

# Chapter IV (alternative)

Eddie  
04/10/24

## Learning & Minimax theory in a nutshell

I Empirical risk minimization: A simple oracle bound

II Minimax lower bounds: A bayesian method applied to regression

### I] Empirical risk minimization

We observe  $(X_1, Y_1), \dots, (X_n, Y_n)$  iid with distribution  $P_{(X,Y)}$  over  $\mathcal{X} \times \mathcal{Y}$ .

The goal is to learn a function  $f: \mathcal{X} \rightarrow \mathcal{Y}$  with small risk

$$R(f) = E_P[\ell(Y, f(X))]$$

where  $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  is a prescribed loss function

Ex: Classification

$$\mathcal{Y} = \{-1, 1\}$$

$$\ell(y, y') = \mathbb{1}_{y \neq y'}$$

Regression

$$\mathcal{Y} = \mathbb{R}$$

$$\ell(y, y') = (y - y')^2$$

We choose to search  $f$  in a prescribed set  $\mathcal{F}$  of (measurable) functions<sup>(2)</sup>

E.g.  $\mathcal{F} :=$  linear functions on  $X = \mathbb{R}^d$   
 $:=$  A class of neural nets  
 $:=$  sum of  $N$  Fourier functions

We try to minimize  $R$  by the empirical risk

$$\boxed{\hat{R}(f) := \frac{1}{m} \sum_{i=1}^m \ell(Y_i, f(X_i))}$$

The ERM on  $\mathcal{F}$  is then

$$\boxed{\hat{f} \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} \hat{R}(f)}$$

The goal is to compare the (expected) performances of  $\hat{f}$ , compared to the best possible over the class  $\mathcal{F}$

$$\boxed{\hat{f}^* \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} R(f)}$$

We have, almost surely,

$$\underline{R(\hat{f})} = \underline{R(f^*)} + \underline{R(\hat{f}) - R(f^*)}$$

(Random!)  
 performance  
 (Deterministic)  
 approximation error  
 Estimation error

The approximation error depends on  $\mathcal{H}$  and  $P_{X,Y}$ . In regression,  $\mathcal{H}$  has been studied in chapter 2. We focus on the estimation term

$$\begin{aligned} R(\hat{f}) - R(f^*) &= (R(\hat{f}) - \hat{R}(\hat{f})) + (\hat{R}(\hat{f}) - \hat{R}(f^*)) \\ &\quad \xrightarrow{\leq 0 \text{ by definition of } \hat{f}} + (\hat{R}(f^*) - R(f^*)) \\ &\leq \delta \text{ if risk minimization is done approximately} \end{aligned}$$

$$\leq 2 \sup_{f \in \mathcal{H}} |\hat{R}(f) - R(f)|$$

From the central limit theorem we know that for fixed  $f \in \mathcal{H}$ ,

$$\sqrt{n} (\hat{R}(f) - R(f)) \xrightarrow[n \rightarrow \infty]{d} N(0, \text{Var}(\ell(Y, f(X))))$$

so that  $|\hat{R}(f) - R(f)| \leq \frac{1}{\sqrt{n}}$  whp

Here we need a stronger bound holding whp on all  $f \in \mathcal{H}$   
 $\implies$  Depends on the complexity of the class  $\mathcal{H}$

For simplicity, assume that  $\mathcal{Y} = \mathbb{R}^p$

(4)

$$\bullet \forall f \in \mathcal{F}, \|f\|_\infty \leq M$$

(B)

$\bullet \ell$  is  $G$ -Lipschitz in the second variable

(L)

$$|\ell(y, y') - \ell(y, y'')| \leq G \|y' - y''\|$$

From (L), we get that  $\forall f, f' \in \mathcal{F}$ ,

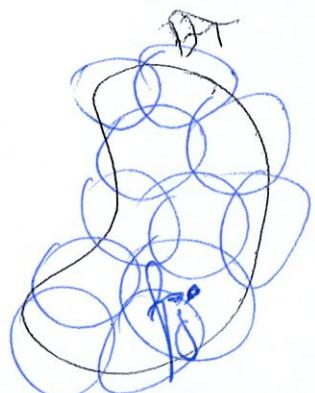
$$\begin{aligned} |\mathcal{R}(f) - \mathcal{R}(f')| &= \left| \mathbb{E} [\ell(Y, f(X)) - \ell(Y, f'(X))] \right| \\ &\leq G \left[ \mathbb{E} |f(X) - f'(X)| \right] \\ &\leq G \|f - f'\|_\infty \end{aligned}$$

and similarly,  $|\hat{\mathcal{R}}(f) - \hat{\mathcal{R}}(f')| \leq G \|f - f'\|_\infty$

This will allow us to discretize the supremum  $\sup_{f \in \mathcal{F}}$  and simplify it to a finite family.

