

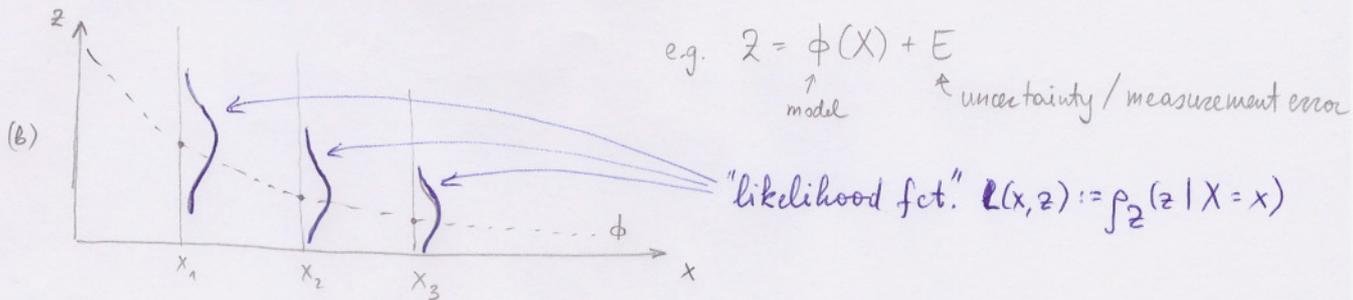
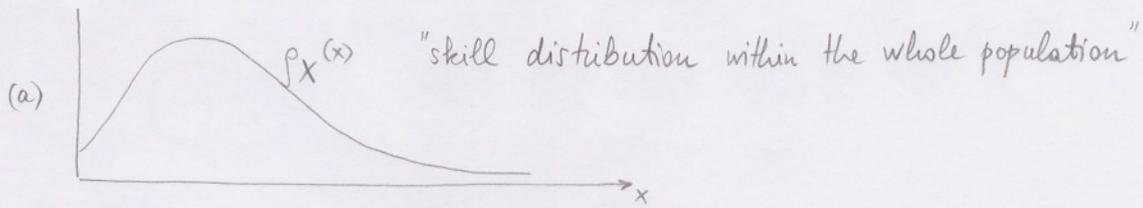
# Frequentist Consistency of Frequentist & Bayesian Methods

- (I) Introductory Example
- (II) Setup & Notation
- (III) What are we going to talk about?
- (IV) Frequentist Consistency of Frequentist Methods (MLE)
- (V) Asymptotic Normality of Frequentist Methods (MLE)
- (VI) Frequentist Consistency of Bayesian Methods

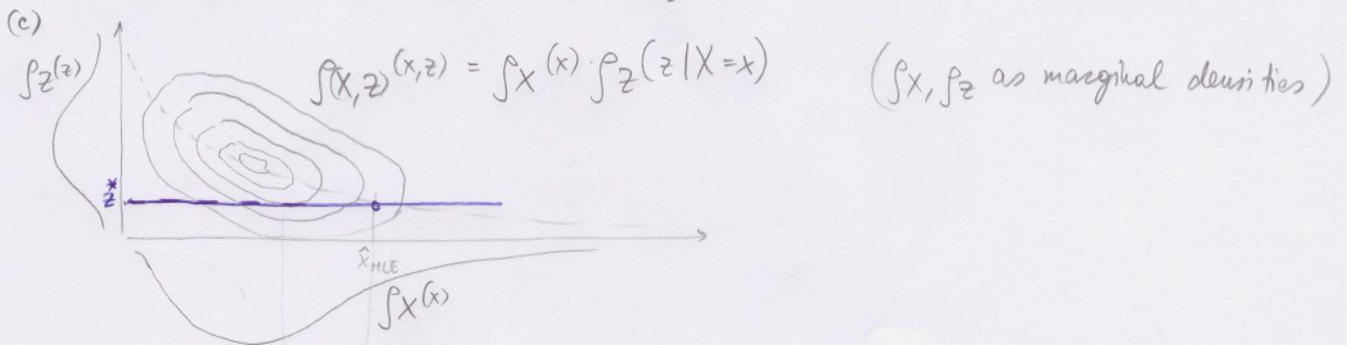
Sources: Tim Sullivan - Introduction to Uncertainty Quantification [S]  
Richard Nickel - Statistical Theory [N]

