

DDA4210/AIR6002 Advanced Machine Learning

Lecture 03 Learning Theory

Tongxin Li

School of Data Science, CUHK-Shenzhen

Spring 2024

Overview

- 1 Introduction
- 2 Empirical Risk Minimization
- 3 Growth Function and VC dimension
- 4 Rademacher Complexity

What is machine learning theory

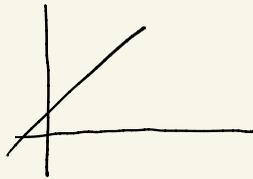
- Machine Learning Theory is also known as *Computational Learning Theory*.
- It aims to understand the fundamental principles of learning as a computational process and combines tools from Computer Science and Statistics.
 - Create mathematical models of machine learning and analyze the inherent ease or difficulty of different types of learning problems.
 - Proving guarantees for algorithms (e.g., under what conditions will they succeed, how much data and computation time is needed)
 - Developing machine learning algorithms that provably meet desired criteria.
 - Mathematically analyzing general issues (e.g., "When can one be confident about predictions made from limited data?", "What kinds of methods can learn even in the presence of large quantities of distracting information?")

A toy example: $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ (Classification Task)

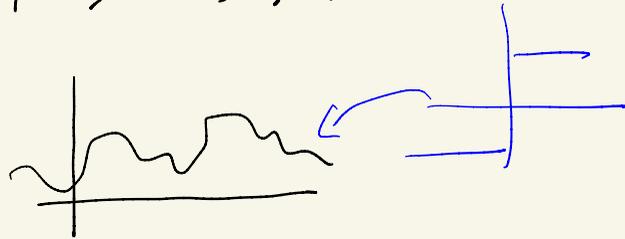
e.g. $\{(2, +1), (1, +1), (0.5, +1), (-2, -1), (-4, -1), (-0.5, -1)\}$

$$(y = \text{sign}(x))$$

- Which function class does y belong to? e.g. linear y ?
- How much data do we need? (sample complexity)
- How to measure the "complexity" of y ?

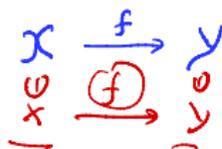


simple?



complicated?

Basic notation



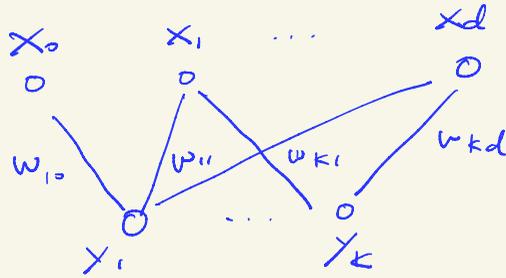
- Input space/feature space : \mathcal{X}
 - Feature is a numerical description for a sample or object.
 - Feature extraction is an art.
 - Output space/label space: \mathcal{Y}
 - E.g.: $\{+1, -1\}$, $\{1, 2, \dots, K\}$, \mathbb{R} -valued output, structured output.
 - Loss function: $\ell : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$
 - E.g.: 0 - 1 loss $\ell(y, \hat{y}) = 1\{y \neq \hat{y}\}$, square loss $\ell(y, \hat{y}) = (y - \hat{y})^2$, absolute loss $\ell(y, \hat{y}) = |y - \hat{y}|$, cross-entropy loss $\ell(y, \hat{y}) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$.
 - It measures performance/cost per instance (e.g., inaccuracy or error of prediction).
 - Model class/hypothesis class: $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$ (or \mathcal{H} or \mathbb{H})
 - E.g.: $\mathcal{F} = \{x \mapsto f^T x : \|f\|_2 \leq 1\}$, $\mathcal{F} = \{x \mapsto \text{sign}(f^T x)\}$
 $\mathcal{F} = \{ \text{continuous } f \}$
- (optional) $\left\{ \begin{array}{l} \mathcal{F} = \text{RKHS (Reproducing kernel Hilbert space)} \\ \mathcal{F} = \text{NTK (Neural Tangent Kernel)} \end{array} \right. \text{ etc}$

Basic notation

$$x \xrightarrow{f} y$$

- Input space/feature space : \mathcal{X}
 - Feature is a numerical description for a sample or object.
 - Feature extraction is an art.
- Output space/label space: \mathcal{Y}
 - E.g.: $\{+1, -1\}$, $\{1, 2, \dots, K\}$, \mathbb{R} -valued output, structured output.
- Loss function: $\ell : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$
 - E.g.: 0 - 1 loss $\ell(y, \hat{y}) = 1\{y \neq \hat{y}\}$, square loss $\ell(y, \hat{y}) = (y - \hat{y})^2$, absolute loss $\ell(y, \hat{y}) = |y - \hat{y}|$, cross-entropy loss $\ell(y, \hat{y}) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$.
 - It measures performance/cost per instance (e.g., inaccuracy or error of prediction).
- Model class/hypothesis class: $\mathcal{F} \subset \mathcal{Y}^{\mathcal{X}}$ (or \mathcal{H} or \mathbb{H})
 - E.g.: $\mathcal{F} = \{x \mapsto f^T x : \|f\|_2 \leq 1\}$, $\mathcal{F} = \{x \mapsto \text{sign}(f^T x)\}$
 $\mathcal{F} = \{ \text{continuous } f \}$
 $\mathcal{F} = \text{RKHS}$ (Reproducing kernel Hilbert space)
 $\mathcal{F} = \text{NTK}$ (Neural Tangent Kernel) etc

$F = \{ \text{a set of single-layer NNs} \}$



$$y = W \cdot x$$

↓

a K -by- $(d+1)$ matrix

(Recall: Hypothesis set H)

Next: Three learning frameworks.

Probably approximately correct (PAC) learning

- Learner only observes training samples

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

• $x_1, x_2, \dots, x_n \sim D_X$, $y_i = f^*(x_i)$, $i = 1, 2, \dots, n$, where $f^* \in \mathcal{F}$.

- Goal: find $\hat{f} \in \mathcal{Y}^X$ to minimize

$$\mathbb{P}_{x \sim D_X} [\hat{f}(x) \neq f^*(x)]$$

generalization error

- **Probably approximately correct (PAC)** [Valiant 1984] learning is a framework for mathematical analysis of machine learning.

Probably Approximately Correct (PAC) Learning

- In PAC learning, the learner receives samples and must select a generalization function (called the hypothesis) from a certain class of possible functions. The goal is that, with high probability ("probably"), the selected function will have low generalization error ("approximately correct"). The learner must be able to learn the concept given any arbitrary approximation ratio, probability of success, or distribution of the samples.

- Sample complexity (definition):

Given $\delta > 0$, $\epsilon > 0$, and sample complexity $n(\epsilon, \delta)$ is the smallest n such that we can always find forecaster \hat{f} s.t. with probability at least $1 - \delta$,

$$\mathbb{P}_{x \sim D_X} \left[\hat{f}(x) \neq \hat{a} \underline{f^*(x)} \right] \leq \epsilon$$

• f^* is optimal!

* [The learner knows that there exists a perfect f^* that generates the label.]

