf{x}))^2 > 0$.

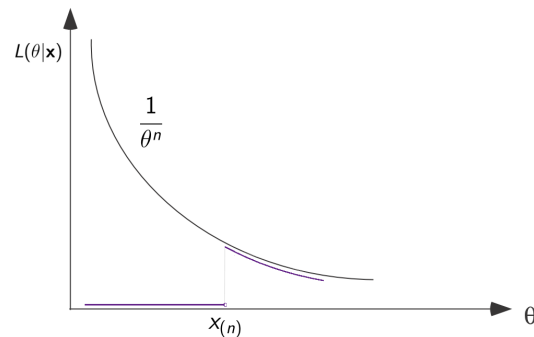
The figure below shows some log-likelihood functions and their extremum values/points in the case of one and two parameters.



Example 5.3.

- Let $X_i, i = 1, \dots, n$, be an iid sample from $Unif[0, \theta], \theta > 0$.

$$L_n(\theta|\mathbf{x}) = \frac{1}{\theta^n} I(x_{(n)} \leq \theta).$$



$\implies X_{(n)}$ is the MLE of θ . We have seen that $X_{(n)}$ is biased, so it is not efficient.

- Let $X_i, i = 1, \dots, n$, be an iid sample from $Unif[\theta, \theta + 1], \theta > 0$.

$$L_n(\theta|\mathbf{x}) = I(x_{(n)} - 1 \leq \theta \leq x_{(1)}).$$

So any $\hat{\theta}$ in $[X_{(n)} - 1, X_{(1)}]$ is a MLE estimator of θ .

- Let $X_i, i = 1, \dots, n$, be an iid sample from $Ber(\pi)$, $\pi \in (0, 1)$.

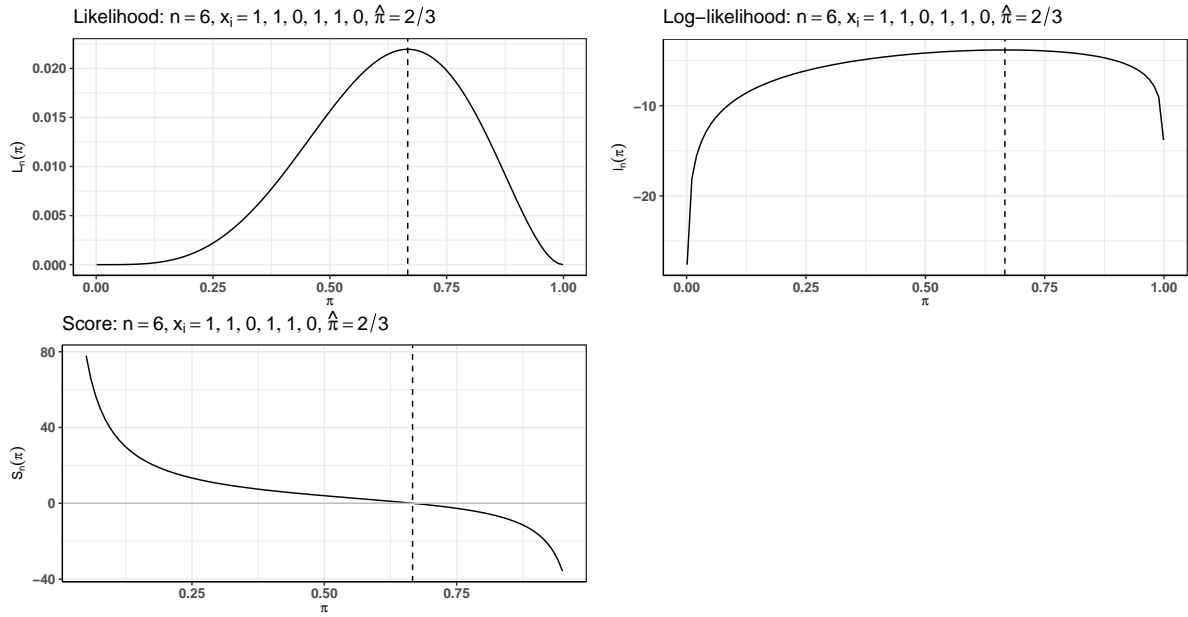
$$L_n(\pi|\mathbf{x}) = \pi^{\sum_{i=1}^n x_i} (1 - \pi)^{n - \sum_{i=1}^n x_i}$$

$$\ell_n(\pi|\mathbf{x}) = \left(\sum_{i=1}^n x_i \right) \log(\pi) + \left(n - \sum_{i=1}^n x_i \right) \log(1 - \pi)$$

$$S_n(\pi, \mathbf{x}) = \ell'_n(\pi|\mathbf{x}) = \frac{\sum_i x_i}{\pi} - \frac{n - \sum_i x_i}{1 - \pi}$$

$$S'_n(\pi, \mathbf{x}) = \ell''_n(\pi|\mathbf{x}) = - \left(\frac{\sum_i x_i}{\pi^2} + \frac{n - \sum_i x_i}{(1 - \pi)^2} \right).$$

$S_n = 0 \Leftrightarrow \hat{\pi} = \frac{1}{n} \sum_{i=1}^n x_i$, and ℓ_n is strictly concave ($\ell''_n < 0, \forall \pi$), so $\hat{\pi} = \bar{X}_n$ is the MLE of π .



- Let $X_i, i = 1, \dots, n$, be an iid sample from $N(\mu, \sigma^2)$, $\mu \in (-\infty, \infty)$, and $\sigma^2 \in (0, \infty)$.

$$L_n(\mu, \sigma^2|\mathbf{x}) = (2\pi\sigma^2)^{-n/2} \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right)$$

$$\ell_n(\mu, \sigma^2|\mathbf{x}) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

- Let's first consider the case of an unknown μ and a known σ^2 . The score for μ and its derivative (with respect to μ) are

$$S_{1n} = \partial_\mu \ell_n(\mu, \sigma^2|\mathbf{x}) = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu)$$

$$\partial_\mu S_{1n} = \partial_\mu^2 \ell_n(\mu, \sigma^2|\mathbf{x}) = -n/\sigma^2.$$

As a function of μ , ℓ_n is strictly concave. Hence, setting the first derivative equal to zero (i.e. $S_{1n} = 0$) and solving for μ we get \bar{X}_n as the MLE of μ .

- Now, let's consider the case of a known μ and an unknown σ^2 . The score for σ^2 and its derivative (with respect to σ^2) are

$$S_{2n} = \partial_{\sigma^2} \ell_n(\mu, \sigma^2 | \mathbf{x}) = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (x_i - \mu)^2$$

$$\partial_{\sigma^2} S_{2n} = \partial_{\sigma^2}^2 \ell_n(\mu, \sigma^2 | \mathbf{x}) = \frac{n}{2\sigma^4} - \frac{1}{\sigma^6} \sum_{i=1}^n (x_i - \mu)^2.$$

As a function of σ^2 , ℓ_n is not strictly concave, because $\partial_{\sigma^2}^2 \ell_n(\mu, \sigma^2 | \mathbf{x})$ is not negative for all possible parameter values. Now, by setting the first derivative equal to zero (i.e. $S_{2n} = 0$) and solving for σ^2 we obtain as the *unique* solution $\tilde{\sigma}_n^2 = n^{-1} \sum_{i=1}^n (x_i - \mu)^2$, and it is easy to see that

$$\partial_{\sigma^2}^2 \ell_n(\mu, \tilde{\sigma}_n^2 | \mathbf{x}) = -\frac{n}{2\tilde{\sigma}_n^4} < 0.$$

We conclude that $\tilde{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ is the MLE of σ^2 .

- Finally, let's consider the case of unknown $\theta = (\mu, \sigma^2)$. The score (vector) and the Hessian (matrix) are

$$S_n(\theta, \mathbf{x}) = (S_{1n}, S_{2n})$$

$$\nabla_{\theta}^2 \ell_n(\theta | \mathbf{x}) = \begin{pmatrix} \partial_{\mu} S_{1n} & \partial_{\sigma^2} S_{1n} \\ \partial_{\mu} S_{2n} & \partial_{\sigma^2} S_{2n} \end{pmatrix} = - \begin{pmatrix} \frac{n}{\sigma^2} & \frac{1}{\sigma^4} \sum_i (x_i - \mu) \\ \frac{1}{\sigma^4} \sum_i (x_i - \mu) & -\frac{n}{2\sigma^4} + \frac{1}{\sigma^6} \sum_i (x_i - \mu)^2 \end{pmatrix}.$$

We can see that $S_n(\hat{\theta}, \mathbf{x}) = (0, 0) \Leftrightarrow \hat{\mu} = \bar{x}_n$, and $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}_n)^2$. So, $\hat{\theta} = (\bar{x}_n, \hat{\sigma}_n^2)$ is the unique candidate for the MLE. Now,

$$\nabla_{\theta}^2 \ell_n(\hat{\theta} | \mathbf{x}) = - \begin{pmatrix} \frac{n}{\hat{\sigma}_n^2} & 0 \\ 0 & \frac{n}{2\hat{\sigma}_n^4} \end{pmatrix} \prec 0 \quad (\text{i.e. negative definite}).$$

We conclude that $\hat{\theta} = (\bar{X}_n, n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2)$ is the MLE of $\theta = (\mu, \sigma^2)$.

- Let $\mathbf{X}_i = (X_{1i}, X_{2i}, X_{3i})$, $i = 1, \dots, n$, be iid rve from the trinomial distribution with joint pd

$$f(\mathbf{x}; \boldsymbol{\pi}) = \frac{m!}{x_1! x_2! x_3!} \pi_1^{x_1} \pi_2^{x_2} \pi_3^{x_3},$$

where $m \geq 1$, $\mathbf{x} = (x_1, x_2)$, and $\boldsymbol{\pi} = (\pi_1, \pi_2)$ is unknown; with $x_1 = 0, \dots, m$, $x_2 = 0, \dots, m$, $x_3 = m - x_1 - x_2$, $\pi_1 \in (0, 1)$, $\pi_2 \in (0, 1)$, and $\pi_3 = 1 - \pi_1 - \pi_2$. We

have that

$$\begin{aligned}
L_n(\pi_1, \pi_2 | \mathbf{x}) &= \frac{(m!)^n}{\prod_{i=1}^n (x_{1i}! x_{2i}! x_{3i}!)} \pi_1^{\sum_{i=1}^n x_{1i}} \pi_2^{\sum_{i=1}^n x_{2i}} \pi_3^{\sum_{i=1}^n x_{3i}}, \\
\ell_n(\pi_1, \pi_2 | \mathbf{x}) &= \sum_{i=1}^n x_{1i} \log(\pi_1) + \sum_{i=1}^n x_{2i} \log(\pi_2) + \sum_{i=1}^n x_{3i} \log(\pi_3) + \text{const}, \\
S_{1n} &:= \partial_{\pi_1} \ell_n(\pi_1, \pi_2 | \mathbf{x}) = \frac{\sum_i x_{1i}}{\pi_1} - \frac{\sum_i x_{3i}}{\pi_3} \text{ and} \\
S_{2n} &:= \partial_{\pi_2} \ell_n(\pi_1, \pi_2 | \mathbf{x}) = \frac{\sum_i x_{2i}}{\pi_2} - \frac{\sum_i x_{3i}}{\pi_3}.
\end{aligned}$$

The score equation is equivalent to $\pi_3 \sum_i x_{1i} = \pi_1 \sum_i x_{3i}$ and $\pi_3 \sum_i x_{2i} = \pi_2 \sum_i x_{3i}$. Summing these two equations yields to $nm\hat{\pi}_3 = \sum_i x_{3i}$. So the unique candidate for the MLE of (π_1, π_2) is $(\hat{\pi}_1, \hat{\pi}_2)$, with $\hat{\pi}_k = \sum_i x_{ki} / nm$. The Hessian of the log-likelihood is

$$\nabla_{\boldsymbol{\pi}}^2 \ell_n(\boldsymbol{\pi} | \mathbf{x}) = - \begin{pmatrix} \frac{\sum_i x_{1i}}{\pi_1^2} + \frac{\sum_i x_{3i}}{\pi_3^2} & \frac{\sum_i x_{3i}}{\pi_3^2} \\ \frac{\sum_i x_{3i}}{\pi_3^2} & \frac{\sum_i x_{2i}}{\pi_2^2} + \frac{\sum_i x_{3i}}{\pi_3^2} \end{pmatrix}$$

This matrix is negative-definite, hence ℓ_n is strictly concave. We conclude that $(\hat{\pi}_1, \hat{\pi}_2) = \left(\frac{\bar{X}_{1n}}{m}, \frac{\bar{X}_{2n}}{m} \right)$ is the MLE of (π_1, π_2) . \square

Chapter 6

Finite and large sample properties of MLE

We now turn our attention to some appealing mathematical properties of the maximum likelihood estimators. We will start with key finite sample properties before discussing, in the next section, some asymptotic features.

6.1 Finite sample properties

6.1.1 Efficiency

Without loss of generality, we consider here the one parameter case. Suppose that an efficient estimator $\hat{\delta} \equiv \hat{\delta}(X)$ of θ exists. By the CRLB attainment theorem (assuming the required assumptions are met), $\hat{\delta}$ must satisfy the equation

$$S_n(\theta, \mathbf{X}) = I_n(\theta) (\hat{\delta} - \theta), \quad (6.1)$$

where $I_n(\theta) = \text{Var}(S_n^2(\theta, \mathbf{X})) = -E(\partial_\theta^2 \ell_n(\theta|x))$ is the Fisher information, contained in the sample, about θ . Now, let $\hat{\theta}$ be the MLE of θ , then $\hat{\theta}$ satisfies $S_n(\hat{\theta}, \mathbf{X}) = 0$. Thus, provided that $I_n(\hat{\theta}) > 0$, the MLE $\hat{\theta}$ coincides with the efficient estimator $\hat{\delta}$.

Theorem 6.1 (MLE and efficiency). *If an efficient estimator exists, then the maximum likelihood method of estimation will produce it.*

In other words, $\hat{\theta}$ is efficient $\implies \hat{\theta}$ is the MLE. The opposite of this statement is not true; as a counter-example, see the example above with $\text{Unif}[0, \theta]$.

6.1.2 Invariance to re-parameterization

Theorem 6.2. *A MLE is invariant with respect to any bijective transformation. That is, if $g : \Theta \longrightarrow \Lambda$ is bijective, then*

$$\hat{\theta} \text{ is a MLE of } \theta \Leftrightarrow g(\hat{\theta}) \text{ is a MLE of } g(\theta).$$

To see why this is the case, let $L_n(\theta)$ be the likelihood with the parametrization θ , i.e. $L_n(\theta) = \prod_i f(x_i, \theta)$, and $L_n^*(\eta)$ be the likelihood with the parametrization $\eta = g(\theta)$, i.e. $L_n^*(\eta) = \prod_i f^*(x_i, \eta)$, where $f^*(x, \eta) := f(x, g^{-1}(\eta))$. Put $\hat{\eta} = g(\hat{\theta})$, and let η be any point in Λ and $\theta = g^{-1}(\eta)$. Since $\hat{\theta} = \arg \max_{\theta \in \Theta} L_n(\theta)$, we have that,

$$L_n^*(\hat{\eta}) = L_n(g^{-1}(\hat{\eta})) = L_n(\hat{\theta}) \geq L_n(\theta) = L_n(g^{-1}(\eta)) = L_n^*(\eta).$$

We conclude that $L_n^*(\hat{\eta}) \geq L_n^*(\eta)$, $\forall \eta \in \Lambda$, and so $\hat{\eta}$ is the MLE of η .

Example 6.1. Let X_i , $i = 1, \dots, n$, be an iid sample from $Ber(\pi)$, $\pi \in (0, 1)$. We have seen that the likelihood function is given by $L_n(\pi|\mathbf{x}) = \pi^{\sum_{i=1}^n x_i} (1 - \pi)^{n - \sum_{i=1}^n x_i}$, and the MLE of π is $\hat{\pi} = \bar{X}_n$.

Now, put $\eta = \log(\pi/(1 - \pi))$; with this parametrization, the likelihood function becomes

$$L_n^*(\eta|\mathbf{x}) = \left(\frac{1}{1 + e^{-\eta}} \right)^{\sum_{i=1}^n x_i} \left(\frac{1}{1 + e^{\eta}} \right)^{n - \sum_{i=1}^n x_i}.$$

One can check, by taking the logarithm and then differentiate, that the maximizer of this likelihood is $\hat{\eta} = \log(\sum_i X_i / (n - \sum_i X_i))$. Hence, the MLE of $\eta = \log(\pi/(1 - \pi))$ is $\hat{\eta} = \log(\hat{\pi}/(1 - \hat{\pi}))$. The same result can be obtained directly (without any calculation) by simply applying the above theorem. \square

This simple version of the invariance of MLE is not always useful because many of the functions we are interested in are not bijective. For example, we cannot apply the above result to obtain the MLE of $\pi(1 - \pi)$ in the case of $Ber(\pi)$, or to obtain the MLE of $cv := \sigma/\mu$ in the case of $N(\mu, \sigma^2)$. Now, whether g is bijective or not, if $\hat{\theta}$ is the MLE of θ then it is reasonable to use $g(\hat{\theta})$ as an estimator for $g(\theta)$. *Even if g is not bejective, some authors refer to $g(\hat{\theta})$ as the MLE of $g(\theta)$* , although, according to the (classical) definition of MLE, we can't really call it so because we can't necessarily express the pd as a genuine function of $g(\theta)$.