# (I) Introductory Example

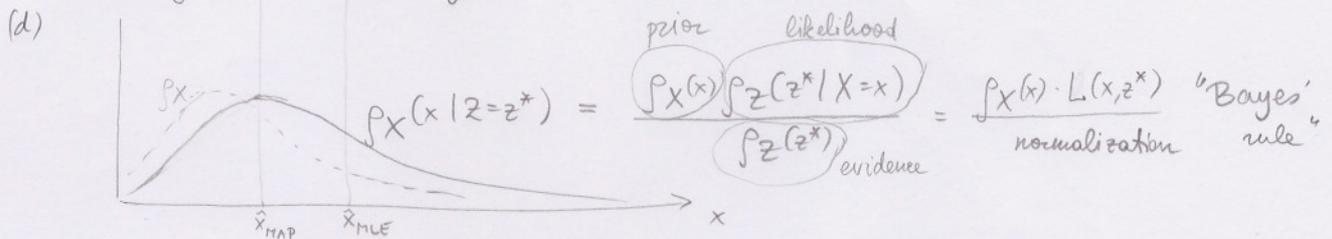
- $X$  = skiing skill of a person - "value of interest"
- $Z$  = time measurement on some racing track



(a) and (b) give the joint probability density of  $X$  and  $Z$ :



Restricting (c) to  $Z = z^*$  yields the posterior density  $f_X(x | Z = z^*)$



If we have many <sup>independent</sup> measurements for the same person,  $Z_1 = z_1^*, \dots, Z_M = z_M^*$ ,

Bayes' rule becomes:  $f_X(x | z_1 = z_1^*, \dots, z_M = z_M^*) = \frac{f_X(x) \cdot \prod_{m=1}^M f_Z(z_m^* | X=x)}{\text{norm.}} = \frac{f_X(x) \cdot f_{Z_1, \dots, Z_M}(z_1^*, \dots, z_M^* | X=x)}{\text{joint likelihood}}$

## (II) Setup & Notation (5)

(2)

- $X \in \mathbb{R}^d$  unknown parameter / value of interest
- $Z \in \mathbb{R}^n$  measurement with prob. density  $f_Z$
- Input: - prob. density  $f_X \in C^1(\mathbb{R}^d)$  of  $X$ 
  - likelihood fct.  $L(x, z) = f_Z(z|X=x)$
  - one or more realization(s)  $z_1 = z_1^*, \dots, z_M = z_M^*$
- Output: - Estimation  $\hat{x}_M$  of  $X$ 
  - its uncertainty: - [the distribution of] the error  $\hat{x}_M - X$
- We will always assume that there exists an (unknown) "true" parameter value  $x_T \in \text{supp}(f_X)$  which "induces" the measurements, i.e.

$$z_1, \dots, z_M \sim Z_T, \quad Z_T \text{ given by the density } f_Z(z|X=x_T)$$

Def: A parametric model  $\left( \underbrace{\{f_Z(\cdot|X=x) | x \in \mathbb{R}^d\}}_{\text{likelihood model}}, \underbrace{\mathbb{P}}_{\text{data-generating distribution}} \right)$  is said to be

"correctly or well-specified" if there exist an  $x_T \in \mathbb{R}^d$  s.t.  $\mathbb{P}$  has the prob. density  $f_Z(z|X=x_T)$ .