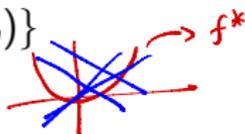
Statistical Learning (agnostic PAC)

- Learner only observes training samples

$$(x_i, y_i) \sim D \quad \underline{\underline{\text{IID}}}$$

$$S = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

drawn iid from joint distribution D on $\mathcal{X} \times \mathcal{Y}$



- Goal: find \hat{f} to minimize expected loss over future instances

$$\mathbb{E}_{(x,y) \sim D}[\ell(\hat{f}(x), y)] - \underbrace{\inf_{f \in \mathcal{F}} \mathbb{E}_{(x,y) \sim D}[\ell(f(x), y)]}_{\text{the best one can do}} \neq f^*$$

- Sample complexity** (definition, denote $L(g) = \mathbb{E}[\ell(g, \cdot)]$):
Given $\delta > 0$, $\epsilon > 0$, and sample complexity $n(\epsilon, \delta)$ is the smallest n such that we can always find forecaster \hat{f} s.t. with probability at least $1 - \delta$,

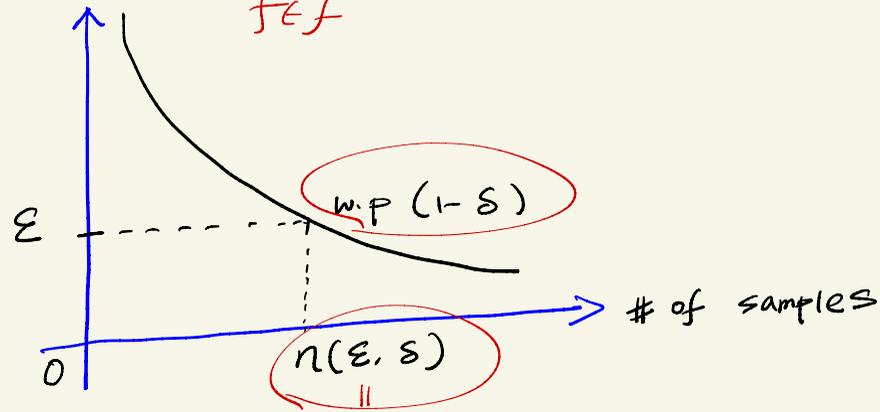
$$L_D(\hat{f}) - \inf_{f \in \mathcal{F}} L_D(f) \leq \epsilon$$

- * The learner doesn't assume that \mathcal{F} contains an error free hypothesis f .

Key difference: Agnostic to (f^*)?

(δ, ϵ) - notion for error analysis

$$L_D(\hat{f}) - \inf_{f \in \mathcal{F}} L_D(f)$$



Online learning

(sequential input)

x_t
 t

- Online learning

For $t = 1$ to n

Learner receives $x_t \in \mathcal{X}$

Learner predicts output $\hat{y}_t \in \mathcal{Y}$, $\hat{y}_t = \hat{f}(x_t)$

True output $y_t \in \mathcal{Y}$ is revealed

EndFor

$(x_1, x_2, \dots, x_t, x_{t+1}, x_t)$
 $(y_1, y_2, \dots, y_{t+1}, ?)$ | x_t

- Goal: minimize regret

$$\text{Reg}_n(\mathcal{F}) := \frac{1}{n} \sum_{t=1}^n \ell(\hat{y}_t, y_t) - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t)$$

fixed

This course will only introduce the learning theory of offline and supervised learning.

Online learning

(sequential input)

- Online learning

For $t = 1$ to n

Learner receives $x_t \in \mathcal{X}$

Learner predicts output $\hat{y}_t \in \mathcal{Y}$, $\hat{y}_t = \hat{f}(x_t)$

True output $y_t \in \mathcal{Y}$ is revealed

EndFor

- Goal: minimize regret

$$\text{Reg}_n(\mathcal{F}) := \frac{1}{n} \sum_{t=1}^n \ell(\hat{y}_t, y_t) - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t)$$

Handwritten notes: "fixed" with arrows pointing to the x_t and y_t terms in the second sum.

Q: What is the underlying assumption on $((x_t, y_t) : t=1 \dots n)$?

A: At each time t , (x_t, y_t) is sampled from a (stationary) distribution.

This course will only introduce the learning theory of offline and supervised learning.

what if D_t is changing over time?

(x_t, y_t) (Dynamic regret) adaptive regret

Minimax Rate

Previously, w.h.p. ... (stochastic models)

- How well does the **best** learning algorithm do in the **worst** case scenario? *Imagine: There is an adversary attacking my model.*

Minimax Rate = "Best Possible Guarantee"

- PAC** framework

$$\mathcal{V}_n^{\text{PAC}}(\mathcal{F}) := \inf_{\hat{f}} \sup_{D_X, f^* \in \mathcal{F}} \mathbb{E}_{S: |S|=n} \left[\mathbb{P}_{X \sim D_X} \left(\hat{f}(X) \neq f^*(X) \right) \right] \quad (1)$$

Handwritten notes: $\{(x_i, y_i) : i=1 \dots n\}$ (with an arrow pointing to the sample set S in the equation); "there exists f^* !" (with an arrow pointing to $f^*(X)$ in the equation).

A problem is "PAC learnable" if $\mathcal{V}_n^{\text{PAC}} \rightarrow 0$ as $n \rightarrow \infty$.

- Statistical learning**

Q. Can we bound $\mathcal{V}_n^{\text{stat}}(\mathcal{F})$? (under some assumptions)

$$\mathcal{V}_n^{\text{stat}}(\mathcal{F}) := \inf_{\hat{f}} \sup_D \mathbb{E}_{S: |S|=n} \left[L_D(\hat{f}) - \inf_{f \in \mathcal{F}} L_D(f) \right] \quad (2)$$

A problem is "statistically learnable" if $\mathcal{V}_n^{\text{stat}} \rightarrow 0$ as $n \rightarrow \infty$.

Empirical Risk Minimization

- Empirical Risk Minimization (ERM): pick the hypothesis from model class \mathcal{F} that best fits the sample, i.e.,

$$\hat{f}_{\text{erm}} \leftarrow \underset{f \in \mathcal{F}}{\text{argmin}} \underbrace{\frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t)}_{\text{empirical risk}} \triangleq R_{\text{emp}}(f) \quad (3)$$

- For a fixed function f , according to the law of large numbers, we have

$$R_{\text{emp}}(f) \rightarrow R_f = \mathbb{E}[\underbrace{\ell(f(x), y)}_{\text{true risk}}] \quad \text{for } n \rightarrow \infty$$

randomness w.r.t. data distribution

Empirical Risk Minimization

- Empirical Risk Minimization (ERM): pick the hypothesis from model class \mathcal{F} that best fits the sample, i.e.,

$$\hat{f}_{\text{erm}} \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \triangleq R_{\text{emp}}(f) \quad (3)$$

- For a fixed function f , according to the law of large numbers, we have

$$R_{\text{emp}}(f) \longrightarrow R_f = \underbrace{\mathbb{E}[\ell(f(x), y)]}_{\text{true risk}} \quad \text{for } n \longrightarrow \infty$$

randomness w.r.t. data distribution

- Bayes optimal function

$$f' := \underset{f}{\operatorname{argmin}} \mathbb{E}[\ell(f(x), y)]$$

$$\Delta := \mathbb{E}[\ell(\hat{f}_{\text{erm}}(x), y)] - \mathbb{E}[\ell(f'(x), y)] \quad \text{"excess risk"}$$

In practice, f' is hard to get.

- 1 Introduction
- 2 Empirical Risk Minimization
- 3 Growth Function and VC dimension
- 4 Rademacher Complexity

→ Just some error benchmarks so far.

