

## parameter estimation

1. maximum likelihood estimation (regularization)
  2. empirical risk minimization (structural)
- 

### MLE

Def:  $x_1, \dots, x_n \sim P_\theta(x)$

likelihood  $L(x_i, \theta) \triangleq P_\theta(x_i)$

Def: likelihood for the entire dataset

$$L_n(\theta) \triangleq P_\theta(x_1, \dots, x_n)$$

Def: MLE

$$\hat{\theta}_n = \underset{\theta \in \Theta}{\operatorname{argmax}} L_n(\theta)$$

MLE maps  $x_1, \dots, x_n \rightarrow \hat{\theta}_n$

Def: log-likelihood  $l_n(\theta) \triangleq \log L_n(\theta)$

Def: bias of  $\hat{\theta}_n$   $\mathbb{E} \hat{\theta}_n - \underbrace{\theta}_{\text{truth}}$

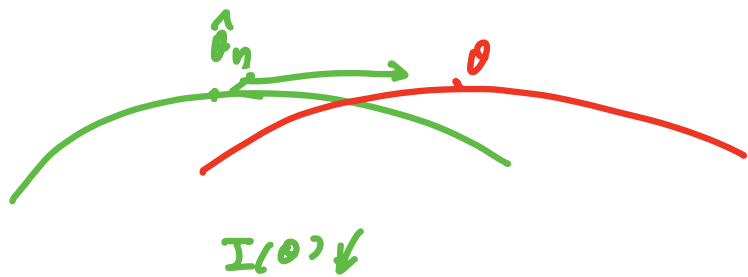
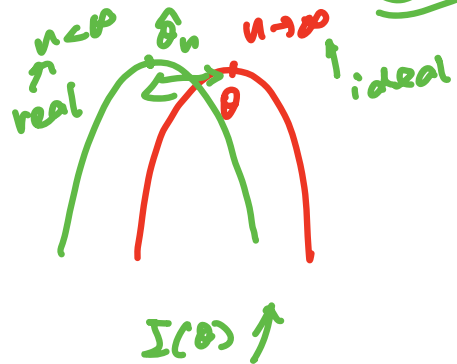
Def: if  $\text{bias}(\hat{\theta}_n) = 0$  then  $\hat{\theta}_n$  is the unbiased estimator

### Theory of MLE

MLE is asymptotically normal and efficient

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow \infty]{d} N(0, I^{-1}(\theta))$$

$I(\theta) \triangleq -\mathbb{E}_\theta \left[ \frac{\partial^2}{\partial \theta^2} \log P_\theta(x) \right]$  Fisher info...



$$f \in \mathcal{F} \quad z_1, \dots, z_n \in \mathcal{Z}$$

Empirical risk minimization.

$$l: \mathcal{F} \times \mathcal{Z} \rightarrow \mathbb{R}$$

$l$ : loss

$$L(f) \triangleq \mathbb{E}_P[l(f, z)] \quad \text{risk}$$

goal:  $\min_f L(f)$  ← the population version

$$\text{reality} \quad \min_f \hat{L}_n(f) \triangleq \frac{1}{n} \sum_{i=1}^n l(f, z_i) \quad -\log p_{\theta}(x_i)$$

ERM

Example 1 (Binary classification problem)

$$\underbrace{\mathbb{E}[l(f, z)]}_{\text{risk}} = \mathbb{E}[\underbrace{1_{f(x)y \leq 0}}_{\text{indicator}}]$$

Example 2 (multiclass classification problem)

$$y \in \{1, \dots, K\} \quad f: \mathcal{X} \rightarrow \mathbb{R}^K$$

$$l(f, x, y) = \mathbb{1}_{\{\max_{i \neq y} f_i(x) \geq f_y(x)\}}$$

Example 3: (hinge or logistic)  $\mathcal{F} = \{f_\theta: f_\theta(x) = \underline{\underline{\sigma^T x}}\}$

$$l_{\text{hinge}}(f_\theta, (x, y)) = [1 - y \underline{\underline{x^T \theta}}]_+$$

$$l_{\text{logit}}(f_\theta, (x, y)) = \log(1 + \exp(-y \underline{\underline{x^T \theta}}))$$

Given  $z_1, \dots, z_n$  pick  $\hat{f}_n \in \mathcal{F}$  Generalization

$$\underline{\underline{\hat{L}_n(\hat{f}_n)}} = \frac{1}{n} \sum_{i=1}^n l(\hat{f}_n, z_i) \leq \underline{\underline{R(\hat{f}_n)}} + \underline{\underline{\varepsilon_{\text{error}}}}$$

$\hat{f}_n(x) = y$  if  $(x, y)$  in the training dataset  
 $\hat{f}_n(x) = -1$  otherwise

$$l(f, (x, y)) = (y - f(x))^2$$



$$\hat{L}_n(\hat{f}_n) = 0$$

$$R(\hat{f}_n) \uparrow$$