

## Lecture 5 : Basics on Generalization Error

- Last lecture:
- we constructed linear approximation spaces in high, infinite-dimensional inputs with kernels.
  - efficient learning algorithms (regularised least-squares)

- Today:
- Basic tools to control generalisation error.

Reminder: Next week guest lecture by Alberto Bietti

→ Focus on ERM (Empirical Risk Minimisation). Let  $\nu = \nu(x, y)$  be the joint data distribution over  $X \times Y$ .

Goal: Learn  $f: X \rightarrow Y$  that minimises a risk

$$R(f) = E_{\nu} [l(y, f(x))] \quad l: \text{_squared-loss, logistic loss, etc.}$$

→ ERM: Consider a normed function class  $F \subset \{f: X \rightarrow Y\}$  and consider an estimator (ERM) of the form

$$\hat{f} \in \arg \min_{\|f\|_F \leq \delta} \hat{R}(f) = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i)).$$

→ Recall the fundamental decomposition of the risk:

$$R(\hat{f}) - \inf_g R(g) = \boxed{R(\hat{f}) - \inf_{f \in F_{\delta} = F} R(f)} + \boxed{\inf_{f \in F} R(f) - \inf_{g \in G} R(g)}$$

from now on.

↳ estimation error      ↳ approximation error

→ Estimation error decomposed in terms of ERN:

$$\begin{aligned} R(\hat{f}) - \inf_{f \in \mathcal{F}} R(f) &= R(\hat{f}) - R(g) = R(\hat{f}) - \hat{R}(f) \\ &\quad + \hat{R}(f) - \hat{R}(g) + \\ g &\in \arg \min_{f \in \mathcal{F}} R(f) \\ &\quad + \hat{R}(g) - R(g) \\ &\leq \sup_{f \in \mathcal{F}} |R(f) - \hat{R}(f)| + \sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)| \end{aligned}$$

→ When  $\hat{f}$  is not the exact minimizer, then we need to add optim. error.

→ Q: What is the expected / intuitive behavior of this estimation error term?

→ Recall that at fixed  $f$ , we have

$\hat{R}(f)$  is unbiased estimator of  $R(f)$  (Law of Large Numbers)  
 $(\hat{R}(f) - R(f))$  are  $O(\sqrt{n})$  from CLT: approximate

Q: How to go from the (CLN) view at fixed  $f$  to Gaussianity.  
(i) a uniform deviation bound?

(ii) How to give non-asymptotic guarantees?

Finite Function Classes. Assume that the loss is

bounded a.s.  $\begin{cases} \hat{R}(f) \leq A & \text{with proba 1} \\ R(f) \leq A \end{cases}$

The simplest tool is the union bound:

$$P(|R(\hat{f}) - R(f)| > t) \leq P\left(2 \sup_{f \in F} |R(\hat{f}) - R(f)| \geq t\right)$$

$$\left(\sum_{f \in F} P(|R(\hat{f}) - R(f)| > t/2)\right)$$

union bound.

$(y_i, x_i) \sim$

→ For any fixed  $f \in F$ ,

$\ell(y_i, f(x_i))$  is a bounded RV.

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

$\hat{R}(f)$  is an average of bounded, iid random variables.

Prop (Hoeffding Ineq): For  $z_1 \dots z_n$  iid random variables in  $[0, c]$ , then  $\forall t \geq 0$ ,

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n z_i - E[z]\right| \geq t\right) \leq 2 \exp\left(-2nt^2/c^2\right)$$

so using Hoeffding

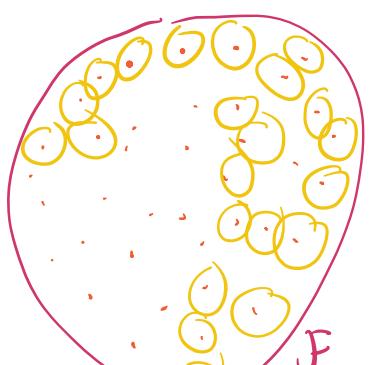
$$P\left(|\hat{R}(f) - R(f)| \geq t/2\right) \leq 2 \exp\left(-\frac{nt^2}{2A^2}\right)$$

$(\ell(y_i, f(x_i)) \leq A \text{ a.s.})$  for each  $f \in F$ .

