

Lecture 7: Dudley's integral and chaining

1 / 18

Announcements and plan

- Project proposals due next Tuesday **24.10.**, send to Konstantin and supervisor
- One page is enough, instructions on project website (plan how you split up work among the group)

Plan today

- Pollard: One-step discretization → Finer argument via Dudley's integral: Chaining
- Moving from classification to (non-parametric) regression

2 / 18

Recap: Metric entropy to bound excess risk

- Excess risk $R(\hat{f}_n) - R(f^*)$ bounded by generalization gap and standard concentration terms.
- For bounded losses, generalization gap $R(\hat{f}_n) - R_n(\hat{f}_n)$ is bounded by Rademacher complexity w.h.p.
- Can bound (population) R.C. via sup of empirical R.C.
- View the empirical R.C. as expected supremum of **subgaussian process** $X_\theta := \frac{1}{\sqrt{n}} \langle \epsilon, \theta \rangle$ for Rademacher vector ϵ and $\theta \in \mathcal{H}(x_1^n) = \{(h(x_1), \dots, h(x_n)) \mid h \in \mathcal{H}\}$
- Bounded this expectation using the covering number (Pollard's bound)

3 / 18

Recap: Covering number

Proposition (using Pollard's bound - MW Prop 5.17)

Let $\delta > 0$. If a set of points $\theta^1, \dots, \theta^N$ is a covering of \mathbb{T} in the metric $\rho = \frac{\|\cdot\|_2}{\sqrt{n}}$, i.e. it satisfies $\min_j \rho(\theta, \theta^j) \leq \delta$ for all $\theta \in \mathbb{T}$ and $\sup_{\theta, \theta' \in \mathbb{T}} \rho(\theta, \theta') \leq \sigma$, then we have

$$\tilde{\mathcal{R}}_n(\mathbb{T}) \leq \frac{1}{\sqrt{n}} \mathbb{E} \sup_{\theta, \theta' \in \mathbb{T}} X_\theta - X_{\theta'} \leq 2\left[\delta + 2\sigma \sqrt{\frac{\log N(\delta)}{n}}\right]$$

This bound holds in particular for the covering number

Definition (covering number, metric entropy)

For a metric ρ let the ϵ -covering number $\mathcal{N}(\epsilon; \mathbb{T}, \rho)$ be the smallest N such that a set of N points $S = \{\theta_i\}_{i=1}^N$ satisfies $\max_{\theta \in \mathbb{T}} \min_i \rho(\theta_i, \theta) \leq \epsilon$ (S is ϵ -cover). The *metric entropy* is $\log \mathcal{N}(\epsilon; \mathbb{T}, \rho)$.

4 / 18

Recap: Examples

Example I: Smoothly parameterized function class \mathcal{H}_1 with h s.t.

$$\sup_z |h(z; u) - h(z; u')| \leq L \|u - u'\|_2$$

where $u \in \mathbb{B}_2(1) \subset \mathbb{R}^d$ is the 2-norm ball of radius 1.

Covering number: order $\log(1 + \frac{L}{\delta})$ and $\mathcal{R}_n(\mathcal{H}_1) \leq O(\sqrt{\frac{d \log n}{n}})$.

Example II: Smooth non-parametric function classes \mathcal{H}_2^α with $h : [0, 1] \rightarrow \mathbb{R}$ s.t. $|h^{(\alpha)}(x) - h^{(\alpha)}(x')| \leq L|x - x'|$

For $\alpha = 0$, covering number: order $\frac{L}{\delta}$ and $\mathcal{R}_n(\mathcal{H}_2^0) \leq O(n^{-1/3})$.

For general α we have $\mathcal{R}_n(\mathcal{H}_2^\alpha) \leq O(n^{-\frac{1}{2} \frac{(2\alpha+2)}{(2\alpha+3)}})$ (MW Ex. 5.10., 5.11. and 5.21).

