# M1 Data Science, University of Lille
# Statistics 2 - Lecture notes

Emilie Kaufmann (CNRS, Univ. Lille)
emilie.kaufmann@univ-lille.fr

February 27, 2025

2

**Pre-requisite**   In statistics 1, you have seen:

- classical distributions

- examples of estimators

- confidence intervals

- the statistical testing protocol (type I and type II error)

- example of classical tests

In statistics 2, we will revisit statistical and testing with a focus on *optimality*.
We will notably discuss:

- different performance measure for estimators

- generic estimation strategies, notably the maximum likelihood principle

- asymptotic properties of estimators

- likelihood-ratio based testing procedures

Several examples will come from an important family of distributions called exponential families.

# Chapter 1

# Estimation

## 1.1  Statistical inference

In statistical inference, we observe a realization of some random variable (or random vector) $X$, called the observation, whose distribution over some space $\mathcal{X}$ is $P_X$. The goal is to discover ("infer") some properties of this underlying distribution, assuming that $P_X$ belongs to some set of possible distributions, called the *statistical model*.

Depending on the situation, we may make assumptions on the cumulative distribution function (cdf) of $X$, $F_X$ or on its density $f_X$ with respect to some reference measure and the statistical model may be a set of distribution, a set of cdfs or a set of pdfs parameterized by some parameter $\theta$:

$$\mathcal{M} = \{P_\theta, \theta \in \Theta\}, \quad \mathcal{M} = \{F_\theta, \theta \in \Theta\} \quad \text{or} \quad \mathcal{M} = \{f_\theta, \theta \in \Theta\}.$$

When the parameter space $\Theta \subseteq \mathbb{R}^d$, the model is called parametric, otherwise it is non-parametric. Given the "true" parameter $\theta$ (i.e. $\theta$ such that $P_X = P_\theta$), the probability space on which $X$ is defined is denoted by $(\Omega, \mathcal{F}, \mathbb{P}_\theta)$, and the corresponding expectation is denoted by $\mathbb{E}_\theta$.

**The $n$-sample example**   The observation $X$ will often be a random vector of the form $X = (X_1, \ldots, X_n)$ where the $X_i$ are assumed to be iid realizations of the same distribution. These iid copies represent the repetition of some random experiment (for example the vote expressed by one individual in a population, or the effect of a treatment on one patient). These random variables $X_i$ are defined on some common probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and will most of the time take values in $\mathbb{R}$ (we could consider some multi-dimensional outcomes in, e.g. two-sample testing problems).

In the $n$-sample setting, we denote by $P$ the distribution of $X_1$ (which is the common distribution of all $X_i$'s), by $F$ the cdf of this distribution and by $f$ its density (with respect to some reference measure $\nu$), if it admits one. We will write indifferently

$$X_1, \ldots, X_n \overset{\text{i.i.d.}}{\sim} P, \quad X_1, \ldots, X_n \overset{\text{i.i.d.}}{\sim} F \quad \text{or} \quad X_1, \ldots, X_n \overset{\text{i.i.d.}}{\sim} f.$$

In that case, the statistical model is typically expressed as possible candidates for $P$, $F$ or $f$ directly, instad of possible candidates for $P^{\otimes n}$, $F^{\otimes n}$ and $f^{\otimes n}$. By a slight abuse of notation, we will also denote by $P_\theta$, $F_\theta$ and $f_\theta$ the possible candidate for the distribution of $X_1$, for $\theta$ in some parameter space $\Theta$.

**Example 1.1.** *We consider a Gaussian $n$-sample with unit variance and unknown mean. That is, we have $X_1, \ldots, X_n \overset{iid}{\sim} \mathcal{N}(\theta, 1)$ for $\theta \in \Theta = \mathbb{R}$. Let $f_\theta$ be the density of a $\mathcal{N}(\theta, 1)$ variable with respect to*

*the Lebesgue measure (in $\mathbb{R}$):*

$$f_\theta(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(x-\theta)^2}{2}\right) .$$

*If we look at the observation $X = (X_1, \ldots, X_n)$, the statistical model $\mathcal{M}$ for $X$ is a set of multivariate Gaussian distributions whose densities with respect to the Lebesgue measure in $\mathbb{R}^n$ is*

$$f_\theta(x_1, \ldots, x_n) = \prod_{i=1}^{n} f_\theta(x_i)$$

*for some parameter $\theta \in \mathbb{R}$.*

In statistical inference, we are interested in making statements about the "true" parameter $\theta$ generating the data or about some *parameter of interest* which can be some function of $\theta$, denoted by $g(\theta)$. This statement can be a guess for its value (estimation), an interval to which it belongs (confidence interval) or the answer to some question about this parameter (statistical test).

**Example 1.2.** *$X_1, \ldots, X_n \overset{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$. The parameter of the model is $\theta = (\mu, \sigma)$ and the parameter space is $\Theta = \{(\mu, \sigma), \mu \in \mathbb{R}, \sigma > 0\}$. If we are solely interested in estimating the mean, the parameter of interest is $\mu$ and $\sigma$ may be called a nuisance parameter.*

*In some situations, we may be interested in estimating more complex functions of $\theta$. For example, assume that $X_i$ models the amount of antibodies produced 15 days after receiving a vaccine. For a given disease, the vaccine is considered efficient if this amount exceeds some threshold $v$. A possible parameter of interest is the probability of efficacy of the vaccine, $p = p(\mu, \sigma)$, which can be expressed as*

$$p = \mathbb{P}(X_1 \geq v) = 1 - \mathbb{P}(X_1 < v) = 1 - \mathbb{P}\left(\frac{X_1 - \mu}{\sigma} < \frac{v - \mu}{\sigma}\right) = 1 - \Phi\left(\frac{v - \mu}{\sigma}\right)$$

*where $\Phi$ is the cdf of a $\mathcal{N}(0, 1)$ random variable.*

**Example 1.3** (regression model). *$Z_1, \ldots, Z_n \overset{iid}{\sim} P$. $X_i = (Z_i, Y_i) \in \mathcal{X} \times \mathbb{R}$ such that*

$$Y_i = h(Z_i) + \varepsilon_i$$

*where $\varepsilon_i \overset{iid}{\sim} \mathcal{N}(0, 1)$ and $h : \mathcal{X} \to \mathcal{Y}$ is the regression function. The observation is $X = (X_1, \ldots, X_n)$ and the parameters of the model are $P$ (that could belong to some parametric class of probability distributions) and the regression function $h$ (that could belong to a parametric families of functions, e.g. linear functions). In that case the "parameter" of interest is usually the regression function.*

## 1.2   Performance of an estimator

An estimator of $g(\theta)$ is any function of the observation $\widehat{g} = h(X)$ that is supposed to be "close" to the parameter of interest $g(\theta)$. When $X = (X_1, \ldots, X_n)$ has the $n$-sample structure, we will materialize the dependency in $n$ of the estimator by writing $\widehat{g}_n = h(X_1, \ldots, X_n)$.

From its definition, $\widehat{g}$ is a random variable (or a random vector, when we estimate a multi-dimensional parameter), hence its quality will be expressed in terms of some properties of its distribution, which should ideally be concentrated around $g(\theta)$. Two important characteristics of this distributions are its mean and its variance, both expressed with expectations.

### 1.2.1 Recap: Densities and Expectations

In general, if $Z$ is a random variable taking values in $\mathcal{Z}$ whose distribution $P$ has a density $f$ with respect to some reference measure $\nu$, we have, for all function $\phi$,

$$\mathbb{E}_\theta[\phi(Z)] = \int_{\mathcal{Z}} \phi(x)f(x)d\nu(x).$$

We will mostly see examples of random variables defined on $\mathcal{Z} = \mathbb{R}^d$ whose distributions have a density with respect to the Lebesgue measure in $\mathbb{R}^d$, or of discrete random random variables (i.e. for which $\mathcal{Z}$ is discrete) that have a density with respect to the counting measure. In the discrete case, the density is simply defined, for all $z \in \mathcal{Z}$, by

$$f(z) = P(\{z\}) = \mathbb{P}_{Z \sim P}(Z = z) \ .$$

Back to our statistical model, in the most common $n$-sample case in which $X = (X_1, \ldots, X_n) \overset{iid}{\sim} P_\theta$, we will often uncounter two cases. Either $X_i \in \mathbb{R}$ and $P_\theta$ has a density with respect to the Lebesgue measure. Then

- for any $\phi : \mathbb{R}^n \to \mathbb{R}$, $\mathbb{E}_\theta[\phi(X)] = \int_{\mathbb{R}} \phi(x_1, \ldots, x_n)f_\theta(x_1, \ldots, x_n)dx_1 \ldots dx_n$
- for any $\phi : \mathbb{R} \to \mathbb{R}$, $\mathbb{E}_\theta[\phi(X_1)] = \int_{\mathbb{R}} \phi(u)f_\theta(u)du$

Or $X_i \in \mathcal{S}$ for some discrete set $\mathcal{S}$ (typically a subset of $\mathbb{N}$) and we have

- for any $\phi : \mathcal{S}^n \to \mathbb{R}$, $\mathbb{E}_\theta[\phi(X)] = \sum_{x \in \mathcal{S}^n} \phi(x_1, \ldots, x_n)f_\theta(x_1, \ldots, x_n)$
- for any $\phi : \mathcal{S} \to \mathbb{R}$, $\mathbb{E}_\theta[\phi(X_1)] = \sum_{u \in \mathcal{S}} \phi(u)f_\theta(u)$

### 1.2.2 Bias, Variance and Quadratic Risk

**Definition 1.4.** *The* bias *of estimator $\widehat{g}$ of $g(\theta)$ is defined as* $\mathrm{b}_\theta(\widehat{g}) = \mathbb{E}_\theta[\widehat{g}] - g(\theta)$.
*When $\mathrm{b}_\theta(\widehat{g}) = 0$, the estimator is called* unbiased.

**Definition 1.5.** *The* variance *of a real-valued estimator $\widehat{g}$ is* $\mathrm{Var}_\theta[\widehat{g}] := \mathbb{E}_\theta[(\widehat{g} - \mathbb{E}_\theta[\widehat{g}])^2]$.

A good (real-valued) estimator has ideally a small bias and a small variance, which indicates that on average, its value is close to $g(\theta)$ and that under different realizations of the experiments, its value would not change too much. The closeness from $\widehat{g}$ to $g(\theta)$ can also directly be measured using their average distance, a notion that can also be meaningful in the multi-dimensional setting.

**Definition 1.6.** *The* quadratic risk *of an estimator $\widehat{g}$ of $g(\theta) \in \mathbb{R}^p$ is*

$$\mathrm{R}_\theta(\widehat{g}) = \mathbb{E}_\theta\left[\|\widehat{g} - g(\theta)\|^2\right],$$

*where $\|u\|$ is the Euclidian norm in $\mathbb{R}^p$, such that $\|u\|^2 = u^\top u$. In the one-dimensional case ($p = 1$), this quantity is sometimes called the* mean-squared error, *and denoted by* $\mathrm{MSE}_\theta(\widehat{g})$.