Def: The maximum likelihood estimator (MLE)  $\hat{x}_M$  based on  $z_1, \dots, z_M$  is given by

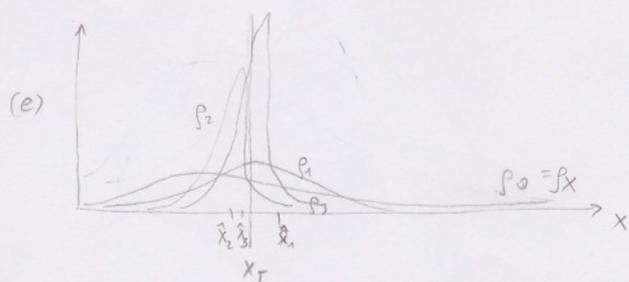
$$\hat{x}_M = \arg\max_x \prod_{m=1}^M L(x, z_m) = \arg\max_x \sum_{m=1}^M \log L(x, z_m) = \arg\min_x \sum_{m=1}^M q(x, z_m)$$

where  $q(x, z) := -\log L(x, z)$  is the negative log-likelihood

(III) What are we going to talk about?

Example 2(a): Same as introductory example, but with many measurements

Let  $f_M(x) := \int_X(x | z_1=z_1^*, \dots, z_M=z_M^*)$  and  $\hat{x}_M$  be the MLE based on  $z_1, \dots, z_M$

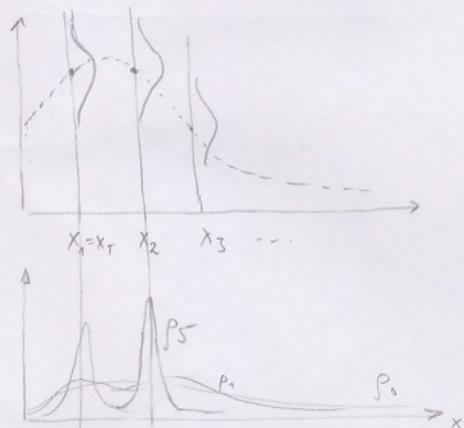


Natural Questions: What happens to  $\hat{x}_M$  and  $f_M$  for  $M \rightarrow \infty$ ?

- (i)  $\hat{x}_M \rightarrow x_T$ ? In what sense?
- (ii)  $f_M$  concentrates around  $x_T$ ? How exactly?
- (iii) What are sufficient conditions for (i) and (ii)?

Let's start with (iii)

Example 2(b): Same as 2(a), but with different "model"  $\phi$  (more general, different likelihood model  $\{L(x, \cdot) | x \in \mathbb{R}^d\}$ )



(i) and (ii) go wrong!

So, a necessary condition is the so-called identifiability assumption:

Def: A likelihood model  $\{L(x, \cdot) | x \in \mathbb{R}^d\}$  fulfills the identifiability assumption (IA),

$$\text{if } \forall x_1, x_2; x_1 \neq x_2 : \int_{\mathbb{Z}} (\cdot | X=x_1) \neq \int_{\mathbb{Z}} (\cdot | X=x_2) \text{ (in } L^1) \quad \text{(IA)}$$

Further criteria we are going to use:

(C1)  $L(x, z) > 0 \forall x, z$ ,  $L$  is continuous in  $x$  (and  $\int L(x, z) dz = 1 \forall x$ )

(C2)  $E[\sup_x |L(x, z_T)|] < \infty$

(C3)  $L(\cdot, z) \in C^2(\mathcal{U}) \forall z$  on some nbh.  $x_T \in \mathcal{U} \subset \mathbb{R}^d$ , s.t.

•  $E[D_x^2 q(x_T, z_T)] \in GL(d, \mathbb{R})$  and  $E[\|\partial_x q(x_T, z_T)\|^2] < \infty$

•  $\exists K = \mathcal{B}_{x_T}(\tau)$  compact ball around  $x_T$  with radius  $\tau > 0$  s.t.

$$E[\sup_{x \in K} |D_x^2 q(x, z_T)|] < \infty, \int \sup_{x \in K} \|\partial_x L(x, z)\| dz < \infty, \int \sup_{x \in K} \|D^2 L(x, z)\| dz < \infty$$

#### (IV) Frequentist Consistency of the MLE:

(4)

Thm: Assuming Setup (S), <sup>(IA)</sup> and (C1), (C2), the MLE  $\hat{x}_M$  based on  $Z_1, \dots, Z_M$  is consistent,  
[N, Thm 2] i.e.  $\hat{x}_M \xrightarrow{\mathbb{P}} x_T$   
[S, Thm 6.13]