Given  $\delta > 0$ , write  $f_1, \dots, f_N$  for a minimal  $\delta$ -covering of  $\mathcal{F}$  in  $\ell^\infty$  norm. That is  $N = N(\mathcal{F}, \|\cdot\|_\infty, \delta)$

$$\mathcal{F} \subset \bigcup_{j=1}^N B_{\|\cdot\|_\infty}(f_j, \delta)$$



For all  $f \in \mathcal{H}$ , there exists  $f_j$  with  $\|f - f_j\|_\infty \leq \delta$ , and hence

$$\begin{aligned} |\mathcal{R}(f) - \hat{\mathcal{R}}(f)| &\leq |\mathcal{R}(f) - \mathcal{R}(f_j)| + |\mathcal{R}(f_j) - \hat{\mathcal{R}}(f_j)| + |\hat{\mathcal{R}}(f_j) - \hat{\mathcal{R}}(f)| \\ &\leq 2G\|f - f_j\|_\infty + |\mathcal{R}(f_j) - \hat{\mathcal{R}}(f_j)| \\ &\leq 2G\delta + |\mathcal{R}(f_j) - \hat{\mathcal{R}}(f_j)| \end{aligned}$$

Hence,  $\sup_{f \in \mathcal{B}} |\mathcal{R}(f) - \hat{\mathcal{R}}(f)| \leq 2G\delta + \max_{j \leq N_f} |\mathcal{R}(f_j) - \hat{\mathcal{R}}(f_j)|$

Let now  $f_j$  be fixed. We assume that the loss function is bounded

$$\boxed{\forall y, y' \in \mathcal{Y}, |\ell(y, y')| \leq \ell_\infty}$$

Then  $\mathcal{R}(f_j) - \hat{\mathcal{R}}(f_j) = \frac{1}{m} \sum_{i=1}^m (\mathbb{E}[\ell(Y_i, f_j(X_i))] - \ell(Y_i, f_j(X_i)))$

is the sum of  $m$  centered variables bounded by  $\frac{\ell_\infty}{m}$ . From Hoeffding's concentration inequality,  $\forall t > 0$ ,

$$\begin{aligned} P(|\mathcal{R}(f_j) - \hat{\mathcal{R}}(f_j)| > t) &\leq 2 \exp\left(-\frac{2t^2}{m \left(\frac{2\ell_\infty}{m}\right)^2}\right) \\ &= 2 \exp\left(-\frac{m t^2}{2\ell_\infty^2}\right) \end{aligned}$$

At the end of the day, we apply a union bound to get

$$\begin{aligned} P\left(\sup_{f \in \mathcal{F}} |\mathcal{R}(f) - \hat{\mathcal{R}}(f)| > t\right) &\leq P\left(\bigcup_{j \leq N_8} |\mathcal{R}(f_j) - \hat{\mathcal{R}}(f_j)| > t\right) \\ &\leq N_8 \cdot \Pr_{j \leq N_8} (|\mathcal{R}(f_j) - \hat{\mathcal{R}}(f_j)| > t) \\ &\leq N_8 \times 2 \exp\left(-\frac{mt^2}{2\ell_\infty}\right) \end{aligned}$$

En posant  $\beta = 2N_8 \exp\left(-\frac{mt^2}{2\ell_\infty}\right)$ , on obtient donc qu'avec proba au moins  $1 - \beta$ ,

$$\Leftrightarrow \frac{mt^2}{2\ell_\infty} = \log\left(\frac{eN_8}{\beta}\right) \Leftrightarrow t = \sqrt{\frac{2\ell_\infty \log\left(\frac{eN_8}{\beta}\right)}{m}}$$

$$\sup_{f \in \mathcal{F}} |\mathcal{R}(f) - \hat{\mathcal{R}}(f)| \leq 2G\delta + \sqrt{\frac{2\ell_\infty \log\left(\frac{eN_8}{\beta}\right)}{m}}$$

Coming back to the initial bound, with proba  $\geq 1 - \beta$ ,

$$\boxed{\mathcal{R}(\hat{f}) \leq \mathcal{R}(f^*) + \sqrt{\frac{8\ell_\infty \log\left(\frac{2N_8}{\beta}\right)}{m}} + 4G\delta}$$

Oracle bound  
in probability

Using that  $E[Z] = \int_0^\infty P(Z \geq t) dt$  if  $Z \geq 0$ , we get a similar expected bound.

## II Minimax lower bound technique

Given a particular estimator  $\hat{\theta}_n$  that we designed over  $n$ -samples  $Z_1, \dots, Z_n \sim P$ , one may wonder whether one couldn't get a more precise one, for instance on average. Minimax lower bounds aim at doing this, by showing some intrinsic information-theoretic limitations of specific problems.

