

MS = MAXIMUM LIKELIHOOD ESTIMATION

Let $(\Omega, \mathcal{F}, \{P_\theta\}_{\theta \in \Theta})$ be a statistical model. Assume that n observations X_1, \dots, X_n are drawn from P_{θ^*} , for some $\theta^* \in \Theta$. The parameter θ^* is referred to as the true parameter, and the aim is to estimate it based solely on the observations X_1, \dots, X_n .

Given our family of suspects $\{P_\theta\}_{\theta \in \Theta}$, the goal is thus to select which member of the family is the most likely to have given rise to X_1, \dots, X_n . The approach discussed in this chapter selects the distribution in $\{P_\theta\}_{\theta \in \Theta}$ that matches the most the empirical distribution of X_1, \dots, X_n .

The discrepancy between the two distributions is measured by the Kullback-Leibler divergence, that we introduce in the next section.

I. MAXIMUM LIKELIHOOD ESTIMATION

I.1. Kullback-Leibler divergence.

The Kullback-Leibler (KL) divergence between two probability measures P_θ and $P_{\theta'}$ is defined by:

$$KL(P_\theta \parallel P_{\theta'}) = \sum_{x \in \Omega} p_\theta(x) \log \left\{ \frac{p_\theta(x)}{p_{\theta'}(x)} \right\} \quad \text{if } \Omega \text{ is at most countable,}$$

with $p_\theta(x) = P_\theta(X=x)$, $X(\omega) = \omega$

and

$$KL(P_\theta \parallel P_{\theta'}) = \int_{\Omega} f_\theta(x) \log \left\{ \frac{f_\theta(x)}{f_{\theta'}(x)} \right\} dx \quad \text{if } \Omega \text{ is continuous} \quad (2)$$

Assuming that the X s are absolutely continuous, with density f_θ : $P_\theta(X \in B) = \int_B f_\theta(x) dx$, $X(\omega) = \omega$

$$\hookrightarrow \text{Note that } KL(P_\theta \parallel P_{\theta'}) = \mathbb{E}_\theta \left\{ \log \frac{f_\theta(X)}{f_{\theta'}(X)} \right\},$$

and similarly in the discrete case.

x Properties of the KL divergence:

(i) KL divergence is *not* symmetric. In general, we have that $KL(P_\theta \parallel P_{\theta'}) \neq KL(P_{\theta'} \parallel P_\theta)$.

\Rightarrow Be careful in which order you write $P_\theta \neq P_{\theta'}$.

Consequence: $KL(\cdot \parallel \cdot)$ is *not* a metric.

(ii) $KL(P_\theta \parallel P_{\theta'}) \geq 0$, and equals zero if and only if $P_\theta = P_{\theta'}$.

Consequence: $KL(\cdot \parallel \cdot)$ can be used as a measure of discrepancy between two probability measures.

Positivity of the KL divergence follows from GIBB'S INEQUALITY: for any two densities f_θ and $f_{\theta'}$,

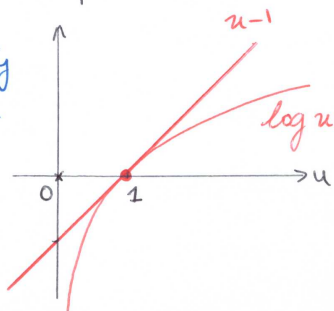
$$\int f_\theta(x) \log \{f_\theta(x)\} dx \geq \int f_\theta(x) \log \{f_{\theta'}(x)\} dx,$$

with equality iff $f_\theta(x) = f_{\theta'}(x)$ almost everywhere. The inequality holds in the discrete case as well,

$$\sum_x p_a(x) \log p_a(x) \geq \sum_x p_a(x) \log p_{a'}(x) \quad (3)$$

proof: we prove Gibb's inequality in the discrete case, the AC case being treated similarly.

Note that $\forall x > 0$, $\log x \leq x-1$, with equality for $x=1$.



Thus

$$\begin{aligned} \sum_x p_a(x) \log \frac{p_{a'}(x)}{p_a(x)} &\leq \sum_x p_a(x) \left\{ \frac{p_{a'}(x)}{p_a(x)} - 1 \right\} \\ &= \sum_x p_{a'}(x) - \sum_x p_a(x) = 0. \end{aligned}$$

Equality occurs at the value $p_{a'}(x)/p_a(x) = 1$, that is when the two distributions coincide.

Note: Alternatively, Gibb's inequality can be proved using Jensen's inequality. ■

Example: KL divergence between two Gaussian densities with common variance: $f_{\mu, \sigma^2}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$

$$\begin{aligned} \text{KL}(\mathbb{P}_{\mu, \sigma^2} \parallel \mathbb{P}_{\mu', \sigma^2}) &= \int f_{\mu, \sigma^2}(x) \log \left\{ \frac{f_{\mu, \sigma^2}(x)}{f_{\mu', \sigma^2}(x)} \right\} dx \\ &= \frac{1}{2\sigma^2} \int \left\{ (x-\mu')^2 - (x-\mu)^2 \right\} f_{\mu, \sigma^2}(x) dx \\ &= \frac{1}{2\sigma^2} \left\{ \mathbb{E}_{\mu, \sigma^2} (X-\mu')^2 - \mathbb{E}_{\mu, \sigma^2} (X-\mu)^2 \right\} \\ &= \frac{1}{2\sigma^2} \mathbb{E}_{\mu, \sigma^2} (-2X(\mu'-\mu) + (\mu')^2 - \mu^2) \end{aligned}$$

$$\begin{aligned} &= \frac{1}{2\sigma^2} \left\{ -2(\mathbb{E}_{\mu, \sigma^2} X)(\mu'-\mu) + (\mu')^2 - \mu^2 \right\} \quad (4) \\ &= \frac{(\mu-\mu')^2}{2\sigma^2} \end{aligned}$$