1 Introduction ✓ { PAC learning
Statistical learning
online learning

2 Empirical Risk Minimization

3 Growth Function and VC dimension

4 Rademacher Complexity

Error benchmarks

Sample complexity 1

Sample complexity 2

Regret

- There are many other interesting frameworks/models!
- You can define your own as long as it's consistent
- We will be focusing on the statistical learning setting.

Empirical Risk Minimization

- Empirical Risk Minimization (ERM): pick the hypothesis from model class \mathcal{F} that best fits the sample, i.e.,

$$\hat{f}_{\text{erm}} = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \triangleq R_{\text{emp}}(f) \quad (3)$$

Handwritten notes: "this what we do in practice!" with an arrow pointing to the right side of the equation. "empirical risk." with a bracket under the sum.

- For a fixed function f , according to the law of large numbers, we have

$$R_{\text{emp}}(f) \longrightarrow R_f = \underbrace{\mathbb{E}[\ell(f(x), y)]}_{\text{true risk}} \quad \text{for } n \longrightarrow \infty$$

Handwritten notes: "Remember this!" with an arrow pointing to the true risk term.

- Generalization error bound

$$\left| \underbrace{\mathbb{E}[\ell(f(x), y)]}_{\text{test error}} - \underbrace{\frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t)}_{\text{training error}} \right| \leq ?$$

Handwritten notes: "my model" with an arrow pointing to f in a circle. "generalizability" with an arrow pointing to the right side of the inequality. "It reflects the ability of 'generalizing' from training to testing." with an arrow pointing to the right side of the inequality.

- * Connection with Statistical Learning?

Empirical Risk Minimization

- Hoeffding's inequality

- Let X_1, X_2, \dots, X_n be independent random variables.
- Suppose $S_n = X_1 + X_2 + \dots + X_n$ and $a_i \leq X_i \leq b_i \forall i$.

$$P(|S_n - \mathbb{E}[S_n]| \geq \epsilon) \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

Example: $X_i \sim \text{Bernoulli}(p)$

$$b_i = 1$$

$$a_i = 0.$$

$$X_i \begin{cases} p & 1 \\ 1-p & 0 \end{cases}$$

$$\Pr\left(|\sum X_i - n \cdot p| \geq \epsilon\right) \leq 2 \exp\left(-\frac{2\epsilon^2}{n}\right)$$

Empirical Risk Minimization

• Hoeffding's inequality

- Let X_1, X_2, \dots, X_n be independent random variables.
- Suppose $S_n = X_1 + X_2 + \dots + X_n$ and $a_i \leq X_i \leq b_i \forall i$.

$$P(|S_n - \mathbb{E}[S_n]| \geq \epsilon) \leq 2 \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

• Hoeffding's inequality for ERM

- Suppose $\sup_{y, y' \in \mathcal{Y}} |\ell(y, y')| \leq 1$ assume $X_t \geq 0$. data pt

$$P\left(\left|\underbrace{\mathbb{E}[\ell(f(x), y)]}_{\text{true risk}} - \underbrace{\frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t)}_{\text{empirical risk}}\right| \geq \epsilon\right) \leq 2 \exp\left(-\frac{\epsilon^2 n}{2}\right) \quad (4)$$

* What's the drawback of this bound?

$$\frac{1}{n} \mathbb{E}\left[\sum_{t=1}^n \ell(f(x_t), y_t)\right] = \mathbb{E}[\ell(f(x), y)].$$

since the IID assumption

Empirical Risk Minimization

Cardinality: $\mathcal{F} = \{f(x) = ax, x \in \mathbb{R} \mid a = \{1, 2, 3\}\}$ $|\mathcal{F}| = 3$

- Assume $((x_t, y_t) \sim \mathcal{P} : t=1, \dots, n)$ are generated i.i.d from a distribution \mathcal{P}
 - ERM with finite class

Proposition 1

Consider the case when the hypothesis \mathcal{F} has finite cardinality, that is $|\mathcal{F}| < \infty$. For any loss ℓ satisfies $\sup_{y, y' \in \mathcal{Y}} |\ell(y, y')| \leq 1$, we have

$$V_n^{\text{stat}}(\mathcal{F}) \leq \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \leq 8 \sqrt{\frac{\log(n|\mathcal{F}|^2)}{n}}$$

The minimax rate is $O\left(\sqrt{\frac{\log|\mathcal{F}|}{n}}\right)$.

$|\mathcal{F}|$:= cardinality of a function class \mathcal{F} .

Empirical Risk Minimization

Cardinality: $\mathcal{F} = \{f(x) = ax, x \in \mathbb{R} \mid a = \{1, 2, 3\}\}$ $|\mathcal{F}| = 3$

- Assume $((x_t, y_t) \sim \mathcal{P} : t=1, \dots, n)$ are generated i.i.d from a distribution \mathcal{P}
 - ERM with finite class

Proposition 1

Consider the case when the hypothesis \mathcal{F} has finite cardinality, that is $|\mathcal{F}| < \infty$. For any loss ℓ satisfies $\sup_{y, y' \in \mathcal{Y}} |\ell(y, y')| \leq 1$, we have

$$V_n^{\text{stat}}(\mathcal{F}) \leq \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \leq 8 \sqrt{\frac{\log(n|\mathcal{F}|^2)}{n}}$$

The minimax rate is $O\left(\sqrt{\frac{\log|\mathcal{F}|}{n}}\right)$. *assuming $|\mathcal{F}|$ is large*

minimax rate in statistical learning.

- The iid assumption can be replaced by Martingales.
- In practice, iid assumption is very unrealistic \rightarrow causal learning Week 9.

Empirical Risk Minimization

$$P(A \cup B) = P(A) + P(B).$$

- ERM with finite class

$$\inf_{f \in \mathcal{F}} \sup_{\mathcal{D}} [\quad]$$

Proposition 1

Consider the case when the hypothesis \mathcal{F} has finite cardinality, that is $|\mathcal{F}| < \infty$. For any loss ℓ satisfies $\sup_{y, y' \in \mathcal{Y}} |\ell(y, y')| \leq 1$, we have

$$V_n^{\text{stat}}(\mathcal{F}) \leq \mathbb{E}_{\mathcal{S}} \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \leq 8 \sqrt{\frac{\log n |\mathcal{F}|^2}{n}}$$

The minimax rate is $O\left(\sqrt{\frac{\log |\mathcal{F}|}{n}}\right)$.

proof sketch:

$$V_n^{\text{stat}}(\mathcal{F}) \leq \sup_{f \in \mathcal{F}} (\text{Generalization Gap}(f)) \leq$$

Empirical Risk Minimization

Proof (part I):

$$\begin{aligned} & \mathbb{E}_S \left[L_D(\hat{f}_{erm}) - \inf_{f \in \mathcal{F}} L_D(f) \right] \\ &= \mathbb{E}_S \left[L_D(\hat{f}_{erm}) \right] - \inf_{f \in \mathcal{F}} \mathbb{E}_S \left[\frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \quad \text{def of } L_D(\cdot) \\ &\leq \mathbb{E}_S \left[L_D(\hat{f}_{erm}) - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \end{aligned}$$

Goal: bound V_n^{stat} .