**Theorem 1.7** (bias-variance decomposition)**.** *Assume $g(\theta) \in \mathbb{R}$. We have*

$$\mathrm{R}_\theta(\widehat{g}) = (\mathrm{b}(\widehat{g}))^2 + \mathrm{Var}_\theta[\widehat{g}] \ .$$

*Proof.* When $g(\theta) \in \mathbb{R}$, we write

$$
\begin{aligned}
\mathrm{R}_\theta(\widehat{g}) &= \mathbb{E}_\theta\left[(\widehat{g} - g(\theta))^2\right] = \mathbb{E}_\theta\left[((\widehat{g} - \mathbb{E}_\theta[\widehat{g}]) + (\mathbb{E}_\theta[\widehat{g}] - g(\theta)))^2\right] \\
&= \mathbb{E}_\theta\left[(\widehat{g} - \mathbb{E}_\theta[\widehat{g}])^2\right] + \mathbb{E}_\theta\left[(\mathbb{E}_\theta[\widehat{g}] - g(\theta))^2\right] + 2\mathbb{E}_\theta\left[(\widehat{g} - \mathbb{E}_\theta[\widehat{g}])(\mathbb{E}_\theta[\widehat{g}] - g(\theta))\right] \\
&= \mathbb{E}_\theta\left[(\widehat{g} - \mathbb{E}_\theta[\widehat{g}])^2\right] + (\mathbb{E}_\theta[\widehat{g}] - g(\theta))^2 + 2\left(\mathbb{E}_\theta[\widehat{g}] - g(\theta)\right)\underbrace{\mathbb{E}_\theta\left[\widehat{g} - \mathbb{E}_\theta[\widehat{g}]\right]}_{=0} \\
&= \mathrm{Var}_\theta\left[\widehat{g}\right] + (\mathrm{b}(\widehat{g}))^2 .
\end{aligned}
$$

$\square$

The quadratic risk can be used to compare estimators, and we say that an estimator $\widehat{g}$ is better than an estimator $\widetilde{g}$ if for all $\theta \in \Theta$, $\mathrm{R}_\theta(\widehat{g}) \leq \mathrm{R}_\theta(\widetilde{g})$. However, this relationship is not a total order, as there may exists estimators for which $\mathrm{R}_{\theta_1}(\widehat{g}) \leq \mathrm{R}_{\theta_1}(\widetilde{g})$ for some parameter $\theta_1 \in \Theta$ but $\mathrm{R}_{\theta_2}(\widehat{g}) > \mathrm{R}_{\theta_2}(\widetilde{g})$ for a different parameter $\theta_2 \in \Theta$.

**Definition 1.8.** *An estimator $\widehat{g}$ of $g(\theta)$ is called* admissible *is there exists no estimator $\widetilde{g}$ which is strictly better than $\widehat{g}_n$ i.e. for which*

$$
\forall \theta \in \Theta, \quad \mathrm{R}_\theta(\widetilde{g}) \leq \mathrm{R}_\theta(\widehat{g})
$$

*and the inequality is strict for at least one value $\theta_0$.*

**Influence of the sample size**   When $X$ is a $n$-sample, the performance measures for an estimator $\widehat{g}_n$ are all defined for a fixed sample size $n$, and are not capturing another desirable property of an estimator: $\widehat{g}_n$ should get closer to $g(\theta)$ when the sample size $n$ goes larger. We expect to $\widehat{g}_n$ to get closer to $g(\theta)$, meaning that its distribution concentrates for and more around $g(\theta)$. We will discuss these asymptotic properties in the next chapter.

## 1.3   Estimation procedures

### 1.3.1   The moment method

When $X = (X_1, \ldots, X_n) \overset{iid}{\sim} P_\theta$, the moment method can be used when the parameter of interest $g(\theta)$ can be expressed as a function of the moments of $X_1$.

In the simplest case, the parameter of interest can directly be written as an expectation:

$$
g(\theta) = \mathbb{E}_\theta\left[\phi(X_1)\right]
$$

for some function $\phi$ such that $\mathbb{E}[|\phi(X_1)|] < \infty$. Motivated by the law of large numbers, we define the moment estimator

$$
\widehat{g}_n := \frac{1}{n}\sum_{i=1}^{n}\phi(X_i)
$$

which satisfies $\widehat{g}_n \to g(\theta)$, $\mathbb{P}_\theta$ – a.s.. Hence, this estimator is naturally going to be close to $g(\theta)$ at least for a large sample size $n$.

**Example 1.9** (the empirical cdf). *Given iid samples $X_1 \ldots, X_n$ from some distribution $P$ in $\mathbb{R}$ whose cdf is $F$, we want to estimate the function $F$, that is, for each value $x \in \mathbb{R}$ we want to estimate the quantity $F(x) = \mathbb{P}(X \le x) = \mathbb{E}[\mathbb{1}(X \le x)]$. As this quantity can be written as an expectation, its moment estimator is simply*

$$\widehat{F}_n(x) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}(X_i \le x)$$

*The function $x \mapsto \widehat{F}_n(x)$ is called empirical cdf of $X = (X_1, \ldots, X_n)$.*

*We can further compute the biais and the bias of this estimator*

$$b(\widehat{F}_n(x)) = \mathbb{E}\left[\widehat{F}_n(x)\right] - F(x) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}\left[\mathbb{1}(X_i \le x)\right] - F(x) = F(x) - F(x) = 0$$

*and its variance:*

$$
\begin{aligned}
\mathrm{Var}[\widehat{F}_n(x)] &= \frac{1}{n^2} \mathrm{Var}\left[\sum_{i=1}^{n} \mathbb{1}(X_i \le x)\right] = \frac{1}{n^2} \sum_{i=1}^{n} \mathrm{Var}\left[\mathbb{1}(X_i \le x)\right] \\
&= \frac{1}{n^2} \sum_{i=1}^{n} F(x)(1 - F(x)) = \frac{F(x)(1 - F(x))}{n} \ .
\end{aligned}
$$

*For the variance computation we have used that the variance of the sum of independent random variables is the sum of their variances, and that $\mathbb{1}(X_i \le x)$ which takes value in $\{0, 1\}$ is a Bernoulli distribution with mean $p = \mathbb{P}(X_i \le x) = F(x)$, whose variance if $p(1 - p)$.*

More generally, suppose that we seek to estimate a multi-dimensional parameter $\theta = (\theta_1, \ldots, \theta_k)^\top \in \mathbb{R}^k$ and that for $1 \le j \le k$ the $j$-th moment can be expressed as some function of the parameter $\theta$:

$$\mathbb{E}_\theta[X^j] = \alpha_j(\theta).$$

Letting $\widehat{\alpha}_j = \frac{1}{n} \sum_{i=1}^{n} X_i^j$ the $j$-th sample moment, the moment estimator is defined as the solution $\widehat{\theta}_n$ of the system of equations

$$\alpha_1(\theta) = \widehat{\alpha}_1, \quad \ldots \quad, \alpha_k(\theta) = \widehat{\alpha}_k.$$

**Example 1.10.** $X_1, \ldots, X_n \overset{iid}{\sim} \mathcal{N}(\mu, \sigma^2)$. *We can find the moment estimator for the parameter $\theta = (\mu, \sigma^2)$. There are two parameters so we can look at the first two moments.*

$$
\begin{aligned}
\mathbb{E}_\theta[X_1] &= \mu \\
\mathbb{E}_\theta[X_1^2] &= \mathrm{Var}_\theta[X_1] + (\mathbb{E}_\theta[X_1])^2 = \sigma^2 + \mu^2
\end{aligned}
$$

*The empirical first and second moments are $\frac{1}{n} \sum_{i=1}^{n} X_i$ and $\frac{1}{n} \sum_{i=1}^{n} X_i^2$ so we get the system of equations*

$$
\begin{cases}
\mu &= \frac{1}{n} \sum_{i=1}^{n} X_i \\
\mu^2 + \sigma^2 &= \frac{1}{n} \sum_{i=1}^{n} X_i^2
\end{cases}
$$

*from which we get $\widehat{\theta}_n = (\widehat{\mu}_n, \widehat{\sigma}_n^2)$ where*

$$\widehat{\mu}_n = \frac{1}{n} \sum_{i=1}^{n} X_i \quad and \quad \widehat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^{n} X_i^2 - \left(\frac{1}{n} \sum_{i=1}^{n} X_i\right)^2$$

*We recognize the well-known empirical mean and (unadjusted) empirical variance, which can also be rewritten*

$$\widehat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i - \widehat{\mu}_n)^2 \ .$$

### 1.3.2   The plug-in method

The plug-in method is also suited for the $n$-sample setting, when the parameter of interest can be expressed as some functional of the distribution of $X_1$ (for example some moment of this distribution, or some quantile), that is we have

$$g(\theta) = H(P)$$

where $P$ is the cdf of $X_1$.

A plug-in estimator replaces ("plugs in") the unknown distribution $P$ by an empirical variant of this distribution

$$\widehat{g}_n := H(\widehat{P}_n)$$

where the empirical distribution $\widehat{P}_n$ is defined below.

**Definition 1.11.** *The empirical distribution of $X = (X_1, \ldots, X_n)$ is defined as*

$$\widehat{P}_n = \frac{1}{n} \sum_{i=1}^{n} \delta_{X_i},$$

*where the Dirac measure in $x$ is defined as, $\delta_x(A) = 1$ if $x \in A$, $\delta_x(A) = 0$ otherwise, for all event $A$.*

*$\widehat{P}_n$ is a discrete distribution supported on $\mathcal{S} = \{X_1, \ldots, X_n\}$, the set of distinct values in our sample, and for all $x$,*

$$\widehat{P}_n(\{x\}) = \frac{1}{n} \sum_{i=1}^{n} \delta_{X_i}(x) = \frac{\#\{i : X_i = x\}}{n}.$$

*The cdf of the distribution $\widehat{P}_n$ can be shown to be the empirical cdf $\widehat{F}_n$, presented in Example 1.9.*

For any function $\phi$, the expectation of $\phi(Z)$ when $Z$ is distributed according to the empirical distribution $\widehat{P}_n$ is given by

$$\mathbb{E}_{Z \sim \widehat{P}_n}[\phi(Z)] = \sum_{x \in S} \phi(x)\widehat{P}_n(x) = \frac{1}{n} \sum_{i=1}^{n} \phi(X_i) .$$

**Remark 1.12.** *When the functional $H(P)$ is defined as some expectation under $P$, the moment method and the plug-in method actually concide. Indeed, if*

$$g(\theta) = \mathbb{E}_{X \sim P}[\phi(X)]$$

*the plug-in method yields*

$$\widehat{g}_n = \mathbb{E}_{X \sim \widehat{P}_n}[\phi(X)] = \frac{1}{n} \sum_{i=1}^{n} \phi(X_i).$$

*Such estimators are also called "empirical estimators".*

*But plug-in estimator can be more general when $H$ is not defined as some expectation. For example when we want to estimate some quantile of a distribution $P$, e.g. its median, the plug-in estimator, also called empirical quantile is the corresponding quantile of the empirical distribution $\widehat{P}_n$.*

**Example 1.13** (variance estimation)**.** *Given that the variance of a distribution can be written*

$$\mathrm{Var}[X] = \mathbb{E}_{X \sim P}[(X - \mathbb{E}_{X \sim P}[X])^2]$$

*the plug-in method yields the estimator*

$$
\begin{aligned}
\widehat{\sigma}_n^2 &= \mathbb{E}_{X \sim \widehat{P}_n}[(X - \mathbb{E}_{X \sim \widehat{P}_n}[X])^2] \\
&= \mathbb{E}_{X \sim \widehat{P}_n}[(X - \widehat{\mu}_n)^2] \\
&= \frac{1}{n} \sum_{i=1}^n (X_i - \widehat{\mu}_n)^2
\end{aligned}
$$

*where we introduce the empirical mean $\widehat{\mu}_n = \mathbb{E}_{X \sim \widehat{P}_n}[X] = \frac{1}{n} \sum_{i=1}^n X_i$. We recover the same estimator as the one derived before the Gaussian distributions.*

*As for its properties, it is well known that the empirical variance is biased. Indeed one can check that $\mathbb{E}[\widehat{\sigma}_n^2] = \frac{n-1}{n} \sigma^2$ if $\sigma^2 = \mathrm{Var}[X_1]$. Hence an unbiased estimator of the variance is the adjusted empirical variance*

$$
\widetilde{\sigma}_n^2 = \frac{n}{n-1} \widehat{\sigma}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \widehat{\mu}_n)^2
$$

*which is often preferred in practise. Still, comparing $\widehat{\sigma}_n^2$ and $\widetilde{\sigma}_n^2$ in terms of mean-squared error for certain distributions (e.g. Gaussian) shows that the unadjusted estimator can have a smaller MSE.*

### 1.3.3 Maximum Likelihood Estimation (MLE)

The maximum likelihood approach can be used to estimate $g(\theta) = \theta$ when the statistical model is of the form

$$
\mathcal{M} = \{P_\theta : P_\theta \text{ has a density } f_\theta \text{ with respect to } \nu, \theta \in \Theta\}
$$

where $\nu$ is a fixed reference measure (which is the same for all the distributions in the model). Such a model is called *dominated* (by the reference measure $\nu$).

In most practical cases, this reference measure will be the Lebesgue measure in $\mathbb{R}^d$ (when the distributions are continuous) or the counting measure on discrete set (when the distributions are discrete). In that case, the density is given by $f_\theta(x) = \mathbb{P}_\theta(X = x)$.