Example: Similarly, we compute the KL divergence between two multivariate Gaussian distributions $f_{\mu, \Sigma}(x)$ and $f_{\mu', \Sigma'}(x)$ as:

$$\text{KL}(\mathbb{P}_{\mu, \Sigma} \parallel \mathbb{P}_{\mu', \Sigma'}) = \int_{\mathbb{R}^d} f_{\mu, \Sigma}(x) \log \left\{ \frac{f_{\mu, \Sigma}(x)}{f_{\mu', \Sigma'}(x)} \right\} dx$$

$$= \frac{1}{2} \log \left\{ \frac{|\Sigma'|}{|\Sigma|} \right\} + \mathbb{E}_{\mu, \Sigma} \left\{ -\frac{1}{2}(X-\mu)^t \Sigma^{-1} (X-\mu) + \frac{1}{2}(X-\mu')^t (\Sigma')^{-1} (X-\mu') \right\}$$

$|\Sigma|$ = determinant of Σ .

$$= \frac{1}{2} \log \left\{ \frac{|\Sigma'|}{|\Sigma|} \right\} - \frac{1}{2} \mathbb{E}_{\mu, \Sigma} \left\{ \text{Tr}(\Sigma^{-1} (X-\mu)(X-\mu)^t) \right\} + \frac{1}{2} \mathbb{E}_{\mu, \Sigma} \left\{ (X-\mu')^t (\Sigma')^{-1} (X-\mu') \right\}$$

Use: $\mathbb{E}(X^t A X) = \mu^t A \mu + \text{Tr}(A \Sigma)$, for X with mean μ & covariance matrix Σ .

$$= \frac{1}{2} \log \left\{ \frac{|\Sigma'|}{|\Sigma|} \right\} - \frac{1}{2} \text{Tr} \left\{ \Sigma^{-1} \mathbb{E}_{\mu, \Sigma} (X-\mu)(X-\mu)^t \right\} + \frac{1}{2} \left[(\mu-\mu')^t (\Sigma')^{-1} (\mu-\mu') + \text{Tr}((\Sigma')^{-1} \Sigma) \right]$$

$$\begin{aligned} \text{Use: } \mathbb{E}_{\mu, \Sigma} (X X^t - 2X \mu^t + \mu \mu^t) &= \Sigma + \mu \mu^t - 2\mu \mu^t + \mu \mu^t = \Sigma. \end{aligned}$$

$$= \frac{1}{2} \log \left\{ \frac{|\Sigma'|}{|\Sigma|} \right\} - \frac{1}{2} \text{Tr}(\Sigma^{-1}\Sigma) = d$$

$$+ \frac{1}{2} \left[(\mu - \mu')^t (\Sigma')^{-1} (\mu - \mu') + \text{Tr} \{ (\Sigma')^{-1} \Sigma \} \right]$$

5

$$\Rightarrow \text{KL}(\mathbb{P}_{\mu, \Sigma} \parallel \mathbb{P}_{\mu', \Sigma'}) = \frac{1}{2} \left(\log \left\{ \frac{|\Sigma'|}{|\Sigma|} \right\} - d + \text{Tr} \{ (\Sigma')^{-1} \Sigma \} + (\mu - \mu')^t (\Sigma')^{-1} (\mu - \mu') \right)$$

I.2. The Maximum likelihood principle.

Recall: true model is \mathbb{P}_{θ^*} for some $\theta^* \in \Theta$. Consider $X \sim \mathbb{P}_{\theta^*}$.

Idea: find the value of θ which minimizes the KL divergence between \mathbb{P}_{θ^*} and \mathbb{P}_{θ} :

$$\text{KL}(\mathbb{P}_{\theta^*} \parallel \mathbb{P}_{\theta}) = \mathbb{E}_{\theta^*} \{ \log f_{\theta^*}(X) \} - \mathbb{E}_{\theta^*} \{ \log f_{\theta}(X) \}$$

Assuming that X is AC with density f_{θ^*} (the discrete case is treated similarly).

Minimization is performed with respect to θ . The first term $\mathbb{E}_{\theta^*} \{ \log f_{\theta^*}(X) \}$ is independent of θ .
 \Rightarrow Minimization of the KL divergence is equivalent to the...

$$\text{Maximization of } \mathbb{E}_{\theta^*} \{ \log f_{\theta}(X) \}.$$

in practice, the expectation $\mathbb{E}_{\theta^*} \{ \dots \}$ is estimated from a random sample X_1, \dots, X_n , where $X_i \sim \mathbb{P}_{\theta^*}$, iid.

The goal is now to maximize $\frac{1}{n} \sum_{i=1}^n \log f_{\theta}(X_i)$

6

Empirical mean; converges a.s. to $\mathbb{E}_{\theta^*} \{ \log f_{\theta}(X_i) \}$, SLLN.

$$\frac{1}{n} \log \left(\prod_{i=1}^n f_{\theta}(X_i) \right)$$

$$\Leftrightarrow \max_{\theta} \log \prod_{i=1}^n f_{\theta}(X_i)$$

The Maximum likelihood principle.