6.2 Large sample properties

We will show that the MLE is *often*: (1) consistent, (2) asymptotically normal, and (3) asymptotically efficient.

For (1) to hold, we need some regularity conditions:

- We observe X_i , $i = 1, \dots, n$, an iid sample from a pd $f(x; \theta)$, $\theta \in \Theta$.
- The model is identifiable; that is, if $f(x; \theta) = f(x; \tilde{\theta}) \forall x$, then $\theta = \tilde{\theta}$.
- The model is correctly specified. *We denote by θ_0 is the true parameter value.*
- Θ is a compact set (closed and bounded) and $\theta \mapsto f(x; \theta)$ is continuous on Θ .
- $|f(X; \theta)| \leq d(X)$, $\forall \theta \in \Theta$, and $E_{\theta_0}(d(X)) < \infty$.

For (2)-(3) to hold, we need in addition to the above conditions the next assumptions:

- The support of f is independent of θ .

- θ_0 is in the interior of Θ .
- f is twice continuously differentiable in θ and $\int f(x, \theta) dx$ can be differentiated two times under the integral sign.
- The Fisher information satisfies $0 < I(\theta_0) < \infty$.
- In a neighborhood of θ_0 , $|\partial_\theta^3 \log f(x; \theta)| \leq M(x) \forall$; and $E_{\theta_0}(M(X)) < \infty$.

These assumptions are sufficient (to prove consistency and asymptotic normality) but not necessary. More general and weaker conditions can be found in the literature.

6.2.1 Consistency

Theorem 6.3 (Consistency of MLEs). *Under the regularity assumptions stated above, $\hat{\theta}_n \xrightarrow{p} \theta_0$.*

We will not prove this result here, but will only sketch out why this is happening.

First, by definition, the MLE $\hat{\theta}_n$ is the maximizer of $\bar{\ell}_n(\theta) = n^{-1} \sum_{i=1}^n \log f(X_i; \theta)$ which is the log-likelihood function normalized by $1/n$ (of course, this does not affect maximization). Second, by the law of large numbers (WLLN),

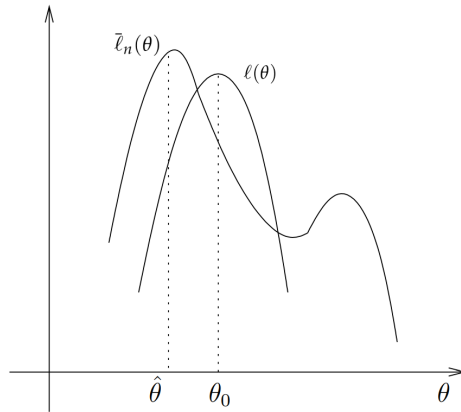
$$\bar{\ell}_n(\theta) \xrightarrow{p} \ell(\theta) := E_{\theta_0}(\log f(X; \theta)), \forall \theta.$$

The expectation operator E is indexed by θ_0 to explicitly point out that the expectation is evaluated using the true parameter θ_0 , i.e. acknowledging that the pd of X is $f(x; \theta_0)$. Third, by Jensen's inequality, for any $\theta \neq \theta_0$,

$$\ell(\theta) - \ell(\theta_0) = E_{\theta_0} \left(\log \frac{f(X; \theta)}{f(X; \theta_0)} \right) < \log E_{\theta_0} \left(\frac{f(X; \theta)}{f(X; \theta_0)} \right) = 0.$$

So, θ_0 is the maximizer of $\ell(\theta)$.

To sum up, we know that $\theta_0 = \arg \max \ell(\theta)$, $\hat{\theta}_n = \arg \max \bar{\ell}_n(\theta)$, and $\bar{\ell}_n(\theta) \xrightarrow{p} \ell(\theta)$, $\forall \theta$. So, we (intuitively) expect $\hat{\theta}_n$ to approach θ_0 as the sample size increases.



In fact, the assumed hypothesis guarantees that this is indeed the case. Thus, $\hat{\theta}_n \xrightarrow{p} \theta_0$.

6.2.2 Asymptotic normality and asymptotic efficiency

Theorem 6.4 (Asymptotic normality of MLEs). *Under the regularity assumptions stated above,*

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N\left(0, \frac{1}{I(\theta_0)}\right),$$

where $I(\theta_0)$ is the Fisher information evaluated at θ_0 .

The above result can be expressed as $\hat{\theta} \sim_a N(\theta_0, I_n^{-1}(\theta_0))$, where $I_n(\theta) = nI(\theta)$. And thus, the MLE $\hat{\theta}$ is asymptotically efficient for θ_0 (according to the definition given previously).

The proof uses Taylor's theorem, CLT, and Slutsky's theorem. In fact, Taylor expansion of $\theta \mapsto \bar{\ell}'_n(\theta)$ around θ_0 yields to $0 = \bar{\ell}'_n(\hat{\theta}) \approx \bar{\ell}'_n(\theta_0) + (\hat{\theta} - \theta_0)\bar{\ell}''_n(\theta_0)$. So,

$$\sqrt{n}(\hat{\theta} - \theta_0) \approx -\frac{1}{\bar{\ell}''_n(\theta_0)}\sqrt{n}\bar{\ell}'_n(\theta_0).$$

$\sqrt{n}\bar{\ell}'_n(\theta_0) = \sqrt{n}(n^{-1}\sum_i \partial_\theta \log f(X_i; \theta_0) - 0) \xrightarrow{d} N(0, I(\theta_0))$ by the CLT and the fact that $\partial_\theta \log f(X_i; \theta_0)$ has mean zero and variance $I(\theta_0)$. For the denominator, by WLLN, we have that $-\bar{\ell}''_n(\theta_0) = -n^{-1}\sum_i \partial_{\theta^2}^2 \log f(X_i; \theta_0) \xrightarrow{p} I(\theta_0)$. The proof is completed by applying Slutsky's Theorem.

Attention: To simplify the notations, henceforth, we suppress the subscript 0 in θ_0 , and write $\hat{\theta} \xrightarrow{p} \theta$, which must be understood as $\hat{\theta}$ converges to the true value of θ whatever this one is. In the same way, we write $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, I^{-1}(\theta))$.

The results stated above (consistency, asymptotic normality and asymptotic efficiency) can be extended to the multiparameter case.

If $\hat{\theta}$ is d -dimensional MLE of θ , then, under some regularity assumptions (similar to those stated above),

- $\hat{\theta} \xrightarrow{p} \theta$,
- $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N_d(\mathbf{0}, I^{-1}(\theta))$, where I^{-1} is the inverse of the Fisher information matrix, and
- $\hat{\theta}$ is asymptotically efficient for θ .

Example 6.2.

- Let $X_i, i = 1, \dots, n$, be an iid sample from $\text{Bin}(m, \pi)$, $\pi \in (0, 1)$. We have that

$$\begin{aligned} L_n(\pi|\mathbf{x}) &= \left(\prod_{i=1}^n C_m^{x_i} \right) \pi^{\sum_{i=1}^n x_i} (1 - \pi)^{\sum_{i=1}^n (m - x_i)} \\ \ell_n(\pi|\mathbf{x}) &= \left(\sum_{i=1}^n x_i \right) \log(\pi) + \left(nm - \sum_{i=1}^n x_i \right) \log(1 - \pi) + \text{const}, \\ S_n(\pi, \mathbf{x}) &= \partial_\pi \ell_n(\pi|\mathbf{x}) = \frac{\sum_i x_i}{\pi} - \frac{nm - \sum_i x_i}{1 - \pi}, \\ \partial_\pi^2 \ell_n(\pi|\mathbf{x}) &= \partial_\pi S_n(\pi, \mathbf{x}) = -\frac{\sum_i x_i}{\pi^2} - \frac{nm - \sum_i x_i}{(1 - \pi)^2}. \end{aligned}$$

Hence, the MLE of π is $\hat{\pi} = \frac{\sum_i X_i}{nm} = m^{-1} \bar{X}_n$. This later is consistent and asymptotically normal. More precisely

$$\hat{\pi} \sim_a N(\pi, I_n^{-1}(\pi)),$$

where $I_n(\pi) := -E(\partial_{\pi^2} \ell_n(\pi|\mathbf{x})) = \frac{nm}{\pi(1 - \pi)}$.

Note that the asymptotic distribution of $\hat{\pi}$, as given above, can also be obtained by applying the CLT directly to \bar{X}_n .

- Let $\mathbf{X}_i = (X_{1i}, X_{2i}, X_{3i})$, $i = 1, \dots, n$, be iid rve from the trinomial distribution with joint pd

$$f(\mathbf{x}; \boldsymbol{\pi}) = \frac{m!}{x_1! x_2! x_3!} \pi_1^{x_1} \pi_2^{x_2} \pi_3^{x_3},$$

We have seen the the MLE of (π_1, π_2) is $(\bar{X}_{1n}/m, \bar{X}_{2n}/m)$. We have also seen that

$$I^{-1}(\pi_1, \pi_2) = m^{-1} \begin{pmatrix} \pi_1(1 - \pi_1) & -\pi_1\pi_2 \\ -\pi_1\pi_2 & \pi_2(1 - \pi_2) \end{pmatrix}.$$

Hence,

$$\sqrt{n} \begin{pmatrix} \hat{\pi}_1 - \pi_1 \\ \hat{\pi}_2 - \pi_2 \end{pmatrix} \xrightarrow{d} N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, m^{-1} \begin{pmatrix} \pi_1(1 - \pi_1) & -\pi_1\pi_2 \\ -\pi_1\pi_2 & \pi_2(1 - \pi_2) \end{pmatrix} \right).$$

Again, the same result can be obtained by applying the CLT directly on $\bar{\mathbf{X}}_n = (\bar{X}_{1n}, \bar{X}_{2n})$.

□

The theoretical results presented above can be extended to the situation where the parameter of interest is $\mathbf{g}(\boldsymbol{\theta})$ rather than $\boldsymbol{\theta}$. In fact, let $\hat{\boldsymbol{\theta}}$ be the MLE of $\boldsymbol{\theta}$. Then, *whether* $\mathbf{g} : \mathbb{R}^d \mapsto \mathbb{R}^p$ is bijective or not,

- By the continuous mapping theorem, $\mathbf{g}(\hat{\boldsymbol{\theta}}) \xrightarrow{p} \mathbf{g}(\boldsymbol{\theta})$;
- By the Delta method, $\sqrt{n}(\mathbf{g}(\hat{\boldsymbol{\theta}}) - \mathbf{g}(\boldsymbol{\theta})) \xrightarrow{d} N_p(\mathbf{0}, \mathbf{I}^{-1}(\mathbf{g}(\boldsymbol{\theta})))$, where $\mathbf{I}^{-1}(\mathbf{g}(\boldsymbol{\theta})) = \dot{\mathbf{g}}(\boldsymbol{\theta}) \mathbf{I}^{-1}(\boldsymbol{\theta}) \dot{\mathbf{g}}^t(\boldsymbol{\theta})$;

- $g(\hat{\theta})$ is asymptotically efficient for $g(\theta)$.

6.2.3 Observed Fisher information

We have seen that the MLE $\hat{\theta}$ of θ is asymptotically normal, i.e. $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N_d(\mathbf{0}, I^{-1}(\theta))$. Equivalently, we can write that

$$\sqrt{I_n(\theta)}(\hat{\theta} - \theta) \xrightarrow{d} N_d(\mathbf{0}, \mathbb{1}) ,$$

where $\mathbb{1}$ is the identity matrix and $\sqrt{I_n(\theta)}$ is the *square-root matrix* of the FI

$$I_n(\theta) = -E \left\{ \left[\partial_{\theta_j \theta_k} \ell_n(\theta) \right]_{j,k} \right\} = nI(\theta) , \text{ with } I(\theta) = -E \left\{ \left[\partial_{\theta_j \theta_k} \log f(X; \theta) \right]_{j,k} \right\} .$$

To use this asymptotic normality in practical inference, $I(\theta)$ must be estimated. The most obvious estimator of $I(\theta)$ is $I(\hat{\theta})$. Since $\hat{\theta} \xrightarrow{p} \theta$, $I(\hat{\theta}) \xrightarrow{p} I(\theta)$, provided that $\theta \mapsto I(\theta)$ is a continuous function. In this case, Slutsky's theorem ensures that

$$\sqrt{I_n(\hat{\theta})}(\hat{\theta} - \theta) \xrightarrow{d} N_d(\mathbf{0}, \mathbb{1}) .$$

The disadvantage of this approach is that, in practice, it is not always easy to calculate $I(\hat{\theta})$ because of the difficulties of working out the expectations.

Let's define the matrix

$$J_n(\theta) := -\nabla^2 \ell_n(\theta) = - \left[\partial_{\theta_j \theta_k} \ell_n(\theta) \right]_{j,k} = - \left[\sum_{i=1}^n \partial_{\theta_j \theta_k} \log f(X_i; \theta) \right]_{j,k} .$$

$J_n(\theta)$ is the Hessian of the negative log-likelihood. This matrix is known as *the sample Fisher information* or *the observed Fisher information*. It can always be calculated as long as the second partial derivatives can be calculated. Observe that $I_n(\theta) = E(J_n(\theta))$. The law of large numbers guarantees that $n^{-1}J_n(\theta) \xrightarrow{p} I(\theta)$. So,

$$\sqrt{J_n(\theta)}(\hat{\theta} - \theta) \xrightarrow{d} N_d(\mathbf{0}, \mathbb{1}) .$$

Again Slutsky's theorem can be applied to show that

$$\sqrt{J_n(\hat{\theta})}(\hat{\theta} - \theta) \xrightarrow{d} N_d(\mathbf{0}, \mathbb{1}) .$$

To sum up, *for large sample sizes*, one can use $I_n(\theta)$, $I_n(\hat{\theta})$, $J_n(\theta)$ and $J_n(\hat{\theta})$ interchangeably.

Example 6.3. Let X_i , $i = 1, \dots, n$, be an iid sample from the pd

$$f(x; \theta) = \frac{1 + \theta x}{2} I(-1 \leq x \leq 1); -1 < \theta < 1 .$$

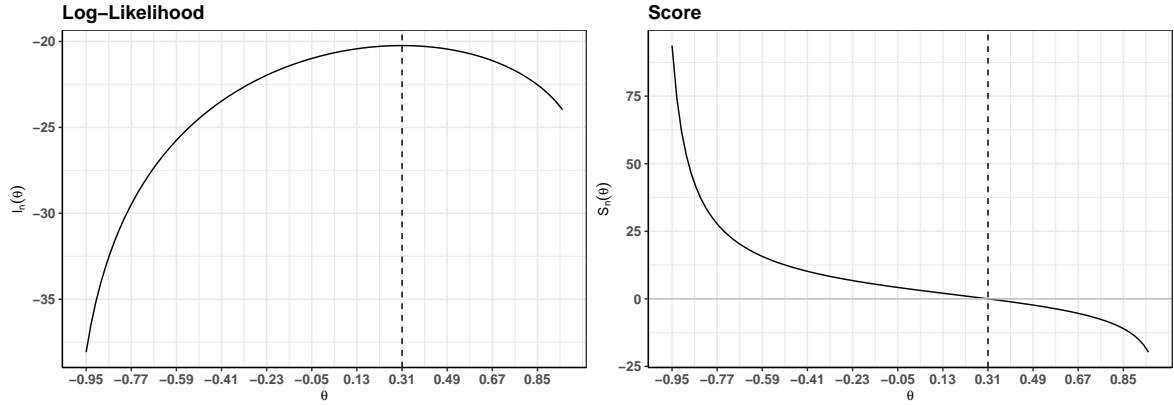
The log-likelihood, the Score and the observed FI are given by

$$\begin{aligned}\ell_n(\theta|x) &= \sum_{i=1}^n \log(1 + \theta x_i) - n \log(2), \\ S_n(\theta, x) &= \partial_\theta \ell_n(\theta|x) = \sum_{i=1}^n \frac{x_i}{1 + \theta x_i}, \\ J_n(\theta) &= -\partial_\theta S_n(\theta, x) = -\partial_\theta^2 \ell_n(\theta|x) = \sum_{i=1}^n \frac{x_i^2}{(1 + \theta x_i)^2}.\end{aligned}$$

As a function of θ , the log-likelihood is *continuous and strictly concave* ($J_n > 0$), so there is a unique MLE, say $\hat{\theta}$, of θ . Now, although the likelihood equation of this model *cannot be solved explicitly* to get the analytic expression of $\hat{\theta}$, the theory tells us that $\sqrt{I_n(\theta)}(\hat{\theta} - \theta) \xrightarrow{d} N(0, 1)$, where $I_n(\theta) = E(J_n(\theta))$.