Framework: Model  $P$

- Parameter of interest  $\Theta: \mathcal{P} \rightarrow \mathbb{H}$
- Metric  $d: \mathbb{H} \times \mathbb{H} \rightarrow \mathbb{R}_{\geq 0}$

The minimax risk for estimating  $\Theta(P)$  over  $\mathcal{P}$  is

$$R_m(\mathcal{P}) := \inf_{\hat{\theta}_m} \sup_{P \in \mathcal{P}} E_{P^{\otimes m}} [d(\Theta(P), \hat{\theta}_m)]$$

where  $\hat{\theta}_m = \hat{\theta}_m(Z_1, \dots, Z_m)$  ranges among all the estimators over  $m$  points  $Z_1, \dots, Z_m$

$$\text{E.g.: } \mathcal{P} = \mathcal{D}(X, Y) \left| \begin{array}{l} \cdot X \sim \text{Unif}[0, 1]^d \\ \cdot Y = f(X) + \varepsilon \\ \cdot X \perp \varepsilon \\ \cdot \text{Var}(\varepsilon) \leq \sigma^2, \mathbb{E}[\varepsilon] = 0 \\ \cdot f \in \mathcal{C}^1 \end{array} \right. \quad \left| \begin{array}{l} \cdot \Theta(P) = f^* \\ \cdot d(f, f') = \|f - f'\|_2 \end{array} \right. /$$

Studying  $R_m(P)$  is twofold

Upper bounds: Exhibit a particular  $\hat{\theta}_m$  which has good  
 $R_m(P) \leq \dots$  estimation rate uniformly over  $P$

Lower bounds: Show that one cannot do better; at least up  
 $R_m(P) \geq \dots$  to constants

For lower bounds, the main idea is to replace the  $\sup_{P \in \mathcal{P}}$  by an average  $\int_P -\pi(dP)$ .

Indeed, if  $\pi$  is a probability distribution over  $\mathcal{P}$ , then for all  $\hat{\theta}_m$ ,

↳ Prior distribution

$$\sup_{P \in \mathcal{P}} E_{P^m} [d(\theta(P), \hat{\theta}_m)] \geq \int_P E_{P^m} [d(\theta(P), \hat{\theta}_m)] \pi(dP)$$

The right-hand term can usually be studied relatively easily,

since • WE choose  $\pi$

• The best possible estimator  $\hat{\theta}_\pi$  (bayes estimator) can be explicit

To cover the main ideas, we will present the simplest case where

$$\boxed{\overline{P} = \frac{1}{2}(P_0 + P_1)}$$

$\curvearrowleft$  two "hypotheses"

Thm: (Le Cam's lemma)

For all  $P_0, P_1 \in \mathcal{P}$ ,

$$R_n(P) \geq \frac{1}{2} d(\sigma(P_0), \sigma(P_1)) \quad \left. \begin{array}{l} P_0^{\otimes n}, P_1^{\otimes n}, d\Gamma^{\otimes n} \\ \text{where } P_0, P_1 < 0 \end{array} \right\}$$

$$\left. \begin{array}{l} P_f = \frac{dP}{d\Gamma} \\ \text{where } P_0, P_1 < 0 \end{array} \right\}$$

Proof:  $R_n(P) \geq \inf_{\hat{\theta}_n} \frac{1}{2} \left( E_{P_0^{\otimes n}} [d(\sigma(P_0), \hat{\theta}_n)] + E_{P_1^{\otimes n}} [d(\sigma(P_1), \hat{\theta}_n)] \right)$

$$= \inf_{\hat{\theta}_n} \frac{1}{2} \int_{\mathbb{Z}^n} \left( d(\sigma(P_0), \hat{\theta}_n(y_1, \dots, y_n)) P_0^{\otimes n}(\Gamma) + d(\sigma(P_1), \hat{\theta}_n(y_1, \dots, y_n)) P_1^{\otimes n}(\Gamma) \right) d\Gamma$$

$$\geq \frac{1}{2} d(\sigma(P_0), \sigma(P_1)) \int_{\overline{P}} P_0^{\otimes n}, P_1^{\otimes n}, d\Gamma^{\otimes n}$$

Rmk: Actually,  $\int P_0 \wedge P_1 = 1 - \frac{1}{2} \int |P_0 - P_1|$   $\curvearrowright$  Total Variation

We can further bound the right-hand term

Lemma:  $\int P_0^{\otimes n} \wedge P_1^{\otimes n} \geq \frac{1}{2} \text{Bcp}\left(-n \text{KL}(P_0, P_1)\right),$

where  $\text{KL}(P_0, P_1) = \begin{cases} \int \log\left(\frac{P_0}{P_1}\right) P_0 & \text{if } P_0 \ll P_1, \\ \infty & \text{otherwise} \end{cases}$

Prof: Jensen, See lemma 2.6 in Tsybakov [Introduction to Nonparametric estimation]