**Definition 1.14.** *The likelihood of the observation $X$ given a parameter $\theta$ is defined by*

$$
L(X; \theta) = f_\theta(X).
$$

*In the $n$-sample case, due to independence, the log-likelihood can be written*

$$
L(X_1, \dots, X_n; \theta) = \prod_{i=1}^n f_\theta(X_i). \tag{1.1}
$$

**Example 1.15.** *If $X_1, \dots, X_n \sim \mathcal{B}(\theta)$. The density of a Bernoulli distribution with parameter $\theta$ can be written*

$$
f_\theta(x) = \theta \mathbb{1}(x = 1) + (1 - \theta)\mathbb{1}(x = 0) = \theta^x (1 - \theta)^{1-x} \mathbb{1}(x \in \{0, 1\})
$$

*hence we have*

$$
L(X_1, \dots, X_n; \theta) = \prod_{i=1}^n \theta^{X_i} (1 - \theta)^{1-X_i} = \theta^{\sum_{i=1}^n X_i} (1 - \theta)^{n - \sum_{i=1}^n X_i}
$$

*If $X_1, \dots, X_n \sim \mathcal{N}(\theta, \sigma^2)$, we get*

$$
L(X_1, \dots, X_n; \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X_i - \theta)^2}{2\sigma^2}\right).
$$

The likelihood can be interpreted as the probability of making observation $X$ if the underlying parameter is $\theta$. Indeed, if $\mathcal{M}$ is a set of discrete distributions (i.e. when $\nu$ is the counting measure), we have $f_\theta(x_i) = \mathbb{P}_\theta(X_i = x_i)$. Due to independence, we have

$$\mathbb{P}_\theta(X_1 = x_1, \ldots, X_n = x_n) = \prod_{i=1}^{n} f_\theta(x_i) = L(x; \theta) \quad \text{where} \ x = (x_1, \ldots, x_n)$$

In the continuous case (i.e. when $\nu$ is the Lebesgue measure), the probability of a given $x = (x_1, \ldots, x_n)$ is zero and we replace it by the value of the (joint) density in the point. This interpretation motivates the maximum likelihood estimator as the estimator of $g(\theta) = \theta$ seeking the parameter $\theta$ for which the actual observation $X$ is the most likely (i.e. which as it the largest "probability").

**Definition 1.16.** *A maximum likelihood estimator (MLE) of a parameter $\theta$ is an estimator satisfying*

$$\widehat{\theta} \in \underset{\theta \in \Theta}{argmax} \ L(X; \theta).$$

As we will see in some exercises, the maximum likelihood is not always unique. From a computational perspective (and due to the common product form of the likelihood, see (1.1)) it is often more convenient to maximize the logarithm of the likelihood, which then becomes a sum.

**Definition 1.17.** *The log-likelihood of the observation $X$ given a parameter $\theta$ is denoted by*

$$\ell(X; \theta) = \log L(X; \theta).$$

**Example 1.18** (Bernoulli distributions)**.** *As written above, the likelihood of an $n$-sample from a Bernoulli distribution with parameter $\theta$ is*

$$L(X_1, \ldots, X_n; \theta) = \prod_{i=1}^{n} \theta^{\sum_{i=1}^{n} X_i} (1 - \theta)^{n - \sum_{i=1}^{n} X_i}$$

*The log-likelihood takes the simple form*

$$\ell(X; \theta) = \left( \sum_{i=1}^{n} X_i \right) \log(\theta) + \left( n - \sum_{i=1}^{n} X_i \right) \log(1 - \theta) \coloneqq g(\theta)$$

*In order to find the maximizer of $g(\theta)$, we compute the derivative*

$$g'(\theta) = \frac{\sum_{i=1}^{n} X_i}{\theta} - \frac{n - \sum_{i=1}^{n} X_i}{1 - \theta} = \frac{\sum_{i=1}^{n} X_i - \theta n}{\theta(1 - \theta)}.$$

*There is a unique solution to $g'(\theta) = 0$ given by $\widehat{\theta}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$. Moreover $g'(\theta) > 0$ iff $\theta < \widehat{\theta}_n$, hence this solution is a maximizer, the MLE.*

*We remark that the same estimator could also have been obtained using the moment method by remarking that the parameter $\theta$ is also the mean of the Bernoulli distribution: $\theta = \mathbb{E}_{Z \sim \mathcal{B}(\theta)}[Z]$.*

**Example 1.19** (linear regression)**.** *We collect pairs of independent samples $(X_i, Y_i)$ such that $X_i \in \mathbb{R}^d$ comes from distribution with density $f$ and $Y_i = \theta^\top X_i + \varepsilon_i$. Assuming that the noise $\varepsilon_i$ is Gaussian with known variance $\sigma^2$ allows to write the likelihood of $n$ independent observations:*

$$L((X_1, Y_1), \ldots, (X_n, Y_n); \theta) = \prod_{i=1}^{n} f(X_i) \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left( -\frac{(Y_i - \theta^\top X_i)^2}{2\sigma^2} \right)$$

*The MLE is a minimizer of the function $g : \mathbb{R}^d \to \mathbb{R}$ given by*

$$g(\theta) = \sum_{i=1}^{n}(Y_i - \theta^\top X_i)^2,$$

*also known as the least-squares estimate. The gradient of $g$ is given by*

$$\nabla g(\theta) = -2\sum_{i=1}^{n}(Y_i - \theta^\top X_i)X_i$$

*and solving $\nabla g(\theta) = 0$ (and checking that the Hessian matrix is negative definite in this critical point) yields that the MLE is given by*

$$\widehat{\theta}_n = \left(\sum_{i=1}^{n} X_i X_i^\top\right)^{-1} \sum_{i=1}^{n} Y_i X_i,$$

*provided that the design matrix $\sum_{i=1}^{n} X_i X_i^\top$ is invertible.*

**Computational considerations** In the two examples above, the Maximum Likelihood Estimator can be computed explicitly, by finding the critical point (for which the derivative, or the gradient is zero) and proving that it is indeed a maximizer (e.g., by checking that the second derivative, or the Hessian is negative in the critical point). In more complex cases, the maximizer in the definition of the MLE can only be approximated using some optimization algorithm converging towards the maximizer (e.g. a gradient ascent).

In complex models involving latent variables, i.e. variable that are not actually observed (e.g. the membership of some individual in some cluster, which we also try to infer) as in mixture models, more fancy approximation scheme are needed, like the Expectation Maximization (EM) iterative algorithm.

**Example 1.20.** *In the logistic regression model, there are iid pairs of observations $(X_i, Y_i)$ where $X_i$ comes from some distribution on $\mathbb{R}^d$ that is assumed to have some density and $Y_i \in \{-1, 1\}$ is such that*

$$\mathbb{P}\left(Y_i = 1 | X_i = x\right) = \frac{1}{1 + e^{-x^\top\theta}}$$

*where $\theta \in \mathbb{R}^d$ is a regression parameter.*

*To define the likelihood of the data, we admit that the density of $(X_1, Y_1) \in \mathbb{R}^d \times \{0, 1\}$ is*

$$f_\theta(x, y) = \mathbb{P}(Y_1 = y | X_1 = x) f(x).$$

*You can verify that for all $x \in \mathbb{R}^d$ and all $y \in \{-1, 1\}$, $\mathbb{P}(Y_1 = y | X_1 = x) = \frac{1}{1 + e^{-yx^\top\theta}}$. The likelihood can therefore be written*

$$L((X_1, Y_1), \ldots, (X_n, Y_n)) = \prod_{i=1}^{n} f(X_i)\left(\frac{1}{1 + e^{-Y_i(X_i^\top\theta)}}\right)$$

*and a maximum likelihood estimator $\widehat{\theta}_n$ satisfies*

$$\widehat{\theta}_n \in \underset{\theta \in \mathbb{R}^d}{argmin} \sum_{i=1}^{n} \log\left(1 + e^{-Y_i(X_i^\top\theta)}\right).$$

*In this example, no closed-form expression exists for the MLE and we have to resort to an optimization algorithm.*

**$M$-estimators**    The MLE estimator is actually an example of a more general family of estimators called $M$-estimators, that are obtained as the minimization of some cumulative loss function of the data. A $M$ estimator is of the form

$$\widehat{\theta}_n \in \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \, M_n(\theta) \quad \text{where} \quad M_n(\theta) = \sum_{i=1}^{n} m(X_i; \theta).$$

In the particular case of the MLE, we have $m(X; \theta) = -\log f_\theta(X)$.

## 1.4    Beyond the likelihood

Under some additional regularity conditions on some dominated model it is possible to define an important quantity called the Fisher information, which is useful to provide a lower bound on the quality of an (unbiased) estimator (see Section 1.5). The Fisher information will also be useful in the next chapter to characterize the asymptotic distribution of the maximum likelihood estimator.

To ease the presentation, we define everything in the single-parameter setting, that is when the parameter space $\Theta$ is a subset of $\mathbb{R}$. All this concepts can be extended to the multi-dimensional setting by replacing derivative with gradients, variances with covariances, and second derivative with Hessian. We will briefly discuss this extension afterwards.

**Definition 1.21.** *A (uni-dimensional) parameteric model $\mathcal{M} = \{P_\theta, \theta \in \Theta \subseteq \mathbb{R}\}$ is regular if*

1. *it is dominated by some reference measure $\nu$ and for all $\theta$, the support of $f_\theta$, $S = \{x \in \mathcal{X} : f_\theta(x) > 0\}$ is independent of $\theta$*

2. *for all $x \in S$, $\theta \mapsto f_\theta(x)$ is twice differentiable on $\Theta$ and its second derivative is continuous*

3. *for any event $\mathcal{E}$, we have*

$$\frac{\partial}{\partial \theta} \int_{\mathcal{E}} f_\theta(x) d\nu(x) = \int_{\mathcal{E}} \frac{\partial}{\partial \theta} f_\theta(x) d\nu(x)$$
$$\frac{\partial^2}{\partial^2 \theta} \int_{\mathcal{E}} f_\theta(x) d\nu(x) = \int_{\mathcal{E}} \frac{\partial^2}{\partial^2 \theta} f_\theta(x) d\nu(x)$$

**Example 1.22.** *We can show that many classical parametric model satisfy this assumption (e.g. Bernoulli models, Gaussian model, Poisson model). A counter-example that will be studied in an exercise is the family of uniform distributions on $[0, \theta]$ for $\theta \in \mathbb{R}^+$, which already violates assumption 1.*

### 1.4.1    The Fisher information

**Definition 1.23.** *The score function is defined as the derivative of the log-likelihood.*

$$s(X; \theta) = \frac{\partial \ell(X; \theta)}{\partial \theta} = \frac{1}{f_\theta(X)} \frac{\partial f_\theta(X)}{\partial \theta}$$

An important property of the score under a regular model is the following.

**Lemma 1.24.** *Under a regular model, for all $\theta \in \Theta$, $\mathbb{E}_\theta[s(X; \theta)] = 0$.*

*Proof.*

$$\mathbb{E}_\theta[s(X;\theta)] = \int \frac{\partial \ell(x;\theta)}{\partial \theta} f_\theta(x) d\nu(x) = \int \frac{\frac{\partial}{\partial \theta} f_\theta(x)}{f_\theta(x)} f_\theta(x) d\nu(x) = \int \frac{\partial}{\partial \theta} f_\theta(x) d\nu(x)$$

$$\underset{(a)}{=} \int_S \frac{\partial}{\partial \theta} f_\theta(x) d\nu(x) \underset{(b)}{=} \frac{\partial}{\partial \theta} \left( \int_S f_\theta(x) d\nu(x) \right) \underset{(c)}{=} \frac{\partial}{\partial \theta}(1) = 0$$

where $(a)$ uses property 1. of a regular model, $(b)$ uses property 3 and $(c)$ uses that $f_\theta$ is a density.

$\square$

The Fisher information matrix is defined as the variance of the score, which is equal to its second moment as the score is centered.

