

# Chapter 4: Information Geometry

## Fisher Information

Why "Inform" geometry?

↳ about distances in the space of probability measures variations  
norms

$P$  defined on  $X$   $P \rightarrow P + \delta P$   
 $P'$   
 $d(P, P')$ ?  
 $\hookrightarrow KL(P || P')$

↳ motivation:

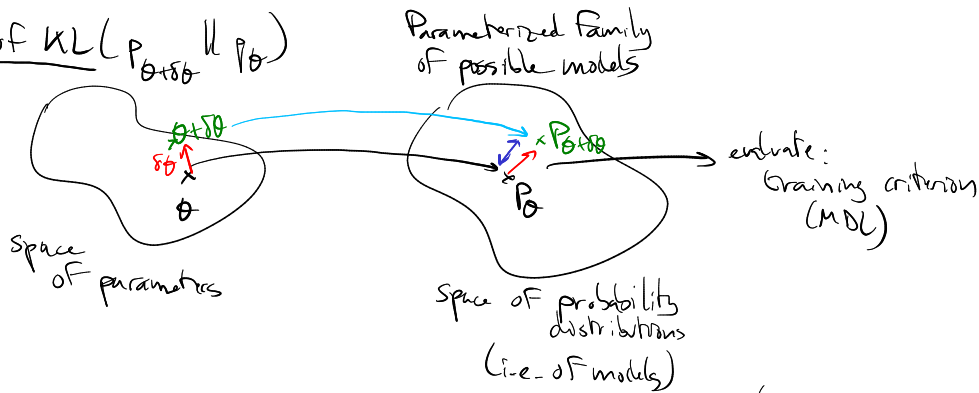
ML problem:  $\mathcal{D} = \{ (x_1, y_1), (x_2, y_2), \dots \}$

↳ fit a  $P_\theta$   $y = P_\theta(x) \rightarrow$  non-deterministic:  $P_\theta(y|x)$

finding the best parameters  $\theta$   $\rightarrow$  optimize  $\theta$   
i.e. the best  $P_\theta$

↳  $\nabla$  descent  $\rightarrow$  consider  $P_{\theta+\delta\theta}$

### Second order of $KL(P_{\theta+\delta\theta} || P_\theta)$



$KL(P_{\theta+\delta\theta} || P_\theta)$  = infinitesimal distance between  $P_\theta$  &  $P_{\theta+\delta\theta}$  ( $\delta\theta$ : small parameter variation)

$$= \int_{x \in X} P_{\theta+\delta\theta}(x) \log \frac{P_{\theta+\delta\theta}(x)}{P_\theta(x)} dx$$

Tricks:

$$\cdot \log(1+z) = z - \frac{1}{2}z^2 + O(z^3)$$

$$\cdot \frac{d \log F(z)}{dz} = \frac{1}{F(z)} \frac{dF}{dz}$$

$$\cdot \int_{x \in X} P_\theta(x) = 1 \quad \forall \theta$$

$$\Rightarrow \frac{d^k}{d\theta^k} (1) = 0$$

$$\Rightarrow \int_{x \in X} \frac{d^k}{d\theta^k} P_\theta(x) dx = 0$$

$$\frac{P_{\theta+\delta\theta}(x)}{P_\theta(x)} = \underbrace{\frac{P_\theta(x)}{P_\theta(x)}}_{\in \mathbb{R}} + \underbrace{\frac{dP_\theta(x)}{d\theta}}_{\text{matrix of size } |X| \times |Y|} \times \underbrace{\delta\theta}_{\text{vector } |Y|} + \frac{1}{2} \underbrace{\delta\theta}_{\text{matrix of size } |X| \times |Y|} \underbrace{\frac{d^2 P_\theta(x)}{d\theta^2}}_{\text{matrix of size } |X| \times |Y| \times |Y|} \underbrace{\delta\theta}_{\text{vector } |Y|} + O(\delta\theta^3)$$

$$\frac{P_{\theta+\delta\theta}(x)}{P_\theta(x)} = 1 + \frac{1}{P_\theta(x)} \frac{dP_\theta(x)}{d\theta} \delta\theta + \frac{1}{2} \delta\theta \frac{1}{P_\theta(x)} \frac{d^2 P_\theta(x)}{d\theta^2} \delta\theta + O(\delta\theta^3)$$