Note that $\prod_{i=1}^n f_{\theta}(X_i) = f_{\theta}(X) =$ joint density of X_1, \dots, X_n .
 The function $l(X; \theta) = \log f_{\theta}(X)$ is known as the log-likelihood function.

In the discrete case, replace f_{θ} by the probability masses.

Summary: The maximum likelihood principle returns the value of θ maximizing $\begin{cases} \log \prod_{i=1}^n f_{\theta}(x_i) & \text{in the AC case} \\ \log \prod_{i=1}^n p_{\theta}(x_i) & \text{in the discrete case,} \end{cases}$ having observed $X_1 = x_1, \dots, X_n = x_n$.

We write $\hat{\theta}_{ML} \in \underset{\theta \in \Theta}{\text{argmax}} l(X; \theta)$,

where

$$l(X; \theta) = \begin{cases} \sum_{i=1}^n f_{\theta}(X_i) & \text{in the AC case} \\ \sum_{i=1}^n p_{\theta}(X_i) & \text{in the discrete case.} \end{cases}$$

$\Rightarrow \hat{\theta}_{ML}$ is a Random Variable. We need to investigate its properties.

Remark: Denoting \hat{P}_n the empirical distribution placing a mass $1/n$ on each observation X_1, \dots, X_n , we see that the maximum likelihood estimator $\hat{\theta}_{ML}$ corresponds to the value of $\theta \in \Theta$ minimizing the KL divergence between \hat{P}_n and P_θ :

$$\forall \theta \in \Theta, \quad KL(\hat{P}_n \parallel P_{\hat{\theta}_{ML}}) \leq KL(\hat{P}_n \parallel P_\theta).$$

↑
since we replaced P_{θ^*} (true distribution) with \hat{P}_n on the top of page 6, when estimating $E_{\theta^*} \{ \dots \}$ using the empirical mean.

I.3. Alternatives to KL divergence.

• The KL divergence is a tool to quantify how far away two probability measures are from each other. However, there are many alternatives, and we introduce a selected few in this section.

• One potential drawback to the KL divergence is its asymmetry. Therefore, one often work with the HELLINGER METRIC

$$H(P_\theta, P_{\theta'}) = \begin{cases} \left[\int (\sqrt{f_\theta(x)} - \sqrt{f_{\theta'}(x)})^2 dx \right]^{1/2} & \text{if AC cont.} \\ \left[\sum_x (\sqrt{p_\theta(x)} - \sqrt{p_{\theta'}(x)})^2 \right]^{1/2} & \text{if discrete.} \end{cases}$$

It clearly satisfies $H(P_\theta, P_{\theta'}) = H(P_{\theta'}, P_\theta)$, and in fact satisfies all the properties of a metric.

Moreover, $H(P_\theta, P_{\theta'}) \leq \sqrt{KL(P_\theta \parallel P_{\theta'})}$ (8)

↳ an so a small KL divergence between two distributions implies that the two distributions are close as measured by the Hellinger metric.

$$\begin{aligned} \text{proof} = H^2(P_\theta, P_{\theta'}) &= \int (\sqrt{f_\theta(x)} - \sqrt{f_{\theta'}(x)})^2 dx \\ &= 2 - 2 \int \sqrt{f_\theta(x) f_{\theta'}(x)} dx \\ (*) &= 2 \left(1 - \int \sqrt{\frac{f_{\theta'}(x)}{f_\theta(x)}} f_\theta(x) dx \right) \\ &= 2 \left(1 - E_\theta \left\{ \sqrt{\frac{f_{\theta'}(x)}{f_\theta(x)}} \right\} \right) \\ &\stackrel{1-x \leq -\log x}{\leq} 2 \log \left(E_\theta \left\{ \sqrt{\frac{f_{\theta'}(x)}{f_\theta(x)}} \right\} \right) \\ &\stackrel{\text{Jensen!}}{\leq} -2 E_\theta \left\{ \log \sqrt{\frac{f_{\theta'}(x)}{f_\theta(x)}} \right\} \\ &= E_\theta \left\{ \log \left(\frac{f_\theta(x)}{f_{\theta'}(x)} \right) \right\} \\ &= KL(P_\theta \parallel P_{\theta'}). \end{aligned}$$

We can repeat the previous steps, replacing equality (*) with $2 \left(1 - \int \sqrt{\frac{f_\theta(x)}{f_{\theta'}(x)}} f_{\theta'}(x) dx \right)$, so that we obtain as well that $H^2(P_\theta, P_{\theta'}) \leq KL(P_{\theta'} \parallel P_\theta)$.

x Example: Consider again two univariate Gaussian densities with different means, and common variance. (9)

$$H(\mathbb{P}_{\mu, \sigma^2}, \mathbb{P}_{\mu', \sigma^2}) = 2 \left(1 - \int \sqrt{f_{\mu, \sigma^2}(x) f_{\mu', \sigma^2}(x)} dx \right)$$

Again, using
 $1-x \leq -\log x, x > 0$

$$\leq -2 \log \left\{ \int \sqrt{f_{\mu, \sigma^2}(x) f_{\mu', \sigma^2}(x)} dx \right\}$$

↑
 This quantity inside the log is known as the AFFINITY between $\mathbb{P}_{\mu, \sigma^2}$ and $\mathbb{P}_{\mu', \sigma^2}$.

$$= -2 \log \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \int \exp\left(-\frac{(x-\mu)^2}{4\sigma^2} - \frac{(x-\mu')^2}{4\sigma^2}\right) dx \right\}$$

$$= -2 \log \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} \int \exp\left(-\frac{(\mu-\mu')^2}{8\sigma^2}\right) \exp\left(-\frac{(x-\frac{\mu+\mu'}{2})^2}{2\sigma^2}\right) dx \right\}$$

$$= -2 \log \exp\left\{-\frac{(\mu-\mu')^2}{8\sigma^2}\right\}$$

$$= \frac{(\mu-\mu')^2}{4\sigma^2} \quad \blacksquare$$

• The TOTAL VARIATION DISTANCE is also a common choice. It is defined as

$$TV(\mathbb{P}_Q, \mathbb{P}_{Q'}) = \sup_{A \in \mathcal{F}} |\mathbb{P}_Q(A) - \mathbb{P}_{Q'}(A)|$$

↑ One can show that it is a distance / metric.