Note that $V_n^{stat}(\mathcal{F}) = \inf_{\hat{f}} \sup_D \mathbb{E}_S \left[L_D(\hat{f}) - \inf_{f \in \mathcal{F}} L_D(f) \right]$

Let's focus on this!

[Let's consider using \hat{f}_{erm} !]

Empirical Risk Minimization

Proof (part I):

$$\begin{aligned} & \mathbb{E}_S \left[L_D(\hat{f}_{erm}) - \inf_{f \in \mathcal{F}} L_D(f) \right] \\ &= \mathbb{E}_S \left[L_D(\hat{f}_{erm}) \right] - \inf_{f \in \mathcal{F}} \mathbb{E}_S \left[\frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \\ & \text{(why?) } \leq \mathbb{E}_S \left[L_D(\hat{f}_{erm}) - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \end{aligned}$$

def of $L_D(\cdot)$

Empirical Risk Minimization

Proof (part I):

$$\begin{aligned} & \mathbb{E}_S \left[L_D(\hat{f}_{erm}) - \inf_{f \in \mathcal{F}} L_D(f) \right] \\ &= \mathbb{E}_S \left[L_D(\hat{f}_{erm}) \right] - \inf_{f \in \mathcal{F}} \mathbb{E}_S \left[\frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \\ &\leq \mathbb{E}_S \left[L_D(\hat{f}_{erm}) - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \quad \leftarrow \text{from last slide} \\ &\leq \mathbb{E}_S \left[\underbrace{\mathbb{E}[\ell(\hat{f}_{erm}(x), y)]}_{\text{def of } L_D} - \frac{1}{n} \sum_{t=1}^n \ell(\underbrace{\hat{f}_{erm}(x_t), y_t}_{\text{def of } \hat{f}_{erm}}) \right] \\ &\stackrel{(\text{why?})}{\leq} \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left[\mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \right] \\ &\leq \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \end{aligned}$$

Empirical Risk Minimization

Proof (part I):

$$\begin{aligned} & \mathbb{E}_S \left[L_D(\hat{f}_{erm}) - \inf_{f \in \mathcal{F}} L_D(f) \right] \\ &= \mathbb{E}_S \left[L_D(\hat{f}_{erm}) \right] - \inf_{f \in \mathcal{F}} \mathbb{E}_S \left[\frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \\ &\leq \mathbb{E}_S \left[L_D(\hat{f}_{erm}) - \inf_{f \in \mathcal{F}} \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \quad \leftarrow \text{from last slide} \\ &\leq \mathbb{E}_S \left[\underbrace{\mathbb{E}[\ell(\hat{f}_{erm}(x), y)]}_{\text{def of } L_D} - \frac{1}{n} \sum_{t=1}^n \ell(\underbrace{\hat{f}_{erm}(x_t), y_t}_{\text{def of } \hat{f}_{erm}}) \right] \\ &\leq \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left[\mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \right] \\ &\leq \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left[\mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right] \right] \end{aligned}$$

since $\hat{f}_{erm} \in \mathcal{F}$

\mathcal{Z}

Empirical Risk Minimization

Proof (part II):

$$\begin{aligned} \mathcal{V}_n^{\text{stat}}(\mathcal{F}) &= \inf_{\hat{f}} \sup_D \mathbb{E}_S \left[L_D(\hat{f}) - \inf_{f \in \mathcal{F}} L_D(f) \right] \\ &\leq \sup_D \mathbb{E}_S \left[L_D(\hat{f}_{\text{erm}}) - \inf_{f \in \mathcal{F}} L_D(f) \right] \\ &\stackrel{\text{from last slide}}{\leq} \sup_D \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left| \underbrace{\mathbb{E}[\ell(f(x), y)]}_{\Sigma} - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \\ &\stackrel{D.}{\leq} \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left[\underbrace{\mathbb{E}[\ell(f(x), y)]}_{\text{test error}} - \underbrace{\frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t)}_{\text{training error}} \right] \right] \end{aligned}$$

Empirical Risk Minimization

Proof (part II):

$$\begin{aligned} \mathcal{V}_n^{\text{stat}}(\mathcal{F}) &= \inf_{\hat{f}} \sup_D \mathbb{E}_S \left[L_D(\hat{f}) - \inf_{f \in \mathcal{F}} L_D(f) \right] \\ &\leq \sup_D \mathbb{E}_S \left[L_D(\hat{f}_{\text{erm}}) - \inf_{f \in \mathcal{F}} L_D(f) \right] \\ &\stackrel{\text{from last slide}}{\leq} \sup_D \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left| \underbrace{\mathbb{E}[\ell(f(x), y)]}_{\Sigma} - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \\ &\leq \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left| \underbrace{\mathbb{E}[\ell(f(x), y)]}_{\text{test error}} - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \end{aligned}$$

generalization gap!

Jensen's inequality

$\sup \mathbb{E}[X] \leq \mathbb{E}[\sup(X)]$
 x is a random variable.

• It remains to bound this generalization gap!

Empirical Risk Minimization

Proof (part III): Idea: using the Hoeffding's Inequality!

(from last lecture)

$$\begin{aligned} & \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \\ = & \mathbb{E}_S \left[\mathbb{1}_{\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \leq \epsilon} \sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \\ & \left(+ \mathbb{E}_S \left[\mathbb{1}_{\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| > \epsilon} \sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \right) \end{aligned}$$

What are we doing?

Why do we need this?

Empirical Risk Minimization

Proof (part III): Idea: using the Hoeffding's Inequality!

(from last lecture)

$$\begin{aligned} & \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \\ = & \mathbb{E}_S \left[\mathbb{1}_{\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \leq \epsilon} \sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \\ & + \mathbb{E}_S \left[\mathbb{1}_{\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| > \epsilon} \sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \end{aligned}$$

- Truncate the value of the random variable

$$\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) \right| =: \sup_{f \in \mathcal{F}} \mathbb{Z} = \mathbb{Z}'$$

into two parts:

①. $\mathbb{Z}' \leq \epsilon$

②. $\mathbb{Z}' > \epsilon$

- $\epsilon^{>0}$ is a free parameter that's tunable
- This is a standard "truncation technique" in probability theory, which turns a concentration inequality into expectation bounds
- May not be TIGHT! (depending on the upper bounds of \mathbb{Z})

Empirical Risk Minimization

Proof (part III): Let's consider the two parts one-by-one:

$$\begin{aligned} & \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \\ &= \mathbb{E}_S \left[\mathbb{1}_{\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \leq \epsilon} \sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \\ & \quad + \mathbb{E}_S \left[\mathbb{1}_{\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| > \epsilon} \sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \\ &\leq \epsilon + 2P \left(\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| > \epsilon \right) \rightarrow \text{Can we bound this probability?} \end{aligned}$$

$$\leq \epsilon + 2|\mathcal{F}|P \left(\left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| > \epsilon \right) \rightarrow \text{We use the union bound. It's only bounded when}$$

$$\leq \epsilon + 4|\mathcal{F}| \exp\left(-\frac{\epsilon^2 n}{2}\right) \text{ How to bound this?}$$

- $|\mathcal{F}|$ is finite
↓
function class

$$\exp\left(-\frac{\log n |\mathcal{F}|^2}{2}\right) \cdot 4|\mathcal{F}|$$

Let $\epsilon = \sqrt{\log(n|\mathcal{F}|^2)/n}$, we have $\mathcal{V}_n^{\text{stat}}(\mathcal{F}) \leq 8\sqrt{\frac{\log(n|\mathcal{F}|^2)}{n}}$. This finished the proof.