The following data are simulated from the f above with $\theta = 0.5$ (hereafter, we'll pretend that we don't know the θ that generated the data and see how maximum likelihood behaves).

```
x <- c(0.9852, 0.0450, -0.6123, -0.7518, -0.2824, 0.7085, -0.0711, 0.9625,
      -0.4746, 0.1617, -0.4592, -0.3113, 0.6800, -0.6694, 0.1512, -0.7048,
      0.3421, -0.9658, 0.9809, -0.1205, 0.4730, -0.1665, 0.9956, 0.8720,
      0.9849, -0.7650, 0.4528, 0.2190, 0.9611, -0.0257)
```



By inspecting the graphs above, we can see that $\hat{\theta} = 0.31$ (we will see later how to obtain this result numerically).

As we discussed above, we can estimate $I_n(\theta)$ using $I_n(\hat{\theta})$ or $J_n(\hat{\theta})$. To use the former, we must first derive the expression of $I_n(\theta) = nE\left(\frac{X^2}{(1 + \theta X)^2}\right) = \frac{n}{2} \int_{-1}^1 \frac{x^2}{1 + \theta x} dx$. Clearly, $I_n(0) = \frac{n}{2} \int_{-1}^1 x^2 = \frac{n}{3}$, and for $\theta \neq 0$,

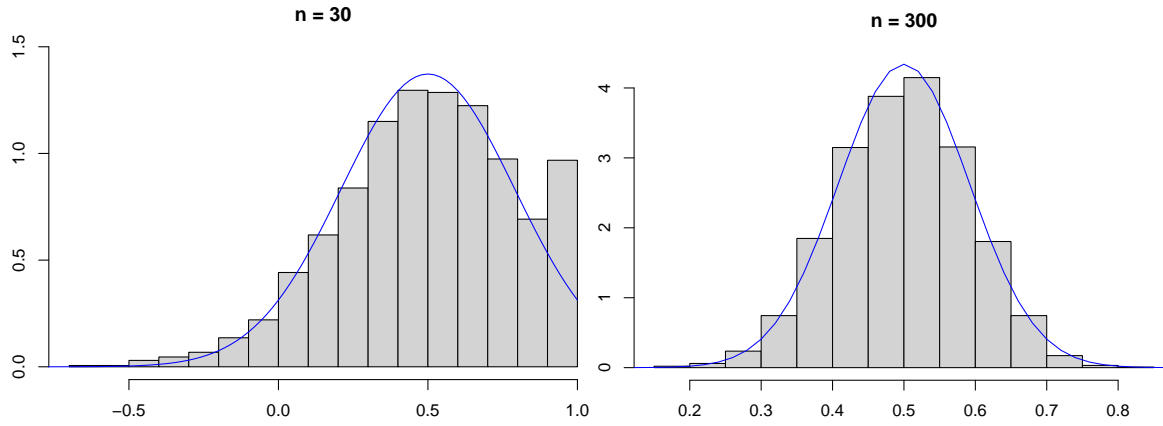
$$\begin{aligned}I_n(\theta) &= \frac{n}{2\theta^3} \int_{-1}^1 \frac{(u-1)^2}{u} du = \frac{n}{2\theta^3} \int_{-1}^1 \left(u - 2 + \frac{1}{u}\right) du \\ &= \frac{n}{2\theta^3} \left[\frac{u^2}{2} - 2u + \log|u| \right]_{-1}^{1+\theta} = -\frac{n}{\theta^2} + \frac{n}{2\theta^3} \log\left(\frac{1+\theta}{1-\theta}\right).\end{aligned}$$

In our case, $n = 30$, $\hat{\theta} = 0.31$, $I_{30}(\hat{\theta}) = 10.619$ and $1/I_{30}(\hat{\theta}) = 0.094$. Hence, based on the observed data, we conclude that $\hat{\theta}_{30} \sim_a N(\theta, 0.094)$.

The second method for estimating $I_n(\theta)$ is much simpler and consists in using the observed FI which is given here by $J_{30}(\hat{\theta}) = 11.846$. On this basis, we can write that $\hat{\theta}_{30} \sim_a N(\theta, 0.084)$.

Note that the true (but *unknown*) asymptotic variance of $\hat{\theta}_{30}$ is actually $1/I_{30}(0.5) = 0.085$.

To see how the asymptotic normal approximation works in practice, we repeat the data generation procedure 5000 times and calculate $\hat{\theta}$ each time. We did this for $n = 30$ and $n = 300$, respectively. The graph below shows the histogram of the simulated $\hat{\theta}$ and the curve of the true asymptotic normal density, i.e. $N(0.5, I_n^{-1}(0.5))$, in blue.



□

Chapter 7

More about likelihood

7.1 Numerical maximization of the likelihood

In several interesting cases, the maximization problem has no analytical solution. In other words, it is not possible to write $\hat{\theta}$ explicitly as a function of the data. In these cases, it is necessary to resort to numerical algorithms for the maximization of the likelihood.

7.1.1 The Newton-Raphson (NR) method

One of the most used method for optimization in statistics is the Newton-Raphson method. It is based on approximating the log-likelihood $\ell_n(\theta)$ by a quadratic function. For a given *starting point* $\theta_0 \in \Theta$, such that $\ell_n''(\theta_0) < 0$, define

$$\tilde{\ell}_n(\theta) = \ell_n(\theta_0) + (\theta - \theta_0)\ell_n'(\theta_0) + (\theta - \theta_0)^2\ell_n''(\theta_0)/2.$$

This is the second order Taylor serie approximation of $\ell_n(\theta)$ around θ_0 .

The solution of the first-order condition for maximizing $\tilde{\ell}_n(\theta)$ is

$$\theta_1 = \theta_0 - \frac{\ell_n'(\theta_0)}{\ell_n''(\theta_0)}.$$

Since $\tilde{\ell}_n$ is an approximation of ℓ_n , θ_1 , defined above, provides a guess value for the MLE.

We can try to improve the approximation by taking θ_1 as the new starting point and *keep repeating the process until convergence*. This suggests the following iterative procedure:

$$\theta_{k+1} = \theta_k - \frac{\ell_n'(\theta_k)}{\ell_n''(\theta_k)}, \quad k = 0, 1, \dots$$

Equivalently, this algorithm can also be written as

$$\theta_{k+1} = \theta_k + \frac{S_n(\theta_k)}{J_n(\theta_k)}, \quad k = 0, 1, \dots,$$

where $S_n(\theta) = \ell'_n(\theta)$ and $J_n(\theta) = -\ell''_n(\theta)$ are the score and the observed FI.

As we said, the procedure should be run until convergence, i.e. until there is no “significant” difference between θ_k and θ_{k+1} . No “significant” difference means that changes between consecutive iterations are less than a user-defined tolerance. For example, we may stop the algorithm whenever the difference $|\theta_{k+1} - \theta_k|$, or the relative difference $|\theta_{k+1} - \theta_k|/|\theta_k|$, is smaller than 10^{-8} .

Note that, $\theta_{k+1} = \theta_k$ is equivalent to $\ell'_n(\theta_k) = 0$. So, at the end of the process (i.e. when the iterations stop), the algorithm converges to (a neighborhood of) a stationary point of ℓ_n . This point could be a maximum point, a minimum point or even a inflection point. However, if $\ell''_n(\theta_k) < 0$, then the convergence point is a **local maximizer**. Moreover, if ℓ_n is strictly concave, then the convergence point is the (unique) global maximizer.

In all cases, when it converges, the NR algorithm reaches the stationary point closest to its starting point θ_0 . *If several stationary points are present, the choice of the starting point becomes critical.* In many situations, the MoM can be used to obtain a reasonable starting point.

In the case where the likelihood is not strictly concave, it is recommended to:

- If possible, start by visually checking the graph of the (log-)likelihood function. Sometimes it's easy to see where an extremum occurs. But sometimes, local fluctuations of a relatively small scale can hide such points.
- Rerun the algorithm from different starting points, and then choose, among all the convergence points, the one that globally maximizes the likelihood.
- Perturb the convergence point by a “small” amount, and use this as a starting point for a new run of the algorithm. Then, see if it converges to a “better point”, or “always” to the same one.

The arguments for deriving the NR algorithm for optimization in one dimension can be directly extended to multi-dimensional problems giving the multi-parameter NR method:

$$\theta_{k+1} = \theta_k + J_n^{-1}(\theta_k)S_n(\theta_k), \quad k = 0, 1, \dots$$

where $S_n(\theta) = \nabla_{\theta} \ell_n(\theta) = \sum_{i=1}^n \nabla_{\theta} \log f(X_i, \theta)$ is the Score vector and $J_n(\theta) = -\nabla_{\theta}^2 \ell_n(\theta) = -\sum_{i=1}^n \nabla_{\theta}^2 \log f(X_i, \theta)$ is the observed FI matrix (i.e. Hessian of negative log-likelihood).

Let's consider our previous example with $f(x; \theta) = \frac{1+\theta x}{2} I(-1 \leq x \leq 1)$; $-1 < \theta < 1$, and the observed data

```
x <- c(0.9852, 0.0450, -0.6123, -0.7518, -0.2824, 0.7085, -0.0711, 0.9625,
      -0.4746, 0.1617, -0.4592, -0.3113, 0.6800, -0.6694, 0.1512, -0.7048,
      0.3421, -0.9658, 0.9809, -0.1205, 0.4730, -0.1665, 0.9956, 0.8720,
      0.9849, -0.7650, 0.4528, 0.2190, 0.9611, -0.0257)
```

The following code gives the R functions needed for running the NR algorithm.

```

LogLik <- function(theta, x) {sum(log((1 + theta * x) / 2))}
LogLikGrad <- function(theta, x) {sum(x / (1 + theta * x))}
LogLikHess <- function(theta, x) {-sum((x / (1 + theta * x))^2)}
NRoptim <- function(theta0, x, eps = 1e-06, trace = FALSE) {
  diff <- Inf
  theta <- theta0
  LL <- LogLik(theta, x)
  grad <- LogLikGrad(theta, x)
  hess <- LogLikHess(theta, x)
  detail <- data.frame(theta = theta, LL = LL, grad = grad, hess = hess, diff = diff)
  while (diff > eps) {
    theta.old <- theta
    theta <- theta.old - grad / hess
    diff <- abs(theta - theta.old)
    LL <- LogLik(theta, x)
    grad <- LogLikGrad(theta, x)
    hess <- LogLikHess(theta, x)
    detail <- rbind(detail, c(theta = theta, LL = LL, grad = grad, hess = hess, diff = diff))
  }
  if(trace) print(detail)
  c(Estimate = theta, Std.Error = sqrt(-1/hess))
}

```

Let's run this function with starting point -0.3 .

```
NRoptim(-0.3, x, trace = TRUE)
```

```

      theta    LL      grad  hess    diff
1 -0.300 -22.5  7.93e+00 -18.1      Inf
2  0.140 -20.4  1.96e+00 -11.5 4.40e-01
3  0.310 -20.2 -1.24e-02 -11.8 1.71e-01
4  0.309 -20.2 -3.66e-06 -11.8 1.05e-03
5  0.309 -20.2 -3.16e-13 -11.8 3.09e-07

Estimate Std.Error
      0.309      0.291

```

Here, with $\text{eps} = 1\text{e-}06$, the NR algorithm reaches its target after only 4 iterations. And since the log-likelihood function is strictly concave, there is no need for further investigation, and the point of convergence (0.309) is certainly the MLE.

The two figures below show how ℓ_n is approximated by $\tilde{\ell}_n$ and how the algorithm moves to the maximum.

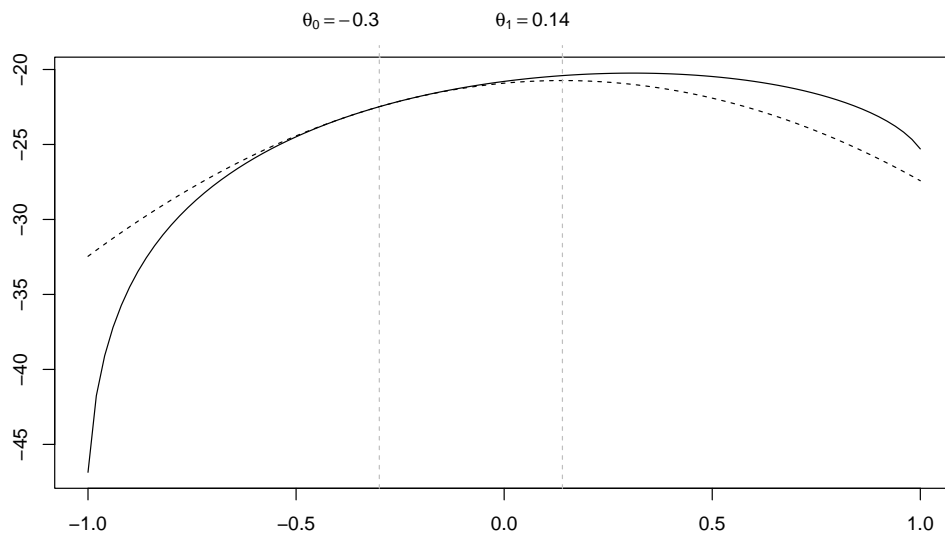


Figure 7.1: ℓ_n (solid line) and its quadratic approximation $\tilde{\ell}_n$ (dashed line) at $\theta_0 = -0.3$

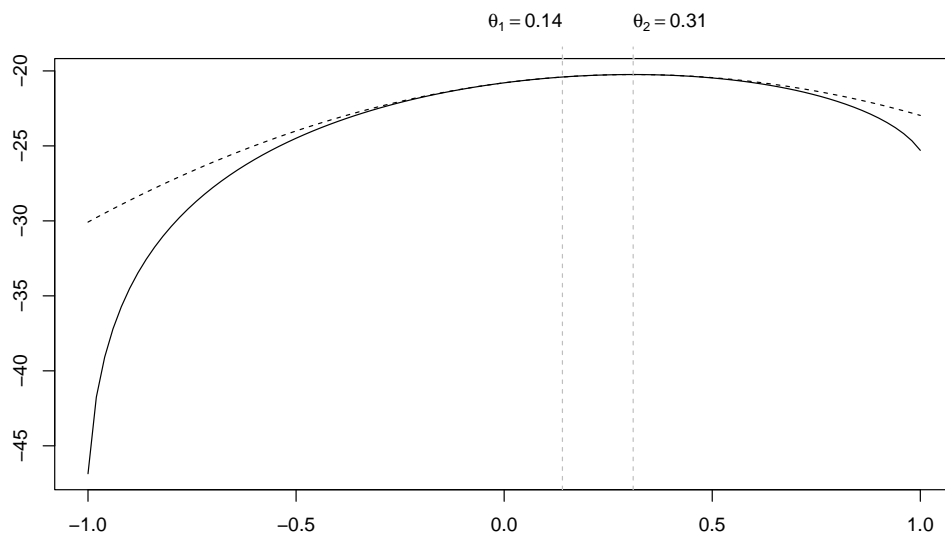


Figure 7.2: ℓ_n (solid line) and its quadratic approximation $\tilde{\ell}_n$ (dashed line) at $\theta_1 = 0.14$

As starting point, we could use the MoM estimator of θ . In fact, it is easy to check that $E(X) = \theta/3$, so the MoM of θ is $\hat{\theta}_0 = 3\bar{x}_n = 0.36$.

```
NRoptim(0.36, x)
```

Estimate	Std.Error
0.309	0.291

When it works, the NR method converges very quickly to the maximum, especially when the starting point is not very far from the maximizer. However, this method has some serious problems, especially if the likelihood is non-concave:

- NR is sensitive to the starting point.
- A NR step may jump far away from the target.

- As the number of parameters increases, the NR method becomes very computationally expensive.

A large number of enhanced alternative methods (as for example the [Quasi-Newton methods](#)) can be found in the literature, but this is beyond the scope of this course.

7.1.2 Maximum Likelihood in R

R provides a function called `optim()` which, by default, *performs minimization*. To maximize the likelihood, provide `optim()` with *the negative of the log-likelihood* (for any function f , minimizing $(-f)$ maximizes f).

As main arguments, `optim()` takes: `par`, the starting vector point, i.e. the initial values for the parameters to be optimized over; `fn`: the function to be minimized, *with first argument the vector of parameters over which minimization is to take place*; and some other optional arguments.

```
optim(par, fn, gr = NULL, ..., # Further arguments to be passed to fn and gr.
      method = c("Nelder-Mead", "BFGS", "CG", "L-BFGS-B", "SANN", "Brent"),
      lower = -Inf, upper = Inf, control = list(), hessian = FALSE)
```

The default algorithm (`method`) for `optim()` is a derivative-free optimization routine called the “Nelder-Mead” simplex algorithm. The other optimization methods are: BFGS, CG, L-BFGS-B, Here we will use the L-BFGS-B method. L-BFGS-B is a variant of BFGS method, an optimization algorithm in the family of quasi-Newton methods. L-BFGS-B is more memory-efficient than BFGS and allows the incorporation of “box” constraints, i.e. constraints of the form $a < \theta < b$; see the the Help of `optim()` for more details.