In the discrete case, the TV distance between \mathbb{P}_Q and $\mathbb{P}_{Q'}$ is a simple function of the probability mass functions $p_Q(x)$ and $p_{Q'}(x)$: (10)

$$TV(\mathbb{P}_Q, \mathbb{P}_{Q'}) = \frac{1}{2} \sum_x |p_Q(x) - p_{Q'}(x)|$$

And in the AC case, one can show that

$$TV(\mathbb{P}_Q, \mathbb{P}_{Q'}) = \frac{1}{2} \int |f_Q(x) - f_{Q'}(x)| dx$$

Moreover, it is known that $TV(\mathbb{P}_Q, \mathbb{P}_{Q'}) \leq \sqrt{\frac{1}{2} KL(\mathbb{P}_Q \| \mathbb{P}_{Q'})}$.

- Finally, note that the KL divergence can be made symmetrical, by considering for instance $KL(\mathbb{P}_Q \| \mathbb{P}_{Q'}) + KL(\mathbb{P}_{Q'} \| \mathbb{P}_Q)$,
 or
 $\min \{ KL(\mathbb{P}_Q \| \mathbb{P}_{Q'}), KL(\mathbb{P}_{Q'} \| \mathbb{P}_Q) \}$.

Alternatively, the JENSEN-SHANNON DIVERGENCE (JS) is sometimes used, and is a symmetrized & smoothed version of the KL divergence. It is defined as

$$JS(\mathbb{P}_Q \| \mathbb{P}_{Q'}) = \frac{1}{2} KL(\mathbb{P}_Q \| \mathbb{Q}_{Q,Q'}) + \frac{1}{2} KL(\mathbb{P}_{Q'} \| \mathbb{Q}_{Q,Q'})$$

where

$$\mathbb{Q}_{Q,Q'} := \frac{1}{2} (\mathbb{P}_Q + \mathbb{P}_{Q'})$$

- Other alternatives: Levy-Prokhorov metric (metrizes weak convergence), Wasserstein-Kantorovich metric, ...

II. PROPERTIES OF THE MLE

(11)

II.1. Introductory examples.

x Example: Let $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ iid. Put $\theta = (\mu, \sigma^2)$

The log-likelihood is

$$l(\underline{X}; \theta) = \log \left\{ \prod_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} \exp \left\{ -\frac{1}{2} \frac{(X_i - \mu)^2}{\sigma^2} \right\} \right\}$$
$$= -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2.$$

Setting the partial derivatives to zero yields:

$$\begin{cases} \frac{\partial l}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \mu) = \frac{n}{\sigma^2} (\bar{X} - \mu) = 0 \\ \frac{\partial l}{\partial \sigma} = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n (X_i - \mu)^2 = 0 \end{cases}$$

The first equation gives $\hat{\mu}_{ML} = \bar{X}$, that we can substitute into the second equation, to obtain

$$\hat{\sigma}_{ML}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Observations: • $\hat{\mu}_{ML} = \bar{X} \xrightarrow{a.s.} \mu$ as $n \rightarrow \infty$ (SLLN) consistency.

Moreover, $\hat{\mu}_{ML}$ is unbiased for μ .

• $\hat{\sigma}_{ML}^2 \xrightarrow{a.s.} \sigma^2$ as $n \rightarrow \infty$ (consistency). However, $\hat{\sigma}_{ML}^2$ is a biased estimator of σ^2 (we need to renormalize by $(n-1)$ to obtain an unbiased estimator). ■

x Example: $X_1, \dots, X_n \sim \mathcal{P}(\lambda)$ iid. The log-lik is (12)

$$l(\underline{X}, \lambda) = \log \left\{ \prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right\}$$
$$= -n\lambda - \sum_{i=1}^n \log(x_i!) + \log \lambda \sum_{i=1}^n x_i.$$

$$\frac{\partial l(\underline{X}; \lambda)}{\partial \lambda} = -n + \frac{1}{\lambda} \sum_{i=1}^n x_i = 0 \Rightarrow \hat{\lambda}_{ML} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Observations: $\hat{\lambda}_{ML}$ is unbiased for λ , consistent, and efficient (it has the smallest possible variance amongst the class of unbiased estimators) (see page 21/22 in MS = PARAMETRIC INFERENCE). ■

II.2. General properties of the MLE.