$\#F \in F$ .  
so, if  $g \in \arg \min_{f \in F} R(f)$ ,

$$P(R(\hat{f}) - R(g) \geq t) \leq \sum_{f \in F} 2 \exp\left(-\frac{nt^2}{2A^2}\right) =$$

$$= |F| \cdot \exp\left(-\frac{nt^2}{2A^2}\right)$$

- Now, we set  $\delta = 2|f| \exp\left(\frac{-nt^2}{2A^2}\right)$  and solve for  $t$ ;
- we conclude that with probability at least  $1-\delta$ ,
- $$R(\hat{f}) - R(g) \leq \frac{2A}{\sqrt{n}} \sqrt{\log \frac{2|f|}{\delta}}$$
- $$\sqrt{\log |f|} + \sqrt{\log \delta}$$
- corresponds to looking at "extremal" statistics
- ↳ term that comes from concentration of measure.
- Remark: Going from expected excess  $q_{ij/c}$   
 $E[R(\hat{f}) - R(g)]$  to high-probability bounds uses "classic" tools; eg MacDiarmid inequality
  - Remark: The term  $\sqrt{\log |f|}$  appears when computing the expectation of the max of  $|f|$  zero-mean, subgaussian indep random variables.
- $$E \max_i |z_i| \leq \Gamma \sqrt{2 \log n}$$
- ① How to extend to infinite function classes?
- idea: let's grid the domain!
- Assume both  $R, \hat{R}$  are  $\beta$ -Lipschitz with respect to a distance metric  $d$  over  $F$ .
- 

$$F. \quad |R(f) - R(\hat{f})| \leq \beta \cdot \text{dist}(f, \hat{f}).$$

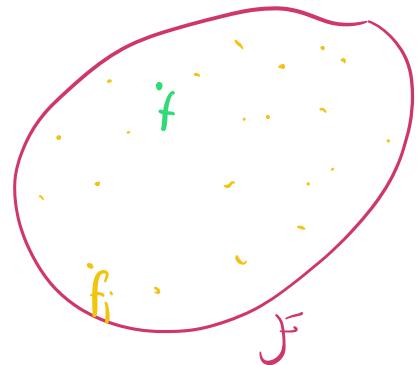
$$\left. \begin{array}{c} f_\theta \\ f_{\tilde{\theta}} \end{array} \right\} \left. \begin{array}{c} \hat{f} \\ f \end{array} \right\} (\theta, \tilde{\theta})$$

→ Fix  $\epsilon > 0$ . Assume there exists  $m_\epsilon$  elements  $f_1 \dots f_{m_\epsilon}$  such that  $\sup_{f \in F} \inf_i d(f, f_i) \leq \epsilon$ . The smallest  $m_\epsilon$  satisfying this property is the covering number of  $F$ .

→ 0:  $m_\epsilon$  blows up as  $\epsilon \rightarrow 0$ . For parametrized function classes  $F = \{f_\theta ; \theta \in \Theta \subseteq \mathbb{R}^d\}$  we expect  $m_\epsilon \sim \epsilon^{-\tilde{d}}$   $\tilde{d}$  is the ""intrinsic"" dimension of  $F$ .

→ Let's proceed trying to reuse our previous finite analysis:

$$\begin{aligned} \sup_{f \in F} |R(f) - \hat{R}(f)| &\leq \\ &\leq \sup_{f \in F} \left\{ |\hat{R}(f) - \hat{R}(f_i)| + \right. \\ &\quad + |\hat{R}(f_i) - R(f_i)| + \\ &\quad \left. + |R(f_i) - R(f)| \right\} \end{aligned}$$



$$(?) 2\beta \cdot \text{dist}(f, f_i) + \sup_i |\hat{R}(f_i) - R(f_i)|$$

$\hookrightarrow R$  and  $\hat{R}$   
are  $\beta$ -Lipschitz.