$$\log(\dots) = \frac{1}{P_\theta(x)} \frac{dP_\theta(x)}{d\theta} \delta\theta + \frac{1}{2} \delta\theta \frac{1}{P_\theta(x)} \frac{d^2 P_\theta(x)}{d\theta^2} \delta\theta - \frac{1}{2} \delta\theta \frac{1}{P_\theta(x)^2} \frac{dP_\theta(x)}{d\theta} \frac{dP_\theta(x)}{d\theta}^T \delta\theta + O(\delta\theta^3)$$

$$KL = \mathbb{E}_{x \sim P_{\theta+\delta\theta}} [\dots] = \mathbb{E}_{x \sim P_\theta} [\dots] + \int_{x \in X} \frac{1}{P_\theta(x)} \delta\theta \frac{dP_\theta(x)}{d\theta} \frac{dP_\theta(x)}{d\theta}^T \delta\theta dx + O(\delta\theta^3)$$

$$= \underbrace{\left( \int_{x \in \mathcal{X}} \frac{dP_\theta(x)}{d\theta} dx \right)}_0 \cdot \delta\theta + 0 - \frac{1}{2} \mathbb{E} \left[ \int_{x \in \mathcal{X}} \frac{1}{P_\theta(x)} \frac{dP_\theta(x)}{d\theta} \frac{dP_\theta(x)^T}{d\theta} dx \right] \delta\theta^2 + O(\delta\theta^3)$$

+ 1 (... Same stuff)

$$KL(P_{\theta+\delta\theta} \| P_\theta) = \frac{1}{2} \mathbb{E} \left[ \frac{d \log P_\theta(x)}{d\theta} \frac{d \log P_\theta(x)^T}{d\theta} \right] \delta\theta^2 + O(\delta\theta^3)$$

matrix of size  $|\Theta| \times |\Theta|$       vector of size  $|\Theta|$

Second order quantity:  $O(\delta\theta^2)$

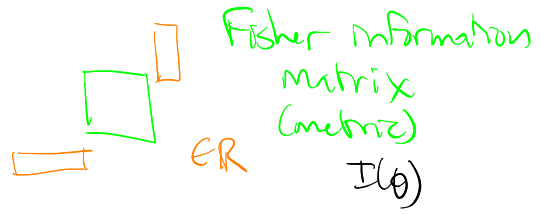
Optimize

$$KL(q \| p)$$

w.r.t  $q$ :

$q=p$  is the global optimum

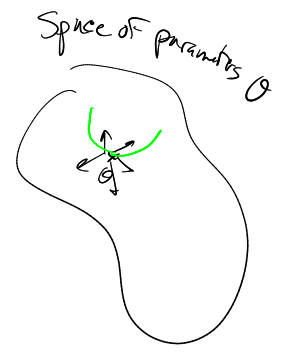
$$\Rightarrow \frac{\partial KL(q \| p)}{\partial q} \Big|_{q=p} = 0$$



Other Formula for Fisher information  $I(\theta)$

$$I(\theta) = \mathbb{E} \left[ \frac{d^2}{d\theta^2} (-\log P_\theta(x)) \right]$$

average of  $\frac{d^2}{d\theta^2}$  of the encoding cost w.r.t  $\theta$  curvature



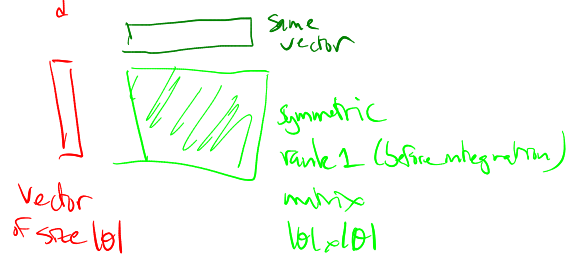
Equivalence between the formulas

$$\log P_\theta(x)$$

$$\frac{d}{d\theta} \log P_\theta(x) = \frac{1}{P_\theta(x)} \frac{dP_\theta(x)}{d\theta}$$

$$\frac{d^2}{d\theta^2} (\dots) = \frac{d}{d\theta} \left( \frac{1}{P_\theta(x)} \right) = \frac{-1}{P_\theta(x)^2} \frac{dP_\theta(x)}{d\theta} \times \frac{dP_\theta(x)^T}{d\theta} + \frac{1}{P_\theta(x)} \frac{d^2 P_\theta(x)}{d\theta^2}$$

$$I(\theta) = \mathbb{E} \left[ - \frac{d \log P_\theta(x)}{d\theta} \frac{d \log P_\theta(x)^T}{d\theta} \right] + \int_{x \in \mathcal{X}} \frac{d^2 P_\theta(x)}{d\theta^2} dx$$