`optim()` returns a list of values of which `par`: the local minimizer, `value`: the target function evaluated at the solution found, `convergence`: an integer code indicating successful convergence (code 0) or a warning or an error code (see the Help), and `hessian`: the Hessian at the solution found (if `hessian = TRUE`).

Let’s consider again our previous example with $f(x; \theta) = \frac{1+\theta x}{2} I(-1 \leq x \leq 1)$; $-1 < \theta < 1$, and the observed data

```
x <- c(0.9852, 0.0450, -0.6123, -0.7518, -0.2824, 0.7085, -0.0711, 0.9625,
      -0.4746, 0.1617, -0.4592, -0.3113, 0.6800, -0.6694, 0.1512, -0.7048,
      0.3421, -0.9658, 0.9809, -0.1205, 0.4730, -0.1665, 0.9956, 0.8720,
      0.9849, -0.7650, 0.4528, 0.2190, 0.9611, -0.0257)
```

```
negLogLik <- function(theta, x) {-sum(log((1 + theta * x) / 2))}
```

```
optim(-0.3, fn = negLogLik, x = x, method = "L-BFGS-B", lower = -1, upper = 1, hessian = TRUE)
```

```
$par
```

```
[1] 0.309
```

```

$value
[1] 20.2

$counts
function gradient
      6      6

$convergence
[1] 0

$message
[1] "CONVERGENCE: REL_REDUCTION_OF_F <= FACTR*EPSMCH"

$hessian
      [,1]
[1,] 11.8

```

There are many packages and functions in R designed to facilitate MLE calculation, most of which call the `optim()` function, in the background, to perform the optimization. This includes the functions `stats4::mle()`, `maxLik::maxLik()` and `MASS::fitdistr()`, to name but a few. In general, these functions make things a little easier to code and perform additional processing to make the results more useful for estimation and/or inference. In the following, we'll look at the `mle()` function.

```

library(stats4)
mle.negLogLik <- mle(\(theta) negLogLik(theta, x = x), # set the arguments that must be held f
                    start = -0.3, lower = -1, upper = 1, method = "L-BFGS-B")
mle.negLogLik

```

Call:

```

mle(minuslogl = function(theta) negLogLik(theta, x = x), start = -0.3,
    method = "L-BFGS-B", lower = -1, upper = 1)

```

Coefficients:

```

theta
0.309

```

The objective returned by the `mle()` function can be used by some useful well-known generic R functions like `summary()`, `logLik()`, `vcov()`. This latter gives the asymptotic variance-covariance matrix of the estimated vector. The following example provides an illustration.

Example 7.1 (Two dimensional case). Here we consider the problem of estimating $\theta = (l, s)$

the location-scale parameters of the [Cauchy distribution](#):

$$f(x;l,s) = \frac{1}{\pi s \left[1 + \left(\frac{x-l}{s} \right)^2 \right]}, \quad l \in \mathbb{R}, s > 0.$$

We start by generating 30 observations with true values $l = 0$ and $s = 1$.

```
set.seed(1)
y <- rcauchy(30)
y

[1]  1.1025  2.3538 -4.2926 -0.2966  0.7346 -0.3305 -0.1756 -1.8082
[9] -2.3286  0.1966  0.7556  0.6195 -1.5015  2.6241 -0.8825 138.3476
[17] -1.2274 -0.0254  2.5265 -0.8409 -0.2081  0.7865 -1.9374  0.4163
[25]  1.1145  2.6747  0.0421  2.5821 -0.4339  1.8237

neglogLi2 <- function(l, s, y) {-sum(dcauchy(y, location = l, scale = s, log = TRUE))}
mle.negLogLik2 <- mle(\(l, s) neglogLi2(l, s, y = y),
                      start = list(l = median(y), s = IQR(y)),
                      lower = c(-Inf, 0), upper = c(Inf, Inf),
                      method = "L-BFGS-B")
summary(mle.negLogLik2)
```

Maximum likelihood estimation

Call:

```
mle(minuslogl = function(l, s) neglogLi2(l, s, y = y), start = list(l = median(y),
s = IQR(y)), method = "L-BFGS-B", lower = c(-Inf, 0), upper = c(Inf,
Inf))
```

Coefficients:

	Estimate	Std. Error
l	0.129	0.266
s	0.962	0.235

-2 log L: 141

```
# the maximum log-likelihood value,
# i.e. the log-likelihood function evaluated at the MLE
logLik(mle.negLogLik2)
```

'log Lik.' -70.6 (df=2)

```
# asymptotic variance-covariance matrix
vcov(mle.negLogLik2)
```

```
      1      s
1 0.07082 0.00425
s 0.00425 0.05520
```

This latter is nothing but the inverse of the negative of the Hessian matrix (of the log-likelihood) evaluated at the maximum likelihood.

```
solve(mle.negLogLik2@details$hessian)
```

```
      1      s
1 0.07082 0.00425
s 0.00425 0.05520
```

Remarks

- In the example above, the same solution can be found directly via `optim()` as follows

```
neglogLi2 <- function(ls, y) {
  -sum(dcauchy(y, location = ls[1], scale = ls[2], log = TRUE))
}
optim(c(l = median(y), s = IQR(y)), fn = neglogLi2, y = y, method = "L-BFGS-B",
      lower = c(-Inf, 0), upper = c(Inf, Inf), hessian = TRUE)
```

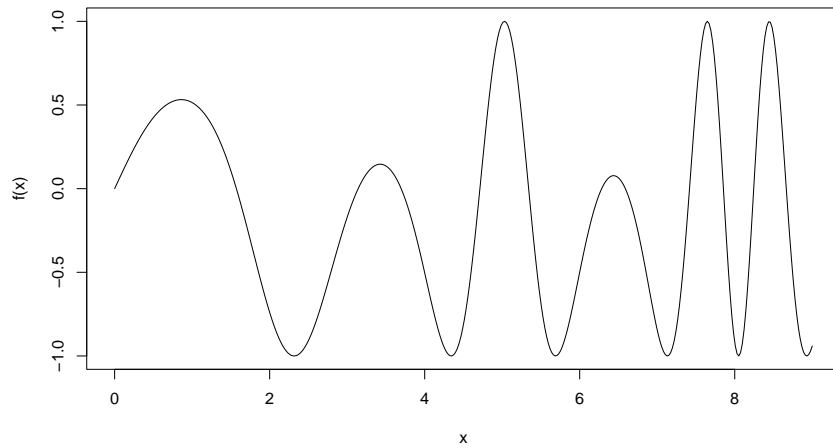
- The same can be done as follows; see the Help for more details

```
MASS::fitdistr(y, densfun = "cauchy")
```

```
location    scale
0.129       0.962
(0.266)    (0.235)
```

- Be careful when performing numerical optimization, as there is no guarantee that the resulting convergence point is actually the global maximizer. Here is an illustration.

```
f <- function(x) sin(x * cos(x))
curve(f, 0, 9, n = 400)
```



```
mle(\(x) -f(x), start = 2)@coef
[1] 0.86
mle(\(x) -f(x), start = 4)@coef
[1] 3.43
mle(\(x) -f(x), start = 5)@coef
[1] 5.03
mle(\(x) -f(x), start = 6)@coef
[1] 8.45
```

The problem becomes more and more difficult as the number of parameters increases.

7.2 Profile likelihood

Although the definition of likelihood covers the *multiparameter case*, the resulting multidimensional likelihood function can be difficult to deal with. Furthermore, in many practical multi-parameter problems, only a subset of the parameters is of interest; in the normal model, we might be interested only in the mean μ , while σ^2 is a “nuisance”, which is there only to make the model correct. And even if we are interested in several parameters, it is always easier to study a single parameter or only a very few at once.

For a given model, let (θ, η) be the full set of parameters (θ and η may be vectors), where both θ and η are unknown. Let’s say that θ is our primary parameter of interest and η is a nuisance parameter.

Definition 7.1 (Profile likelihood). Given a model with (full) likelihood $L(\theta, \eta)$, the *profile likelihood* for θ is

$$L_p(\theta) = \max_{\eta} L(\theta, \eta).$$

Let $\hat{\eta}(\theta) = \arg \max_{\eta} L(\theta, \eta)$, i.e. $\hat{\eta}(\theta)$ is the maximizer of the function $\eta \mapsto L(\theta, \eta)$ when θ is regarded as known. In this way, the definition of the profile likelihood can be made more explicit using

$$L_p(\theta) = L(\theta, \hat{\eta}(\theta)).$$

Let $\hat{\theta} = \arg \max_{\theta} L_p(\theta)$, a bit of logical deduction shows that $(\hat{\theta}, \hat{\eta}(\hat{\theta}))$ is the MLE of (θ, η) . In fact,

$$L(\hat{\theta}, \hat{\eta}(\hat{\theta})) = L_p(\hat{\theta}) \geq L_p(\theta) = \max_{\eta} L(\theta, \eta) \geq L(\theta, \eta), \forall (\theta, \eta).$$

This procedure is illustrated in the following example.

Example 7.2. Let $X_i, i = 1, \dots, n$, be an iid sample from $N(\mu, \sigma^2)$, where $\mu \in (-\infty, \infty)$, and $\sigma^2 \in (0, \infty)$ are unknown. The log-likelihood is

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

We know that if we consider μ as known and maximize this function with respect to σ^2 we get the following explicit expression for the MLE of σ^2

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

So, the profile log-likelihood for μ is

$$\ell_p(\mu) = \ell(\mu, \hat{\sigma}^2(\mu)) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \left(\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \right) - \frac{n}{2}.$$

Clearly, the $\arg \max_{\mu} \ell_p(\mu)$ is \bar{x}_n , and so the MLE of (μ, σ^2) is $(\bar{X}_n, n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2)$. \square

In addition to reducing the dimension of the problem, profile likelihood functions can be used in the same way as (ordinary) likelihood functions not only to estimate efficiently the parameter(s) of interest, but also to obtain their asymptotic variance, which is equal to the inverse of the negative of the second derivative (or the Hessian) of the profile log-likelihood. Profile likelihood is also very useful when it comes to inference; more on this later. Note, however, that in practice it is rarely possible to calculate the profile likelihood explicitly (as we did above for the Normal density); numerical evaluation is typically the way to proceed.

Exercise 7.1. Let $X_i, i = 1, \dots, n$, be an iid sample from the [gamma distribution](#) with pd

$$f(x; \alpha, \sigma) = \frac{1}{\sigma^{\alpha} \Gamma(\alpha)} x^{\alpha-1} \exp(-x/\sigma) I(x > 0),$$

where $\alpha > 0$ (shape), $\sigma > 0$ (scale) are unknown, and $\Gamma(\cdot)$ is the [gamma function](#).

Show that the maximum likelihood estimator of (α, σ) is $(\hat{\alpha}, \hat{\alpha}^{-1} \bar{X}_n)$ where $\hat{\alpha}$ is the arg max of

$$\ell_p(\alpha) = -n \log \Gamma(\alpha) - n\alpha \log(\alpha^{-1} \bar{X}_n) + (\alpha - 1) \sum_{i=1}^n \log(X_i) - n\alpha.$$

Use the above result to calculate the MLE of (α, σ) from the following simulated data. Estimate the asymptotic standard deviation of $\hat{\alpha}$.

```
set.seed(1)
x <- rgamma(n = 35, shape = 5, scale = 1/3)
x
```

```
[1] 1.090 2.588 2.535 1.808 0.609 1.864 2.068 1.935 1.292 0.988 1.093 1.296
[13] 1.468 0.936 2.137 1.949 2.220 2.104 1.553 0.423 1.877 2.366 1.051 1.181
[25] 1.338 1.428 1.076 0.684 2.284 1.212 2.379 2.088 1.386 1.197 1.919
```

7.3 Likelihood for regression models

Consider the simple linear regression model with Gaussian noise, defined as $Y = \beta_0 + \beta_1 x + \epsilon$. ϵ is an unobservable noise/error variable with $N(0, \sigma^2)$ distribution. The aim is to estimate and make inference about $\theta = (\beta_0, \beta_1, \sigma^2)$ on the basis of an iid sample (Y_i, x_i) of (Y, x) . For now, we assume that x_i are *known non-random* quantities (fixed design). According to this model $Y_i \equiv Y|X = x_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$, so the log-likelihood function is given by

$$\ell_n(\theta) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2.$$

It's easy to see that maximizing this function is equivalent to resolve the following optimization problems

$$(1) \ (\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{\beta_0, \beta_1} \sum_{i=1}^n (x_i - (\beta_0 + \beta_1 x_i))^2,$$

$$(2) \ \hat{\sigma}^2 = \arg \max_{\sigma^2} \ell_n(\hat{\beta}_0, \hat{\beta}_1, \sigma^2).$$

Equation (1) above is identical to the optimization equation that defines linear least-squares (LS) regression. So the MLE of (β_0, β_1) is the same the least squares estimator. In fact, setting the derivatives to zero and solving produces

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x}, \ \hat{\beta}_1 = \frac{\bar{xY} - \bar{x}\bar{Y}}{\bar{x^2} - \bar{x}^2}, \text{ and } \hat{\sigma}^2 = n^{-1} \sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2.$$

Maximum likelihood theory can of course be applied in this situation (Fisher information, asymptotic normality/efficiency, etc.).

For a random sampling design in which X_i are also sampled, we have that $Y_i|X_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$ and the (full or “joint”) likelihood is

$$\prod_i f_{Y|X}(y_i|x_i; \theta) f_X(x_i),$$

where $f_{Y|X}$ is the pd of $Y|X$, and f_X is the pd of X . If the latter does not depend on θ , as is commonly supposed, then f_X plays no role in the likelihood function. In this case, it is equivalent to work with the “joint” log-likelihood or with the “conditional” log-likelihood

as given by $\sum_i \log f_{Y|X}(y_i|x_i; \theta) = \ell_n(\theta)$, i.e. the same as for fixed design. Consequently, the estimators and information matrix are the same as before, and the asymptotic results are the same also.

```
set.seed(5)
x <- 1:10
y <- 10 + 20 * x + rnorm(10, sd = 10)

neglogLiReg <- function(a, b, s, y) {-sum(dnorm(y, mean = a + b * x, sd = s, log = TRUE))}
mle(\(a, b, s) neglogLiReg(a, b, s, y = y),
    start = list(a = mean(y), b = 0, s = sd(y))) |> coef()
      a      b      s
10.445 19.775  9.011

lm(y ~ x) |> coef()
(Intercept)          x
    10.45         19.77
```

The same theory applies to multiple linear regression. Without loss of generality, let's consider the case of 2 covariates X_1 and X_2 . The linear equation is $Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon$, with $\epsilon \sim N(0, \sigma^2)$. The log-likelihood (under fixed or random design) is

$$-\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu_i)^2,$$

where $\mu_i = \beta^t x_i$, $\beta^t = (\beta_0, \beta_1, \beta_2)$ and $x_i^t = (x_{i0}, x_{i1}, x_{i2})$, with $x_{i0} = 1, \forall i$. Here again, the maximum likelihood method and the LS method lead to the same estimators. $\hat{\beta}$ is the root of the system of equations

$$\nabla_{\beta} \sum_{i=1}^n (y_i - \mu_i)^2 := \left(\partial_{\beta_1} \sum_{i=1}^n (y_i - \mu_i)^2, \partial_{\beta_2} \sum_{i=1}^n (y_i - \mu_i)^2, \partial_{\beta_3} \sum_{i=1}^n (y_i - \mu_i)^2 \right) = \mathbf{0}.$$

Since, $\partial_{\beta_k} \sum_{i=1}^n (y_i - \mu_i)^2 = -2 \sum_i (y_i - \mu_i) x_{ik}$,

$$\nabla_{\beta} \sum_{i=1}^n (y_i - \mu_i)^2 = -2(X^t \mathbf{y} - X^t X \beta),$$

where $X = \begin{bmatrix} x_1^t \\ \vdots \\ x_n^t \end{bmatrix}$ and $\mathbf{y}^t = (y_1, \dots, y_n)$. It follows that $\hat{\beta} = (X^t X)^{-1} X^t \mathbf{y}$. The maximum

likelihood estimator of σ^2 has the same form as in the one-covariate case, i.e., the average squared residual $n^{-1} \sum_i (Y_i - \hat{\mu}_i)^2$, with $\hat{\mu}_i = \hat{\beta}^t x_i$.

Linear regression is the simplest regression model that can be estimated and inferred using

likelihood theory. This theory can be applied to a very wide range of other regression settings such as, for example, non-linear regression models, generalized linear models (based on the exponential family), generalized linear mixed models, censored regression models, to name only a few.

