

Mathematical Foundations of ML (CS 4785/5783)

Lecture 3

Uniform Convergence, Symmetrization and Rademacher
Complexity

<http://www.cs.cornell.edu/Courses/cs4783/2022sp/notes03.pdf>

STATISTICAL LEARNING FRAMEWORK

D is a distribution on $\mathcal{X} \times \mathcal{Y}$

D captures the idea of this set **U**

Training sample $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ Each $(x_t, y_t) \sim \mathbf{D}$

Risk of a model g defined as $L_{\mathbf{D}}(g) = \mathbb{E}_{(x,y) \sim \mathbf{D}} [\ell(g(x), y)]$

(Future instances drawn from **D**)

Excess risk of model g w.r.t. model class \mathcal{F} defined as

$$L_{\mathbf{D}}(g) - \min_{f \in \mathcal{F}} L_{\mathbf{D}}(f)$$

Goal: provide an algorithm for which excess risk is small

EMPIRICAL RISK MINIMIZATION

Pick a model in class that minimizes training error

$$\hat{f}_{\text{ERM}} \in \arg \min_{f \in \mathcal{F}} \hat{L}_S(f)$$

- When does this succeed?
 - When model class is too complex, we already saw this can fail
 - When model class is say just one function, it succeeds due to law of large numbers (concentration)
 - In general how well does this algorithm do?

ERM OVER FINITE CLASS

If losses are bounded by 1 (in absolute) and $|\mathcal{F}| < \infty$, then, for any $\delta > 0$ with probability at least $1 - \delta$,

$$L_{\mathbf{D}}(\hat{f}_{\text{ERM}}) - \min_{f \in \mathcal{F}} L_{\mathbf{D}}(f) \leq \sqrt{\frac{8 \log(2|\mathcal{F}|/\delta)}{n}}$$

ERM OVER FINITE CLASS

Hoeffding Inequality: Let Z_1, \dots, Z_n be a sequence of n random variables bounded by 1, drawn iid from a fixed distribution. Then:

$$P\left(\left|\frac{1}{n}\sum_{t=1}^n Z_t - \mathbb{E}Z\right| > \epsilon\right) \leq 2 \exp\left(-\frac{n\epsilon^2}{2}\right)$$

ERM OVER FINITE CLASS

Hoeffding Inequality: Let Z_1, \dots, Z_n be a sequence of n random variables bounded by 1, drawn iid from a fixed distribution. Then:

$$P\left(\left|\frac{1}{n}\sum_{t=1}^n Z_t - \mathbb{E}Z\right| > \epsilon\right) \leq 2 \exp\left(-\frac{n\epsilon^2}{2}\right)$$

Proof idea:

For each $f \in \mathcal{F}$ define $Z_t^f = \ell(f(x_t), y_t)$

Apply Hoeffding for each f individually

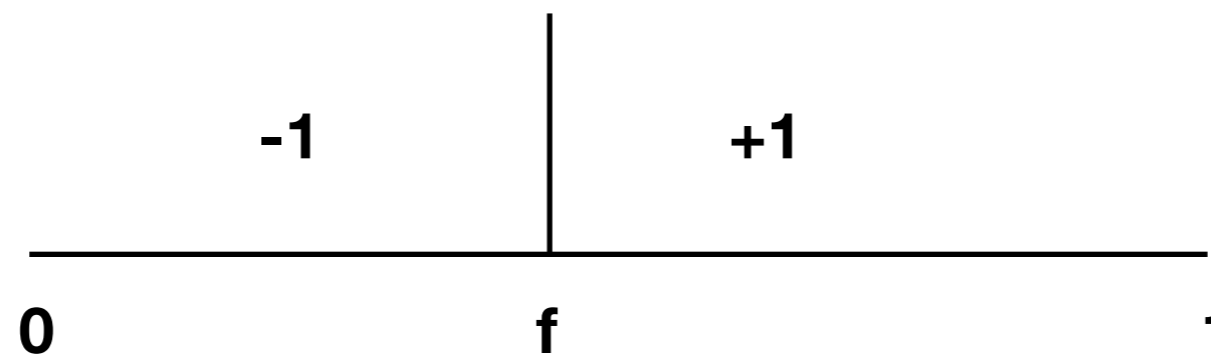
Use union bound to move to uniform deviation

BEYOND FINITE MODEL CLASS

- Idea 1: Find a finite set \mathcal{F}' such that for any $f \in \mathcal{F}$ there exists an $f' \in \mathcal{F}'$ s.t.

$$\forall x, y, \quad |\ell(f'(x), y) - \ell(f(x), y)| < \Delta$$

- But this may not always work, consider the example of learning thresholds:



$$\mathcal{X} = [0, 1] \quad f(x) = \text{sign}(x - f)$$

\mathcal{F} indexed by set $[0, 1]$

For any $\Delta < 1/2$, this class cannot be approximated by a finite set.

UNIFORM CONVERGENCE

We have shown that for any $\epsilon > 0$,

$$P\left(L_{\mathbf{D}}(\hat{f}_{\text{ERM}}) - \min_{f \in \mathcal{F}} L_{\mathbf{D}}(f) > 2\epsilon\right) \leq P\left(\max_{f \in \mathcal{F}} |\hat{L}_S(f) - L_{\mathbf{D}}(f)| > \epsilon\right)$$

Next, we will see that $\max_{f \in \mathcal{F}} |\hat{L}_S(f) - L_{\mathbf{D}}(f)|$ is concentrated near its expectation.