In practice:  $I(\theta) \approx \frac{1}{n} \sum_{i \in \mathcal{D}} \vec{v}_i \otimes \vec{v}_i$  with  $\vec{v}_i = \frac{d \log P_\theta(x_i)}{d\theta}$

number of points analyzed



Cramer-Rao bound (theorem by Amari) → Cover & Thomas ch 11.10

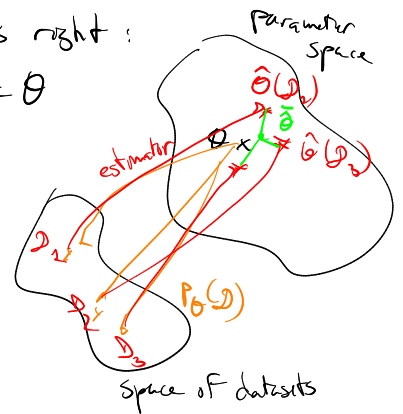
- given observations  $\mathcal{D}_n(x_i)$ , guess the parameter  $\theta$  s.t.  $p_\theta(x)$  is the highest (best fit)  
 ↳ estimator

- an estimator  $\hat{\theta}(\mathcal{D})$  is said to be unbiased if on average it is right:

$$E[\hat{\theta}(\mathcal{D})] = \theta$$

$\mathcal{D} \sim p_\theta$

- question: average of the estimator  
 but what about its variance?  
 ↳ is  $|\theta - \hat{\theta}|$  typically large or small?



Theorem Cramer-Rao inequality

In dimension 1:

$$\text{variance}(\text{any unbiased estimator}) \geq \frac{1}{J(\theta)} \quad \text{E.R.}$$

In dim > 1:

$$\text{covariance}(\dots) \geq J(\theta)^{-1} \quad (\text{in terms of matrix})$$

$J = n I$   
 ↳ number of samples

Proof: online material

II Parameter precision

Yann Ollivier → Jérôme Bessandon

model  $\mu_\theta \rightarrow$  encode  $\theta$  up to which precision?

$$\theta = \theta^* + \epsilon$$

↳ precision

$$\theta^* = 0,32085 \dots$$

$$\theta = 0,32$$

(K) Cost to describe  $\theta$  and the model etc:  $C(\epsilon) = \underbrace{-\log \epsilon}_{\text{cost to describe } \theta \text{ up to } \epsilon} - \log \mu_\theta(\epsilon) = \log \epsilon - \log \mu_{\theta^*}(\epsilon) + \frac{1}{2} \epsilon^2 \underbrace{J(\theta^*)}_{\frac{d^2}{d\theta^2} \log \mu_\theta} + O(\epsilon^3)$

Set  $\epsilon$ ?  $\frac{\partial}{\partial \epsilon} C = 0$

$$\frac{\partial}{\partial \epsilon} C = -\frac{1}{\epsilon} + \epsilon J(\theta^*) = 0$$

$$\boxed{\epsilon = \frac{1}{\sqrt{J(\theta^*)}}} \quad \text{optimal precision!}$$

$$J(\theta^*) = n I(\theta^*)$$

$$\epsilon^* \approx \frac{1}{\sqrt{n}} \frac{1}{\sqrt{I(\theta^*)}}$$

↳ classical statistics

III Misc

Justification of BIC:

$$\text{BIC} : K(\mu) := \frac{1}{2} \text{number of parameters} \times \log(\text{number of observations})$$

↳ encode model with optimal precision

↳ encode the parameters: for each parameter, length of encoding:  $-\log \epsilon = \frac{1}{2} \log n + \frac{1}{2} \log J(\theta)$   
 ↳  $\mathcal{D} \sim \mu$  × number of parameters constant

$$\text{precision: } \epsilon = \frac{1}{\sqrt{n}} \frac{1}{\sqrt{I(\theta)}}$$

→ strong links with natural gradient

→ Jeffrey's prior