Proof (sketch): (1)  $\sup_{x \in \mathbb{R}^d} \left| \underbrace{\mathbb{E}[q(x, Z_T)]}_{=: Q(x) \text{ (fct.)}} - \underbrace{\frac{1}{M} \sum_{m=1}^M q(x, Z_m)}_{=: Q_M(x) \text{ (z.v.)}} \right| \xrightarrow[M \rightarrow \infty]{\mathbb{P}\text{-a.s.}} 0$  "uniform law of large numbers"

[proof see [N., Prop 2]]

(2)  $x_T$  minimizes  $Q$  (uniquely!):  $x_T = \operatorname{argmin} Q(x)$

$$\forall x \neq x_T: Q(x_T) - Q(x) = - \int q(x_T, z) L(x_T, z) dz - \int q(x, z) L(x_T, z) dz = \int \log \left[ \frac{L(x_T, z)}{L(x, z)} \right] L(x_T, z) dz$$

strict Jensen ineq.  
(IA)  $\log \int \frac{L(x, z)}{L(x_T, z)} L(x_T, z) dz = \log 1 = 0$

$-D_{KL}(\pi_{(x_T, \cdot)} \| \pi_{(x, \cdot)}) > 0$   
(1A)

(3) Part (a) implies

$$\operatorname{argmin}_x \underbrace{Q_M(x)}_{\text{|| Def.}} \xrightarrow{\mathbb{P}} \operatorname{argmin}_x \underbrace{Q(x)}_{\text{|| (2)}} \quad \text{[Technical proof see [N., Thm 1]]}$$

Picture

(V) Asymptotic Normality of the MLE

Def: The Fisher information  $I: \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$  of a  $C^1$ -likelihood model  $\{f_z(\cdot | X=x) | x \in \mathbb{R}^d\}$  is defined as  $I(x) = \text{Cov}[\partial_x q(x, z_x)]$  where  $q(x, z) = -\log f_z(z | X=x)$  and  $z_x$  has the prob. density  $f_z(\cdot | X=x)$

In particular,  $I(x_T) = \text{Cov}[\partial_x q(x_T, z_T)] \stackrel{\text{technical detail}}{=} E[(D_x^2 q)(x_T, z_T)]$ .  $\otimes$   
(under regularity assumptions)

Thm: Assuming Setup (S), (IA) and (C1)-(C3), the MLE  $\hat{x}_M$  based on  $z_1, \dots, z_M$  fulfills  $[N, \text{Thm 3}] [S, \text{Thm 6.6}]$   
 $\sqrt{M} (\hat{x}_M - x_T) \xrightarrow[M \rightarrow \infty]{\text{in dist.}} \mathcal{N}(0, I(x_T)^{-1})$ .

Proof (sketch):  $0 \stackrel{\text{Def.}}{\approx} \partial_x Q_M(\hat{x}_M) \stackrel{\text{Taylor}}{=} \partial_x Q_M(x_T) + D_x^2 Q_M(\tilde{x}_M) \cdot (\hat{x}_M - x_T)$   $\otimes \otimes$   
 $e[x_T, \hat{x}_M]$

So,  $\underbrace{\sqrt{M}(\hat{x}_M - x_T)}_{\text{guy of interest}} = - \underbrace{(D_x^2 Q_M(\tilde{x}_M))^{-1}}_{\text{A}} \underbrace{(\frac{1}{M} \partial_x Q_M(x_T))}_{\text{B}}$

$\text{A}$ : It can be shown that (since  $\hat{x}_M \xrightarrow{\mathbb{P}} x_T$  and thereby  $\tilde{x}_M \xrightarrow{\mathbb{P}} x_T$ , law of large numbers)

$D_x^2 Q_M(\tilde{x}_M) \xrightarrow[M \rightarrow \infty]{\mathbb{P}} E[D_x^2 q(x_T, z_T)] = I(x_T)$   
 $\frac{1}{M} \sum_{m=1}^M q(\tilde{x}_M, z_m)$

$\text{B} \cdot E[\partial_x q(x_T, z_T)] = \partial_x \underbrace{E[q(x_T, z_T)]}_{Q(x_T)} = 0$  since  $x_T$  minimizes  $Q$ .