The Maximum likelihood principle returns the value $\hat{\theta}_{ML}$ of $\theta \in \Theta$ maximizing $\log \prod_{i=1}^n f_{\theta}(X_i) = \sum_{i=1}^n l(X_i; \theta)$ (see page 6)

As $n \rightarrow \infty$, expect the MLE $\hat{\theta}_{ML}$ to be close to the value of θ maximizing $\mathbb{E}_{\theta^*} \{ l(X_1; \theta) \} = \mathbb{E}_{\theta^*} \{ \log f_{\theta}(X_1) \}$, since $\frac{1}{n} \sum_{i=1}^n l(X_i; \theta) \xrightarrow{a.s.} \mathbb{E}_{\theta^*} \{ l(X_1; \theta) \}$ as $n \rightarrow \infty$.

Moreover, Gibb's inequality informs us that the value of θ maximizing $\mathbb{E}_{\theta^*} \{ \log f_{\theta}(X_1) \} = \int f_{\theta^*}(x) \log f_{\theta}(x) dx$ is precisely θ^* ; so that when n is large, one expects $\hat{\theta}_{ML}$ to be close (in some sense) to the true parameter θ^* .

x Example: Let $X_1, \dots, X_n \sim \mathcal{P}(\lambda^*)$ iid

Then $l(X_i; \theta) = -\theta + X_i \log \theta - \log(X_i!)$,

so that

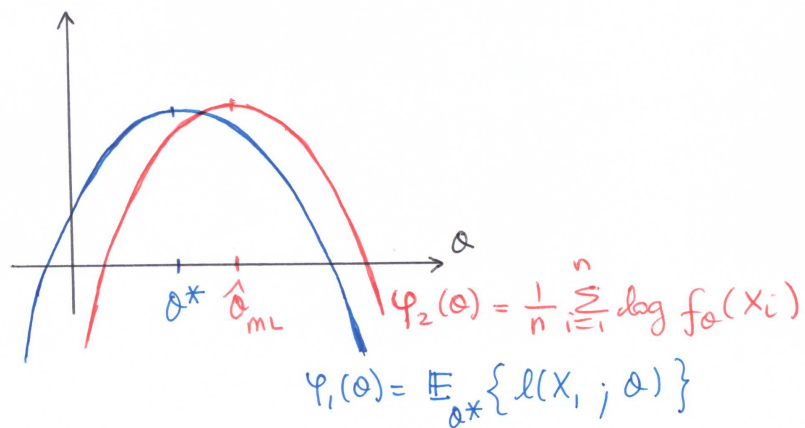
$E_{\lambda^*} \{l(X_i; \theta)\} = E_{\lambda^*} \{-\theta + X_i \log \theta - \log(X_i!)\}$
 $\uparrow = -\theta + \lambda^* \log \theta - E_{\lambda^*} \{\log(X_i!)\}$

The value of θ maximizing $\theta \mapsto E_{\lambda^*} \{l(X_i; \theta)\}$ satisfies $\frac{\partial E_{\lambda^*} \{l(X_i; \theta)\}}{\partial \theta} = -1 + \frac{\lambda^*}{\theta} = 0$

\Rightarrow The maximum occurs at $\theta = \lambda^* =$ the true parameter

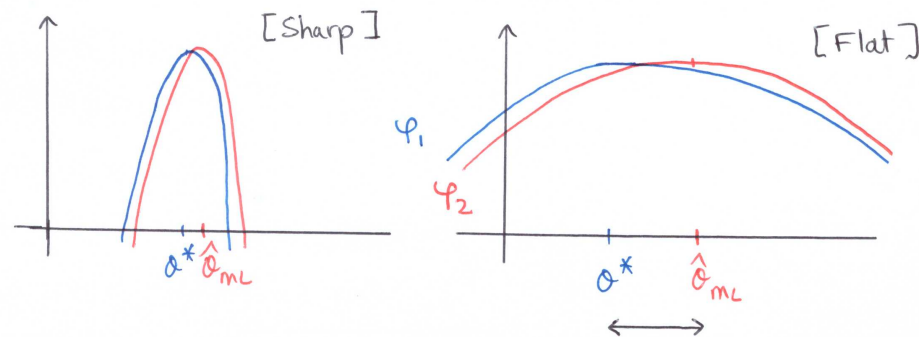
• What we know so far: expect $\hat{\theta}_{MLE}$ to be close to the true parameter $\theta^* =$ maximum point of the function $\varphi_1(\theta) = E_{\theta^*} l(X_i; \theta)$

\hookrightarrow Plot on a same graph $\varphi_1(\theta)$ and $\varphi_2(\theta) = \frac{1}{n} \sum_{i=1}^n l(X_i; \theta)$, where $l(\underline{X}; \theta) = \log \prod_{i=1}^n f_{\theta}(X_i)$.



• Observation: the sharper the peak of $\varphi_1(\theta)$, the closer should be the maxima of φ_1 and φ_2 .

Compare:



\hookrightarrow In both cases, φ_1 and φ_2 are close together $\forall \theta \in \Theta$. However, because of the flat peak of φ_1 , in the second case, the maxima are far from each other. \Rightarrow Estimation of θ^* is related to the curvature of $\varphi_1(\theta)$ at its maximum.

The curvature is defined in terms of the second order derivative, but is not equal to it.

Ex: the function $x \mapsto x^2$ has a constant second order derivative equal to 2, but its curvature is not constant: the graph of $x \mapsto x^2$ is "flatter" as we move away from the origin.

Properties of the MLE $\hat{\theta}_{MLE}$ should be related to the magnitude of $\varphi_1''(\theta)$ evaluated at $\theta = \theta^*$.

However, at the point where the derivative vanishes (i.e. at a local maximum or minimum for example), the curvature is given by $|f''(\cdot)|$ at that point.