Rk: For instance, gaussian satisfy  $\boxed{\text{KL}(N(0, \sigma^2), N(m, \sigma^2)) = \frac{m^2}{2\sigma^2}}$

Application to regression:

Say that  $X = [-1, 1]^d$ , and  $d(f, f') = \|f(0) - f'(0)\|$   
 $Y = \mathbb{R}$

$$Y = f^*(X) + \varepsilon \quad \text{with } \varepsilon \sim N(0, \sigma^2)$$

Assume also that  $f \in \mathcal{C}_L^\beta$ :  $|f(x) - f(y)| \leq L|x-y|^\beta$

$$\hookrightarrow \mathcal{P}_\beta := \left\{ \mathcal{D}(X, Y), \quad \downarrow \right\}$$

Apply Le Cam's lemma to  $P_0 = P_{f_0}$ ,  $P_1 = P_{f_1}$ :

$$R_n(P) \geq \frac{|f_0(0) - f_1(0)|}{2} \exp(-n KL(P_{f_0}, P_{f_1}))$$

When  $\epsilon \sim N(0, \sigma^2)$ , we can show that  $\begin{cases} Y_0 = f_0(X) + \epsilon \\ Y_1 = f_1(X) + \epsilon \end{cases}$  satisfy  $X \sim \text{Unif}([-1, 1]^d)$

$$\begin{aligned} KL(P_{f_0}, P_{f_1}) &\leq \int_{[-1, 1]^d} KL(P_{Y|X}, P_{Y'|X}) dx \\ &= \int_{[-1, 1]^d} KL(N(f_0(x), \sigma^2), N(f_1(x), \sigma^2)) dx \\ &= \frac{\|f_0 - f_1\|_{L^2}^2}{8\sigma^2} \end{aligned}$$

Hence,  $R_n(P) \geq |f_0(0) - f_1(0)| \exp\left(-n \|f_0 - f_1\|_{L^2([-1, 1]^d)}^2\right)$

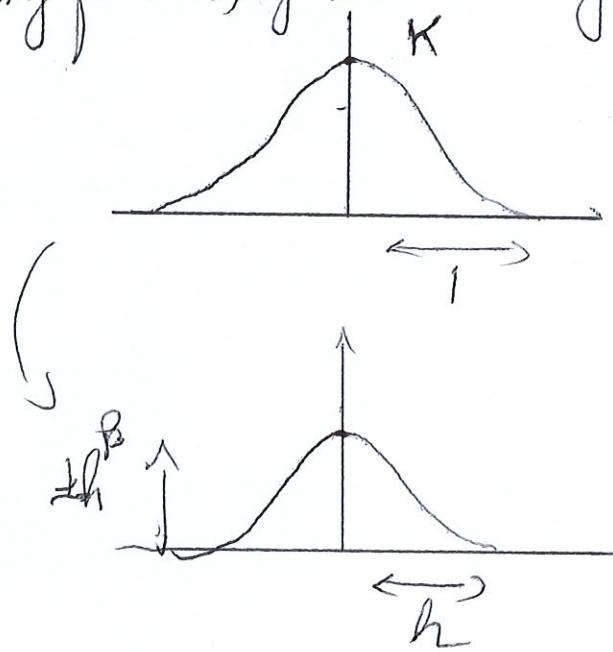
↳ Need to exhibit  $f_0, f_1 \in C^3([-1, 1]^d)$  such that

- $|f_0(0) - f_1(0)|$  is large
- $\|f_0 - f_1\|_{L^2([-1, 1]^d)}^2$  is small.

For this, we let  $h > 0$  to be chosen later, and

$$\begin{cases} \text{. } f_0 = 0 \\ \text{. } f_1(x) = Lh^\beta K\left(\frac{x}{h}\right) \end{cases}$$

where  $K \in C^{\beta}(I)$  is any function, by should be thought of as a localized bump



We have:  $f_0, f_1 \in C^{\beta}(L) \quad \forall h > 0$

$$\cdot |f_0(0) - f_1(0)| = Lh^\beta K(0)$$

$$\cdot \|f_0 - f_1\|_{L^2}^2 = \int_{B(0,h)} Lh^\beta K\left(\frac{x}{h}\right)^2 dx = L^2 h^{2\beta+d} \|K\|_{L^2}^2$$

$$\forall h > 0 \quad R_m(\beta_\beta) \geq Lh^\beta \log\left(-m \frac{Lh^{2\beta+d}}{\sigma^2}\right)$$

Choosing  $\frac{mLh^{2\beta+d}}{\sigma^2} = 1$  (i.e.  $h = \left(\frac{\sigma^2}{L^2 m}\right)^{\frac{1}{2\beta+d}}$ ), we get

$$R_m(\beta_\beta) \geq \left(\frac{\sigma^2}{L^2 m}\right)^{\frac{\beta}{2\beta+d}}$$