Can check for yourself in both cases that the diameter $\sup_{\theta, \theta' \in \mathbb{T}} \frac{\|\theta - \theta'\|_2}{\sqrt{n}}$ is bounded by a constant

5 / 18

Metric entropy refinement: chaining

- Remember Pollard's bound with $D = \sup_{\theta, \tilde{\theta} \in \mathbb{T}} \rho(\theta, \tilde{\theta})$

$$\tilde{\mathcal{R}}_n(\mathbb{T}) \leq \frac{2}{\sqrt{n}} \inf_{\delta > 0} [\delta \sqrt{n} + 2D \sqrt{\log N(\delta)}]$$

- For the last term we're combining a large D with a small δ (hence big $N(\delta)$) \rightarrow lose lose.
- Intuitive question: can we use a finer argument such that small δ is paired with big $N(\delta)$?

Theorem (Dudley's entropy integral - MW Thm 5.22.)

Let $\{X_\theta, \theta \in \mathbb{T}\}$ be a zero-mean subgaussian process wrt some metric ρ . Define $D = \sup_{\theta, \tilde{\theta} \in \mathbb{T}} \rho(\theta, \tilde{\theta})$. Then for any $\delta \in [0, D]$ we have

$$\mathbb{E} \max_{\theta, \tilde{\theta} \in \mathbb{T}} X_\theta - X_{\tilde{\theta}} \leq 2 \mathbb{E} \sup_{\gamma, \gamma' : \rho(\gamma, \gamma') \leq \delta} X_\gamma - X_{\gamma'} + 16 \int_{\delta/4}^D \sqrt{\log \mathcal{N}(t; \mathbb{T}, \rho)} dt$$

Re Tightness: for non-decreasing functions Pollard's bound yields $O((\frac{\log n}{n})^{1/3})$ vs. Dudley: $O((\frac{\log n}{n})^{1/2})$ (exercise, nontrivial)

6 / 18

Example of using Dudley for Lipschitz functions

Remember the examples of the parametric and non-parametric function classes.

Example I: Smoothly parameterized function class \mathcal{H}_1 with h s.t.

$$\sup_z |h(z; u) - h(z; u')| \leq \|u - u'\|_2$$

where $u \in \mathbb{B}_2(1) \subset \mathbb{R}^d$ is the 2-norm ball of radius 1.

The covering number is of order $d \log(\frac{1}{\delta})$.

Example II: Smooth non-parametric function classes \mathcal{H}_2^0 with $h : [0, 1]^d \rightarrow \mathbb{R}$ s.t. $|h(x) - h(x')| \leq \|x - x'\|_\infty$.

The covering number is of order $(\frac{1}{\delta})^d$.

With your neighbor: Use these approximate covering numbers to compute an upper bound for the Rademacher complexity using Dudley's entropy integral and compare the rates obtained using Pollard's bound (focus on $d = 1$ first)

7 / 18

Solution for Example II

Note that we want to find the infimum over δ of the upper bound $\tilde{\mathcal{R}}_n(\mathbb{T}) \leq \frac{2}{\sqrt{n}} \inf_{\delta > 0} [\delta \sqrt{n} + 16 \int_{\delta/4}^D \sqrt{\log \mathcal{N}(t; \mathbb{T}, \rho)} dt]$ where we used the same argument to bound $\mathbb{E} \sup_{\gamma, \gamma': \rho(\gamma, \gamma') \leq \delta} X_\gamma - X_{\gamma'}$ as in Pollard's bound. We are going to ignore constants in almost all steps. . .

Primarily, we need to 1) compute the integral and 2) since the two terms have opposite tendencies when δ decreases, set both terms to be of equal order.

- For $d \leq 2$, it suffices to upper bound the integral by

$$\int_0^D \sqrt{\log \mathcal{N}(t; \mathbb{T}, \rho)} dt = \int_0^D t^{-d/2} dt \leq \begin{cases} 2\sqrt{D} & d = 1 \\ \log D & d = 2 \end{cases}. \text{ No matter}$$

how small we choose δ , we will get a bound of order $\frac{1}{\sqrt{n}}$.

- For $d > 2$, we use a more fine-grained upper bound of $\int_{\delta/4}^D \sqrt{\log \mathcal{N}(t; \mathbb{T}, \rho)} dt = \int_{\delta/4}^D t^{-d/2} dt \leq c(\frac{d}{2} - 1)^{-1} L^{d/2} \delta^{-d/2+1}$ and choosing $\delta = O(n^{-\frac{1}{d}})$ makes both terms of equal order.