→ Since Bounded random variables in particular are subgaussian

$$E \left[ \sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)| \right] \leq 2\beta \varepsilon + |E| \underbrace{\sup_{i=1 \dots m_\varepsilon} |\hat{R}(f_i) - R(f_i)|}_{\text{std is } \frac{A}{\sqrt{n}}}$$

$$\leq 2\beta \varepsilon + A \sqrt{\frac{2 \log(2m_\varepsilon)}{n}}$$

$\rightarrow$  So, if  $m_\varepsilon \sim \varepsilon^{-d}$ , then optimising this upper bound wrt  $\varepsilon$  gives  $\varepsilon \sim 1/\sqrt{n}$  and a bound

$$E \left[ \sup_{f \in \mathcal{F}} |\hat{R}(f) - R(f)| \right] = O\left(\sqrt{\left(\frac{d}{n}\right) \log n}\right)$$

$\star \rightarrow$  Dependency on  $d$ ; is this real or artifact of the proof?

$\rightarrow l(y, f(x))$  bounded for e.g. squared loss.  
+ Lipschitz.

it is Lipschitz because  $\|f\| \leq \tilde{A}$   
and  $(y, x)$  is bounded a.s.

A: we have been a bit sloppy.

$\rightarrow$  VC dimension is one tool that gives a generalised notion of cardinality, via the "shattering" notion.

$\hookrightarrow$  Not very popular these days to analyse large overparametrised models such as NNs.

$\rightarrow$  Instead, we will use Rademacher Complexity

$\rightarrow$  let  $Z = (x, y) \in \mathbb{Z}$  and consider a class  $\mathcal{H}$  of functions from  $\mathbb{Z}$  to  $\mathbb{M}$

functions from  $\mathcal{Z}$  to  $\mathbb{R}$ .

$$H = \{(x, y) \mapsto l(y, f(x)) \text{ for } f \in \mathcal{F}\}.$$

→ we want to control  $\sup_{f \in \mathcal{F}} R(f) - \hat{R}(f) =$

$$= \sup_{h \in H} \left\{ |E[h(z)] - \frac{1}{n} \sum_{i=1}^n h(z_i)| \right\}$$

Def: (Rademacher Complexity)  $R_n(H)$  given data

$D = \{z_1 \dots z_n\}$  is

$$R_n(H) := E_{D, \varepsilon} \left( \sup_{h \in H} \frac{1}{n} \sum_{i=1}^n \varepsilon_i \cdot h(z_i) \right)$$

where  $\varepsilon_1 \dots \varepsilon_n$  are Rademacher independent r.v.

$$\varepsilon_i = \begin{cases} +1 & \text{with prob } 1/2 \\ -1 & 1/2 \end{cases}$$

↳ Measuring how well can functions in  $H$  correlate with random labels.

Proposition (Symmetrisation) we have

$$|E \left[ \sup_{h \in H} \left( \frac{1}{n} \sum_i h_i(z_i) - E[h(z)] \right) \right] \leq 2 \cdot R_n(H)$$

and  $|E \left[ \sup_{h \in H} \left( E[h(z)] - \frac{1}{n} \sum_i h_i(z_i) \right) \right] \leq 2 \cdot R_n(H)$

Proof: let  $D' = \{z'_1 \dots z'_n\}$  independent copy of  $D$ .

$\varepsilon_1 \dots \varepsilon_n$  iid Rademacher r.v.s.

$$\mathbb{P}[h(z'_1, \dots, z'_n) - h(z_1, \dots, z_n)]$$

$$(E(z' | D) - E(z))$$

$$\begin{aligned}
& \mathbb{E} \left[ \sup_{h \in H} \left( \mathbb{E} h(z) - \frac{1}{n} \sum_{i=1}^n h(z_i) \right) \right] = \dots \\
& = \mathbb{E} \left[ \sup_{h \in H} \left( \frac{1}{n} \sum_{i=1}^n (\mathbb{E}(h(z'_i) | D) - h(z'_i)) \right) \right] \\
& \left( \mathbb{E}(h(z'_i) | D) \right) \leftarrow = \mathbb{E} \left[ \sup_{h \in H} \left( \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left( h(z'_i) - h(z_i) \mid D \right) \right) \right] \\
& \sup \mathbb{E} \leq \mathbb{E} \sup \leftarrow \text{(≤)} \mathbb{E} \left[ \mathbb{E} \left( \sup_{h \in H} \left( \frac{1}{n} \sum_{i=1}^n h(z'_i) - h(z_i) \right) \right) \mid D \right] \\
& = \mathbb{E} \sup_{h \in H} \frac{1}{n} \sum_{i=1}^n (h(z'_i) - h(z_i)) \\
& \xrightarrow{\substack{\text{symmetrization} \\ \text{of } P, D \\ \text{and } L_i}} \mathbb{E} \sup_{h \in H} \frac{1}{n} \sum_{i=1}^n \varepsilon_i (h(z'_i) - h(z_i)) \\
& \leq 2 \mathbb{E} \sup_{h \in H} \frac{1}{n} \sum_{i=1}^n \varepsilon_i h(z_i) = 2 R_n(H) \quad \square
\end{aligned}$$