**Definition 1.25.** *In a regular model, the Fisher information of the observation $X$ is defined as*

$$I^X(\theta) = \mathrm{Var}_\theta\left[s(X;\theta)\right] = \mathbb{E}_\theta\left[(s(X,\theta))^2\right] .$$

*In the $n$-sample case, we will write $I_n(\theta)$ to denote the Fisher information of the $n$-sample, and $I(\theta)$ the Fisher information of the observation made of a single realisation $X_1 \sim P_\theta$.*

### 1.4.2 Some properties of the Fisher information

**Lemma 1.26.** *Under a regular model, it holds that $I^X(\theta) = -\mathbb{E}_\theta\left[\frac{\partial^2 \ell(X;\theta)}{\partial^2 \theta}\right]$.*

**Exercise 1.27.** *Prove it. Hint: start by computing the right-hand side, using property 3. of a regular model as in the proof of Lemma 1.24.*

The above lemma can be useful for the computation of the Fisher information. We now present another interesting property which is the additivity of the Fisher information. This property follows from the fact that the density of a couple of independent random variable is the product of their densities, and uses properties of the logarithm.

**Lemma 1.28.** *If $X$ and $Y$ are two independent random variables defined on the same probability space $(\Omega, \mathcal{F}, \mathbb{P}_\theta)$, then*

$$I^{(X,Y)}(\theta) = I^X(\theta) + I^Y(\theta) .$$

*It follows that for a $n$ sample $X = (X_1, \ldots, X_n) \overset{iid}{\sim} P_\theta$,*

$$I_n(\theta) = I^X(\theta) = nI^{X_1}(\theta) = nI(\theta) .$$

**Example 1.29.** *Consider the Bernoulli model $X_1, \ldots, X_n \overset{iid}{\sim} \mathcal{B}(\theta)$. We have seen above that $I_n(\theta) = nI(\theta)$ where $I(\theta)$ is the Fisher information in a model with one Bernoulli observation $X_1$. In this model, we have*

$$\begin{aligned}
L(X_1;\theta) &= \theta^{X_1}(1-\theta)^{1-X_1} \\
\ell(X_1;\theta) &= X_1 \log(\theta) + (1-X_1)\log(1-\theta) \\
\frac{\partial \ell(X_1;\theta)}{\partial \theta} &= \frac{X_1}{\theta} - \frac{1-X_1}{1-\theta} \\
\frac{\partial^2 \ell(X_1;\theta)}{\partial^2 \theta} &= -\frac{X_1}{\theta^2} + \frac{1-X_1}{(1-\theta)^2}
\end{aligned}$$

*hence*

$$I(\theta) = -\mathbb{E}_\theta\left[\frac{\partial^2 \ell(X_1;\theta)}{\partial^2 \theta}\right] = \mathbb{E}_\theta\left[\frac{X_1}{\theta^2} - \frac{1-X_1}{(1-\theta)^2}\right] = \frac{1}{\theta} - \frac{1}{1-\theta} = \frac{1}{\theta(1-\theta)}$$

*Finally, using Lemma 1.28, we get $I_n(\theta) = \frac{n}{\theta(1-\theta)}$.*

**Extension to the multi-dimensional setting**   If $\theta = (\theta_1, \ldots, \theta_d)$, the score is a vector in $\mathbb{R}^d$, defined as

$$s(X;\theta) = \nabla_\theta \ell(X;\theta) = \left(\frac{\partial \ell(X;\theta)}{\partial \theta_1}, \ldots, \frac{\partial \ell(X;\theta)}{\partial \theta_d}\right)^\top.$$

In (an extension of the definition of a) regular model, the score satisfies $\mathbb{E}[s(X;\theta)] = 0$ and the Fisher information is defined as the (covariance) of the score, ie

$$I(\theta) = \mathbb{E}\left[(s(X,\theta))(s(X,\theta))^\top\right].$$

The Fisher information is therefore a $d \times d$ matrix, and a counterpart of Lemma 1.26 can be proved:

$$I(\theta) = -\mathbb{E}\left[\left(\frac{\partial^2 \ell(X;\theta)}{\partial \theta_i \partial \theta_j}\right)_{\substack{1 \le i \le d \\ 1 \le j \le d}}\right].$$

### 1.4.3   Interpretation of the Fisher information *(more advanced)*

The Fisher information will be shortly related to the minimal variance that a unbiased estimator can have. In this section, we give some elements of explanation as to why it can be called "information".

First, due to its additivity property (Lemma 1.28), if we interpret $I(\theta)$ as an amount of "information" brought by one sample, we note that the Fisher information of a $n$-sample is the sum of all the information brought by individual samples. Moreover, another property is that given an observation $X$, any "summary" of this observation in the form of a statistic $S = s(X)$ has a smaller Fisher information.

**Lemma 1.30.** *For any statistic $S = s(X)$ of an observation $X$ defined on a probability space $(\Omega, \mathcal{F}, \mathbb{P}_\theta)$, we have $I^S(\theta) \le I^X(\theta)$.*

*Proof.* Let's write down the proof assuming that $X$ takes values in a discrete space $\mathcal{X}$ (to avoid the concept of conditional density). $X$ and $S = s(X)$ are clearly defined on the same probability space $(\Omega, \mathcal{F}, \mathbb{P}_\theta)$. We can write

$$\mathbb{P}_\theta(X = x) \quad = \quad \mathbb{P}_\theta(X = x, S = s(x)) = \mathbb{P}_\theta(X = x|S = s(x))\mathbb{P}_\theta(S = s(x))$$

Hence, for any $x \in \mathcal{X}$, writing $s = s(x)$, we have

$$f_\theta(x) = f_\theta(x|s)\widetilde{f}_\theta(s)$$

where we introduce $f_\theta$ the density of $X$, $\widetilde{f}_\theta$ the density of $S$ and $f_\theta(x|s) := \mathbb{P}_\theta(X = x|S = s)$. Taking the logarithm and differentiating twice yields

$$\frac{\partial^2 \log f_\theta(x)}{\partial^2 \theta} = \frac{\partial^2 \log \widetilde{f}_\theta(s)}{\partial^2 \theta} + \frac{\partial^2 \log f_\theta(x|s)}{\partial^2 \theta}$$

and in particular

$$\frac{\partial^2 \log f_\theta(X)}{\partial^2 \theta} = \frac{\partial^2 \log \widetilde{f}_\theta(S)}{\partial^2 \theta} + \frac{\partial^2 \log f_\theta(X|S)}{\partial^2 \theta}$$

Taking the expectation and using Lemma 1.26 yields

$$I^X(\theta) = I^S(\theta) - \mathbb{E}_\theta \left[ \frac{\partial^2 \log f_\theta(X|S)}{\partial^2 \theta} \right]$$

We conclude by noting that

$$-\mathbb{E}_\theta \left[ \frac{\partial^2 \log f_\theta(X|S)}{\partial^2 \theta} \right] = \sum_s \mathbb{P}_\theta(S = s) \underbrace{\left[ -\mathbb{E}_\theta \left[ \frac{\partial^2 \log \mathbb{P}_\theta(X|S = s)}{\partial^2 \theta} \right] \right]}_{\geq 0}$$

and the term between brackets is positive as it is the Fisher information of the conditional distribution of $X$ given $(S = s)$.

$\square$

From this result a good statistic $S = s(X)$ is one that doesn't loose information, i.e. for which $I^S(\theta) = I^X(\theta)$. Sufficient statistic have this property, and are defined below.

**Definition 1.31.** *A statistic $S = s(X)$ is called sufficient for $\theta$ if the distribution of $X = (X_1, \ldots, X_n)$ conditionally to $S$ does not depend on $\theta$.*

We admit the following characterization.

**Theorem 1.32** (Neyman-Fisher)**.** *The statistic $S = s(X_1, \ldots, X_n)$ is sufficient for $\theta$ is there exists two positive functions $g$ and $h$ such that the density of $X$ can be written*

$$f_\theta(x_1, \ldots, x_n) = g(x_1, \ldots, x_n) h(s(x_1, \ldots, x_n); \theta) .$$

### 1.4.4 The Kullback-Leibler divergence

We define another information theoretic quantity that is related to the likelihood (or actually rather to a likelihood ratio) and provides some notion of "distance" (although it is not a distance in the topological sense) between probability measures.

**Definition 1.33.** *For two probability measure $P$ and $Q$ that have a densities $f$ and $g$ with respect to the same probability measure $\nu$ and such that $g(x) = 0 \Rightarrow f(x) = 0$, we have*

$$\mathrm{KL}(P, Q) = \mathbb{E}_{X \sim P} \left[ \log \frac{f(x)}{g(x)} \right].$$

*In particular, if $P_\theta$ and $P_{\theta'}$ are two distributions in a regular model (actually assumption 1. in Definition 1.21 is sufficient), we can define*

$$\mathrm{K}(\theta, \theta') := \mathrm{KL}(P_\theta, P_{\theta'}) = \mathbb{E}_\theta \left[ \log \frac{f_\theta(X)}{f_{\theta'}(X)} \right] .$$

**Example 1.34.** *The KL divergence between $\mathcal{N}(\mu, \sigma^2)$ and $\mathcal{N}(\mu', \sigma^2)$ is*

$$\mathrm{K}(\mu, \mu') = \frac{(\mu - \mu')^2}{2\sigma^2} \; .$$

*The KL divergence between two Bernoulli distributions of parameters $\theta$ and $\theta'$ is*

$$\mathrm{K}(\theta, \theta') = \theta \log\left(\frac{\theta}{\theta'}\right) + (1 - \theta) \log\left(\frac{1 - \theta}{1 - \theta'}\right) \; .$$

**Proposition 1.35.** $\mathrm{KL}(P, Q) \geq 0$ *and* $\mathrm{KL}(P, Q) = 0$ *if and only if* $P = Q$.

## 1.5   The Cramer-Rao lower bound

The Fisher information defined in the previous section enables us (in the case of uni-dimensional estimation) to solve the following question: what is the minimal variance of an unbiased estimator? We consider this question for regular models.

**Theorem 1.36.** *Assume the statistical model is regular. Let $\widehat{g}$ be an estimator of $g(\theta) \in \mathbb{R}$ where $g$ is differentiable. We assume that $\widehat{g} = h(X)$ is such that $\mathbb{E}_\theta[\widehat{g}_n] = g(\theta)$ (unbiased estimator) and*

$$\frac{\partial}{\partial \theta} \int h(x) f_\theta(x) d\nu(x) = \int h(x) \left(\frac{\partial}{\partial \theta} f_\theta(x)\right) d\nu(x)$$

*Then, for all $\theta \in \Theta$,*

$$\mathrm{Var}_\theta[\widehat{g}] \geq \frac{(g'(\theta))^2}{I(\theta)}.$$

*Proof.* The idea of the proof is to differentiate $g(\theta) = \mathbb{E}_\theta[h(X)]$ and introduce the score. Using one of the assumptions, we can write

$$
\begin{aligned}
g'(\theta) \quad &= \quad \frac{\partial}{\partial \theta} \int_S h(x) f_\theta(x) d\nu(x) = \int_S h(x) \left(\frac{\partial}{\partial \theta} f_\theta(x)\right) d\nu(x) \\
&= \quad \int_S h(x) \left(\frac{\partial}{\partial \theta} \log f_\theta(x)\right) f_\theta(x) d\nu(x) \\
&\overset{(a)}{=} \quad \int_S (h(x) - \mathbb{E}_\theta[h(X)]) \left(\frac{\partial}{\partial \theta} \log f_\theta(x)\right) f_\theta(x) d\nu(x) + \mathbb{E}_\theta[h(X)] \underbrace{\int_S \left(\frac{\partial}{\partial \theta} \log f_\theta(x)\right) f_\theta(x) d\nu(x)}_{=0} \\
&\overset{(b)}{=} \quad \int_S (h(x) - \mathbb{E}_\theta[h(X)]) \left(\frac{\partial}{\partial \theta} \log f_\theta(x) - \mathbb{E}_\theta\left[\frac{\partial}{\partial \theta} \log f_\theta(X)\right]\right) f_\theta(x) d\nu(x)
\end{aligned}
$$

where both $(a)$ and $(b)$ use that the expected score is zero by Lemma 1.24.