Chapter 8

Hypothesis testing

There are some problems we meet in statistical practice in which estimation of a parameter is not the primary goal of the analysis; rather, we wish to use our data to decide to trust or not a particular claim about a given population. For example, a drug company may claim that its new drug reduces “bad” cholesterol (or LDL cholesterol) level by more than 20 points (after a certain period of treatment). To try to prove this, the drug company may design a clinical trial and collect data on the level of LDL reduction in selected subjects. The observed data can be analyzed to confirm (or refute) the manufacturer’s claim. Properly formulated, such inference problems are called *testing of hypothesis* problems or simply test problems.

8.1 Basic concepts

Before attempting any inferential procedure, the model and all its underlying assumptions must clearly be defined. This is referred to the *maintained hypothesis*. These hypothesis define the framework within which the inferential procedure can be applied and interpreted correctly. In the context of the parametric hypothesis testing problem, the maintained hypothesis are simply the assumption made on the distribution of the observed data. In simple situations, this reduces to assuming that we observe an iid sample $\mathbf{X} = (X_1, \dots, X_n)$ from a pd $f(x; \theta)$, for some $\theta \in \Theta$ (real or a vector parameter space).

Based on the observed data, we have to decide whether $\theta \in \Theta_0$ or $\theta \in \Theta_1$, where Θ_0, Θ_1 are a partition of the parameter space Θ , i.e. $\Theta_0 \cap \Theta_1 = \emptyset$ and $\Theta_0 \cup \Theta_1 = \Theta$. We formulate this in the following manner:

$$H_0 : \theta \in \Theta_0 \text{ vs } H_1 : \theta \in \Theta_1.$$

The hypothesis H_0 is called the *null hypothesis* and H_1 is called the *alternative hypothesis*. In a statistical test procedure, these two hypothesis play an asymmetric role:

- H_0 represents the *current theory*: “null” means statu quo, no change or no effect. Typically, H_0 is also simpler (lower-dimensional) than the alternative hypothesis.

- H_1 is referred as the **researcher's hypothesis**: it is the new claim that one would really like to validate or the question to be answered.

Example 8.1. Take the cholesterol example. Suppose that, in a clinical trial, the observed LDL decrease of n subjects are $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, 25)$, for some $\mu \in (-\infty, \infty)$. Here, iid and normality (with var. = 25) are the maintained hypothesis. The drug company's claim is $\mu > 20$. This claim can be formulated as a test problem with two hypothesis: the null hypothesis $H_0 : \mu \leq 20$, and the alternative hypothesis $H_1 : \mu > 20$. \square

A hypothesis is called **simple** if it completely specifies the underlying distribution; otherwise it is called **composite**. In our example above, we are testing a composite null against a composite alternative.

In addition, a test is called **one-tailed** if the alternative hypothesis is articulated directionally ($H_1 : \theta < \theta_0$ or $H_1 : \theta > \theta_0$, for some given θ_0); otherwise it is called **two-tailed** ($H_1 : \theta \neq \theta_0$). In our example above, we have a one-tailed test.

We will use our data to choose between the two hypothesis H_0 and H_1 . For that, the basic idea is to try to “compare” the observed sample with the model under H_0 and under H_1 . The comparison is typically based on a statistic $T \equiv T_n(\mathbf{X})$, called the **test statistic**, which is designed in order to measure the discrepancy between the data and the models under H_0 and under H_1 . Most of the time, a test statistic is taken to be a (function of a) **sufficient statistic** for the parameter of interest θ , whose distribution (under H_0) is known exactly or asymptotically.

Assume that an appropriate test statistic T has been chosen and, without loss of generality, suppose that small values of T support H_0 , while large values support H_1 . The next step is to select a number, called the **critical value**, k and apply the following rule:

if $T_n(\mathbf{x}) \geq k$, reject H_0 in favor of H_1 ; otherwise, do not reject H_0 .

In this way, our test is nothing but the statistic/function $\varphi_n \equiv \varphi_n(\mathbf{X}) = I(T_n(\mathbf{X}) \geq k) : \mathbb{R}^n \mapsto \{0, 1\}$, called the **test function** or simply test, that takes only values 0 and 1:

- $\varphi_n(\mathbf{x}) = 1 \Leftrightarrow$ reject H_0 in favor of H_1 .
- $\varphi_n(\mathbf{x}) = 0 \Leftrightarrow$ do not reject H_0 .

The subset $\mathcal{R}_n = \{\mathbf{x} : T_n(\mathbf{x}) \geq k\}$ of the sample space for which H_0 will be rejected is called the **rejection region or critical region**. \mathcal{R}_n^c , the complement of the rejection region, is called the **non rejection region**.

Remark The tests used in most practical situations can be expressed as $I(T_n \leq k)$ or $I(T_n \geq k)$. We use $I(T_n \leq k)$ in situations where T_n tends to be small under H_1 and $I(T_n \geq k)$ in situations where T_n tends to be large under H_1 . In these notes, we will focus on the last case, but all the methodology can be applied directly to the first situation by using $-T_n$ as the actual test statistic.

Example 8.2. In the cholesterol example, our hypothesis are $H_0 : \mu \leq 20$ vs $H_1 : \mu > 20$.

As a consistent estimator of μ , \bar{X}_n tends to be larger under H_1 than under H_0 . Hence, it is natural to reject H_0 for large values of \bar{X}_n and so to consider, for example, the test function $\varphi_1 = I(\bar{X}_n \geq 21)$. The same reasoning can be applied to the sample median, say M_n , which leads to the test function $\varphi_2 = I(M_n \geq 21)$. The first test φ_1 uses \bar{X}_n as test statistic, while the second test φ_2 uses M_n as test statistic. The rejection region of φ_1 is $\{(x_1, \dots, x_n) : \bar{x}_n \geq 21\}$ and that of φ_2 is $\{(x_1, \dots, x_n) : m_n \geq 21\}$. \square

This example raises a couple of questions:

- Why 21 is used as the critical value? How to choose a suitable critical value?
- How to choose between φ_1 and φ_2 ? Is there an “optimal” test? If so, how to find it?

8.2 Evaluating a test

In a situation where the analyst has to decide between H_0 and H_1 , and given that in reality one of these hypothesis is true and the other is false, there are four possible states summarized in the following table:

Decision/Reality	H_0 is true	H_1 is true
Do not reject H_0	Correct	Error of type II
Reject H_0	Error of type I	Correct

For a given test φ_n , the function $\pi_n(\theta)$ defined on Θ by

$$\pi_n(\theta) := E_\theta(\varphi_n) = P_\theta(\varphi_n = 1)$$

is called the **power function** of the test φ_n . This function indicates the probability of rejecting H_0 for every possible value of θ ; i.e. $\theta \mapsto P_\theta(\text{Reject } H_0)$. It plays an important role in evaluating the quality of a test and in comparing different test procedures.

Related to the power function we define:

- **Probability of type I error.** This is the probability that the test rejects a true null hypothesis; i.e. $P_\theta(\text{Reject } H_0 | H_0) = P(\varphi_n = 1 | \theta \in \Theta_0)$.
- The **size** is the largest probability of committing a type I error, i.e. $\max_\theta P(\varphi_n = 1 | \theta \in \Theta_0)$. If the size of a test is known not to exceed a given $\alpha \in (0, 1)$, i.e. $\max_\theta P(\varphi_n = 1 | \theta \in \Theta_0) \leq \alpha$, then we say that the test is of (significance) **level** α .
- **Power.** This is the probability that the test rejects a false null hypothesis; i.e. $P_\theta(\text{Reject } H_0 | H_1) = P(\varphi_n = 1 | \theta \in \Theta_1)$.
- **Probability of type II error.** This is the probability that the test does not reject a false null hypothesis; i.e. $P_\theta(\text{Do not reject } H_0 | H_1) = P(\varphi_n = 0 | \theta \in \Theta_1) \equiv 1 - \text{Power}$.

Example 8.3. Recall our cholesterol example, with $H_0 : \mu \leq 20$ vs $H_1 : \mu > 20$. Assuming that $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, 25)$, the power function of the test $\varphi_n = I(\bar{X}_n \geq 21)$ is

$$\pi_n(\mu) := P_\mu(\bar{X}_n \geq 21) = 1 - \Phi\left(\sqrt{n} \frac{21 - \mu}{5}\right),$$

where Φ is the cdf of $N(0, 1)$. As $\mu \mapsto \pi_n(\mu)$ is an increasing function in μ , the size of our test, i.e. $\max_{\mu \leq 20} \pi_n(\mu)$, is $1 - \Phi(\sqrt{n}/5)$. For example, with $n = 10$, the size is 0.2635 and with $n = 100$ the size is 0.0228.

Here is the plot of $\pi_n(\mu)$, for different values of the sample size n .

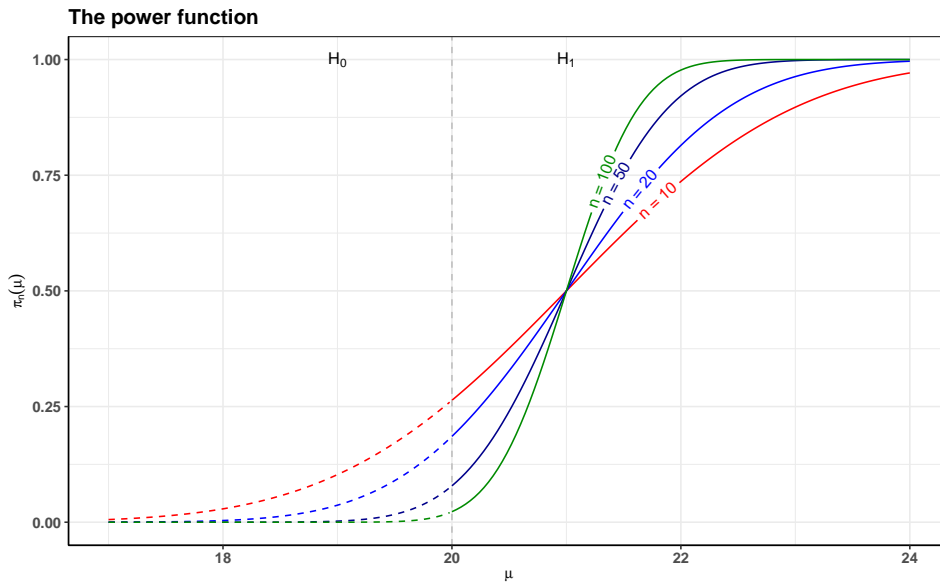


Figure 8.1: Power (solid line) and Type-I-error probability (dashed line)

For n fixed, as the true μ moves away from 20, the power increases to 1 and the probability of type I error decreases to 0. On the other hand, for fixed μ , as n increases, the power increases (for μ larger than 21) and the probability of type I error decreases. \square

Ideally, we would like to have a statistical test with the lowest probability of Type I error (0) and the highest power (1). For a fixed sample size, *it is usually impossible to control both Type I error and power*. In fact, to minimize the probability of type I error, one must not reject H_0 more often and to maximize the power, one must reject H_0 more often: these two goals work against each other. Typically, even in very simple problems, larger power comes at the expense of a larger probability of Type I error. Vice versa, when one tries to reduce the probability of type I error, the power also gets reduced. All that can be done is to control one of these two quantities by adjusting the critical value of the test, or, if possible, by increasing the sample size. Here is an example.

Example 8.4. Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(\pi)$, where $\pi \in \{0.3, 0.5\}$. Consider testing

$$H_0 : \pi = 0.5 \text{ vs } H_1 : \pi = 0.3.$$

Put $T_n = \sum_{i=1}^n X_i$. It is reasonable to reject the null hypothesis if the total number of successes T_n is “too small”. Let’s define the test $\varphi_n = I(T_n \leq k)$, where k is the critical value that needs to be fixed. The power function of φ_n is given by

$$\pi_n(\pi) = P_\pi(T_n \leq k) = \sum_{\ell=0}^k C_n^\ell \pi^\ell (1 - \pi)^{(n-\ell)}.$$

For example, with $n = 15$ we get (using the R function `pbinom(k, size = n, prob = π)`)

	Type-I-Error	Power
k	$P(T_n \leq k \pi = 0.5)$	$P(T_n \leq k \pi = 0.3)$
1	4.88×10^{-4}	0.03
4	0.06	0.52
7	0.50	0.95
10	0.94	≈ 1

For example, we can see that the test $\varphi_n = I(T_n \leq 4)$ gives a low type-I-error probability (0.06), but gives also a low power (0.52) and therefore does not provide adequate protection against a type-II-error (probability 0.48). We can increase the power to 0.95 by taking $\varphi_n = I(T_n \leq 7)$. However, by doing this, we increase the type-I-error probability to 0.5.

This demonstrates the conjunction of type-I-error probability and power of a test. The only way to decrease the probability of type I error and to increase the power simultaneously is to increase the sample size and to choose an appropriate critical value. For example, with $n = 150$ and $k = 60$, we get

k	Type-I-Error	Power
60	0.009	0.996

□

8.3 Testing strategy: the Neyman-Pearson approach

Given that it is not possible to reduce simultaneously the probabilities of both types of errors, [Neyman](#) and [Pearson](#) suggested to hold one of the two error probabilities at a pre-specified level that can be considered tolerable, and then minimize the other error probability. More specifically, Neyman-Pearson framework of statistical hypothesis is to

- (1) choose a small number $\alpha \in (0, 1)$ ($\alpha = 0.01, 0.05$, and 0.1 are commonly used in practice). *Restrict attention to tests with the pre-chosen level α only,*
- (2) *among these tests, try to find a test that has the largest power.* When it exists, such a test is called the **uniformly most powerful test** (UMP) of level α . More precisely, a test φ^*

of level α is UMP if $\forall \varphi$ of level α , $\pi_n^{\varphi^*}(\theta) \geq \pi_n^{\varphi}(\theta)$, $\forall \theta \in \Theta_1$, $\forall n$.

Since reducing type I error also reduces power, to obtain a UMP test (with a pre-chosen α level), it's necessary to restrict attention to tests with a size as close as possible to α (but not exceeding α).

To apply the Neyman-Pearson approach, two questions must be answered:

- (1) how to control the level of a given test and fix it at the pre-chosen α ?
- (2) how to find an UMP test (*when it exists*) ?

We will defer this last question to Section 8.5 and discuss here the first one. Before, observe that the Neyman-Pearson approach makes H_0 **and** H_1 **asymmetric**: we control the type-I-error probability to be at most α ; then, we try to make the type-II-error probability *as small as possible*. However, this latter can still be very big even for an UMP test. If we take this approach, the test must be set up so that H_1 , the alternative hypothesis, is the one we seek to prove (this is why we refer to it as the *researcher's hypothesis*). By using an α level test, with small α , we guard against saying that the data support the research hypothesis (H_1) when it is false.

Now regarding our first question, without loss of generality and to make things clearer, let's consider the case where our test is given by $\varphi_n = I(T_n \geq k)$, i.e. large values of T_n give evidence against H_0 . To fix the size $\max_{\theta} P(\varphi_n = 1 | \theta \in \Theta_0)$ to α , we must choose an appropriate critical value k . More precisely, we must choose $k \equiv k(n, \alpha)$ so that

$$\max_{\theta} P(T_n \geq k | \theta \in \Theta_0) = \alpha.$$

Finding such a k is greatly simplified if **the null distribution of the test statistic T_n** (i.e. the distribution of T_n under H_0) is known; otherwise, k can be obtained by *Monte Carlo methods* (not discussed here). Sometimes it is not possible to set the size exactly to α , in which case one must choose *the smallest possible k* such that

$$\max_{\theta \in \Theta_0} P_{\theta}(T_n \geq k) \leq \alpha.$$

Because $k \mapsto P(T_n \geq k)$ is a decreasing function, choosing the smallest possible k , allows *the maximum admissible type I error (α)* and thus the maximum power.

Once the critical value k is chosen, the next step is to calculate $t_n = T_n(\mathbf{x})$, the observed value of T_n , and then to take one of the following decisions:

- If $t_n \geq k$, then, at level α , the data **reject** H_0 in favor of H_1 because the event $\{T_n \geq k\}$ is quite rare (it occurs with probability at most α) under H_0 . We say that the **test is significant**.
- If $t_n < k$, then, at level α , the data **do not reject** H_0 . In this case, either
 - the power is known to be “very large” \implies **accept** H_0 , i.e reject H_1 in favor of H_0
 - or the power is unknown or known to be “small”. In this case, accepting H_0 is a

risky decision \implies we don't have enough information (i.e. sufficient sample size) to accept or reject H_0 . In this case, we say that the *test is non significant*.