McDIARMID'S INEQUALITY

Let $Z_1, \dots, Z_n \in \mathcal{Z}$ be a sequence of n random variables drawn iid from a fixed distribution. Assume that $\Phi : \mathcal{Z}^n \mapsto \mathbb{R}$ is a function satisfying the condition that: For any $i \in [n]$, and any $z_1, \dots, z_n \in \mathcal{Z}$ and any $z'_i \in \mathcal{Z}$,

$$|\Phi(z_1, \dots, z_i, \dots, z_n) - \Phi(z_1, \dots, z'_i, \dots, z_n)| \leq \frac{C}{n}$$

Then we have the following concentration result :

$$P(|\Phi(Z_1, \dots, Z_n) - \mathbb{E}[\Phi(Z_1, \dots, Z_n)]| > \epsilon) \leq 2 \exp\left(-\frac{2n\epsilon^2}{C^2}\right)$$

UNIFORM CONVERGENCE

Eg: The function $\phi((x_1, y_1), \dots, (x_n, y_n)) = \max_{f \in \mathcal{F}} |\hat{L}_S(f) - L_{\mathbf{D}}(f)|$ satisfies the condition with $C = 2$ when loss is bounded by 1.

UNIFORM CONVERGENCE

Eg: The function $\phi((x_1, y_1), \dots, (x_n, y_n)) = \max_{f \in \mathcal{F}} |\hat{L}_S(f) - L_{\mathbf{D}}(f)|$ satisfies the condition with $C = 2$ when loss is bounded by 1.

Hence we have that for any $\delta > 0$, with probability at least $1 - \delta$,

$$\max_{f \in \mathcal{F}} |\hat{L}_S(f) - L_{\mathbf{D}}(f)| \leq 2 \mathbb{E} \left[\max_{f \in \mathcal{F}} |\hat{L}_S(f) - L_{\mathbf{D}}(f)| \right] + O \left(\sqrt{\frac{\log(1/\delta)}{n}} \right)$$

Complexity Measure

SYMMETRIZATION AND RADEMACHER COMPLEXITY

Let $\epsilon_1, \dots, \epsilon_n \in \{\pm 1\}$ be Rademacher random variables where each ϵ_i is $+1$ with probability $1/2$ and -1 with probability $1/2$.

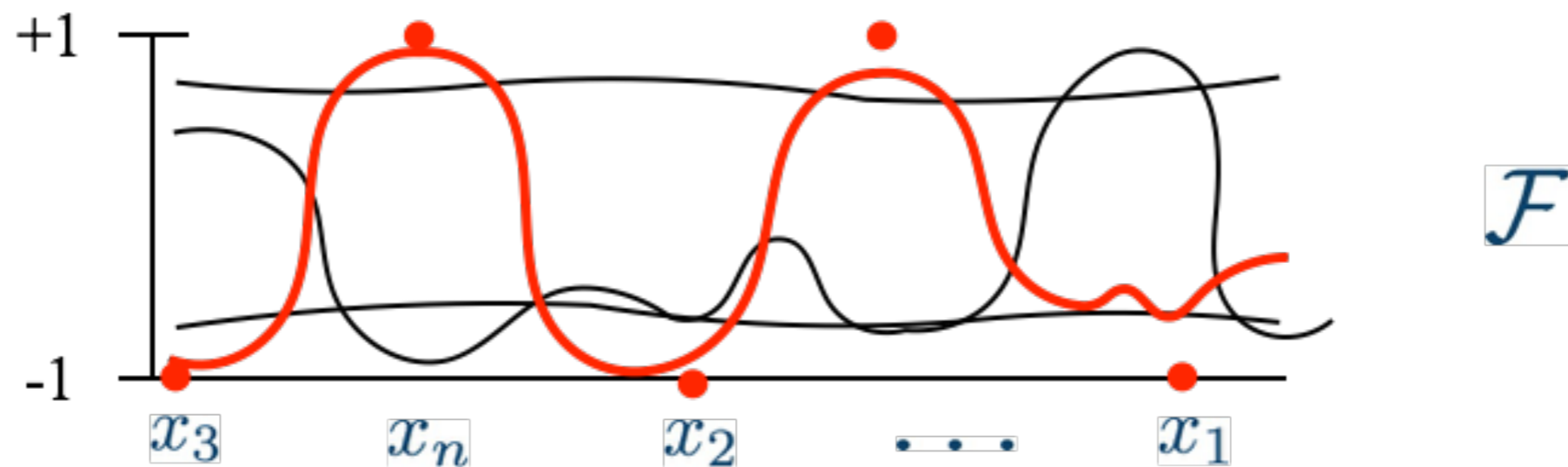
We will see that:

$$\mathbb{E} \left[\max_{f \in \mathcal{F}} |\hat{L}_S(f) - L_{\mathbf{D}}(f)| \right] \leq \frac{2}{n} \mathbb{E}_S \left[\mathbb{E}_{\epsilon} \left[\max_{f \in \mathcal{F}} \left| \sum_{t=1}^n \epsilon_t \ell(f(x_t), y_t) \right| \right] \right]$$

Rademacher Complexity

RADEMACHER COMPLEXITY

Example : $\mathcal{X} = [0, 1]$, $\mathcal{Y} = [-1, 1]$



Proof of the Result

Why is this useful?