Now we assume that $\mathbb{E}_\theta[h^2(X)] < \infty$ (otherwise, the inequality in Theorem 1.36 is trivially true). Then we can use the Cauchy-Schwarz inequality to get

$$
\begin{aligned}
|g'(\theta)| \quad &\leq \quad \sqrt{\mathbb{E}_\theta\left[(h(x) - \mathbb{E}_\theta[h(X)])^2\right]} \sqrt{\mathbb{E}_\theta\left[\left(\frac{\partial}{\partial \theta} \log f_\theta(x) - \mathbb{E}_\theta\left[\frac{\partial}{\partial \theta} \log f_\theta(X)\right]\right)^2\right]} \\
&\leq \quad \sqrt{\mathrm{Var}_\theta[h(X)]} \sqrt{\mathrm{Var}_\theta\left[\frac{\partial}{\partial \theta} \log f_\theta(X)\right]} \\
&\leq \quad \sqrt{\mathrm{Var}_\theta[h(X)]} \sqrt{I(\theta)}
\end{aligned}
$$

where the last step uses the definition of the Fisher information.

$\square$

An unbiased estimator that achieves the Cramer-Rao lower bound for all values of $\theta \in \Theta$ is called efficient (or uniformly efficient). The example below show that there exists efficient estimators.

**Exercise 1.37.** *Verify that in the Bernoulli model $X_1, \ldots X_n \overset{iid}{\sim} \mathcal{B}(p)$ the MLE is an efficient estimator.*

## 1.6 Exponential families

Actually, the reason why in the Bernoulli model we can find an efficient estimator comes from the fact that the set of Bernoulli distributions is a particular example of exponential family. We define exponential families below, and we will see several of their properties in this class.

**Definition 1.38.** *An exponential family is a set of probability distributions on some set $\mathcal{X}$ defined as*

$$\mathcal{P} = \{P_\theta, \theta \in \Theta : \; P_\theta \text{ has a density } \; f_\theta(x) = h(x) \exp\left(a(\theta)^\top T(x) - b(\theta)\right) \text{ wrt to } \nu\}$$

*where $\nu$ is a reference measure (common to all distributions), $h : \mathcal{X} \to \mathbb{R}^+$ is a positive function, $a : \Theta \to \mathbb{R}^d$, $b : \Theta \to \mathbb{R}$ and $T : \mathcal{X} \to \mathbb{R}^d$ are some functions and $u^\top v = \sum_{i=1}^d u_i v_i$ is the scalar product in $\mathbb{R}^d$.*

$T(x) \in \mathbb{R}^d$ is called the canonical statistic and $d$ is the dimension of the exponential family. In a one-dimensional exponential family, the density can simply be expressed

$$f_\theta(x) = h(x) \exp\left(a(\theta)T(x) - b(\theta)\right).$$

**Example 1.39.** *The family of Bernoulli distributions $\mathcal{P} = \{\mathcal{B}(p), p \in (0,1)\}$ form an exponential family (of dimension 1). Indeed, its density with respect to the counting measure is*

$$
\begin{aligned}
f_p(x) &= p^x (1-p)^{1-x} \mathbb{1}(x \in \{0,1\}) \\
&= \exp(x \log(p) + (1-x)\log(1-p))\mathbb{1}(x \in \{0,1\}) \\
&= h(x) \exp\left(x \log \frac{p}{1-p} + \log(1-p)\right)
\end{aligned}
$$

*with $h(x) = \mathbb{1}(x \in \{0,1\})$. Introducing the natural parameter $\theta = \log \frac{p}{1-p}$, we have $p = \frac{e^\theta}{1+e^\theta}$ and $\log(1-p) = -\log(1+e^\theta)$. Hence we have*

$$f_p(x) = h(x) \exp\left(x\theta - b(\theta)\right)$$

*with $b(\theta) = \log(1+e^\theta)$ and the family of Bernoulli distributions can be written as the family of densities*

$$\{f_\theta(x) = h(x) \exp(a(\theta)T(x) - b(\theta)), \theta \in \mathbb{R}\}$$

*where $a(\theta) = \theta$, $T(x) = x$ and $b(\theta) = \log(1+e^\theta)$ and the reference measure is the counting measure.*

We can prove that efficient estimator can only exist in some exponential families, and for a particular parameter to estimate. There are therefore not so much common. In the next chapter, we will define an asymptotic notion of efficiency, which can be easier to attain.

# Chapter 2

# Asymptotic properties of estimators

In this chapter, we focus on the $n$-sample case, in which $X = (X_1, \ldots, X_n) \overset{iid}{\sim} P_\theta$. For each $n$, given an estimator $\widehat{g}_n = h(X_1, \ldots, X_n)$ of a certain parameter of interest $g(\theta)$, we are interested in studying the sequence of estimators $(\widehat{g}_n)_n$ when the sample size $n$ grows large. As the $\widehat{g}_n$ are random variables, we first recap the different notion of convergences, as well as some important results.

## 2.1 Refresher: Convergence of random variables

**Definition 2.1.** *Let $Z_1, Z_2, \ldots$ be a sequence of random variable and let $Z$ be another random variable. Let $F_n$ denote the CDF of $Z_n$ and let $F$ denote the cdf of $Z$.*

1. ***$Z_n$ converges to $Z$ in probability** if, for every $\varepsilon > 0$, $\lim_{n\to\infty} \mathbb{P}(|Z_n - Z| > \varepsilon) = 0$.*

   *We write $Z_n \overset{P}{\longrightarrow} Z$.*

2. ***$Z_n$ converges to $Z$ in distribution** if, $\lim_{n\to\infty} F_n(t) = F(t)$ for all $t$ for which $F$ is continuous.*

   *We write $Z_n \rightsquigarrow Z$.*

3. ***$Z_n$ converges to $Z$ almost surely** if $\mathbb{P}\left(\lim_{n\to\infty} Z_n = Z\right) = 1$. We write $Z_n \overset{a.s.}{\longrightarrow} Z$.*

4. ***$Z_n$ converges to $Z$ in quadratic mean** if $\lim_{n\to\infty} \mathbb{E}\left[(Z_n - Z)^2\right] = 0$. We write $Z_n \overset{L^2}{\longrightarrow} Z$.*

In statistics, the first two notions are the most common, and we will mostly discuss them in the following. The definitions above were all given for real-values random variables, but can be extended to the multi-dimensional setting. For the convergence in probability, the distance between $Z_n$ and $Z$ and $\mathbb{R}^d$ can no longer be measured with the absolute value, but given any distance $d$ on $\mathbb{R}^d$ (for example the Euclidian distance), we define $Z_n \overset{P}{\longrightarrow} Z$ is for all $\varepsilon > 0$, $\lim_{n\to\infty} \mathbb{P}\left(d(Z_n, Z) > \varepsilon\right) = 0$.

The convergence in distribution in $\mathbb{R}^d$ can still be characterized by the cdf, but in this case, the cdf is a multi-variate function and we should have, for all $z = (z_1, \ldots, z_d)$ in which $F$ is continuous,

$$\lim_{n\to\infty} \mathbb{P}\left(Z_n^1 \le z_1, \ldots, Z_n^d \le z_d\right) = \mathbb{P}(Z^1 \le z_1, \ldots, Z^d \le z_d) = 0.$$

**Example 2.2.** *$Z_n \sim \mathcal{N}(0, \frac{1}{n})$. Justify that $Z_n$ converges to $0$ (the random variable that is constant and equal to zero) in distribution and in probability.*

### 2.1.1   Properties

The following relationship between the different convergence notions are useful.

**Lemma 2.3.**      *1. $X_n \xrightarrow{P} X$ implies that $X_n \rightsquigarrow X$*

*2. $X_n \xrightarrow{P} c$ where c is a constant if and only if $X_n \rightsquigarrow X$*

*3. $X_n \xrightarrow{L^2} X$ implies that $X_n \xrightarrow{P} X$*

We note that $(a)$ and $(c)$ are not equivalences. In particular, beyond the case of convergence to constants, the convergence in distribution does not imply the convergence in probability. A (contrived) counter-example is the following: take any symmetric distribution $Y$, that is a distribution for which $Y$ and $-Y$ have the same distribution (for example, a centered Gaussian distribution). Define $Z_n = Y$ for all $n$ and $Z = -Y$. As the cdf and $Z_n$ and that of $Z$ are equal, we have in particular $Z_n \rightsquigarrow Z$. However, $\mathbb{P}\left(|Z_n - Z| > \varepsilon\right) = \mathbb{P}\left(|2Y| > \varepsilon\right)$ does not converge to zero for every $\varepsilon$ (unless $Y = 0$ a.s.).

**Lemma 2.4** (continuous mapping). *Let $g : \mathcal{X} \to \mathbb{R}$ be a continuous function. Then*

- *If $X_n \xrightarrow{P} X$ then $g(X_n) \xrightarrow{P} g(X)$*

- *If $X_n \rightsquigarrow X$ then $g(X_n) \rightsquigarrow g(X)$*

**Lemma 2.5** (Slutsky lemma). *If $X_n \rightsquigarrow X$ and $Y_n \xrightarrow{P} c$ where c is a constant, then, for any continuous function g,*

$$g(X_n, Y_n) \rightsquigarrow g(X, c) .$$

*In particular*

- *$X_n + Y_n \rightsquigarrow X + c$*

- *$X_n Y_n \rightsquigarrow cX$*

Slutsky's lemma is a consequence of the fact that as a couple of random variables $(X_n, Y_n)$ converges in distribution to $(X, c)$ (and the fact that the continuous mapping lemma also applies to multi-variate random variables).

### 2.1.2   Two fundamental theorems

We recall here the two fundamental theorems in statistics: the law of large numbers and the central limit theorem. Given an iid sequence $Z_i$, they provide some convergence results for the empirical average

$$\widehat{Z}_n = \frac{1}{n} \sum_{i=1}^{n} Z_i.$$

**Theorem 2.6** (Law of large numbers). *If $(Z_i)_{i \in \mathbb{N}}$ is an iid sequence with $\mathbb{E}[Z_1] < \infty$, we have*

$$\frac{1}{n} \sum_{i=1}^{n} Z_i \xrightarrow{P} \mathbb{E}[Z_1]$$

Actually, a stronger version of this result (called the strong law of large numbers) holds under the same assumptions, in which the convergence in probability is replaced by an almost sure convergence.

**Theorem 2.7** (Central limit theorem). *If $(Z_i)_{i \in \mathbb{N}}$ is an iid sequence with $\mathbb{E}[Z_1^2] < \infty$, letting $\mu = \mathbb{E}[Z_1]$ and $\sigma^2 = \mathrm{Var}[Z_1]$, we have*

$$\sqrt{\frac{n}{\sigma^2}} \left( \widehat{Z}_n - \mu \right) \rightsquigarrow \mathcal{N}(0, 1)$$

Under the hypotheses of the central limit theorem, $\widehat{Z}_n$ can be written

$$Z_n = \mu + \sqrt{\frac{\sigma^2}{n}} Y_n$$

where $Y_n \rightsquigarrow \mathcal{N}(0, 1)$. Therefore, informally, the distribution of $\widehat{Z}_n$ is close to $\mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$, a Gaussian distribution whose variance decays to zero and is therefore more and more concentrated around $\mu$. We may write $\widehat{Z}_n \approx \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$ and talk about the asymptotic distribution of $\widehat{Z}_n$.

## 2.2 Consistency and asymptotic normality

**Definition 2.8.** *An estimator $\widehat{g}_n$ of $g(\theta)$ is consistent if for every $\theta \in \Theta$, $\widehat{g}_n \xrightarrow{P} g(\theta)$.*

Consistency of estimators will often follow from the law of large numbers. When we further have an almost sure convergence, that is when $\widehat{g}_n \xrightarrow{a.s.} g(\theta)$, we shall say that $\widehat{g}_n$ is strongly consistent.

Lemma 2.3 and Lemma 2.4 also yield the following properties:

- If the quadratic risk $\mathrm{R}_\theta(\widehat{g}_n)$ goes to zero when $n$ goes to infinity, $\widehat{g}_n$ is consistent.

- If $\widehat{\theta}_n$ is a consistent estimator of $\theta$ and $g$ is a continuous mapping, then $\widehat{g}_n = g(\widehat{\theta}_n)$ is a consistent estimator of $g(\theta)$.