Example 8.5. Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, where σ is known. For some given μ_0 , consider testing

$$H_0 : \mu \leq \mu_0 \text{ vs } H_1 : \mu > \mu_0.$$

We have seen that it is natural to reject H_0 when \bar{X}_n is "large". For ease of calculation, it is convenient to replace \bar{X}_n by the test statistic $Z_n = \sqrt{n} \frac{\bar{X}_n - \mu_0}{\sigma}$. So we consider the test $\varphi_n = I(Z_n \geq k)$. All we need now is to choose k . We have that

$$\begin{aligned} P_\mu(Z_n \geq k) &= P\left(Z \geq k + \sqrt{n} \frac{\mu_0 - \mu}{\sigma}\right), \text{ with } Z \sim N(0, 1) \\ \Rightarrow \max_{\mu \leq \mu_0} P_\mu(Z_n \geq k) &= P(Z_n \geq k | \mu = \mu_0) = P(Z \geq k). \end{aligned}$$

$\max_{\mu \leq \mu_0} P_\mu(Z_n \geq k) = \alpha \Leftrightarrow P(Z \geq k) = \alpha \Leftrightarrow k = z_{1-\alpha}$, i.e. the $(1 - \alpha)$ -quantile of $N(0, 1)$ (in R `qnorm(1- α)`). To conclude, the test

$$\varphi_n = I(Z_n \geq z_{1-\alpha}),$$

or equivalently $\varphi_n = I(\bar{X}_n \geq \mu_0 + \frac{\sigma}{\sqrt{n}} z_{1-\alpha})$, is of size α .

The power function of this test is given by

$$\pi_n(\mu) = P\left(Z \geq z_{1-\alpha} + \sqrt{n} \frac{\mu_0 - \mu}{\sigma}\right) = 1 - \Phi\left(z_{1-\alpha} + \sqrt{n} \frac{\mu_0 - \mu}{\sigma}\right).$$

Observe that $\pi_n(\mu) \xrightarrow[n \rightarrow \infty]{} 1, \forall \mu > \mu_0$, such a test is said to be **consistent**. More precisely, *a test of level α is consistent if its power converges to 1 as n tends to infinity*. This property is often considered necessary for a test to be useful in practice.

Here is the plot of $\pi_n(\mu)$, for $\mu_0 = 20$, $\sigma = 5$ and $\alpha = 5\%$, for different values of the sample size n .

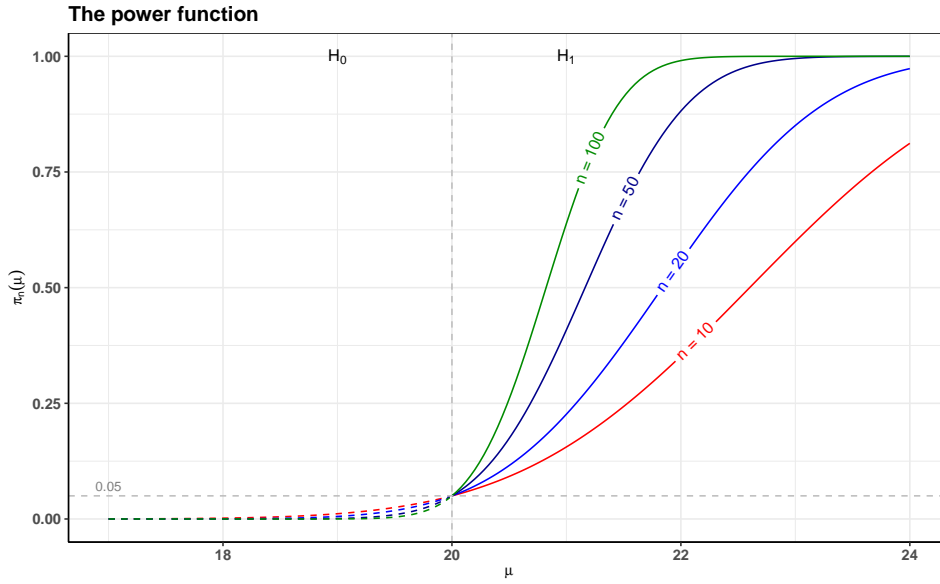


Figure 8.2: Power (solid line) and Type-I-error probability (dashed line)

8.4 The p -Value

The classic approach described above, based on critical values, can be summarized in two main steps:

1. specify the significance level α ,
2. calculate the test statistic and compare it to the critical value, then reject/not reject H_0 .

Simply reporting that a given hypothesis is rejected is not very informative. This says nothing about the fact that the computed value of the test statistic just barely fell into the rejection region or whether it exceeded the critical value by a large amount. What we need is a relevant way to *measure the strength of evidence* against a null hypothesis. And this is exactly what a p -value does.

To make things clearer, let's consider the test we studied in the Example 8.5 above. Recall that our decision is to reject H_0 whenever z_n , the observed value of the test statistic Z_n , is greater than $z_{1-\alpha}$.

Since $z \mapsto P_{\mu_0}(Z_n \geq z) \equiv P(Z_n \geq z | \mu = \mu_0)$ is a decreasing function, we have that

$$\begin{aligned} z_n \geq z_{1-\alpha} &\Leftrightarrow P_{\mu_0}(Z_n \geq z_n) \leq P_{\mu_0}(Z_n \geq z_{1-\alpha}) \\ &\Leftrightarrow P(Z \geq z_n) \leq \alpha, \quad Z \sim N(0, 1). \end{aligned}$$

The quantity $P(Z \geq z_n)$, i.e. $P(Z_n \geq z_n | H_0)$, is called the p -value. From the above calculation, we can reformulate our test decision as follows:

- $p\text{-value} \leq \alpha \Rightarrow$ reject H_0 at level α .
- $p\text{-value} > \alpha \Rightarrow$ do not reject H_0 at level α .

A more general definition of the p -value can be stated as follows. Let $T \equiv T(X)$ be a test statistic such that large values of T give evidence against H_0 . Given $t = T(x)$, the observed sample value of T , the p -value of such a test is

$$\max_{\theta \in \Theta_0} P_{\theta}(T \geq t).$$

In the particular case where $\Theta_0 = \{\theta_0\}$, the above p -value reduces to $P_{\theta_0}(T \geq t) = P(T \geq t | \theta = \theta_0) \equiv P(T \geq t | H_0)$ which can be seen as *the probability of observing a test statistic at least as “extreme” (i.e. at the opposite of H_0 but along the direction of H_1) as the one actually observed, assuming the null hypothesis H_0 is true*. This is a simple and an intuitive definition of the p -value which is used in many textbooks. *The smaller the p -value, the more evidence there is in the observed data against the null hypothesis and in favor of the alternative hypothesis*. Note that, unlike the critical value, the calculation of the p -value does not involve the level of the test α .

Example 8.6. Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, where σ is known. For some given μ_0 , consider testing

$$H_0 : \mu = \mu_0 \text{ vs } H_1 : \mu \neq \mu_0.$$

A natural test function is given by $\varphi_n = I(|Z_n| \geq k)$, where $Z_n = \sqrt{n} \frac{\bar{X}_n - \mu_0}{\sigma}$. The test statistic here is $|Z_n|$.

For this test to be of size α , we need to choose k so that $P_{\mu_0}(|Z_n| \geq k) = \alpha$. Thus, $k = z_{1-\alpha/2}$. So, based on this critical value, one should reject H_0 , at level α , if $|z_n| \geq z_{1-\alpha/2}$.

An alternative and more convenient approach to do this test is to calculate its p -value which, according to the definition above, is given by

$$P_{\mu_0}(|Z_n| \geq |z_n|) = P(|Z| \geq |z_n|) = 2P(Z \geq |z_n|), \text{ with } Z \sim N(0, 1).$$

Based on this, one should reject H_0 , at level α , if $2P(Z \geq |z_n|) \leq \alpha$. It is easy to see that

$$|z_n| \geq z_{1-\alpha/2} \Leftrightarrow 2P(Z \geq |z_n|) \leq \alpha. \square$$

Exercise 8.1. Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, where σ is known. Find a consistent test of size α for

$$H_0 : \mu \geq \mu_0 \text{ vs } H_1 : \mu < \mu_0.$$

Calculate its p -value and its power function.

8.5 Likelihood Ratio Test (LRT)

The likelihood ratio test (LRT) is the most general and the most used test in practice. This method provides a unified approach for developing *optimal test procedures*. In fact, it have

been shown that *whenever the uniformly most powerful test exists, the LRT procedure leads to it*. In addition to hypothesis testing, the LRT procedure can also be used to construct optimal confidence intervals.

Suppose we observe an iid sample $\mathbf{X} = (X_1, \dots, X_n)$ from a pd $f(x; \theta)$, for some $\theta \in \Theta$ (Θ can be a vector space). The **likelihood ratio statistic** for testing $H_0 : \theta \in \Theta_0$ vs $H_1 : \theta \in \Theta_1$, with $\Theta_0 \cap \Theta_1 = \emptyset$ and $\Theta_0 \cup \Theta_1 = \Theta$, is defined by

$$\Lambda_n(\mathbf{x}) := \frac{\max_{\theta \in \Theta_0} L_n(\theta|\mathbf{x})}{\max_{\theta \in \Theta} L_n(\theta|\mathbf{x})} = \frac{\max_{H_0} L_n(\theta|\mathbf{x})}{\max_{H_0 \cup H_1} L_n(\theta|\mathbf{x})} = \frac{L_n(\hat{\theta}_0)}{L_n(\hat{\theta})},$$

where $L_n(\theta|\mathbf{x}) = \prod_i f(x_i; \theta)$ is the likelihood function of θ based on of the observed sample \mathbf{x} ,

$\hat{\theta} \equiv \hat{\theta}(\mathbf{x}) = \arg \max_{\theta \in \Theta} L_n(\theta|\mathbf{x})$ is the MLE of θ obtained by maximizing the likelihood over the whole parameter space Θ , and $\hat{\theta}_0 \equiv \hat{\theta}_0(\mathbf{x}) = \arg \max_{\theta \in \Theta_0} L_n(\theta|\mathbf{x})$ be the **restricted** MLE of θ obtained by maximizing the likelihood over the null parameter space $\Theta_0 \subset \Theta$. Notice that, by definition, $0 \leq \Lambda_n \leq 1$.

To understand the reasoning behind the LRT, consider the discrete case where $f(x; \theta)$ is a probability mass function. In this case, given an observed data \mathbf{x} , Λ_n can be expressed as

$$\frac{P_{\hat{\theta}_0}(\mathbf{X} = \mathbf{x})}{P_{\hat{\theta}}(\mathbf{X} = \mathbf{x})} = \frac{\max_{\theta \in \Theta_0} P_{\theta}(\mathbf{X} = \mathbf{x})}{\max_{\theta \in \Theta_0 \cup \Theta_1} P_{\theta}(\mathbf{X} = \mathbf{x})} = \frac{\max \text{Prob}(\text{obs. data} | H_0)}{\max \text{Prob}(\text{obs. data} | H_0 \cup H_1)}.$$

- A small value of this ratio means that the observed data are less likely to occur under the null hypothesis than in the case where no restrictions (as imposed by the null hypothesis) are applied. This means that there should be a parameter point in Θ_1 for which the observed sample is more likely than for any parameter point in Θ_0 . Consequently, H_0 has to be rejected in favor of H_1 .
- A value of this ratio close to 1 means that the observed data are nearly as likely to occur under the null hypothesis as without it. H_0 is therefore not binding and there is no reason to reject it.

This motivates the definition of LRT as $I(\Lambda_n \leq c)$, for some $c \in [0, 1)$. In this way, the LRT rejects H_0 whenever $\Lambda_n \leq c$. c must be chosen, in $[0, 1)$, so that the test level is α , i.e. $\max_{\theta \in \Theta_0} P_{\theta}(\Lambda_n \leq c) \leq \alpha$. According to the definition given above, the p -value of the LRT is $\max_{\theta \in \Theta_0} P_{\theta}(\Lambda_n \leq \lambda_n)$, where λ_n is the observed value of Λ_n , i.e. $\lambda_n = \Lambda_n(\mathbf{x})$.

To perform LRT, we need to (i) maximize the likelihood function L_n over the full and restricted parameter spaces and (ii) calculate the critical value or the p -value as defined above. To do this, we need to know the “null distribution” of Λ_n (i.e. its distribution under H_0) or to express it as a monotonic function of some statistic T_n whose null distribution is known. The following examples illustrate this.

Example 8.7. Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$. The likelihood function of (μ, σ^2) is given by

$$L_n(\mu, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \right)$$

Assume that $\sigma^2 > 0$ is known.

- Let's find the LRT for testing $H_0 : \mu = \mu_0$ vs $H_1 : \mu \neq \mu_0$.

Since $\max_{\mu=\mu_0} L_n(\mu, \sigma^2) = L_n(\mu_0, \sigma^2)$ and $\max_{\mu} L_n(\mu, \sigma^2) = L_n(\bar{X}_n, \sigma^2)$,

$$\begin{aligned} \Lambda_n &= \frac{L_n(\mu_0, \sigma^2)}{L_n(\bar{X}_n, \sigma^2)} = \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n ((X_i - \mu_0)^2 - (X_i - \bar{X}_n)^2) \right) \\ &= \exp \left(-\frac{n}{2\sigma^2} (\bar{X}_n - \mu_0)^2 \right) = \exp \left(-\frac{1}{2} Z_n^2 \right), \end{aligned}$$

where $Z_n = \sqrt{n} \frac{\bar{X}_n - \mu_0}{\sigma}$.

The LRT rejects H_0 if $\Lambda_n \leq c$ (with $c \in [0, 1)$) $\Leftrightarrow |Z_n| \geq k$, where $k = \sqrt{-2 \log(c)} > 0$. This latter is determined so that the size of the test is α , that is $P_{\mu_0}(|Z_n| \geq k) = \alpha$. This is equivalent to say that $P(|Z| \geq k) = \alpha$, thus $k = z_{1-\alpha/2}$. We conclude that the LRT of level α for the above hypothesis is equivalent to $I(|Z_n| \geq z_{1-\alpha/2})$.

The p -value of this test is given by

$$\max_{\mu: \mu=\mu_0} P(\Lambda_n \leq \lambda_n) = P(\Lambda_n \leq \lambda_n | \mu = \mu_0) = P(|Z_n| \geq |z_n| | \mu = \mu_0) = P(|Z| \geq |z_n|),$$

where $\lambda_n = \exp(-\frac{1}{2} z_n^2)$ is the observed value of Λ_n , and z_n is the observed value of Z_n .

- Let's find the LRT for testing $H_0 : \mu \leq \mu_0$ vs $H_1 : \mu > \mu_0$.

Since $\mu \mapsto \ell_n(\mu, \sigma^2) = \log L_n(\mu, \sigma^2)$ is strictly concave and has a (unique) maximum at \bar{X}_n , we have that

$$\arg \max_{\mu \leq \mu_0} L_n(\mu, \sigma^2) = \arg \max_{\mu \leq \mu_0} \ell_n(\mu, \sigma^2) = \begin{cases} \bar{X}_n & \text{if } \bar{X}_n \leq \mu_0 \\ \mu_0 & \text{if } \bar{X}_n > \mu_0 \end{cases}$$

So,

$$\Lambda_n = \begin{cases} 1 & \text{if } \bar{X}_n \leq \mu_0 \\ \frac{L_n(\mu_0, \sigma^2)}{L_n(\bar{X}_n, \sigma^2)} = \exp(-\frac{1}{2} Z_n^2) & \text{if } \bar{X}_n > \mu_0 \end{cases}$$

The LRT rejects H_0 if $\Lambda_n \leq c$ (with $c \in [0, 1)$) $\Leftrightarrow |Z_n| \geq k$ **and** $\bar{X}_n > \mu_0 \Leftrightarrow Z_n \geq k$ (with $k > 0$). This latter is determined so that the size of the test is α , that is $\max_{\mu \leq \mu_0} P_{\mu}(Z_n \geq k) = \alpha$. This is equivalent to say that $P(Z \geq k) = \alpha$, thus $k = z_{1-\alpha}$. We conclude that the LRT of level α for the above hypothesis is equivalent to $I(Z_n \geq$

$z_{1-\alpha}$). The p -value of this test is

$$\begin{aligned} \max_{\mu: \mu \leq \mu_0} P(\Lambda_n \leq \lambda_n) &= I(\bar{X} \leq \mu_0) + I(\bar{X} > \mu_0) \max_{\mu \leq \mu_0} P(Z_n \geq z_n) \\ &= I(\bar{X} \leq \mu_0) + I(\bar{X} > \mu_0) P(Z \geq z_n). \end{aligned}$$

Assume now that $\sigma^2 > 0$ is *unknown*.