$(2 \text{Var}(X) = \mathbb{E} \|X - X'\|^2 \text{ with } X' \text{ an independent copy of } X)$

→ An important property of Rademacher complexity is that it can be easily composed with Lipschitz maps:

Prop (Contraction Principle of Rademacher complexity)

Given any functions  $b, a_i : \Theta \rightarrow \mathbb{R}$  and  $\ell_i : \mathbb{R} \rightarrow \mathbb{R}$   $\beta$ -Lipschitz, we have

$$\begin{aligned}
& \mathbb{E}_\varepsilon \left[ \sup_{\theta \in \Theta} b(\theta) + \sum_{i=1}^n \varepsilon_i \ell_i(a_i(\theta)) \right] \leq \\
& \leq \beta \mathbb{E}_\varepsilon \left[ \sup_{\theta \in \Theta} b(\theta) + \sum_{i=1}^n \varepsilon_i a_i(\theta) \right]
\end{aligned}$$

→ For us, this means that if  $\tilde{y} \mapsto l(y, \tilde{y})$  is  $\beta$ -Lipschitz

a.s., then

$$\begin{aligned} \mathbb{E}_\epsilon \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i l(y_i, f(x_i)) \mid D \right] &\leq \\ &\leq \beta \cdot \mathbb{E}_\epsilon \left[ \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \epsilon_i f(x_i) \right] \end{aligned}$$

In other words,  $R_n(H) \leq \beta \cdot R_n(\mathcal{F})$ .

Example: Linear models with norm constraint.  $\hookrightarrow L_p$  norm

$$\mathcal{F} = \{ f_\theta(x) = \langle \theta, \phi(x) \rangle ; \text{ with } \|\theta\|_p \leq D \}$$

with  $\phi(x) \in \mathbb{R}^d$  with  $d \gg 1$

$$\text{let's write } \Phi = (\phi(x_1) \dots \phi(x_n)) \in \mathbb{R}^{n \times d}$$

let's now compute  $R_n(\mathcal{F})$ :

$$R_n(\mathcal{F}) = \mathbb{E}_{\epsilon, D} \left[ \sup_{\|\theta\|_p \leq D} \left( \frac{1}{n} \sum_{i=1}^n \epsilon_i \theta^\top \phi(x_i) \right) \right] =$$

$\epsilon = (\epsilon_1 \dots \epsilon_n)$

$$= \mathbb{E} \sup_{\|\theta\|_p \leq D} \frac{1}{n} \theta^\top \Phi \epsilon$$

$$= \frac{D}{n} \underbrace{\mathbb{E} \sup_{\|\theta\|_p \leq 1} \langle \theta, \phi^\top \epsilon \rangle}_{\substack{\text{Dual} \\ \text{Norm}}} = \frac{D}{n} \|\phi^\top \epsilon\|_p^* =$$

$$\left( \|x\|^* = \sup_{\|\theta\| \leq 1} \langle \theta, x \rangle \right)$$

$$= \frac{D}{n} \|\phi^\top \epsilon\|_p \quad \text{with } \frac{1}{p} + \frac{1}{q} = 1$$

$\rightarrow$  When  $p=2$ , then  $q=2$ ,

$$\frac{D}{n} \mathbb{E} \| \phi^T \epsilon \|_2 \stackrel{\text{Jensen}}{\leq} \frac{D}{n} \sqrt{\mathbb{E} \| \phi^T \epsilon \|_2^2}$$

Jensen

$$= \frac{D}{n} \sqrt{\mathbb{E} \langle \phi \phi^T, \epsilon \epsilon^T \rangle}$$

$\text{Cor}(\epsilon) = \text{Id.}$

$$\begin{aligned} \textcircled{=} \quad & \frac{D}{n} \sqrt{\mathbb{E} \| \phi \|^2} = \frac{D}{n} \sqrt{\mathbb{E} \sum_{i=1}^n \| \phi(x_i) \|^2} \\ & = \frac{D}{\sqrt{n}} \sqrt{\mathbb{E}_x \| \phi(x) \|^2} \end{aligned}$$

$\hookrightarrow$  There is no longer a dependency on dimension!