**Example 2.9.** *Using the law of large number directly yields that the empirical mean $\widehat{\mu}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$ of any $n$ sample with a finite expectation $\mu$ is a consistent estimator of $\mu$. As for the empirical variance, which can be written*

$$\widehat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^{n} X_i^2 - \left( \frac{1}{n} \sum_{i=1}^{n} X_i \right)^2$$

*the law of large numbers also yields the convergence almost surely to $\mathbb{E}[X_1^2] - (\mathbb{E}[X_1])^2$, the variance of the distribution, provided that $X_1$ has a second moment. Hence the consistency.*

Given a consistent estimator $\widehat{g}_n$, we may be interested in how fast $\widehat{g}_n - g(\theta)$ converges to zero. To do so, we will look at the limit distribution of (some re-normalization) of this random variable.

If we take the example of the empirical $\widehat{\mu}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$ estimator of the common mean $\mu$ of some $n$ sample $(X_1, \ldots, X_n)$ that has variance $\sigma^2$, the Central Limit Theorem tells us that

$$\sqrt{n} \left( \widehat{\mu}_n - \mu \right) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

Here, the limit distribution is Gaussian, and the convergence speed is $\sqrt{n}$. Due to the generality of the Central Limit Theorem, we expect this Gaussian limit behavior to be a general pattern for estimators, hence the definition of asymptotic normality.

**Definition 2.10.** *An estimator is $\widehat{g}_n$ of $g(\theta)$ is asymptotically normal is it satisfies, for all $\theta \in \Theta$,*

$$\sqrt{n}\left(\widehat{g}_n - g(\theta)\right) \rightsquigarrow \mathcal{N}(0, \sigma_\theta^2)$$

*where $\sigma_\theta^2$ is called the asymptotic variance.*

**Example 2.11.** *Given a $n$ sample $(X_1, \ldots, X_n)$ from some distribution with cdf $F$, we consider the value of the empirical cdf in $x$ as an estimator for $F(x)$ (see Example 1.9). We have*

$$\widehat{F}_n(x) = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}(X_i \le x)$$

*This estimator is the empirical mean of the random variables $Z_i = \mathbb{1}(X_i \le x)$ whose mean is $F(x)$ and whose variance is $F(x)(1 - F(x))$. Using the central limit theorem, we can obtain the asymptotic distribution of $\widehat{F}_n(x)$:*

$$\sqrt{n}\left(\widehat{F}_n(x) - F(x)\right) \rightsquigarrow \mathcal{N}\left(0, F(x)(1 - F(x))\right)$$

*The mean of the estimator is $F(x)$ and it is asymptotically normal with variance $F(x)(1 - F(x))$.*

It is worth mentioning that there exists estimators that are have a limiting distribution, but are not asymptotically normal. It means that they satisfy something like

$$g(n)\left(\widehat{g}_n - g(\theta)\right) \rightsquigarrow Z$$

where $g(n)$ is some convergence speed (that can be different than $\sqrt{n}$) and $Z$ is some fixed distribution (that is not necessarily Gaussian).

**Example 2.12.** *As studied in exercise, in the model $X_1, \ldots, X_n \overset{iid}{\sim} \mathcal{U}([0, \theta])$, the moment estimator is $\widehat{\theta}_n = \frac{2}{n} \sum_{i=1}^{n} X_i$ while the MLE is $\widetilde{\theta}_n = \max_{i=1..n} X_i$.*

*Using the Central Limit Theorem (and the continuous mapping lemma), we can show that*

$$\sqrt{n}\left(\widehat{\theta}_n - \theta\right) \rightsquigarrow \mathcal{N}\left(0, \frac{\theta^2}{3}\right)$$

*hence the moment estimator is asymptotically normal with asymptotic variance $\sigma_\theta = \frac{\theta^2}{3}$.*

*On the other hand, we computed the distribution of $\widetilde{\theta}_n$ in exercise, showing that*

$$\mathbb{P}\left(\widetilde{\theta}_n \le x\right) = \begin{cases} 1 & \text{if } x \ge \theta \\ \frac{x^n}{\theta^n} \mathbb{1}_{[0,\theta]}(x) & \text{else.} \end{cases}$$

*Hence we have, for all $u > 0$,*

$$\mathbb{P}\left(\theta - \widetilde{\theta}_n \ge u\right) = \left(1 - \frac{u}{\theta}\right)^n \mathbb{1}_{[0,\theta]}(u)$$

*and finally, for all $t > 0$,*

$$\mathbb{P}\left(n\left(\theta - \widetilde{\theta}_n\right) \ge t\right) = \left(1 - \frac{t}{n\theta}\right)^n \mathbb{1}_{[0,n\theta]}(t)$$

*The limit of the right-hand side when $n$ goes to infinity is equal to $e^{-\frac{t}{\theta}} = \mathbb{P}(Z > t)$ where $Z$ is an exponential distribution with parameter $1/\theta$. Finally, one can write*

$$n\left(\theta - \widetilde{\theta}_n\right) \rightsquigarrow \mathcal{E}\left(\theta^{-1}\right).$$

*This provide another argument for using the MLE over the moment estimator in this particular case, as its asymptotic convergence is faster.*

*In Section 2.4, we will actually see that for regular models, the MLE is asymptotically Gaussian. The reason for this different behavior stems from the fact that the model considered here is not regular: one can indeed see that all the possible densities do not have the same support.*

**Why asymptotic normality?** Knowing the exact distribution of an estimator would be very useful to derive (exact) confidence region or tests with guaranteed type I error. Knowing it asymptotically allows to derive asymptotic confidence regions, and the corresponding tests, as we will recall at the beginning of the next chapter.

**Comparing asymptotically normal estimators** Between two asymptotically normal estimator, the one with smallest asymptotic variance $\sigma_\theta^2$ is the one that converges "faster" to the parameter $g(\theta)$. This can be measured by the fact that, if we build asymptotic confidence intervals for $g(\theta)$ of level $1 - \alpha$, using the estimator with smallest asymptotic variance will yield the smallest confidence region.

If two asymptotically normal estimators $\widehat{g}_n$ and $\widetilde{g}_n$ have respective asymptotic variances $\sigma_\theta^2$ and $\widetilde{\sigma}_\theta^2$ and that $\sigma_\theta^2 \leq \widetilde{\sigma}_\theta^2$ for all $\theta \in \Theta$ (with at least one strict inequality), we say that $\widehat{g}_n$ is asymptotically more efficient than $\widetilde{g}_n$.

## 2.3 The Delta method

We now present a useful tool to compute asymptotic distributions of some transformation of an asymptotically normal estimator: the so-called Delta method. This result implies that under some mild conditions, if $\widehat{\theta}_n$ is an asymptotically normal estimator of $\theta$, then $g(\widehat{\theta}_n)$ is an asymptotically normal estimator of $g(\widehat{\theta}_n)$.

**Theorem 2.13.** *Suppose that for some sequence of random variance $(Z_n)$,*

$$\sqrt{n}(Z_n - \mu) \rightsquigarrow \mathcal{N}(0, \sigma^2)$$

*and that $g$ is a differentiable function such that $g'(\mu) \neq 0$. Then*

$$\sqrt{n}(g(Z_n) - g(\mu)) \rightsquigarrow \mathcal{N}\left(0, (g'(\mu))^2 \sigma^2\right).$$

*In other words,*

$$Z_n \approx \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right) \quad \text{implies that} \quad g(Z_n) \approx \mathcal{N}\left(g(\mu), (g'(\mu))^2 \frac{\sigma^2}{n}\right).$$

*Proof.* The proof follows from using a Taylor expansion around $\mu$. As $g$ is differentiable, we have that for all $n$, there exists $\mu_n$ in $(Z_n, \mu)$ (if $Z_n < \mu$) or in $(\mu, Z_n)$ (if $Z_n \geq \mu$) such that

$$g(Z_n) = g(\mu) + g'(\mu_n)(Z_n - \mu)$$

hence

$$\sqrt{n}|g(Z_n) - g(\mu)| = g'(\mu_n)\sqrt{n}(Z_n - \mu).$$

As $|\mu_n - \mu| \leq |Z_n - \mu|$ and $Z_n \xrightarrow{P} \mu$, we have that $\mu_n \xrightarrow{P} \mu$. If we assume $g'$ to be continuous[1], it follows from Lemma 2.4 that $g'(\mu_n) \xrightarrow{P} g'(\mu)$.

By assumption, we also have that $\sqrt{n}(Z_n - \mu) \rightsquigarrow Z$ where $Z \sim \mathcal{N}(0,1)$. It follows from Slutsky's lemma (Lemma 2.5) that

$$\sqrt{n}|g(Z_n) - g(\mu)| \rightsquigarrow g'(\mu)Z$$

whose distribution is $\mathcal{N}(0, (g'(\mu))^2 \sigma^2)$.

$\square$

There exists also a multi-variate version of the Delta method, stated below.

**Theorem 2.14.** *Let $Z_n = (Z_{n,1}, \ldots, Z_{n,d})$ be a sequence of random vectors in $\mathbb{R}^d$ such that*

$$\sqrt{n}(Z_n - \mu) \rightsquigarrow \mathcal{N}(0, \Sigma)$$

*where $\mu \in \mathbb{R}^d$ and $\Sigma \in \mathbb{R}^{d \times d}$ is the covariance matrix. Let $g : \mathbb{R}^d \to \mathbb{R}$ be a differentiable function and let $\nabla g(z) = (\frac{\partial g}{\partial z_1}(z), \ldots, \frac{\partial g}{\partial z_d}(z))^\top$ be its gradient. If all the components of $\nabla g(\mu)$ are non-zero, then*

$$\sqrt{n}\left(g(Z_n) - g(\mu)\right) \rightsquigarrow \mathcal{N}\left(0, (\nabla g(\mu))^\top \Sigma \left(\nabla g(\mu)\right)\right)$$

**Example 2.15.** *In a clinical trials involving two treatments, we observe the outcome of treatment 1 (a placebo) on a pool of $n$ patients. For $i \in \{1, \ldots, n\}$, we record $X_i = 1$ if the treatment is a success for patient $i$, $X_i = 0$ is it is a failure. The outcome of treatment 2 (the new drug) is observed on another pool of $n$ patients, with $Y_j \in \{0,1\}$ indicating success of failure for patient $j \in \{1, \ldots n\}$. We assume that all the $X_i$ and $Y_j$ are independent and that $X_i \sim \mathcal{B}(p_1)$ and $Y_j \sim \mathcal{B}(p_2)$ where $p_1$ and $p_2$ are the probability of efficacy of treatment 1 and 2, respectively. We are interested in estimating the treatment effect $\phi := p_2 - p_1$.*

*First, we can derive the MLE estimator of the parameter $p = \begin{pmatrix} p_1 \\ p_2 \end{pmatrix} \in \mathbb{R}^2$. The likelihood is*

$$\ell(X_1, \ldots, X_n, Y_1, \ldots, Y_n; p_1, p_2) = p_1^{\sum_{i=1}^n X_i}(1 - p_1)^{n - \sum_{i=1}^n X_i} p_2^{\sum_{i=1}^n Y_i}(1 - p_2)^{n - \sum_{i=1}^n Y_i}$$

*and it is maximized for $\widehat{p}_n = (\widehat{p}_{1,n}, \widehat{p}_{2,n})^\top$ where $\widehat{p}_{1,n} = \frac{1}{n}\sum_{i=1}^n X_i$ and $\widehat{p}_{2,n} = \frac{1}{n}\sum_{i=1}^n Y_i$. Each of the estimator $\widehat{p}_{i,n}$ for $i = 1, 2$ is an unbiased estimator of $p_i$ and the central limit theorem tells us that $\sqrt{n}(\widehat{p}_{i,n} - p_i) \rightsquigarrow \mathcal{N}(0, p_i(1 - p_i))$. As the estimators $\widehat{p}_{1,n}$ and $\widehat{p}_{2,n}$ are independent, we get that*

$$\sqrt{n}\left(\widehat{p}_n - p\right) \sim \mathcal{N}\left(0, \Sigma\right) \quad \text{with } \Sigma = \begin{pmatrix} p_1(1 - p_1) & 0 \\ 0 & p_2(1 - p_2) \end{pmatrix}$$