Remarks: We have used the following fact (check yourself!)

$$\forall f \in \mathcal{F}, \Pr \left(\left| \underbrace{\mathbb{E}[\ell((x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t)}_{\bar{\Sigma}} \right| > \varepsilon \right)$$

$$\leq 2 \exp \left(-\frac{\varepsilon^2 n}{2} \right)$$

The Hoeffding's Inequality

• Standard Form: $\left\{ \begin{array}{l} \cdot \Pr(|S_n - \mathbb{E}[S_n]| \geq \varepsilon) \leq 2 \exp \left(-\frac{2\varepsilon^2}{\sum_{i=1}^n (b_i - a_i)^2} \right) \\ \cdot x_1, \dots, x_n \text{ iid} \\ \cdot S_n = \sum_{i=1}^n x_i \\ \cdot a_i \leq x_i \leq b_i, \quad \forall i=1, \dots, n. \end{array} \right.$

• Our Form: Define $x_i = \frac{1}{n} \ell(f(x_i), y_i)$

Empirical Risk Minimization

(x, y)
 $\uparrow \quad \uparrow$
 $\mathbb{R} \quad \{0, +1\}$

$$\mathcal{F} := \{ y = \mathbb{1}(x > a) \mid a \in \mathbb{R} \}$$

$$\mathcal{V}_n^{\text{stat}}(\mathcal{F}) \leq \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \right] \leq 8 \sqrt{\frac{\log n |\mathcal{F}|^2}{n}}$$

- It shows the connection to

$$\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right|$$

- It requires that \mathcal{F} is finite, i.e., $|\mathcal{F}| < \infty$
- How about $|\mathcal{F}| = \infty$?

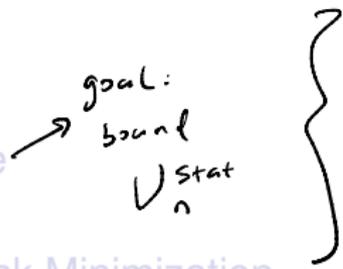
How to measure the complexity of \mathcal{F} ?

- When \mathcal{F} is finite: $|\mathcal{F}|$

A deeper question:

- When \mathcal{F} is an infinite set: ?

Next slide

- 1 Introduction
- 2 Minimax rate 
- 3 Empirical Risk Minimization
- 4 Growth Function and VC dimension
- 5 Rademacher Complexity

Growth Function

- **Growth function** (also known as shattering coefficient)

Given $\{(x_i, y_i)\}_{1 \leq i \leq n}$ and define $S = \{x_1, x_2, \dots, x_n\}$. Let

$\mathcal{F}_S = \mathcal{F}_{x_1, \dots, x_n} = \{f(x_1), \dots, f(x_n) : f \in \mathcal{F}\}$ and suppose

$f(x) \in \{0, 1\}$. The growth function is the maximum number of ways into which n points can be classified by the function class:

$$G(\mathcal{F}, n) = \sup_{x_1, \dots, x_n} |\mathcal{F}_S|$$

$(x_1, x_2) \quad (0, 1)$
 $(x'_1, x'_2) \quad (1, 0)$

- When \mathcal{F} is finite, $G(\mathcal{F}, n) \leq |\mathcal{F}|$.
- It always holds that $G(\mathcal{F}, n) \leq 2^n$.
- We say \mathcal{F} shatters S if $|\mathcal{F}_S| = 2^{|S|}$.

• In other words, "shattering is the ability of a model to classify a set of points perfectly".

• **Idea:** \mathcal{F} is infinite, but can be evaluated by finite samples!

Growth Function

- **Growth function** (also known as shattering coefficient)
Given $\{(x_i, y_i)\}_{1 \leq i \leq n}$ and define $S = \{x_1, x_2, \dots, x_n\}$. Let $\mathcal{F}_S = \mathcal{F}_{x_1, \dots, x_n} = \{f(x_1), \dots, f(x_n) : f \in \mathcal{F}\}$ and suppose $f(x) \in \{0, 1\}$. The growth function is the maximum number of ways into which n points can be classified by the function class:

$$G(\mathcal{F}, n) = \sup_{x_1, \dots, x_n} |\mathcal{F}_S|$$

- When \mathcal{F} is finite, $G(\mathcal{F}, n) \leq |\mathcal{F}|$.
- It always holds that $G(\mathcal{F}, n) \leq 2^n$.
- We say \mathcal{F} shatters S if $|\mathcal{F}_S| = 2^{|S|}$.
- Uniform convergence bound

all possible realizations by \mathcal{F}

claim: $P \left(\sup_{f \in \mathcal{F}} \left| \underbrace{\mathbb{E}[\ell(f(x), y)]}_{\Sigma} - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \geq \epsilon \right) \leq 2G(\mathcal{F}, 2n) \exp \left(-\frac{\epsilon^2 n}{4} \right) \quad (5)$

- * Connection with bound of \mathcal{V}_n^{stat} ?

Derivation of

$$\Pr\left(\sup_{f \in \mathcal{F}} \left| \underbrace{\mathbb{E}[l(f(x), y)]}_{L(f)} - \frac{1}{n} \sum_{t=1}^n \underbrace{l(f(x_t), y_t)}_{L_{\text{emp}}(f)} \right| \geq \varepsilon \right) \leq 2G(\mathcal{F}, 2n) \cdot \exp\left(-\frac{\varepsilon^2 n}{4}\right). \quad (1)$$

proof: by Vapnik & Chervonenkis

Lemma (Symmetrization Lemma)

If $n\varepsilon^2 \geq 2$, we have

$$\Pr\left(\sup_{f \in \mathcal{F}} |L(f) - L_{\text{emp}}(f)| \geq \varepsilon\right) \leq 2\Pr\left(\sup_{f \in \mathcal{F}} |L_{\text{emp}}(f) - L'_{\text{emp}}(f)| > \varepsilon/2\right)$$

an iid dummy copy of the original samples

$(x_1, y_1) \dots (x_n, y_n)$

[similar to the dummy dataset we created in the Radomacher complexity proof]

Proof sketch of the Symmetrization Lemma

Use the truncation technique!

Therefore, if $n\varepsilon^2 \geq 2$

$$\Pr\left(\sup_{f \in \mathcal{F}} |\mathcal{L}(f) - \mathcal{L}_{\text{emp}}(f)| > \varepsilon\right)$$

$$\leq 2\Pr\left(\sup_{f \in \mathcal{F}} |\mathcal{L}_{\text{emp}}(f) - \mathcal{L}'_{\text{emp}}(f)| > \frac{\varepsilon}{2}\right)$$

[Apply the symmetrization Lemma]

$$= 2\Pr\left(\sup_{f \in \mathcal{F}_{2n}} |\mathcal{L}_{\text{emp}}(f) - \mathcal{L}'_{\text{emp}}(f)| > \frac{\varepsilon}{2}\right)$$

[only functions in \mathcal{F}_{2n} are important]

$$\leq 2G(\mathcal{F}, 2n) \cdot \exp\left(-\frac{n\varepsilon^2}{4}\right)$$



since we have $\frac{\varepsilon}{2}$

since \mathcal{F}_{2n} contains
at most $G(\mathcal{F}, 2n)$ functions, i.e.
there are at most $G(\mathcal{F}, 2n)$ possible
realizations!