- Let's find the LRT for testing $H_0 : \mu = \mu_0$ vs $H_1 : \mu \neq \mu_0$.
We know that

$$\max_{\mu=\mu_0, \sigma^2} L_n(\mu, \sigma^2) = L_n(\mu_0, \tilde{\sigma}_0^2) \text{ and } \max_{\mu, \sigma^2} L_n(\mu, \sigma^2) = L_n(\bar{X}_n, \hat{\sigma}_n^2),$$

where $\tilde{\sigma}_0^2 = n^{-1} \sum_{i=1}^n (X_i - \mu_0)^2$ and $\hat{\sigma}_n^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$. So,

$$\begin{aligned} \Lambda_n &= \frac{L_n(\mu_0, \tilde{\sigma}_0^2)}{L_n(\bar{X}_n, \hat{\sigma}_n^2)} = \left(\frac{\tilde{\sigma}_0^2}{\hat{\sigma}_n^2} \right)^{-n/2}, \text{ with} \\ \frac{\tilde{\sigma}_0^2}{\hat{\sigma}_n^2} &= \frac{\sum_i (X_i - \mu_0)^2}{\sum_i (X_i - \bar{X}_n)^2} = \frac{\sum_i (X_i - \bar{X}_n)^2 + n(\bar{X}_n - \mu_0)^2}{\sum_i (X_i - \bar{X}_n)^2} = 1 + \frac{n(\bar{X}_n - \mu_0)^2}{\sum_i (X_i - \bar{X}_n)^2} \\ &= 1 + \frac{n(\bar{X}_n - \mu_0)^2}{(n-1)S_n^2} = 1 + \frac{T_n^2}{n-1}, \end{aligned}$$

where $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ and $T_n = \sqrt{n} \frac{\bar{X}_n - \mu_0}{S_n}$. We see that $\Lambda_n \leq c \iff |T_n| \geq k$. And since under H_0 , $T_n \sim t_{n-1}$, where t_{n-1} is the Student's t-distribution with $n-1$ degrees of freedom, the LRT of size α for the above hypothesis is equivalent to $I(|T_n| \geq t_{n-1; 1-\frac{\alpha}{2}})$.

- Let's find the LRT for testing $H_0 : \mu \leq \mu_0$ vs $H_1 : \mu > \mu_0$.
Under H_0 , the MLE of (μ, σ^2) is

$$\begin{cases} (\bar{X}_n, \hat{\sigma}_n^2) & \text{if } \bar{X}_n \leq \mu_0 \\ (\mu_0, \tilde{\sigma}_0^2) & \text{if } \bar{X}_n > \mu_0 \end{cases}$$

So,

$$\Lambda_n = \begin{cases} 1 & \text{if } \bar{X}_n \leq \mu_0 \\ \frac{L_n(\mu_0, \tilde{\sigma}_0^2)}{L_n(\bar{X}_n, \hat{\sigma}_n^2)} = \left(1 + \frac{T_n^2}{n-1}\right)^{-n/2} & \text{if } \bar{X}_n > \mu_0 \end{cases}$$

The LRT rejects H_0 if $\Lambda_n \leq c$ (with $c \in [0, 1]$) $\iff |T_n| \geq k$ **and** $\bar{X}_n > \mu_0 \iff T_n \geq k$. By the same reasoning as before (see the case of known σ), we conclude that the LRT of level α is equivalent to $I(T_n \geq t_{n-1; 1-\alpha})$. \square

Exercise 8.2. Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$. Find the LRT of level α for testing $H_0 : \mu \geq \mu_0$ vs $H_1 : \mu < \mu_0$, assuming first that σ is known and then without this assumption.

8.6 Asymptotic Tests

LRT construction, as discussed in the previous section, is a complicated process that requires a case-by-case analysis. The main difficulty is to find the *exact* distribution, under H_0 , of Λ_n (or some function of it), so that the appropriate critical (or p -) value can be calculated.

A way to overcome this difficulty is to rely on asymptotic theory, which, as we will see, allows the development of unified test procedures that can be applied to a wide range of problems. Unlike exact tests, asymptotic tests (those based on asymptotic theory) can, by definition, only be applied in situations where *the sample size is large enough*. We say that a test with power function $\pi_n(\theta)$ is of *asymptotic level* α if

$$\lim_{n \rightarrow \infty} \pi_n(\theta) \leq \alpha, \forall \theta \in \Theta_0.$$

There are three classical asymptotic test methods based on the likelihood function: the *Wald* tests, the *Score* (or *Rao*) tests, and the (log-) *likelihood ratio* tests. Under some regularity conditions, these tests can be shown to be asymptotically equivalent. As such, *for large samples*, no test is uniformly more powerful than the others, yet each test has its strengths and weaknesses (see below for details). That said, broadly speaking, LRT offers the best possible power for small to medium-sized samples.

Let $\hat{\theta}$ denote the MLE of $\theta \in \Theta \subset \mathbb{R}^d$. In what follows, the *regularity assumptions* required to ensure that $\sqrt{I_n(\hat{\theta})}(\hat{\theta} - \theta) \xrightarrow{d} N_d(\mathbf{0}, \mathbb{I})$ are supposed to be fulfilled.

8.6.1 Simple null hypothesis

To aid intuition, we first consider the simplest case of one parameter ($d = 1$). Let

$$H_0 : \theta = \theta_0 \text{ vs } H_1 : \theta \neq \theta_0.$$

The three test statistics of interest are:

$$W_n \equiv W_n(\theta_0) := I_n(\hat{\theta})(\hat{\theta} - \theta_0)^2 \quad (\text{Wald})$$

$$R_n \equiv R_n(\theta_0) := \frac{S_n^2(\theta_0)}{I_n(\theta_0)} \quad (\text{Score})$$

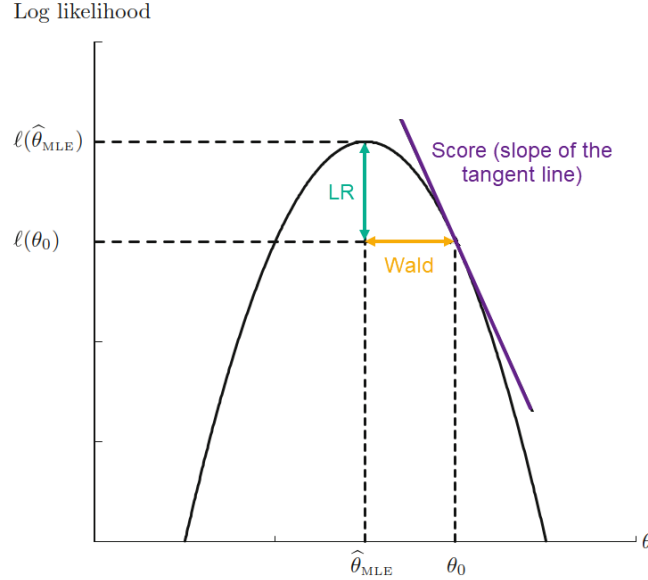
$$T_n \equiv T_n(\theta_0) := 2(\ell_n(\hat{\theta}) - \ell_n(\theta_0)) \quad (\text{Likelihood ratio})$$

In the above formulas $\ell_n(\theta) = \sum_i \log f(X_i, \theta)$ is the log-likelihood and S_n is the corresponding Score function, i.e. $S_n(\theta) = \ell'_n(\theta|X) = \sum_i \partial_\theta \log f(X_i, \theta)$.

Note that $T_n = -2 \log \frac{L_n(\theta_0)}{L_n(\hat{\theta})} = -2 \log \Lambda_n \in [0, \infty)$. This is the so-called *log-likelihood ratio statistic*, but for simplicity's sake we'll just call it the likelihood ratio (LR) statistic.

Although the above formulas defining the three statistics are very different, they are based

on the same principle of assessing the distance between the MLE $\hat{\theta}$ and the null value θ_0 : the greater the difference between these two, the stronger the evidence against the null hypothesis H_0 . Wald uses $(\hat{\theta} - \theta_0)^2$, which we can describe as a “horizontal” difference (see the figure below). LRT uses $(\ell_n(\hat{\theta}) - \ell_n(\theta_0))$, i.e. the “vertical” or likelihood difference. Score uses $(S_n(\hat{\theta}) - S_n(\theta_0))^2 = S_n^2(\theta_0)$, i.e. the “slope” (of the likelihood) difference. In any case, once the “distance” measurement method has been chosen, H_0 must be rejected if it is found to be “too far” away from the data (i.e. the test statistic is “too large”).



In the Wald statistic, $(\hat{\theta} - \theta_0)^2$ is multiplied by the Fisher information I_n to obtain a “standardized distance” with a fixed/known distribution. In fact, we have seen that $\sqrt{I_n(\hat{\theta})}(\hat{\theta} - \theta) \xrightarrow{d} N(0, 1)$, so, by Slutsky’s theorem, $W_n(\theta) \xrightarrow{d} \chi_1^2$ (here θ refers to the true parameter value). Thus, **under H_0** , $W_n \xrightarrow{d} \chi_1^2$.

Note that the Wald statistic can be expressed as $W_n = Z_n^2$, with

$$Z_n \equiv Z_n(\theta_0) := \frac{\hat{\theta} - \theta_0}{\sqrt{\widehat{Avar}(\hat{\theta})}},$$

with $\widehat{Avar}(\hat{\theta}) = I_n^{-1}(\hat{\theta})$. $I_n(\hat{\theta})$ can be replaced by the observed Fisher information $J_n(\hat{\theta}) = -\ell_n''(\hat{\theta})$, or by any asymptotically equivalent quantity, without altering the asymptotic distribution of W_n .

As for the LRT, observe that, by second-order Taylor’s expansion of $\theta \mapsto \ell_n(\theta)$ around $\hat{\theta}$,

$$\begin{aligned} \ell_n(\theta) - \ell_n(\hat{\theta}) &\approx \frac{1}{2}(\theta - \hat{\theta})^2 \ell_n''(\hat{\theta}) \\ \Rightarrow 2(\ell_n(\hat{\theta}) - \ell_n(\theta)) &\approx \frac{-n^{-1} \ell_n''(\hat{\theta})}{I(\theta)} \left(\sqrt{nI(\theta)}(\hat{\theta} - \theta) \right)^2 \end{aligned}$$

Since, $\sqrt{I_n(\theta)}(\hat{\theta} - \theta) \xrightarrow{d} N(0, 1)$ and $-n^{-1} \ell_n''(\hat{\theta}) \xrightarrow{p} I(\theta)$, we conclude that $T_n(\theta) \xrightarrow{d} \chi_1^2$. Thus,

under H_0 , $T_n \xrightarrow{d} \chi_1^2$.

As for the Score test, using the first-order Taylor expansion of $\theta \mapsto S_n(\theta)$ around θ_0 , and following the same line of reasoning as above shows that, under H_0 , $S_n \xrightarrow{d} \chi_1^2$. Similarly to the Wald statistic, in the Score statistic, one could replace $I_n(\theta_0)$ by any consistent estimator of it like for example $J_n(\theta_0)$. Note also that the Score statistic can be expressed as

$$\left(\frac{S_n(\theta_0)}{\sqrt{\text{Var}(S_n(\theta_0))}} \right)^2.$$

To conclude, we can say that, under H_0 , the three statistics (Wald, Score, LR) share the same asymptotic distribution, namely the χ_1^2 . Based on this result, the Wald test, the asymptotic LRT and the Score test reject H_0 , in favor of H_1 , if the observed value of their statistics (i.e. W_n for Wald, R_n for the Score and T_n for LR) is $\geq \chi_{1;1-\alpha}^2$, where $\chi_{1;1-\alpha}^2$ is the $(1 - \alpha)$ -quantile of χ_1^2 . These tests have an *asymptotic size* equal to α . We can check this easily for the LRT, for example, by observing that its type I error probability is $P_{\theta_0}(T_n \geq \chi_{1;1-\alpha}^2) \rightarrow P(\chi_1^2 \geq \chi_{1;1-\alpha}^2) = \alpha$.

In practice, it is common to use the *asymptotic p-value* $P(\chi_1^2 \geq t_n)$ (t_n is the observed value of T_n) and to reject H_0 whenever this quantity is less than α . The exact same considerations applies to the other two statistics (Wald and Score).

Wald, Score, and LR tests are asymptotically equivalent: they reach the same decision with probability approaching 1 as $n \rightarrow \infty$. However, their performance can be *quite different for a finite sample size*. Each method has its strengths and limitations:

- Based on the statistic $Z_n = \sqrt{I_n(\hat{\theta})}(\hat{\theta} - \theta_0)$, it is straightforward to create one-sided Wald tests (e.g. tests of $H_0 : \theta = \theta_0$ vs $H_1 : \theta < \theta_0$ or $H_1 : \theta > \theta_0$), but this is more difficult with Score and likelihood ratio statistics.
- The Wald statistic is by far the simplest and the easiest to interpret. It yields immediate confidence intervals and it is largely available in standard computing packages.
- The Wald test is not limited to MLE estimation, one just need to know the asymptotic distribution of the estimator under studied.
- The Score test does not require the MLE $\hat{\theta}$ whereas the other two tests do. The Score test also tends to give the best Type I error rates for small sample sizes.
- The Score test and likelihood ratio test are invariant under reparameterization, whereas the Wald test is not. For example, the Wald test about a scale parameter σ depends on whether the null hypothesis is expressed as $H_0 : \sigma = \sigma_0$ or $H_0 : \sigma^2 = \sigma_0^2$.
- In general, the likelihood ratio is more difficult to compute than Wald or Score tests, but it tends to have the greatest power.

Example 8.8. Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$. Consider $H_0 : \mu = \mu_0$ vs $H_1 : \mu \neq \mu_0$. Recall the log-likelihood function

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

Assume that $\sigma^2 > 0$ is known. Also recall that

$$S_n(\mu) = \partial_\mu \ell_n(\mu, \sigma^2) = \frac{n}{\sigma^2}(\bar{X}_n - \mu),$$

$$I_n(\mu) = -E\left(S'_n(\mu)\right) = \frac{n}{\sigma^2},$$

and \bar{X}_n is the MLE of μ . It follows that

$$W_n = \frac{n}{\sigma^2}(\bar{X}_n - \mu_0)^2 = \left(\sqrt{n} \frac{\bar{X}_n - \mu_0}{\sigma}\right)^2,$$

$$R_n = \frac{\left(\frac{n}{\sigma^2}(\bar{X}_n - \mu_0)\right)^2}{\frac{n}{\sigma^2}} = \frac{n}{\sigma^2}(\bar{X}_n - \mu_0)^2,$$

$$T_n = \frac{1}{\sigma^2} \sum_{i=1}^n ((X_i - \mu_0)^2 - (X_i - \bar{X}_n)^2) = \frac{n}{\sigma^2}(\bar{X}_n - \mu_0)^2.$$

The three test statistics are identical in this model. \square

Example 8.9. Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(\pi)$. Consider $H_0 : \pi = \pi_0$ vs $H_1 : \pi \neq \pi_0$. Recall that

$$\ell_n(\pi) = n\hat{\pi} \log(\pi) + n(1 - \hat{\pi}) \log(1 - \pi),$$

$$S_n(\pi) = n \frac{\hat{\pi} - \pi}{\pi(1 - \pi)},$$

$$I_n(\pi) = \frac{n}{\pi(1 - \pi)},$$

where $\hat{\pi} = n^{-1} \sum_i X_i$ is the MLE of π . It follows that

$$W_n = n \frac{(\hat{\pi} - \pi_0)^2}{\hat{\pi}(1 - \hat{\pi})} = \left(\sqrt{n} \frac{\hat{\pi} - \pi_0}{\sqrt{\hat{\pi}(1 - \hat{\pi})}}\right)^2,$$

$$R_n = \frac{\left(n \frac{\hat{\pi} - \pi_0}{\pi_0(1 - \pi_0)}\right)^2}{\frac{n}{\pi_0(1 - \pi_0)}} = n \frac{(\hat{\pi} - \pi_0)^2}{\pi_0(1 - \pi_0)},$$

$$T_n = 2 \left\{ n\hat{\pi} \log \frac{\hat{\pi}}{\pi_0} + n(1 - \hat{\pi}) \log \frac{1 - \hat{\pi}}{1 - \pi_0} \right\}. \square$$

Let's check the Type I error and the power of these tests using some simulations.

```
# p-value calculation via simulation
```

```
p.value <- function(p0, tru.p, n) {
  x <- rbinom(n, size = 1, prob = tru.p)
  hat.p <- mean(x)
  Wald <- n * ((hat.p - p0)^2 / (hat.p * (1 - hat.p)))
  Score <- n * ((hat.p - p0)^2 / (p0 * (1 - p0)))
  LR <- 2 * (n * hat.p * ifelse(hat.p == 0, 0, log(hat.p / p0)) +
    n * (1 - hat.p) * ifelse(hat.p == 1, 0, log((1 - hat.p) / (1 - p0))))
```

```

p.valueW <- pchisq(Wald, df = 1, lower.tail = FALSE)
p.valueS <- pchisq(Score, df = 1, lower.tail = FALSE)
p.valueL <- pchisq(LR, df = 1, lower.tail = FALSE)
return(c(wald = p.valueW, Score = p.valueS, LR = p.valueL))
}

```

```

# Type I error ( $p_0 = \text{tru.p} = 0.5$ );  $n = 10, 50, 1000$ 
rbind((replicate(5000, p.value(0.5, 0.5, 10)) <= 0.05) |> rowMeans(),
(replicate(5000, p.value(0.5, 0.5, 50)) <= 0.05) |> rowMeans(),
(replicate(5000, p.value(0.5, 0.5, 1000)) <= 0.05) |> rowMeans())

```

```

      wald  Score    LR
[1,] 0.1088 0.0222 0.1088
[2,] 0.0628 0.0628 0.0628
[3,] 0.0564 0.0564 0.0564

```

```

# Power ( $p_0 = 0.4$  and  $\text{tru.p} = 0.5$ );  $n = 10, 50, 1000$ 
rbind((replicate(5000, p.value(0.4, 0.5, 10)) <= 0.05) |> rowMeans(),
(replicate(5000, p.value(0.4, 0.5, 50)) <= 0.05) |> rowMeans(),
(replicate(5000, p.value(0.4, 0.5, 1000)) <= 0.05) |> rowMeans())

```

```

      wald  Score    LR
[1,] 0.1890 0.0620 0.0718
[2,] 0.3422 0.3422 0.3422
[3,] 1.0000 1.0000 1.0000

```

The generalization of the three tests to the *multiparameter case* is not very difficult. Suppose we wish to test $H_0 : \theta = \theta_0$ vs $H_1 : \theta \neq \theta_0$, where $\theta \in \mathbb{R}^d$. Then

$$W_n \equiv W_n(\theta_0) = (\hat{\theta} - \theta_0)^t I_n(\hat{\theta}) (\hat{\theta} - \theta_0) \quad (\text{Wald})$$

$$R_n \equiv R_n(\theta_0) = S_n^t(\theta_0) I_n^{-1}(\theta_0) S_n(\theta_0) \quad (\text{Score})$$

$$T_n \equiv T_n(\theta_0) = 2 (\ell_n(\hat{\theta}) - \ell_n(\theta_0)) \quad (\text{Likelihood ratio})$$

Under H_0 , these three statistics converge to χ_d^2 , i.e. chi-squared distribution with d degrees of freedom. As consequence, the likelihood ratio test, for example, rejects H_0 when $T_n \geq \chi_{d;1-\alpha}^2$, or, equivalently, when $P(\chi_d^2 \geq t_n) \leq \alpha$. The same applies to the other two statistics.