*As a treatment effect estimator, we propose $\widehat{\phi}_n = \widehat{p}_{2,n} - \widehat{p}_{1,n}$. It can be written $\widehat{\phi}_n = g(\widehat{p}_n)$ here $g : \mathbb{R}^2 \to \mathbb{R}$ is a simple linear function $g(p_1, p_2) = p_2 - p_1$ whose gradient is $\nabla g(p) = (-1, 1)^\top$. Using the multi-variate Delta method, we get that*

$$\sqrt{n}\left(\widehat{\phi}_n - \phi\right) \rightsquigarrow \mathcal{N}\left(0, p_1(1 - p_1) + p_2(1 - p_2)\right)$$

*Using further Slutsky lemma yields*

$$\frac{\widehat{\phi}_n - \phi}{\sqrt{\frac{\widehat{p}_{1,n}(1 - \widehat{p}_{1,n})}{n} + \frac{\widehat{p}_{2,n}(1 - \widehat{p}_{2,n})}{n}}} \rightsquigarrow \mathcal{N}(0,1) \ .$$

---

[1] A slightly more complicated proof can also be given when $g$ is not continuous, see e.g. [Rivoirard and Stoltz, 2009]

## 2.4 Asymptotic properties of the Maximum Likelihood Estimator

Given an iid sample $X_1, \ldots, X_n \sim P_\theta$, we recall that the maximum likelihood estimator of the parameter $\theta$ is defined

$$\widehat{\theta}_n \in \underset{\theta \in \Theta}{\operatorname{argmax}} \, L(X_1, \ldots, X_n; \theta).$$

Despite its implicit definition, as the maximizer of some function, we will see that this estimator enjoys strong asymptotic performance guarantees, when the model satisfies some assumptions. In particular, we will assume that the model is identifiable, that is

$$\forall (\theta, \theta') \in \Theta, P_\theta = P_{\theta'} \quad \text{if and only if} \quad \theta = \theta'.$$

### 2.4.1 Rationale

First, let us try to understand why the MLE is a good estimator. Let us denote by $\theta_\star$ the true parameter from which the data is generated. The maximum likelihood can be rewritten as follows, introducing artificially the likelihood under $\theta_\star$. Indeed, one can write

$$
\begin{aligned}
\widehat{\theta}_n &\in \underset{\theta \in \Theta}{\operatorname{argmax}} \, \frac{f_\theta(X_1) \ldots f_\theta(X_n)}{f_{\theta_\star}(X_1) \ldots f_{\theta_\star}(X_n)} \\
\widehat{\theta}_n &\in \underset{\theta \in \Theta}{\operatorname{argmin}} \, \frac{f_{\theta_\star}(X_1) \ldots f_{\theta_\star}(X_n)}{f_\theta(X_1) \ldots f_\theta(X_n)} \\
\widehat{\theta}_n &\in \underset{\theta \in \Theta}{\operatorname{argmin}} \, \log \frac{f_{\theta_\star}(X_1) \ldots f_{\theta_\star}(X_n)}{f_\theta(X_1) \ldots f_\theta(X_n)} \\
\widehat{\theta}_n &\in \underset{\theta \in \Theta}{\operatorname{argmin}} \, \sum_{i=1}^n \log \left( \frac{f_{\theta_\star}(X_i)}{f_\theta(X_i)} \right) \\
\widehat{\theta}_n &\in \underset{\theta \in \Theta}{\operatorname{argmin}} \, \frac{1}{n} \sum_{i=1}^n \log \left( \frac{f_{\theta_\star}(X_i)}{f_\theta(X_i)} \right)
\end{aligned}
$$

Hence, $\widehat{\theta}_n$ can be rewritten as the minimizer of some empirical average. Introducing the notation

$$M_n(\theta, \theta_\star) = \frac{1}{n} \sum_{i=1}^n \log \left( \frac{f_{\theta_\star}(X_i)}{f_\theta(X_i)} \right)$$

we know by the law of large number that, for all $\theta \in \Theta$,

$$M_n(\theta, \theta_\star) \xrightarrow{P} \mathbb{E}_{\theta_\star} \left[ \log \left( \frac{f_{\theta_\star}(X_1)}{f_\theta(X_1)} \right) \right] = \mathrm{KL} \left( P_{\theta_\star}, P_\theta \right) \tag{2.1}$$

where $\mathrm{KL}(P, P')$ is the KL divergence between distributions, introduced in Definition 1.33. The KL divergence is not a distance but it still satisfies the following important property: $\mathrm{KL}(P, P') = 0$ if an only if $P = P'$. In particular, $\mathrm{KL}(P_{\theta_\star}, P_\theta) = 0$ if and only if $P_{\theta_\star} = P_\theta$, i.e. $\theta = \theta_\star$ as the model is identifiable. Thus we have

$$\underset{\theta \in \Theta}{\operatorname{argmin}} \, \mathrm{KL} \left( P_{\theta_\star}, P_\theta \right) = \theta_\star.$$

Hence, our hope is to prove that, under the model $\mathbb{P}_{\theta_\star}$,

$$\widehat{\theta}_n = \underset{\theta \in \Theta}{\operatorname{argmin}} \, M_n(\theta, \theta_\star) \xrightarrow{P} \underset{\theta \in \Theta}{\operatorname{argmin}} \, \mathrm{KL} \left( P_{\theta_\star}, P_\theta \right) = \theta_\star$$

This will require slightly more sophisticated arguments than the convergence of the objective function to minimize given in (2.1). We present them in the next section for more general $M$-estimators, that are also expressed as minimizer of empirical averages.

### 2.4.2   Consistency of M-estimators

A M-estimator is any estimator defined as a minimizer of some empirical average:

$$\widehat{\theta}_n = \operatorname*{argmin}_{\theta \in \Theta} M_n(\theta) \quad \text{with} \quad M_n(\theta) = \frac{1}{n} \sum_{i=1}^{n} m_\theta(X_i)$$

Letting $M(\theta) = \mathbb{E}[m_\theta(X_1)]$, if this expectation is finite we have $M_n(\theta) \xrightarrow{P} M(\theta)$ for all $\theta \in \Theta$, and we hope that $\widehat{\theta}_n$, a minimizer of $M_n(\theta)$, converges to $\theta_0 = \operatorname*{argmin}_{\theta \in \Theta} M(\theta)$.

**Example 2.16.** *In supervised learning, we observe iid pairs of the form $(X_i, Y_i)$ coming from some unknown distribution $P$ on $\mathcal{X} \times \mathcal{Y}$ where $\mathcal{X}$ is the feature space, often $\mathbb{R}^d$ and label space which is either finite (classification) or continuous (regression). The goal is to produce a predictor $\widehat{f}_n : \mathcal{X} \to \mathcal{Y}$, which is a data-dependent function mapping the feature to the label. Due to the generic empirical risk minimization principle, many predictor can be expressed as $M$-estimators.*

*Given a class of function $\mathcal{F}$, and some loss function $L : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+$, we can define*

$$\widehat{f}_n \in \operatorname*{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} L(f(X_i), Y_i).$$

*In this general (non-parametric) setting, the "parameter" is a function $f$ (possible predictor), and we have $m_f((X_i, Y_i)) = L(f(X_i), Y_i)$. We hope that $\widehat{f}_n$ converges to $f_0 \in \operatorname*{argmin}_{f \in \mathcal{F}} M(f)$ where $M(f) = \mathbb{E}_{(X,Y) \sim P}[L(f(X), Y)]$, that is to a predictor that minimizes the risk associated to the loss function $L$.*

*Sometimes, the class of function $\mathcal{F}$ can be described by a small set of parameters (e.g. a set of linear functions) and the regressor obtained by ridge regression can be defined as $\widehat{f}_n(x) = \widehat{\theta}_n^\top x$ for $x \in \mathbb{R}^d$ where*

$$\widehat{\theta}_n = \operatorname*{argmin}_{\theta \in \mathbb{R}^d : \|\theta\|_2 \leq C} \frac{1}{n} \sum_{i=1}^{n} (Y_i - \theta^\top X_i)^2 .$$

*In this case we hope that $\widehat{\theta}_n$ is close to $\theta_0 \in \operatorname*{argmin}_{\theta \in \mathbb{R}^d, \|\theta\| \leq C} \mathbb{E}_{(X,Y) \sim P}[(Y - \theta^\top X)^2]$.*

To establish the convergence of $\widehat{\theta}_n$ to $\theta_0$, we need two properties. The first one is a property of the minimizer $\theta_0$, which has to be a strict local minima, and the second is about the convergence from $M_n(\theta)$ to $M(\theta)$, which needs to be uniform.

**Theorem 2.17.** *Let $\widehat{\theta}_n = \operatorname*{argmin}_{\theta \in \Theta} M_n(\theta)$ and $\theta_0 = \operatorname*{argmin}_{\theta \in \Theta} M(\theta)$. For $\Theta \subseteq \mathbb{R}^d$, let $d$ a distance on $\mathbb{R}^d$. Assume that the following two properties hold:*

   *1. For all $\varepsilon > 0$, $\sup_{d(\theta, \theta_0) \geq \varepsilon} M(\theta) > M(\theta_0)$.*

   *2. $\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{P} 0$.*

*Then $\widehat{\theta}_n \xrightarrow{P} \theta_0$.*

*Proof.* From assumption 1., for every $\varepsilon > 0$, there exists $\eta_\varepsilon$ such $d(\theta, \theta_0) \geq \varepsilon$ implies that $M(\theta) \geq M(\theta_0) + \eta_\varepsilon$. One can write

$$
\begin{aligned}
\mathbb{P}\left(d(\widehat{\theta}_n, \theta_0) \geq \varepsilon\right) &\leq \mathbb{P}\left(M(\widehat{\theta}_n) \geq M(\theta_0) + \eta_\varepsilon\right) \\
&= \mathbb{P}\left(\eta_\varepsilon \leq M(\widehat{\theta}_n) - M(\theta_0)\right) \\
&= \mathbb{P}\left(\eta_\varepsilon \leq M(\widehat{\theta}_n) - M_n(\widehat{\theta}_n) + M_n(\widehat{\theta}_n) - M_n(\theta_0) + M_n(\theta_0) - M(\theta_0)\right)
\end{aligned}
$$

As $\widehat{\theta}_n$ is a minimizer of $M_n$, we have $M_n(\widehat{\theta}_n) - M_n(\theta_0) \leq 0$ and

$$
\begin{aligned}
\mathbb{P}\left(d(\widehat{\theta}_n, \theta_0) \geq \varepsilon\right) &\leq \mathbb{P}\left(\eta_\varepsilon \leq M(\widehat{\theta}_n) - M_n(\widehat{\theta}_n) + M_n(\theta_0) - M(\theta_0)\right) \\
&\leq \mathbb{P}\left(\eta_\varepsilon \leq 2 \sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)|\right) \\
&= \mathbb{P}\left(\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \geq \frac{\eta_\varepsilon}{2}\right),
\end{aligned}
$$

and the right-hand side tends to zero by assumption 2., which concludes the proof.

$\square$

**Remark 2.18.** *Consistency also holds if $\widehat{\theta}_n$ is not an exact minimizer of $M_n(\theta)$ (which can be hard to compute in some practical cases), as long as its approximation error converges to zero (in probability). A sufficient condition to obtain consistency for an approximate minimizer is to further assume that*

$$
M_n(\widehat{\theta}_n) \leq M_n(\theta_0) + E_n
$$

*for some random variable $E_n \xrightarrow{P} 0$.*

**Application to the MLE estimator**    Using Theorem 2.17, we can propose some sufficient condition for the MLE to be a consistent estimator of $\theta_\star$ when $(X_1, \ldots, X_n) \sim P_{\theta_\star}$.