8 / 18

Proof of Dudley's integral: Part I

Define shorthand $N_{\mathbb{T}}(\delta) := \mathcal{N}(\delta; \mathbb{T}, \rho)$

- Define $L = \lceil \log_2 \frac{D}{\delta} \rceil$ sets of $\delta_i = D2^{-i}$ covers \mathcal{C}_i of \mathbb{T} with $|\mathcal{C}_i| = N_{\mathbb{T}}(\delta_i)$. The finest cover (original/smallest δ) is \mathcal{C}_L .

- Remember the one-step discretization for Pollard's bound:

$$\begin{aligned} X_{\theta} - X_{\tilde{\theta}} &= X_{\theta} - X_{\theta_{\star}^{(L)}} + X_{\theta_{\star}^{(L)}} - X_{\tilde{\theta}_{\star}^{(L)}} + X_{\tilde{\theta}_{\star}^{(L)}} - X_{\tilde{\theta}} \\ &= 2 \sup_{\rho(\gamma, \gamma') \leq \delta} X_{\gamma} - X_{\gamma'} + \max_{\theta, \theta' \in \mathcal{C}_L} X_{\theta} - X_{\theta'} \end{aligned}$$

where $\theta_{\star}^{(i)}$ denotes closest point of θ in \mathcal{C}_i .

- We can now “recursively” act on $\max_{\theta, \theta' \in \mathcal{C}_L} X_{\theta} - X_{\theta'}$ by using the same argument on the set \mathcal{C}_L with the coarser cover \mathcal{C}_{L-1} .

More generally for any two $\theta, \tilde{\theta} \in \mathcal{C}_i$ we have:

$$\begin{aligned} X_{\theta} - X_{\tilde{\theta}} &\leq X_{\theta} - X_{\theta_{\star}^{(i-1)}} + X_{\theta_{\star}^{(i-1)}} - X_{\tilde{\theta}_{\star}^{(i-1)}} + X_{\tilde{\theta}_{\star}^{(i-1)}} - X_{\tilde{\theta}} \\ &\leq 2 \max_{\theta \in \mathcal{C}_i} X_{\theta} - X_{\theta_{\star}^{(i-1)}} + \max_{\theta, \theta' \in \mathcal{C}_{i-1}} X_{\theta} - X_{\theta'} \end{aligned}$$

9 / 18

Proof of Dudley's integral: Part II

- note that in $\max_{\theta \in \mathcal{C}_i} X_{\theta} - X_{\theta_{\star}^{(i-1)}}$, for each $\theta \in \mathcal{C}_i$ we have $\theta_{\star}^{(i-1)}$ be **its** closest point, not of the “original” θ in \mathbb{T}
- “Rolling out” the induction, we obtain

$$\max_{\theta, \tilde{\theta} \in \mathcal{C}_L} X_{\theta} - X_{\tilde{\theta}} \leq 2 \sum_{i=2}^L \max_{\theta \in \mathcal{C}_i} X_{\theta} - X_{\theta_{\star}^{(i-1)}} + \max_{\theta, \theta' \in \mathcal{C}_1} X_{\theta} - X_{\theta'}$$

Rolling out from $L \rightarrow 1$ or going from \mathcal{C}_L to \mathcal{C}_1 , we iteratively

- reduced the cover cardinality until only one element is left (with large diameter),
- while all the intermediate terms (in sum) are δ_{i-1} -subgaussian (instead of fixed D)
- with increasing δ but decreasing corresponding cover cardinality

10 / 18

Proof of Dudley's integral: Part III

In order to compute the final expectation observe that

1. max of subgaussians: $X_\theta - X_{\theta_\star^{(i-1)}}$ is a δ_{i-1} -subgaussian process \rightarrow

$$\mathbb{E} \max_{\theta \in \mathcal{C}_i} X_\theta - X_{\theta_\star^{(i-1)}} \leq 2\delta_{i-1} \sqrt{\log |\mathcal{C}_i|}$$

2. Covering number non-increasing as δ increases and interval $[D2^{-(i+1)}, D2^{-i}]$ is of length $D2^{-(i+1)} = D2^{-(i-1)} \frac{1}{4}$:

$$\delta_{i-1} \sqrt{\log |\mathcal{C}_i|} = D2^{-(i-1)} \sqrt{\log N_{\mathbb{T}}(D2^{-i})} \leq 4 \int_{D2^{-(i+1)}}^{D2^{-i}} \sqrt{\log N_{\mathbb{T}}(t)} dt$$