#

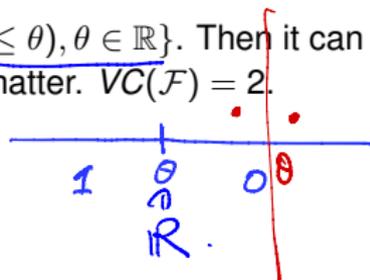
VC dimension

- VC (Vapnik-Chervonenkis) dimension

The VC dimension of a class \mathcal{F} is the largest n such that $G(\mathcal{F}, n) = 2^n$. In other words, VC dimension of a function class F is the cardinality of the largest set that it can shatter. It is a measure of the capacity (complexity, expressive power, richness, or flexibility) of a set of functions.

- Examples

- $\mathcal{F} = \{f(x) = I(x \leq \theta), \theta \in \mathbb{R}\}$. Then it can shatter 2 points but for any three points it cannot shatter. $VC(\mathcal{F}) = 2$.



$$2^1 = 2.$$

VC dimension

- VC (Vapnik-Chervonenkis) dimension

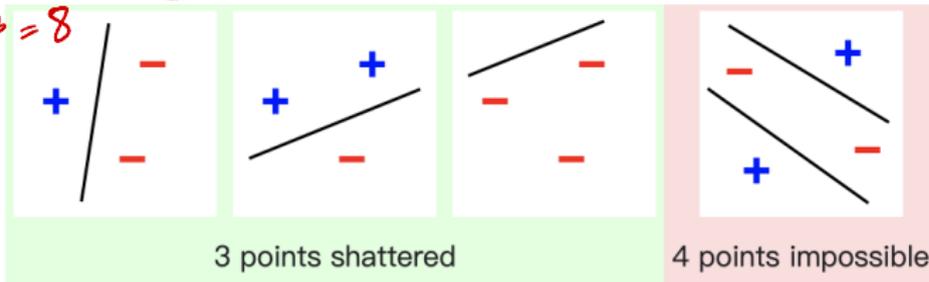
The VC dimension of a class \mathcal{F} is the largest n such that $G(\mathcal{F}, n) = 2^n$. In other words, VC dimension of a function class \mathcal{F} is the cardinality of the largest set that it can shatter. It is a measure of the capacity (complexity, expressive power, richness, or flexibility) of a set of functions.

- Examples

$$|\mathcal{F}| = \infty$$

- $\mathcal{F} = \{f(x) = I(x \leq \theta), \theta \in \mathbb{R}\}$. Then it can shatter 2 points but for any three points it cannot shatter. $VC(\mathcal{F}) = 2$.
- \mathcal{F} is a set of lines in 2-D space: $VC(\mathcal{F}) = 3$. $f(x) = w_0 + w_1 x_1 + w_2 x_2$

$$2^3 = 8$$



- Linear function in \mathbb{R}^d : $VC(\mathcal{F}) = ?$ $d+1$ (Maybe one question in your assignments)
- How about rectangles and circles in 2-D space?

- Sauer's lemma

Lemma 1 (Vapnik, Chervonenkis, Sauer, Shelah)

Let \mathcal{F} be a function class with finite VC dimension d . Then

$$G(\mathcal{F}, n) \leq \sum_{i=0}^d \binom{n}{i} \quad G(\mathcal{F}, d) = 2^d$$

for all $n \in \mathbb{N}$. In particular, for all $n \geq d$, we have

$$G(\mathcal{F}, n) \leq \left(\frac{en}{d}\right)^d.$$

- This bound is tight!
- This lemma can be used to derive lower bounds on VC dimensions

VC generalization bound

- Recall that

Claim. $P\left(\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \geq \epsilon\right) \leq 2G(\mathcal{F}, 2n) \exp\left(-\frac{\epsilon^2 n}{4}\right)$

Let the RHS be some $\delta > 0$ and then solve it for ϵ . We have

$$\mathbb{E}[\ell(f(x), y)] \leq \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) + \sqrt{\frac{4((\log(2G(\mathcal{F}, 2n)) - \log \delta))}{n}}$$

VC generalization bound

- Recall that

$$P \left(\sup_{f \in \mathcal{F}} \left| \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right| \geq \epsilon \right) \leq 2G(\mathcal{F}, 2n) \exp \left(-\frac{\epsilon^2 n}{4} \right)$$

Let the RHS be some $\delta > 0$ and then solve it for ϵ . We have

$$\mathbb{E}[\ell(f(x), y)] \leq \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) + \sqrt{\frac{4((\log(2G(\mathcal{F}, 2n)) - \log \delta))}{n}}$$

VC Inequality

- Using Lemma 1 (suppose $n \geq d$), we have

$$\mathbb{E}[\ell(f(x), y)] \leq \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) + \sqrt{\frac{4 \left(d_{VC} \log\left(\frac{2en}{d_{VC}}\right) - \log \delta \right)}{n}}$$

The bound is very general (loose) since VC dimension only depends function space but not the dataset.

Can we tighten the bound?

- 1 Introduction
- 2 Minimax rate
- 3 Empirical Risk Minimization
- 4 Growth Function and VC dimension
- 5 Rademacher Complexity**

Rademacher complexity

- Rademacher variable σ_i : $P(\sigma_i = 1) = P(\sigma_i = -1) = \frac{1}{2}$
- Empirical Rademacher complexity

$$\mathcal{R}(\mathcal{F}) := \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right]$$

$\sigma_1 = 1 \quad 1$
 $\sigma_2 = -1 \quad -1$

- It is a measure of the capacity of a function space and depends on both dataset and \mathcal{F}

Rademacher complexity

$$f(x_i) \in \{-1, 1\}$$

- Rademacher variable σ_i : $P(\sigma_i = 1) = P(\sigma_i = -1) = \frac{1}{2}$
- Empirical Rademacher complexity

$$\mathcal{R}(\mathcal{F}) := \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right] \quad \mathbb{E}_{\sigma} \left[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \ell(f(x_i), y_i) \right]$$

- It is a measure of the capacity of a function space and depends on both dataset and \mathcal{F}
- Uniform convergence bound

This measures the biggest difference of the losses measured over the whole domain and the sample set.

Lemma 2

$$\mathbb{E}_{\mathcal{S}} \left[\underbrace{\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\}}_{\Sigma} \right] \leq 2 \mathbb{E}_{\mathcal{S}} \underbrace{\mathcal{R}(\ell \circ \mathcal{F})}_{(\text{abuse notation})}$$

Rademacher complexity

Proof (part I):

$$\begin{aligned} & \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \\ &= \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}_{S'} \left[\frac{1}{n} \sum_{t=1}^n \ell(f(x'_t), y'_t) \right] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \\ &\leq \mathbb{E}_S \left[\mathbb{E}_{S'} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \ell(f(x'_t), y'_t) - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \right] \\ &= \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \ell(f(x'_t), y'_t) - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \end{aligned}$$

$(x'_t, y'_t)_t$ are samples
from S'

$(x_t, y_t)_t$ are samples from S .
We have introduced a dummy dataset S' .
What does this inequality mean?