$\sqrt{M} \cdot \frac{\partial_x Q_M(x_T)}{\frac{1}{M} \sum_{m=1}^M \partial_x q(x_T, z_m)} \xrightarrow[M \rightarrow \infty]{\text{in dist. (CLT)}} \mathcal{N}(0, \text{Cov}[\partial_x q(x_T, z_T)]) = \mathcal{N}(0, I(x_T))$

So,  $\sqrt{M}(\hat{x}_M - x_T) = - \underbrace{(D_x^2 Q_M(\tilde{x}_M))^{-1}}_{\xrightarrow{\mathbb{P}} I(x_T)} \cdot \underbrace{\frac{1}{M} \partial_x Q_M(x_T)}_{\xrightarrow[\text{Gau\ss}]{\text{in dist.}} \mathcal{N}(0, I(x_T))} \xrightarrow[M \rightarrow \infty]{\text{in dist.}} \mathcal{N}(0, I(x_T)^{-1} I(x_T) I(x_T)^{-1}) = \mathcal{N}(0, I(x_T)^{-1})$   $\square$

Reminder: (CLT)  $y_1, \dots, y_M \stackrel{\text{iid}}{\sim} y \Rightarrow \sqrt{M} \left( \frac{1}{M} \sum_{m=1}^M y_m - E[y] \right) \xrightarrow[M \rightarrow \infty]{\text{in dist.}} \mathcal{N}(0, \text{Cov}(y))$   
(Gau\ss)  $y \sim \mathcal{N}(\mu, \Sigma), A \in GL(d, \mathbb{R}) \Rightarrow A \cdot y \sim \mathcal{N}(A \cdot \mu, A \Sigma A^T)$

Remark: In dimension  $d > 1$ , the Taylor-expansion  $\otimes \otimes$  has to be performed in each dimension separately:  $\forall j=1, \dots, d: 0 = [\partial_x Q_M(\hat{x}_M)]_j + \underbrace{[D_x^2 Q_M(\tilde{x}_M^j)]_{j, \cdot}}_{\in [x_T, \hat{x}_M]} \cdot (\hat{x}_M - x_T)$

Then replace  $D_x^2 Q_M(\tilde{x}_M)$  by  $H^M$  defined by  $H_{j, \cdot}^M = [D_x^2 Q_M(\tilde{x}_M^j)]_{j, \cdot} \forall j$ :  
 $0 = \partial_x Q_M(\hat{x}_M) = \partial_x Q_M(x_T) + H^M \cdot (\hat{x}_M - x_T)$

Since  $D_x^2 Q_M(\tilde{x}_M^j) \xrightarrow[M \rightarrow \infty]{\mathbb{P}} I(x_T)$ , also  $H^M \xrightarrow[M \rightarrow \infty]{\mathbb{P}} I(x_T)$ , so this detail does not affect the rest of the proof.

### (VI) Frequentist Consistency of Bayesian Methods

Thm: ([N, Thm. 5], [S, Thm. 6.17], Bernstein-von Mises Theorem) :

Assuming Setup (S), (IA) and (C1)-(C3) and denoting

- by  $\hat{x}_M$  the MLE based on  $Z_1, \dots, Z_M$ ,
- by  $\mathbb{P}_M$  the posterior distribution with density  $f_M(x) = f_X(x | Z_1 = z_1^*, \dots, Z_M = z_M^*)$ ,
- by  $\|P-Q\|_{TV} := \sup_{B \in \mathcal{B}_n(\mathbb{R}^d)} |P(B) - Q(B)|$  the total variation distance between  $P, Q$ ,

we have: 
$$\| \mathbb{P}_M - \mathcal{N}(\hat{x}_M, \frac{1}{M} I(x_M^{-1})) \|_{TV} \xrightarrow{M \rightarrow \infty} 0$$