3. Putting things together and because $\delta_L = D2^{-L} \leq \delta$

$$\begin{aligned} \mathbb{E} \max_{\theta, \tilde{\theta} \in \mathcal{C}_L} X_\theta - X_{\tilde{\theta}} &\leq 4 \sum_{i=2}^L D2^{-(i-1)} \sqrt{\log N_{\mathbb{T}}(D2^{-i})} + 2D \sqrt{\log N_{\mathbb{T}}(D/2)} \\ &\leq 16 \int_{\delta/4}^D \sqrt{\log N_{\mathbb{T}}(t)} dt \end{aligned} \quad \square$$

11 / 18

Short navigation slide

Whole topic of this class: For each \mathcal{F} define $f^\star = \arg \min_{f \in \mathcal{F}} R(f)$.
Interested in bounding **excess risk** w.h.p.

$$R(\hat{f}_n) - R(f^\star) = R(\hat{f}_n) - R_n(\hat{f}_n) + \overbrace{R_n(\hat{f}_n) - R_n(f^\star)}^{\leq 0 \text{ by optimality}} + R_n(f^\star) - R(f^\star)$$

- so far: via **uniform convergence** and **Rademacher complexity** using

$$\mathbb{P}(\sup_{h \in \mathcal{H}} \mathbb{E} h(Z) - \frac{1}{n} \sum_{i=1}^n h(Z_i) \geq 2\mathcal{R}_n(\mathcal{H}) + t) \leq e^{-\frac{nt^2}{2b^2}}$$

for $\mathcal{H} = \ell \circ \mathcal{F}$ and bounding empirical Rademacher complexity for finite classes, more generally w/ **metric entropy** and **chaining** (today)

This line of reasoning was useful for **classification**, for the second half of lectures, we'll switch to **regression**. Can we just continue to use this uniform convergence technique to obtain bounds?

12 / 18

(Non-)parametric regression setting - fixed design

- Square loss and constrained regression
- Fixed design, i.e. only care about prediction on training inputs x_1, \dots, x_n
- Gaussian observation noise, i.e. $W = Y - f^*(X) \sim \mathcal{N}(0, \sigma^2)$
- Analyze minimizer $\hat{f} = \arg \min_{f \in \mathcal{F}} R_n(f) := \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2$ or with penalty $\hat{f} = \arg \min_{f \in \mathcal{F}} R_n(f) := \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \|f\|_{\mathcal{F}}$
- Evaluation: Prediction error of some f on fixed design points

$$\|f - f^*\|_n^2 = \frac{1}{n} \sum_{i=1}^n (f(x_i) - f^*(x_i))^2 = \mathbb{E}_Y R_n(f) - \sigma^2 = R(f) - R(f^*)$$

Partner-Q: Derive a h.p. upper bound for $\|f - f^*\|_n^2$ for linear functions $f(x) = \langle w, x \rangle$ with $\|x\|_2 \leq D$, $\|w\|_2 \leq B$. Further assume the noise is bounded. Compare a closed-form vs. a uniform law approach - where might the difference come from? For solution see [Lecture 10](#)

13 / 18

Warm-up using closed-form solution - linear regression

For linear/kernel regression, can directly analyze closed-form solution of both ridge and min-norm interpolator. For linear:

- first recall $y = X\theta^* + w$ and solution $\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^d} \|y - X\theta\|_2^2$
- minimizer $\hat{f}(x) = \hat{\theta}^\top x$ with $\hat{\theta} = (X^\top X)^{-1} X^\top (X\theta^* + w)$
- $\|\hat{f} - f^*\|_n^2 = \frac{1}{n} \|X(\hat{\theta} - \theta^*)\|^2 = \frac{1}{n} w^\top X (X^\top X)^{-1} X^\top w$
- only need to bound $\frac{1}{n} w^\top X (X^\top X)^{-1} X^\top w \rightarrow$ use that the norm of a Gaussian is a Lipschitz function of Gaussian for concentration (here with Lipschitz constant $\sqrt{\frac{\text{rank}(X)}{n}}$ via SVD) and MW Thm 2.26
- Further $\mathbb{E} \frac{1}{n} w^\top X (X^\top X)^{-1} X^\top w = \sigma^2 \frac{\text{rank}(X)}{n}$