Rademacher complexity

Proof (part II):

$$\begin{aligned} & \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \ell(f(x'_t), y'_t) - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \\ &= \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \left(\ell(f(x'_j), y'_j) - \ell(f(x_j), y_j) + \sum_{i \neq j} (\ell(f(x'_i), y'_i) - \ell(f(x_i), y_i)) \right) \right) \right\} \right] \\ &= \mathbb{E}_{S, S', \sigma_j} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \left(\sigma_j \left(\overbrace{\ell(f(x'_j), y'_j)}^{(1)} - \overbrace{\ell(f(x_j), y_j)}^{(2)} \right) + \sum_{i \neq j} (\ell(f(x'_i), y'_i) - \ell(f(x_i), y_i)) \right) \right) \right\} \right] \end{aligned}$$

changes -f signs will only switch ① and ② -
but ① and ② are iid

Rademacher complexity

Proof (part II):

$$\begin{aligned} & \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \ell(f(x'_t), y'_t) - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \\ &= \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \left(\ell(f(x'_j), y'_j) - \ell(f(x_j), y_j) + \sum_{i \neq j} (\ell(f(x'_i), y'_i) - \ell(f(x_i), y_i)) \right) \right\} \right] \\ &= \mathbb{E}_{S, S', \sigma_j} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \left(\sigma_j (\ell(f(x'_j), y'_j) - \ell(f(x_j), y_j)) + \sum_{i \neq j} (\ell(f(x'_i), y'_i) - \ell(f(x_i), y_i)) \right) \right\} \right] \\ &= \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{j=1}^n \sigma_j (\ell(f(x'_j), y'_j) - \ell(f(x_j), y_j)) \right\} \right] \quad \sigma = \{\sigma_1, \sigma_2, \dots, \sigma_n\} \\ &\leq \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{j=1}^n \sigma_j \ell(f(x'_j), y'_j) \right\} + \sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{j=1}^n (-\sigma_j) \ell(f(x_j), y_j) \right\} \right] \\ &= \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{j=1}^n \sigma_j \ell(f(x'_j), y'_j) \right\} + \sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{j=1}^n \sigma_j \ell(f(x_j), y_j) \right\} \right] \end{aligned}$$

Proof (part III):

$$\begin{aligned} & \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \\ & \leq \mathbb{E}_{S, S'} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{t=1}^n \ell(f(x'_t), y'_t) - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \\ & \leq \mathbb{E}_{S, S', \sigma} \left[\sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{j=1}^n \sigma_j \ell(f(x'_j), y'_j) \right\} + \sup_{f \in \mathcal{F}} \left\{ \frac{1}{n} \sum_{j=1}^n \sigma_j \ell(f(x_j), y_j) \right\} \right] \\ & = \mathbb{E}_{S'} \mathcal{R}_{S'}(\ell \circ \mathcal{F}) + \mathbb{E}_S \mathcal{R}_S(\ell \circ \mathcal{F}) \\ & = 2\mathbb{E}_S \mathcal{R}_S(\ell \circ \mathcal{F}) \end{aligned}$$

This finished the proof.

Rademacher complexity bound

Combining $\mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \leq 2\mathbb{E}_S \mathcal{R}_S(\mathcal{F})$ with

Lemma 3 (McDiarmid's Inequality)

Let x_1, \dots, x_n be independent random variables taking on values in a set A and let c_1, \dots, c_n be positive real constants. If $\varphi : A^n \rightarrow \mathbb{R}$ satisfies

$$\sup_{x_1, \dots, x_n, x_i' \in A} |\varphi(x_1, \dots, x_i, \dots, x_n) - \varphi(x_1, \dots, x_i', \dots, x_n)| \leq c_i,$$

for $1 \leq i \leq n$, then

$$P(\varphi(x_1, \dots, x_n) - \mathbb{E}[\varphi(x_1, \dots, x_n)] \geq \epsilon) \leq e^{-2\epsilon^2 / \sum_{i=1}^n c_i^2}$$

Rademacher complexity bound

Combining $\mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left\{ \underbrace{\mathbb{E}[\ell(f(x), y)]}_{\mathcal{Z}(f)} - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \leq 2\mathbb{E}_S \mathcal{R}_S(\ell \circ \mathcal{F})$ with

Lemma 3 (McDiarmid's Inequality)

Let x_1, \dots, x_n be independent random variables taking on values in a set A and let c_1, \dots, c_n be positive real constants. If $\varphi : A^n \rightarrow \mathbb{R}$ satisfies

$$\sup_{x_1, \dots, x_n, x'_i \in A} |\varphi(x_1, \dots, x_i, \dots, x_n) - \varphi(x_1, \dots, x'_i, \dots, x_n)| \leq c_i,$$

for $1 \leq i \leq n$, then

$$P(\varphi(x_1, \dots, x_n) - \mathbb{E}[\varphi(x_1, \dots, x_n)] \geq \epsilon) \leq e^{-2\epsilon^2 / \sum_{i=1}^n c_i^2} \quad // ?$$

• $\mathcal{Z}(f)$ is a random variable that depends on the dataset S and f .

• Idea: $\varphi(x_1, \dots, x_n) \leftrightarrow \sup_{f \in \mathcal{F}} \left(\mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right)$

$$\text{Hence, } \sup_{f \in \mathcal{F}} |\varphi(x_1, \dots, x_i, \dots, x_n) - \varphi(x_1, \dots, x'_i, \dots, x_n)| \leq \frac{1}{n}$$

since $\ell(\cdot, \cdot) \in [0, 1]$.

$$\Rightarrow P_r(\varphi - \mathbb{E}[\varphi] \geq \epsilon) \leq \exp(-2\epsilon^2 n)$$

So, w.p. at least $1 - \exp(-2\epsilon^2 n)$, $\varphi - \mathbb{E}[\varphi] \leq \epsilon$.

Rademacher complexity bound

Combining $\mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] \leq 2\mathbb{E}_S \mathcal{R}_S(\mathcal{F})$ with

Lemma 3 (McDiarmid's Inequality)

Let x_1, \dots, x_n be independent random variables taking on values in a set A and let c_1, \dots, c_n be positive real constants. If $\varphi : A^n \rightarrow \mathbb{R}$ satisfies

$$\sup_{x_1, \dots, x_n, x'_i \in A} |\varphi(x_1, \dots, x_i, \dots, x_n) - \varphi(x_1, \dots, x'_i, \dots, x_n)| \leq c_i,$$

for $1 \leq i \leq n$, then

$$P(\varphi(x_1, \dots, x_n) - \mathbb{E}[\varphi(x_1, \dots, x_n)] \geq \epsilon) \leq e^{-2\epsilon^2 / \sum_{i=1}^n c_i^2}$$

Assume $0 \leq \ell \leq 1$, thus with probability at least $1 - \delta$, we have

Using the standard truncation technique

$$\begin{aligned} & \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \\ & \leq \mathbb{E}_S \left[\sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \right] + \sqrt{\frac{\log(1/\delta)}{2n}} \end{aligned}$$

$$(why?) \leq 2\mathbb{E}_S \mathcal{R}_S(\mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{2n}}$$

Rademacher complexity bound

We have got

$$\begin{aligned} & \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \\ & \leq 2\mathbb{E}_S \mathcal{R}_S(\ell \circ \mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{2n}} \end{aligned}$$

Apply McDiarmid's inequality again on Rademacher complexity itself. The bounded difference of $\mathcal{R}_S(\ell \circ \mathcal{F}) := \mathbb{E}_\sigma \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i)$ is still $1/n$. Then with probability of at least $1 - \delta$, we have