$\hookrightarrow$  Norm-constrained hypothesis spaces are very general.

Summary so far :  $\left\{ \begin{array}{l} \text{generalization bounds that are} \\ \text{(i) non-asymptotic} \\ \text{(ii) in high-probability (not just in} \\ \text{expectation).} \end{array} \right.$

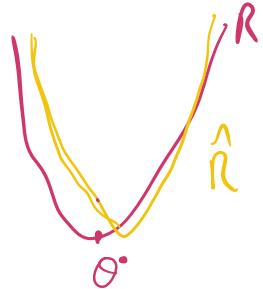
Q: What about asymptotic analysis?

$$\theta \in \mathbb{R}^d \quad R(\theta) = \mathbb{E} [\ell(y, f_\theta(x))], \quad \hat{R}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_\theta(x_i))$$

ensure  $R, \hat{R}$  are convex wrt  $\theta$   
(+  $R$  is strictly convex)

$\rightarrow \theta^* = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} R(\theta)$ , with Hessian  $D^2 R(\theta^*) > 0$   
 $(DR(\theta^*) = 0)$ .

$$\rightarrow \boxed{\hat{\theta}_n \in \underset{\theta}{\operatorname{argmin}} \hat{R}(\theta)}$$



since  $DR(\theta^*) = 0$

$$D\hat{R}_n(\theta^*) = \frac{1}{n} \sum_{i=1}^n D\ell(y_i, f_\theta(x_i)) \Big|_{\theta=\theta^*}$$

} by the LLN  $D\hat{R}(\theta^*) \rightarrow 0$

} by the CLT  $D\hat{R}(\theta^*) \sim N(0, \frac{1}{n} G(\theta^*))$

$$G(\theta^*) = E[D\ell(y, f_\theta(x))] \cdot D\ell(y, f_\theta(x)) \Big|_{\theta=\theta^*}$$

$\hat{\theta}_n \rightarrow \theta^*$ . thanks to strict convexity.

$\rightarrow$  Taylor expansion of  $D\hat{R}$  around  $\theta^*$ :

$$\left. \begin{aligned} 0 &= D\hat{R}(\theta_n) \equiv D\hat{R}(\theta^*) + D^2\hat{R}(\theta^*)(\theta_n - \theta^*) \\ &\text{because } \theta_n \text{ is a minm} \end{aligned} \right]$$

by LLN, we also have

$$D^2\hat{R}(\theta^*) \rightarrow D^2 R(\theta^*)$$

and thus

$$\underline{\theta_n - \theta^*} \approx [D^2 R(\theta^*)]^{-1} \cdot \underline{D\hat{R}(\theta^*)}$$

$\rightarrow$  since  $D\hat{R}(\theta^*)$  is asymp: Normal, we also have

that  $\underline{\theta_n - \theta^*}$  is approx normal, with mean 0 and

Gaussian:

$$\frac{1}{n} D^2 R(\theta^*)^{-1} G(\theta^*) D^2 R(\theta^*)^{-1}$$

→ We can use this asymptotic Gaussianity to measure excess risk:

$$R(\hat{\theta}_n) = R(\theta^*) + \langle DR(\theta^*), \hat{\theta}_n - \theta^* \rangle + \frac{1}{2} (\hat{\theta}_n - \theta^*)^T D^2 R(\theta^*) (\hat{\theta}_n - \theta^*)$$

$$[E(R(\hat{\theta}_n) - R(\theta^*))] \approx \frac{1}{n} \text{Tr}(D^2 R(\theta^*)^{-1} \cdot G(\theta^*))$$

→ NOT an upper bound as before, but the actual expect. value!

→ Fast rate is  $O(1/n)$  versus  $O(1/\sqrt{n})$  !!

↳ In essence, this is a local analysis.