Example 8.10. Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, where μ and σ are unknown. Consider testing

$$H_0 : \mu = \mu_0 \text{ and } \sigma^2 = \sigma_0^2 \text{ vs } H_1 : \mu \neq \mu_0 \text{ or } \sigma^2 \neq \sigma_0^2.$$

We know that

$$\begin{aligned}\ell_n(\mu, \sigma^2) &= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2, \\ \mathbf{S}_n(\boldsymbol{\theta}) &= \left(n \frac{\bar{X}_n - \mu}{\sigma^2}, -\frac{n}{2\sigma^2} + \frac{\sum_i (X_i - \mu)^2}{2\sigma^4} \right)^t, \\ \mathbf{I}_n(\boldsymbol{\theta}) &= \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix},\end{aligned}$$

and that $(\bar{X}_n, \hat{\sigma}_n^2)$, with $\hat{\sigma}_n^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$, is the MLE of (μ, σ^2) . Put $\tilde{\sigma}_0^2 = n^{-1} \sum_{i=1}^n (X_i - \mu_0)^2$.

The Wald statistic is given by

$$\begin{pmatrix} \bar{X}_n - \mu_0 & \hat{\sigma}_n^2 - \sigma_0^2 \end{pmatrix} \begin{pmatrix} \frac{n}{\hat{\sigma}_n^2} & 0 \\ 0 & \frac{n}{2\hat{\sigma}_n^4} \end{pmatrix} \begin{pmatrix} \bar{X}_n - \mu_0 \\ \hat{\sigma}_n^2 - \sigma_0^2 \end{pmatrix} = n \frac{(\bar{X}_n - \mu_0)^2}{\hat{\sigma}_n^2} + \frac{n}{2} \left(1 - \frac{\sigma_0^2}{\hat{\sigma}_n^2} \right)^2 \xrightarrow{d} \chi_2^2, \text{ if } H_0 \text{ is true.}$$

The Score statistic is given by

$$\begin{aligned}& \begin{pmatrix} n \frac{\bar{X}_n - \mu_0}{\sigma_0^2} & -\frac{n}{2\sigma_0^2} + \frac{\sum_i (X_i - \mu_0)^2}{2\sigma_0^4} \end{pmatrix} \begin{pmatrix} \frac{\sigma_0^2}{n} & 0 \\ 0 & \frac{2\sigma_0^4}{n} \end{pmatrix} \begin{pmatrix} n \frac{\bar{X}_n - \mu_0}{\sigma_0^2} \\ -\frac{n}{2\sigma_0^2} + \frac{\sum_i (X_i - \mu_0)^2}{2\sigma_0^4} \end{pmatrix} \\ &= n \frac{(\bar{X}_n - \mu_0)^2}{\sigma_0^2} + \frac{n}{2} \left(1 - \frac{\tilde{\sigma}_0^2}{\sigma_0^2} \right)^2 \xrightarrow{d} \chi_2^2, \text{ if } H_0 \text{ is true.}\end{aligned}$$

The LRT statistic is given by

$$\begin{aligned}T_n &= 2(\ell_n(\bar{X}_n, \hat{\sigma}_n^2) - \ell_n(\mu_0, \sigma_0^2)) \\ &= n \frac{(\bar{X}_n - \mu_0)^2}{\sigma_0^2} - n \log \left(\frac{\hat{\sigma}_n^2}{\sigma_0^2} \right) + n \left(\frac{\hat{\sigma}_n^2}{\sigma_0^2} - 1 \right) \xrightarrow{d} \chi_2^2, \text{ if } H_0 \text{ is true.} \square\end{aligned}$$

8.6.2 Composite null hypothesis

In practice, we are rarely interested in a simple null hypothesis in which all the d components of the vector parameter $\boldsymbol{\theta} \in \mathbb{R}^d$ are assigned, as done previously. Often, we're only interested in a subset of $\boldsymbol{\theta}$, or in some specific constraints on its components. This can be formulated, in most cases, as follows

$$H_0 : \mathbf{A}\boldsymbol{\theta} = \mathbf{a} \text{ vs } H_1 : \mathbf{A}\boldsymbol{\theta} \neq \mathbf{a},$$

where \mathbf{A} is an $r \times d$ full rank matrix ($r \leq d$ and the r rows of \mathbf{A} are linearly independent) and \mathbf{a} is a given vector in \mathbb{R}^r .

For example, suppose that $d = 5$, then

- $\theta_1 = \dots = \theta_5 = 0 \Leftrightarrow \mathbf{A}\boldsymbol{\theta} = \mathbf{0}$, with $\mathbf{A} = \mathbf{1}$ (5×5 identity matrix), and $\mathbf{0} = (0, \dots, 0)^t$.
- $\theta_1 = 1, \theta_2 = 2, \dots, \theta_5 = 5 \Leftrightarrow \mathbf{A}\boldsymbol{\theta} = \mathbf{a}$, with $\mathbf{A} = \mathbf{1}$ and $\mathbf{a} = (1, 2, \dots, 5)^t$.

- $\theta_2 = \theta_4 = 0 \Leftrightarrow A\theta = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$, with $A = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$
- $\theta_2 + \theta_4 = 1 \Leftrightarrow A\theta = 1$, with $A = \begin{pmatrix} 0 & 1 & 0 & 1 & 0 \end{pmatrix}$
- $\theta_1 = 2\theta_2$ and $\theta_3 = \theta_4 \Leftrightarrow A\theta = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$, with $A = \begin{pmatrix} 1 & -2 & 0 & 0 & 0 \\ 0 & 0 & 1 & -1 & 0 \end{pmatrix}$

The generalized Wald test statistic, the generalized Score test statistic and the generalized LRT statistic are given, respectively, by

$$\begin{aligned} W_n &= (A\hat{\theta} - a)^t (A I_n^{-1}(\hat{\theta}) A^t)^{-1} (A\hat{\theta} - a), \\ R_n &= S_n^t(\tilde{\theta}_a) I_n^{-1}(\tilde{\theta}_a) S_n(\tilde{\theta}_a), \\ T_n &= 2(\ell_n(\hat{\theta}) - \ell_n(\tilde{\theta}_a)), \end{aligned}$$

where $\hat{\theta}$ is the MLE of θ , and $\tilde{\theta}_a$ is the *restricted* MLE of θ obtained by maximizing the (log-)likelihood under H_0 , i.e. under the constraint that $A\theta = a$.

Similar to the simple null hypothesis case, *these three statistics converge, Under H_0 , to χ_r^2* , i.e. chi-squared distribution with r degrees of freedom. This latter corresponds to the number of restrictions imposed by H_0 .

In the remainder of this section, in order to simplify calculations and make it easier to follow-up, we'll concentrate solely on the Wald and the LR tests.

Example 8.11. Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, where μ and σ are unknown. Consider testing

$$H_0 : \mu = \mu_0 \text{ vs } H_1 : \mu \neq \mu_0.$$

The Wald statistic is given by

$$\begin{aligned} W_n &= \left(\begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} \bar{X}_n \\ \hat{\sigma}_n^2 \end{pmatrix} - \mu_0 \right)^t \left(\begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} \frac{\hat{\sigma}_n^2}{n} & 0 \\ 0 & \frac{2\hat{\sigma}_n^4}{n} \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right)^{-1} \left(\begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} \bar{X}_n \\ \hat{\sigma}_n^2 \end{pmatrix} - \mu_0 \right) \\ &= n \frac{(\bar{X}_n - \mu_0)^2}{\hat{\sigma}_n^2} \xrightarrow{d} \chi_1^2, \text{ if } H_0 \text{ is true.} \end{aligned}$$

Since the MLE of (μ, σ^2) under H_0 is $(\mu_0, \tilde{\sigma}_0^2)$, where $\tilde{\sigma}_0^2 = n^{-1} \sum_i (X_i - \mu_0)^2$, the LRT statistic is given by

$$\begin{aligned} T_n &= 2(\ell_n(\bar{X}_n, \hat{\sigma}_n^2) - \ell_n(\mu_0, \tilde{\sigma}_0^2)) \\ &= n \log \left(\frac{\tilde{\sigma}_0^2}{\hat{\sigma}_n^2} \right) = n \log \left(1 + \frac{(\bar{X}_n - \mu_0)^2}{\hat{\sigma}_n^2} \right) \xrightarrow{d} \chi_1^2, \text{ if } H_0 \text{ is true.} \square \end{aligned}$$

Example 8.12. Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Ber}(\pi_1)$ and $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Ber}(\pi_2)$ be two independent

samples. Consider testing

$$H_0 : \pi_1 = \pi_2 \text{ vs } H_1 : \pi_1 \neq \pi_2.$$

Based on the sample (X_i, Y_i) , $i = 1, \dots, n$, it's easy to see that the Log-likelihood for (π_1, π_2) is

$$\ell_n(\pi_1, \pi_2) = n\hat{\pi}_1 \log(\pi_1) + n(1 - \hat{\pi}_1) \log(1 - \pi_1) + n\hat{\pi}_2 \log(\pi_2) + n(1 - \hat{\pi}_2) \log(1 - \pi_2),$$

where $(\hat{\pi}_1, \hat{\pi}_2) := (n^{-1} \sum_i X_i, n^{-1} \sum_i Y_i)$ is the MLE of (π_1, π_2) . The FI matrix is

$$I_n(\pi_1, \pi_2) = \begin{pmatrix} \frac{n}{\pi_1(1-\pi_1)} & 0 \\ 0 & \frac{n}{\pi_2(1-\pi_2)} \end{pmatrix}.$$

The Wald statistic is given by

$$\begin{aligned} W_n &= \left(\begin{pmatrix} 1 & -1 \end{pmatrix} \begin{pmatrix} \hat{\pi}_1 \\ \hat{\pi}_2 \end{pmatrix} - 0 \right)^t \left(\begin{pmatrix} 1 & -1 \end{pmatrix} \begin{pmatrix} \frac{\hat{\pi}_1(1-\hat{\pi}_1)}{n} & 0 \\ 0 & \frac{\hat{\pi}_2(1-\hat{\pi}_2)}{n} \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right)^{-1} \left(\begin{pmatrix} 1 & -1 \end{pmatrix} \begin{pmatrix} \hat{\pi}_1 \\ \hat{\pi}_2 \end{pmatrix} - 0 \right) \\ &= n \frac{(\hat{\pi}_1 - \hat{\pi}_2)^2}{\hat{\pi}_1(1 - \hat{\pi}_1) + \hat{\pi}_2(1 - \hat{\pi}_2)} \xrightarrow{d} \chi_1^2, \text{ if } H_0 \text{ is true.} \end{aligned}$$

Under H_0 , the log-likelihood reduces to

$$\ell_n(\pi_1, \pi_1) = n(\hat{\pi}_1 + \hat{\pi}_2) \log(\pi_1) + n(2 - \hat{\pi}_1 - \hat{\pi}_2) \log(1 - \pi_1).$$

It's easy to see that, in this case, the MLE of $\pi_1 (= \pi_2)$ is $\tilde{\pi}_0 = \frac{\sum_{i=1}^n (X_i + Y_i)}{2n}$. As a result, the LRT statistic is given by

$$\begin{aligned} T_n &= 2(\ell_n(\hat{\pi}_1, \hat{\pi}_2) - \ell_n(\tilde{\pi}_0, \tilde{\pi}_0)) \\ &= 2n \sum_{j=1}^2 \left(\hat{\pi}_j \log \left(\frac{\hat{\pi}_j}{\tilde{\pi}_0} \right) + (1 - \hat{\pi}_j) \log \left(\frac{1 - \hat{\pi}_j}{1 - \tilde{\pi}_0} \right) \right) \xrightarrow{d} \chi_1^2, \text{ if } H_0 \text{ is true.} \square \end{aligned}$$

We now turn to a special, but very useful, case of the general formulation of the LR statistic presented above. Consider the case of a parametric model indexed by (θ, η) , where θ and η are unknown. Let's say that θ is our parameter of interest and η is a nuisance parameter (θ and η can be vectors). In this context, our objective is to test $H_0 : \theta = \theta_0$ vs $H_1 : \theta \neq \theta_0$, for some given θ_0 . Let L denote the likelihood function of the model under study, $(\hat{\theta}, \hat{\eta}) = \arg \max_{\theta, \eta} L(\theta, \eta)$ be the MLE of (θ, η) and $\hat{\eta}_0 \equiv \hat{\eta}(\theta_0) := \arg \max_{\eta} L(\theta_0, \eta)$ be the restricted MLE of η under H_0 . According to the general definition given above, the LRT statistic is

$$T_n \equiv T_n(\theta_0) := -2 \log \frac{L(\theta_0, \hat{\eta}_0)}{L(\hat{\theta}, \hat{\eta})} = -2 \log \frac{L_p(\theta_0)}{L_p(\hat{\theta})},$$

where $L_p(\theta) = \max_{\eta} L(\theta, \eta)$ is the profile likelihood for θ , i.e. $L_p(\theta) = L(\theta, \hat{\eta}(\theta))$, with

$\hat{\eta}(\theta) = \arg \max_{\eta} L(\theta, \eta)$. In this case, $T_n \xrightarrow{d} \chi_r^2$, if H_0 is true, where r is the length of θ . T_n is sometimes referred to as the profile likelihood ratio statistic.

Example 8.13. Let's consider the case of the linear model $Y = \beta_0 + \beta_1 x + \epsilon$, where, for example, $\beta_0 = 10$, $\beta_1 = 20$, $x = 1, \dots, 10$ and $\epsilon \sim N(0, 10^2)$. Let's say we observe $n = 10$ observations from (Y, x) . In this context, we'd like to test $H_0 : \beta_1 = 0$ vs $H_1 : \beta_1 \neq 0$. Under the linear model assumption, the LRT statistic is

$$2(\ell_n(\hat{\beta}_0, \hat{\beta}_1, \hat{\sigma}^2) - \ell_n(\hat{\beta}_{00}, 0, \hat{\sigma}_0^2)),$$

where $\hat{\beta}_0, \hat{\beta}_1, \hat{\sigma}^2$ are the unrestricted MLEs (i.e. those we saw in Section 7.3) and $\hat{\beta}_{00}$ and $\hat{\sigma}_0^2$ are the restricted MLEs (i.e. those that maximize likelihood assuming $\beta_1 = 0$).

```
library(stats4)

set.seed(5)
x <- 1:10
y <- 10 + 20 * x + rnorm(10, sd = 10)

neglogLiReg <- function(a, b, s, y) {
  -sum(dnorm(y, mean = a + b * x, sd = s, log = TRUE))
}

loglik <- mle(\(a, b, s) neglogLiReg(a, b, s, y = y),
             start = list(a = mean(y), b = 0, s = sd(y))) |> logLik()
loglik0 <- mle(\(a, s) neglogLiReg(a, b = 0, s, y = y),
              start = list(a = mean(y), s = sd(y))) |> logLik()

LogLR <- (2 * (loglik - loglik0))[1] |> print()
pchisq(LogLR, df = 1, lower.tail = FALSE)

[1] 37.07
[1] 1.14e-09
```

A simple way to perform such a test in R is to use the `drop1()` function.

```
glm(y ~ x) |> drop1(test = "LRT")
```

Single term deletions

Model:

y ~ x

	Df	Deviance	AIC	scaled dev.	Pr(>Chi)
<none>		812	78.3		
x	1	33072	113.4	37.1	1.1e-09 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1