Rademacher complexity bound

We have got

$$\begin{aligned} & \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \\ & \leq 2\mathbb{E}_S \mathcal{R}_S(\ell \circ \mathcal{F}) + \sqrt{\frac{\log(1/\delta)}{2n}} \end{aligned}$$

Apply McDiarmid's inequality again on Rademacher complexity itself. The bounded difference of $\mathcal{R}_S(\ell \circ \mathcal{F}) := \mathbb{E}_\sigma \sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i)$ is still $1/n$. Then with probability of at least $1 - \delta$, we have

$$\begin{aligned} & \sup_{f \in \mathcal{F}} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \\ & \leq 2\mathcal{R}_S(\ell \circ \mathcal{F}) + 3\sqrt{\frac{\log(2/\delta)}{2n}} \leftarrow \text{Apply the McDiarmid's inequality twice} \end{aligned}$$

*Note that $\mathbb{E}_S \mathcal{R}_S(\ell \circ \mathcal{F}) \leq \sqrt{\frac{2 \log G(\mathcal{F}, n)}{n}}$. (why?) The Massart's lemma

Rademacher complexity of linear function class

Examples.

Linear function space: $\mathcal{F}_2 = \{x \rightarrow \langle \mathbf{w}, x \rangle : \|\mathbf{w}\|_2 \leq 1\}$

Lemma 4

Let $S = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$ be vectors in a Hilbert space. Suppose $\|\mathbf{x}_i\| \leq B, i = 1, 2, \dots, n$. Define:

$$\mathcal{F}_2 \circ S = \{(\langle \mathbf{w}, \mathbf{x}_1 \rangle, \dots, \langle \mathbf{w}, \mathbf{x}_n \rangle) : \|\mathbf{w}\|_2 \leq \omega\}.$$

Then $\mathcal{R}(\mathcal{F}_2 \circ S) \leq \frac{\omega B}{\sqrt{n}}$.

Rademacher complexity of linear function class

Proof (part I):

$$\begin{aligned}\mathcal{R}(\mathcal{F}_2 \circ \mathcal{S}) &= \mathbb{E}_\sigma \left[\sup_{\mathbf{a} \in \mathcal{F}_2 \circ \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \sigma_i a_i \right] \\ &= \frac{1}{n} \mathbb{E}_\sigma \left[\sup_{\mathbf{w}: \|\mathbf{w}\| \leq \omega} \sum_{i=1}^n \sigma_i \langle \mathbf{w}, \mathbf{x}_i \rangle \right] && a_i := \langle \mathbf{w}, \mathbf{x}_i \rangle \\ &= \frac{1}{n} \mathbb{E}_\sigma \left[\sup_{\mathbf{w}: \|\mathbf{w}\| \leq \omega} \left\langle \mathbf{w}, \sum_{i=1}^n \sigma_i \mathbf{x}_i \right\rangle \right] && \text{linearity of inner products} \\ &\leq \frac{1}{n} \mathbb{E}_\sigma \left[\sup_{\mathbf{w}: \|\mathbf{w}\| \leq \omega} \|\mathbf{w}\| \left\| \sum_{i=1}^n \sigma_i \mathbf{x}_i \right\| \right] && (\text{Cauchy-Schwartz inequality}) \\ &\stackrel{\text{why?}}{\leq} \frac{\omega}{n} \mathbb{E}_\sigma \left[\left\| \sum_{i=1}^n \sigma_i \mathbf{x}_i \right\| \right] = \frac{\omega}{n} \mathbb{E}_\sigma \left[\left(\left\| \sum_{i=1}^n \sigma_i \mathbf{x}_i \right\|^2 \right)^{1/2} \right] \\ &\leq \frac{\omega}{n} \left(\mathbb{E}_\sigma \left[\left\| \sum_{i=1}^n \sigma_i \mathbf{x}_i \right\|^2 \right] \right)^{1/2} && (\text{Jensen's inequality})\end{aligned}$$

Rademacher complexity of linear function class

Proof (part II):

$$\begin{aligned}\mathcal{R}(\mathcal{F}_2 \circ \mathcal{S}) &= \mathbb{E}_{\sigma} \left[\sup_{\mathbf{a} \in \mathcal{F}_2 \circ \mathcal{S}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbf{a}_i \right] \\ &\leq \frac{\omega}{n} \left(\mathbb{E}_{\sigma} \left[\left\| \sum_{i=1}^n \sigma_i \mathbf{x}_i \right\|^2 \right] \right)^{1/2} \\ &= \frac{\omega}{n} \sqrt{\mathbb{E}_{\sigma} \left[\sum_{i,j} \sigma_i \sigma_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right]} \quad \sigma_i \text{'s are Rademacher variables} \\ &= \frac{\omega}{n} \sqrt{\left(\sum_{i \neq j} \langle \mathbf{x}_i, \mathbf{x}_j \rangle \underbrace{\mathbb{E}_{\sigma} [\sigma_i \sigma_j]}_{\substack{? \\ \mathbb{E}[\sigma_i] \mathbb{E}[\sigma_j] \\ = 0}} \right) + \sum_{i=1}^n \langle \mathbf{x}_i, \mathbf{x}_i \rangle \underbrace{\mathbb{E}_{\sigma} [\sigma_i^2]}_{=1}} \right)} \\ &= \frac{\omega}{n} \sqrt{\sum_{i=1}^n \|\mathbf{x}_i\|^2} \leq \frac{\omega B}{\sqrt{n}}\end{aligned}$$

This finished the proof.

Generalization bound of linear models

Lemma 5

If the loss function ℓ is η -Lipschitz, we have

$$\mathcal{R}(\ell \circ \mathcal{F}) \leq \eta \mathcal{R}(\mathcal{F})$$

$$|\ell(x') - \ell(x)| \leq \eta \|x - x'\|, \forall x, x' \in \mathcal{X}$$

Generalization bound of linear models

Lemma 5

If the loss function ℓ is η -Lipschitz, we have

$$\mathcal{R}(\ell \circ \mathcal{F}) \leq \eta \mathcal{R}(\mathcal{F})$$

Linear function space: $\mathcal{F}_2 = \{x \rightarrow \langle w, x \rangle : \|w\| \leq \omega\}$. Suppose $\|x_i\| \leq B, i = 1, 2, \dots, n$. Then with probability of at least $1 - \delta$, we have

$$\begin{aligned} & \sup_{f \in \mathcal{F}_2} \left\{ \mathbb{E}[\ell(f(x), y)] - \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) \right\} \\ & \leq \frac{2\eta\omega B}{\sqrt{n}} + 3\sqrt{\frac{\log(2/\delta)}{2n}} \end{aligned}$$

Or equivalently, suppose $f \in \mathcal{F}_2$, then with probability of at least $1 - \delta$,

$$\mathbb{E}[\ell(f(x), y)] \leq \frac{1}{n} \sum_{t=1}^n \ell(f(x_t), y_t) + \frac{2\eta\omega B}{\sqrt{n}} + 3\sqrt{\frac{\log(2/\delta)}{2n}}$$

Learning outcomes

- Understand the concepts of PAC, agnostic PAC, generalization bound, growth function, VC dimension, and Rademacher complexity.
- Understand the properties of the **three generalization error bounds** we have learned.
- Be able to compute the Rademacher complexities for some simple function classes.
- Be able to derive the generalization bounds for some simple machine learning models.