This stands in contrast to the uniform law approach where you can use contraction to obtain a bound using Rademacher complexity of linear function classes and at most get a $\frac{1}{\sqrt{n}}$ bound

14 / 18

Beyond closed-form solutions

- First of all, notice the “slow” uniform excess risk bound holds for any \mathcal{F} , including ones for which $f^* \notin \mathcal{F}$!
- Further, in our argument using uniform law, we used optimality of \hat{f}_n only once

$$R(\hat{f}_n) - R(f^*) = R(\hat{f}_n) - R_n(\hat{f}_n) + \overbrace{R_n(\hat{f}_n) - R_n(f^*)}^{\leq 0 \text{ by optimality}} + R_n(f^*) - R(f^*)$$

Next few classes: using *localized complexities* to prove tighter bounds for particular estimator: global minimizer of square loss for regression!

- Idea: By using **optimality of \hat{f}** instead of uniform bound
 1. circumvent uniform boundedness
 2. can get more restricted function space

15 / 18

Basic inequality circumventing boundedness and more

Optimality of \hat{f} yields the *basic inequality*

$$R_n(\hat{f}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \leq \frac{1}{n} \sum_{i=1}^n (y_i - f^*(x_i))^2 = R_n(f^*) \tag{1}$$

$$\|\hat{f} - f^*\|_n^2 \leq \frac{2\sigma}{n} \sum_{i=1}^n w_i (\hat{f}(x_i) - f^*(x_i))$$

- Taking expectations defining $\mathcal{F}^* = \mathcal{F} - f^*$
 $\rightarrow \mathbb{E} \|\hat{f} - f^*\|_n^2 \leq 2\sigma \tilde{\mathcal{G}}_n(\mathcal{F}^*(x_1^n)) := \mathbb{E}_w \sup_{g \in \mathcal{F}^*} \frac{2\sigma}{n} \sum_{i=1}^n w_i g(x_i)$
- Gaussian complexity popped out without needing uniform boundedness (same “order” as Radmacher, satisfies sandwich relationship, proved in HW 2, for each \mathbb{T})

$$\frac{1}{2 \log n} \tilde{\mathcal{G}}_n(\mathbb{T}) \leq \tilde{\mathcal{R}}_n(\mathbb{T}) \leq \sqrt{\frac{\pi}{2}} \tilde{\mathcal{G}}_n(\mathbb{T})$$

- **But still stuck with a huge function space \mathcal{F} !**

The trick is to notice eq. 1 restricts function space!

16 / 18

Motivation for localized Gaussian complexity

- Define $\hat{\Delta} = \hat{f} - f^*$ for simplicity and the space $\mathcal{F}^* = \{f - f^* : f \in \mathcal{F}\}$
 - Furthermore we assume that \mathcal{F}^* is **star-shaped**, i.e. for any $f \in \mathcal{F}^*$, we have $\alpha f \in \mathcal{F}^*$ for all $\alpha \in [0, 1]$
1. Space to control is smaller than all of \mathcal{F}^* since either
 - (i) $\|\hat{\Delta}\|_n \leq \delta_n$ or
 - (ii) if $\|\hat{\Delta}\|_n \geq \delta_n$ then still $\|\hat{\Delta}\|_n^2 \leq \frac{2\sigma}{n} \sum_{i=1}^n w_i \hat{\Delta}(x_i)$ by basic inequality
 2. Further for case (ii), if can show w.h.p.

$$\frac{2\sigma}{n} \sum_{i=1}^n w_i \hat{\Delta}(x_i) \leq 4\|\hat{\Delta}\|_n \delta_n \quad (2)$$

for all $\|\hat{\Delta}\|_n \geq \delta_n$ then we can plug that into RHS of (ii) to obtain $\|\hat{\Delta}\|_n \leq 4\delta_n$ w.h.p.

to be continued...

17 / 18

References

Dudley's integral

- MW Chapter 5

Non-parametric regression

- MW Chapter 13

18 / 18