Provided we can exchange $\mathbb{E}_{\theta^*} \{ \dots \}$ and derivatives (15)
 (in other words, under technical conditions), the second order partial derivative of φ_i is given by

$$\varphi_i''(\theta) = \mathbb{E}_{\theta^*} \{ l''(X; \theta) \} = -I_1(\theta)$$

↑
 Fisher information for one observation
 (see p.20 in chapter MS: PARAMETRIC ESTIMATION)

⇒ Evaluated at $\theta = \theta^*$ = the true parameter, our intuition tells us that the higher $I_1(\theta^*)$, the better the properties of $\hat{\theta}_{ML}$ (hence the name Fisher "information") (since the higher the curvature at the maximum).
 And indeed, we have the following result:

Theorem: Let X_1, \dots, X_n iid $\sim P_{\theta^*}$ for some $\theta^* \in \Theta$.

- $\mathcal{P} = \{ P_{\theta} \}_{\theta \in \Theta}$ a family of suspects
- $\forall \theta \in \Theta$, the support of P_{θ} does not depend on θ , and θ^* is not on the boundary of Θ
- $I_1(\theta) = -\mathbb{E}_{\theta^*} \{ l''(X; \theta) \}$ is invertible in a neighborhood of θ^* , where

$$l(X; \theta) = \begin{cases} \log f_{\theta}(X) & \text{in the AC case} \\ \log p_{\theta}(X) & \text{in the discrete case} \end{cases}$$
- Some additional technical conditions.

Then $\hat{\theta}_{ML}$ satisfies:

(i) $\hat{\theta}_{ML} \xrightarrow{P} \theta^*$ as $n \rightarrow \infty$, with respect to P_{θ^*} .
 (consistency).

(ii) The MLE is asymptotically normal: (16)

$$n^{1/2} (\hat{\theta}_{ML} - \theta^*) \xrightarrow{d} \mathcal{N}(0, I_1^{-1}(\theta^*))$$

(iii) Convergence of moments:

$$\downarrow \mathbb{E}_{\theta^*} \hat{\theta}_{ML} = \theta^* + o(n^{-1/2}) \quad (\text{asympt. unbiased})$$

$$\downarrow \mathbb{E}_{\theta^*} (\hat{\theta}_{ML} - \theta^*)^2 = \frac{1 + o(1)}{n I_1(\theta^*)}$$

↑

The higher $I_1(\theta^*)$, the smaller the MSE.
 The theorem agrees with our intuition.
 Moreover, recall the Cramer-Rao lower bound for unbiased estimators: $1/I_n(\theta)$
 (page 21 in MS: PARAMETRIC ESTIMATION)
 ⇒ $\hat{\theta}_{ML}$ is asymptotically efficient.

proof: We admit consistency. We outline the proof for the asymptotic normal distribution, omitting the technical details.

Taylor expanding $l'(X; \theta^*)$ around the MLE yields:

$$l'(X; \theta^*) = \underbrace{l'(X; \hat{\theta}_{ML})}_{=0} + (\theta^* - \hat{\theta}_{ML}) l''(X; \tilde{\theta}),$$

where $\tilde{\theta}$ is a point between θ^* and $\hat{\theta}_{ML}$. Thus:

$$\sum_{i=1}^n l'(X_i; \theta^*) = -(\hat{\theta}_{ML} - \theta^*) \sum_{i=1}^n l''(X_i; \tilde{\theta})$$

$\times n^{-1/2}$ ↓

$$n^{-1/2} \sum_{i=1}^n l'(X_i; \theta^*) = \underbrace{n^{-1/2} (\hat{\theta}_{ML} - \theta^*)}_{\text{quantity of interest}} \underbrace{(-1) \sum_{i=1}^n l''(X_i; \tilde{\theta})}_{\text{II}}$$

I

$$\downarrow \textcircled{I} = n^{1/2} \sum_{i=1}^n l'(X_i; \theta^*) \xrightarrow{d} \mathcal{N}(0, \mathbb{I}_1(\theta^*)), \quad (17)$$

since $\mathbb{E}_{\theta^*} \{l'(X_i; \theta^*)\} = 0$

& $\mathbb{E}_{\theta^*} \{(l'(X_i; \theta^*))^2\} = \mathbb{I}_1(\theta^*)$

+ Central Limit Theorem.

$\downarrow \textcircled{II}$ Consistency of the MLE : $\hat{\theta}_{ML} \xrightarrow{P} \theta^*$.

Thus, $\tilde{\theta} \xrightarrow{P} \theta^*$, as $n \rightarrow \infty$.

The WLLN + Slutsky theorem gives

$$n^{-1} \sum_{i=1}^n l''(X_i; \tilde{\theta}) \xrightarrow{P} \mathbb{E}_{\theta^*} \{l''(X_i; \theta^*)\} = -\mathbb{I}_1(\theta^*)$$

Putting pieces together, we get $n^{1/2} (\hat{\theta}_{ML} - \theta^*) \xrightarrow{d} \mathcal{N}(0, \mathbb{I}_1^{-1}(\theta^*))$ as required. ■

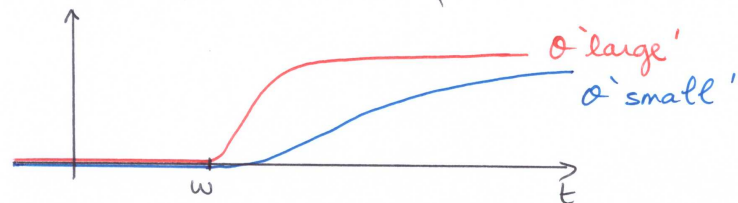
I.3. Illustration.