**Theorem 2.19** (Consistency of the MLE)**.** *Assume that the model $\mathcal{M} = \{f_\theta, \theta \in \Theta\}$ satisfies the following properties:*

1. *$\mathcal{M}$ is identifiable, i.e., $f_\theta = f_{\theta'}$ implies $\theta = \theta'$ for all $(\theta, \theta') \in \Theta$.*

2. *$\Theta$ is compact and for all $x \in \mathcal{X}$, $\theta \mapsto f_\theta(x)$ is continuous.*

3. *For all $\theta \in \Theta$, $\mathbb{E}_\theta\left[\sup_{\theta' \in \Theta} |\log f_{\theta'}(X_1)|\right] < \infty$.*

*Then for all $\theta_\star \in \Theta$, the MLE estimator built from a $n$ sample $X_1, \ldots, X_n \sim f_{\theta_\star}$ satisfies $\widehat{\theta}_n \xrightarrow{P} \theta_\star$ (where the convergence is under the model $\mathbb{P}_{\theta_\star}$).*

     For theses assumption to be satisfied in simple models such as Bernoulli and Gaussian, we would need to restrict the set of possible values for the means (to $[p_0, 1 - p_0]$ for $p_0 > 0$ in the Bernoulli case, or to some bounded interval $[a, b]$ in the Gaussian case). But in these two cases, the consistency of the MLE (which coincides with the empirical means) can already easily be established directly using the law of large number. Still a result such as Theorem 2.19 provide some generic guarantees for the MLE in potentially more complex models, under some restriction on the parameter space.

### 2.4.3    Asymptotic normality of the MLE estimator

Under stronger assumptions, it is also possible to further exhibit the limiting distribution of the MLE estimator. We start by presenting the result and a sketch of proof for the estimation of a one-dimensional parameter $\theta \in \mathbb{R}$. Given a $n$ sample $X_1, \ldots, X_n \sim P_\theta$, we recall that $I(\theta)$ denotes the Fisher information obtained from one sample $X_1$.

**Theorem 2.20.** *Let $\widehat{\theta}_n$ be the MLE of a parameter $\theta \in \mathbb{R}$ computed on a $n$ sample $X_1, \ldots, X_n \sim P_\theta$. If $\widehat{\theta}_n$ is consistent and if the model is regular (according to Definition 1.21) then if the Fisher information satisfies $I(\theta) > 0$, $\sqrt{n}(\widehat{\theta}_n - \theta)$ converges in distribution under $\mathbb{P}_\theta$ towards a Gaussian distribution:*

$$\sqrt{n}(\widehat{\theta}_n - \theta) \rightsquigarrow \mathcal{N}\left(0, \frac{1}{I(\theta)}\right) .$$

*Proof.* Let $\ell(\theta) = \log f_\theta(X_1, \ldots, X_n)$ be a simplified notation for the log-likelihood, and let $\ell'(\theta)$ be its derivative (in $\theta$). As a minimizer of $\ell$, the MLE estimator satisfies $\ell'(\widehat{\theta}_n) = 0$. Using a Taylor expansion of $\ell'$ in the true parameter $\theta$, we can write

$$0 = \ell'(\widehat{\theta}_n) = \ell'(\theta) + (\widehat{\theta}_n - \theta)\ell''(\widetilde{\theta}_n)$$

for some $\widetilde{\theta}_n \in (\theta, \widehat{\theta}_n)$ (or $(\widehat{\theta}_n, \theta)$). Hence, one can write

$$
\begin{aligned}
\widehat{\theta}_n - \theta &= -\frac{\ell'(\theta)}{\ell''(\widetilde{\theta}_n)} \\
\sqrt{n}\left(\widehat{\theta}_n - \theta\right) &= \frac{\frac{1}{\sqrt{n}}\ell'(\theta)}{-\frac{1}{n}\ell''(\widetilde{\theta}_n)}
\end{aligned}
\tag{2.2}
$$

The numerator in (2.2) can be written

$$\frac{1}{\sqrt{n}}\ell'(\theta) = \sqrt{n}\left(\frac{1}{n}\sum_{i=1}^{n} s(X_i, \theta)\right) .$$

Using that under a regular model $\mathbb{E}_\theta\left[s(X_1; \theta)\right] = 0$ and $\operatorname{Var}_\theta\left[s(X_1; \theta)\right] = I(\theta)$, one gets using the Central Limit Theorem that

$$\frac{1}{\sqrt{n}}\ell'(\theta) \rightsquigarrow \mathcal{N}\left(0, I(\theta)\right)$$

From the consistency of $\widehat{\theta}_n$ we know that $\widehat{\theta}_n \xrightarrow{P} \theta$, from which we deduce that $\widetilde{\theta}_n \xrightarrow{P} \theta$. The denominator of (2.2) can be written (this is the part where the proof becomes approximately correct)

$$-\frac{1}{n}\ell''(\widetilde{\theta}_n) = \frac{1}{n}\sum_{i=1}^{n} -\left(\frac{\partial^2 \log f_\theta(X_i)}{\partial^2 \theta}\right)_{\widetilde{\theta}_n} \simeq \frac{1}{n}\sum_{i=1}^{n} -\left(\frac{\partial^2 \log f_\theta(X_i)}{\partial^2 \theta}\right)_{\theta}$$

By the law of large numbers, under the model $\mathbb{P}_\theta$, this empirical average converges in probability to $\mathbb{E}_\theta\left[-\frac{\partial^2 \log f_\theta(X_1)}{\partial^2 \theta}\right]$ which is equal to the Fisher information $I(\theta)$ (using Lemma 1.26).

As $I(\theta) \neq 0$, we can use Slutsky's lemma to get that

$$\sqrt{n}\left(\widehat{\theta}_n - \theta\right) \sim \frac{1}{I(\theta)}\mathcal{N}\left(0, I(\theta)\right) = \mathcal{N}\left(0, \frac{1}{I(\theta)}\right) .$$

$\square$

A useful consequence of Theorem 2.17 is that it allows us to build asymptotic confidence regions around the MLE estimator, by replacing the (unknown) quantity $I(\theta)$ by its empirical version $I(\widehat{\theta}_n)$.

**Corollary 2.21.** *Under the assumptions of Theorem 2.20, if $\theta \mapsto I(\theta)$ is continuous in $\theta$ then under the model $\mathbb{P}_\theta$,*

$$\sqrt{nI(\widehat{\theta}_n)}\left(\widehat{\theta}_n - \theta\right) \rightsquigarrow \mathcal{N}(0,1) \,.$$

*Proof.* Using the continuous mapping lemma and the consistency of $\widehat{\theta}_n$ yields that, under $\mathbb{P}_\theta$, $I(\widehat{\theta}_n) \xrightarrow{P} I(\theta)$. From Theorem 2.20 we have that under $\mathbb{P}_\theta$, $\sqrt{nI(\theta)}\left(\widehat{\theta}_n - \theta\right) \rightsquigarrow \mathcal{N}(0,1)$. Using Slutsky's lemma,

$$\sqrt{nI(\widehat{\theta}_n)}\left(\widehat{\theta}_n - \theta\right) = \sqrt{\frac{I(\widehat{\theta}_n)}{I(\theta)}} \times \sqrt{nI(\theta)}\left(\widehat{\theta}_n - \theta\right)$$

converges in distribution to $\mathcal{N}(0,1)$.

$\square$

Hence, if our model is regular enough, we can build asymptotic confidence intervals of level $1 - \alpha$ around the MLE estimator (and resulting tests, see the next chapter) that are of the form

$$\left[\widehat{\theta}_n - \sqrt{\frac{1}{nI(\widehat{\theta}_n)}}q_{\alpha/2}; \widehat{\theta}_n - \sqrt{\frac{1}{nI(\widehat{\theta}_n)}}q_{\alpha/2}\right]$$

where $q_\alpha$ is such that $\mathbb{P}_{Z \sim \mathcal{N}(0,1)}\left(Z \le q_\alpha\right) = 1 - \alpha$. We provide an example below.

**Example 2.22.** *Consider the Poisson model $X_1, \ldots, X_n \overset{iid}{\sim} \mathcal{P}(\lambda)$ where we recall that*

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

*for all $k \in \mathbb{N}$. The MLE is given by $\widehat{\lambda}_n = \frac{1}{n}\sum_{i=1}^n X_i$ and the Fisher information (of one sample) satisfies*

$$I(\lambda) = \mathbb{E}_\lambda\left[-\frac{\partial^2}{\partial^2\lambda}\log f_\lambda(X_1)\right]$$

*We have*

$$\begin{aligned}
\frac{\partial \log f_\lambda(X_1)}{\partial\lambda} &= \frac{X_1}{\lambda} - 1 \\
\frac{\partial \log f_\lambda(X_1)}{\partial\lambda} &= -\frac{X_1}{\lambda^2}
\end{aligned}$$

*hence $I(\lambda) = \mathbb{E}_\lambda\left[\frac{X_1}{\lambda^2}\right] = \frac{1}{\lambda}$ as the mean of a Poisson distribution with parameter $\lambda$ is $\lambda$. Applying Corollary 2.21 yields that, under $\mathbb{P}_\lambda$,*

$$\sqrt{\frac{n}{\widehat{\lambda}_n}}\left(\widehat{\lambda}_n - \lambda\right) \rightsquigarrow \mathcal{N}(0,1) \,.$$

*Now, let's use this information to build an asymptotic confidence interval on $\lambda$. We have that*

$$\mathbb{P}_\lambda \left( -q_{\alpha/2} \le \sqrt{\frac{n}{\widehat{\lambda}_n}} \left( \widehat{\lambda_n} - \lambda \right) \le q_{\alpha/2} \right) \xrightarrow[n\to\infty]{} \mathbb{P}_{Z\sim\mathcal{N}(0,1)} \left( -q_{\alpha/2} \le Z \le q_{\alpha/2} \right)$$

$$= \quad \mathbb{P}(Z \le q_{\alpha/2}) - \mathbb{P}\left( Z \le -q_{\alpha/2} \right)$$

$$= \quad \mathbb{P}(Z \le q_{\alpha/2}) - \mathbb{P}\left( Z > q_{\alpha/2} \right)$$

$$= \quad 1 - \frac{\alpha}{2} - \frac{\alpha}{2}$$

$$= \quad 1 - \alpha$$

*Putting $\lambda$ in the center of the interval, we have*

$$\mathbb{P}_\lambda \left( \widehat{\lambda}_n - \sqrt{\frac{\widehat{\lambda}_n}{n}} q_{\alpha/2} \le \lambda \le \widehat{\lambda}_n + \sqrt{\frac{\widehat{\lambda}_n}{n}} q_{\alpha/2} \right) \xrightarrow[n\to\infty]{} 1 - \alpha$$

*which provides an asymptotic confidence interval of level $1 - \alpha$.*

**Extensions of Theorem 2.20**    First, this result is also true for the estimation of a multi-dimensional parameter $\theta \in \mathbb{R}^d$ using the MLE. Under similar assumptions, we obtain that under $\mathbb{P}_\theta$,

$$\sqrt{n}(\widehat{\theta}_n - \theta) \rightsquigarrow \mathcal{N}\left( 0, I(\theta)^{-1} \right)$$

but this time the Fisher information $I(\theta)$ is a $d \times d$ matrix, assumed to be invertible (see its definition in Section 1.4.2).

Then, while we presented the consistency results for the general family of M-estimators, we sticked to the MLE case for the asymptotic normality result. A counterpart of Theorem 2.20 also exists for M-estimators, see e.g. the book [Van der Vaart, 1998].

## 2.5   Asymptotic efficiency

In light of the Cramer-Rao lower bound given in Chapter 1, any estimator of a parameter $g(\theta) \in \mathbb{R}$ whose limit distribution satisfies

$$\widehat{g}_n \approx \mathcal{N}\left( g(\theta), \frac{(g'(\theta))^2}{nI(\theta)} \right)$$

is called *asymptotically efficient*. The reason is that, asymptotically, is is unbiased with a variance that is the minimal variance prescribed by the Cramer-Rao lower bound.

For estimating the parameter $\theta$, (under appropriate regularity conditions) the MLE is an example of asymptotically efficient estimator, as we just saw that it satisfies

$$\widehat{g}_n \approx \mathcal{N}\left( \theta, \frac{1}{nI(\theta)} \right) .$$

However, we can find examples of MLE that are not efficient. Take for instance the MLE estimator of the variance of the Gaussian distribution, which is biased. We shall see other examples in exercises.

# Chapter 3

# (Optimal) Testing

# Bibliography

[Rivoirard and Stoltz, 2009]  Rivoirard, V. and Stoltz, G. (2009). *Statistique en Action*. Vuibert.

[Van der Vaart, 1998]  Van der Vaart, A. (1998). *Asymptotic Statistics*. Cambridge University Press.

[Wasserman, 2004]  Wasserman, L. (2004). *All of Statistics, A Concise Course in Statistical Inference*. Springer.