The time T of occurrence of cancer in a tissue, assumed random, can be modelled as the time of occurrence of cancer in one individual cell composing the tissue. Given that a tissue contains a large number of cells, asymptotic theory of extreme values indicate that possible distributions for T can be of three different types (Gumbel, '58, Statistics of Extremes). Following the work of (Pike, '66, A Method of Analysis of a Certain Class of Experiments in Carcinogenesis, Biometrics), we consider that the time of

appearance of cancer, in days, is modelled as $F_{\theta}(t) =$ (18)

$$P(T \leq t) = 1 - e^{-\theta(t-w)^k} \quad \text{for } t \geq w, \quad (*)$$

with $k > 0$, $\theta > 0$. The parameter θ has a huge effect on the distribution of T . For k fixed, we have



A large value of θ indicates that T is more likely to be small. The goal is to quantify the effect of a treatment on the time of appearance of cancer; i.e. on the value of θ .

Pike considers a population of mice, on which we are testing the effect of a pre-treatment on a toxic product. The time t (after the start of the experiment) at which the first effects of the toxic product were observed is recorded for each mouse. In this experiment, we assume k and w known, where
 $w \approx$ time it takes for the toxic product to be active
 $k \approx$ death rate of the mice, function of their age.
 Observations t_1, \dots, t_n ($n =$ number of mice) are assumed to be n independent observations of T_1, \dots, T_n , each following the distribution $(*)$, with density
 $f_{\theta}(t) = k \theta (t-w)^{k-1} e^{-\theta(t-w)^k} \mathbb{1}(t > w)$.

The joint distribution of T_1, \dots, T_n is (19)
 $f_{\Theta}(\underline{t}) = \Theta^n k^n \left(\prod_{i=1}^n (t_i - \omega)^{k-1} \mathbb{1}(t_i > \omega) \right) e^{-\Theta \sum_{i=1}^n (t_i - \omega)}$,

and the log-likelihood is

$$l(\underline{T}; \Theta) = n \log \Theta - \Theta \sum_{i=1}^n (T_i - \omega)^k + \text{something independent of } \Theta.$$

$$\hookrightarrow \frac{\partial l(\underline{T}; \Theta)}{\partial \Theta} = \frac{n}{\Theta} - \sum_{i=1}^n (T_i - \omega)^k,$$

so that the ML estimator of Θ is $\hat{\Theta}_{ML} = \frac{n}{\sum_{i=1}^n (T_i - \omega)^k}$

Fisher information associated with one observation is

$$I_1(\Theta) = -\mathbb{E}_{\Theta} \{ l''(T; \Theta) \} = \frac{1}{\Theta^2}.$$

$\Rightarrow \hat{\Theta}_{ML}$ is asymptotically normally distributed, and

$$n^{1/2} (\hat{\Theta}_{ML} - \Theta) \xrightarrow{d} \mathcal{N}(0, \Theta^2) \quad (**)$$

↑ We can use this result to construct an asymptotic confidence interval for Θ : replace the asymptotic variance Θ^2 by the consistent estimator $\hat{\Theta}_{ML}^2$. Slutsky theorem and **(**)** ensure that

$$\frac{n^{1/2} (\hat{\Theta}_{ML} - \Theta)}{\hat{\Theta}_{ML}} \xrightarrow{d} \mathcal{N}(0, 1) \text{ as } n \rightarrow \infty$$

↑
A pivotal statistic.

Select $z_{1-\alpha/2} = (1 - \frac{\alpha}{2})$ quantile of the standard normal distribution, so that (20)

$$\mathbb{P} \left(-z_{1-\frac{\alpha}{2}} \leq \frac{n^{1/2} (\hat{\Theta}_{ML} - \Theta)}{\hat{\Theta}_{ML}} \leq z_{1-\frac{\alpha}{2}} \right) \approx 1 - \alpha,$$

and we deduce that $\left[\hat{\Theta}_{ML} - z_{1-\frac{\alpha}{2}} n^{-1/2} \hat{\Theta}_{ML}, \hat{\Theta}_{ML} + z_{1-\frac{\alpha}{2}} n^{-1/2} \hat{\Theta}_{ML} \right]$ is an asymptotic CI for Θ , with nominal level $(1 - \alpha)$

• Numerical application: $\omega = 100$, $k = 3$, $n = 17$

$$\sum_{i=1}^n (t_i - 100)^3 = 33\,175\,533$$

observations =

143	164	188	188	190	192	206	209	213	216
220	227	230	234	246	265	304			

We obtain $\hat{\Theta}_{ML} = 5.124 \cdot 10^{-7}$

$$CI = [2.7 \cdot 10^{-7}, 7.6 \cdot 10^{-7}]$$

It turns out that some mice die before the end of the experiment, for reasons that are independent of the injection of the toxic product (such as disease, wounds, ...). These mice are removed from the experiment, and the time s of removal is recorded. For a given mouse, we do not observe the time of appearance t of the first effects directly, but $r = \min(s, t)$. The value s corresponds to the observation of a random variable S with distribution $G_s(s) = \mathbb{P}(S \leq s) = \int_{-\infty}^s g_s(u) du$. The functions G_s and g_s are unknown, and we are not trying

to estimate them. We observe $X := \mathbb{1}(T \leq S)$ as well, $\textcircled{21}$
 where $X=0$ if the mouse is removed from the experiment
 before the first effects of the toxic product appear, and
 $X=1$ otherwise. We write $p = P(X=1)$, so that
 $X \sim B(p)$.

Random sample collected is $(R_1, X_1), \dots, (R_n, X_n)$.
 \Rightarrow We need to derive the joint distribution of (R, X) .
 We have

$$\begin{aligned} \bullet P(R > r, X=1) &= P(S \geq T > r) \\ &= \int \mathbb{1}(s \geq t > r) g_S(s) f_\theta(t) ds dt \\ &= \int \mathbb{1}(t > r) f_\theta(t) \\ &\quad \times \left[\int \mathbb{1}(s \geq t) g_S(s) ds \right] dt \\ &= \int_r^{+\infty} f_\theta(t) (1 - G_S(t)) dt \end{aligned}$$

$$\begin{aligned} \bullet P(R > r, X=0) &= P(T > S > r) \\ &= \int \mathbb{1}(t > s > r) g_S(s) f_\theta(t) ds dt \\ &= \int \mathbb{1}(s > r) g_S(s) \\ &\quad \times \left[\int \mathbb{1}(t > s) f_\theta(t) dt \right] ds \\ &= \int_r^{+\infty} g_S(s) (1 - F_\theta(t)) dt \end{aligned}$$

Differentiating with respect to r , we obtain $\textcircled{22}$

$$\begin{aligned} f(r, x=1) &= k \theta (r-w)^{k-1} e^{-\theta(r-w)_+^k} \mathbb{1}(r > w) (1 - G_S(r)) \\ &= \theta e^{-\theta(r-w)_+^k} \underbrace{k (r-w)_+^{k-1} (1 - G_S(r))}_{\text{independent of } \theta} \\ &=: c(r, 1) \end{aligned}$$

where \uparrow
 $(u)_+ = \max(u, 0)$
 $= \theta c(r, 1) \exp(-\theta(r-w)_+^k)$

likewise,
 $f(r, x=0) = g(r) e^{-\theta(r-w)_+^k} \rightarrow c(r, 0) := g(r)$
 $= c(r, 0) \exp(-\theta(r-w)_+^k)$

Summarizing, $f(r, x) = c(r, x) \theta^x e^{-\theta(r-w)_+^k}$

The joint density of the random sample $(R_1, X_1) \dots (R_n, X_n)$
 is

$$\begin{aligned} f(r, x; \theta) &= \prod_{i=1}^n f(r_i, x_i) \\ &= \theta^{\sum x_i} e^{-\theta \sum (r_i - w)_+^k} \prod_{i=1}^n c(r_i, x_i), \end{aligned}$$

and the log-likelihood is

$$\begin{aligned} \tilde{\ell}(r, x; \theta) &= (\log \theta) \sum x_i - \theta \sum (r_i - w)_+^k \\ &\quad + \text{something independent of } \theta. \end{aligned}$$

$$\begin{aligned} \frac{\partial \tilde{\ell}(r, x; \theta)}{\partial \theta} &= \frac{1}{\theta} \sum x_i - \sum (r_i - w)_+^k \\ \Rightarrow \tilde{\theta}_{ML} &= \frac{\sum_{i=1}^n x_i}{\sum_{i=1}^n (R_i - w)_+^k} \end{aligned}$$

Fisher information, associated with one observation, (23)

is

$$\tilde{I}_1(\theta) = -\mathbb{E}_{\theta} \{ \tilde{\ell}''(R, X; \theta) \} = \frac{\mathbb{E}_{\theta} X_1}{\theta^2} = \frac{P}{\theta^2}.$$

⇒ Asymptotic properties of $\tilde{\theta}_{ML}$:

$$n^{1/2} (\tilde{\theta}_{ML} - \theta) \xrightarrow{d} \mathcal{N}\left(0, \frac{\theta^2}{P}\right)$$

As before, to construct a CI for θ , based on its limiting distribution, we need a consistent estimator of the asymptotic variance.

Put $N_n := \sum_{i=1}^n X_i$ = number of mice on which toxic effects were observed.

$$\text{Then } \left. \begin{array}{l} \frac{N_n}{n} \xrightarrow{a.s.} P \\ \tilde{\theta}_{ML} \xrightarrow{P} \theta \end{array} \right\} \frac{N_n}{n} \frac{1}{\tilde{\theta}_{ML}^2} \xrightarrow{P} \tilde{I}_1(\theta) \text{ i.e. consistent.}$$

Using Slutsky, we conclude that

$$\frac{n^{1/2} (\tilde{\theta}_{ML} - \theta)}{\tilde{\theta}_{ML} / \sqrt{N_n/n}} = \frac{N_n^{1/2} (\tilde{\theta}_{ML} - \theta)}{\tilde{\theta}_{ML}} \xrightarrow{d} \mathcal{N}(0, 1).$$

$$\Rightarrow \left[\tilde{\theta}_{ML} - z_{1-\frac{\alpha}{2}} N_n^{-1/2} \tilde{\theta}_{ML}, \tilde{\theta}_{ML} + z_{1-\frac{\alpha}{2}} N_n^{-1/2} \tilde{\theta}_{ML} \right]$$

is an asymptotic CI for θ , with nominal level $1-\alpha$.

• Numerical Application. Use data page 20, plus two extra observations 216, 244, corresponding to the time these two mice were removed from the experiment. We obtain $\tilde{\theta}_{ML} = 4.5 \cdot 10^{-7}$, and $[2.4 \cdot 10^{-7}, 6.6 \cdot 10^{-7}]$.