

# Chapter6. Empirical Risk Minimization: Abstract risk bounds and Rademacher averages

May 31, 2021

# Last chapter

## Theorem 5.3

$$\hat{f}_n = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \mathbb{E}_{P_n} \ell_f(z)$$

Empirical Risk Minimization (ERM) algorithm is a PAC algorithm if

$$\lim_{n \rightarrow \infty} \sup_{P \in \mathcal{P}} P^n \left( \underbrace{\sup_{f \in \mathcal{F}} |P_n(\ell_f) - P(\ell_f)|}_{\substack{\uparrow \\ \uparrow}} \geq \varepsilon \right) = 0, \quad \forall \varepsilon > 0,$$

for every  $\varepsilon > 0$ .

## Agnostic (model-free) learning

- Sets  $\mathcal{X}$ ,  $\mathcal{Y}$ , and  $\mathcal{U}$
- A class  $\mathcal{P}$  of probability distributions on  $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$
- A class  $\mathcal{F}$  of functions  $f : \mathcal{X} \rightarrow \mathcal{U}$  (the hypothesis space)
- A loss function  $\ell : \mathcal{Y} \times \mathcal{U} \rightarrow [0, 1]$

# Notations

## Agnostic (model-free) learning

- Sets  $\mathcal{X}$ ,  $\mathcal{Y}$ , and  $\mathcal{U}$
- A class  $\mathcal{P}$  of probability distributions on  $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$
- A class  $\mathcal{F}$  of functions  $f : \mathcal{X} \rightarrow \mathcal{U}$  (the hypothesis space)
- A loss function  $\ell : \mathcal{Y} \times \mathcal{U} \rightarrow [0, 1]$

## An abstract framework for ERM

- Set  $\mathcal{Z}$
- A class  $\mathcal{P}$  of probability distributions on  $\mathcal{Z}$
- A class  $\mathcal{F}$  of functions  $f : \mathcal{Z} \rightarrow [0, 1]$  (induced losses)

$$\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$$

$$z = (x, y) \quad g \in \mathcal{H}$$

$$\ell(g(x), y) = f(z)$$

# Notations

- The **expected risk** of any  $f \in \mathcal{F}$ :

$$L_P(f) \quad P(f) := \mathbf{E}_P f(Z)$$

- The **minimum risk** :

$$L_P^*(\mathcal{F}) := \inf_{f \in \mathcal{F}} P(f)$$

- ERM algorithm**:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} P_n(f) = \arg \min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n f(Z_i)$$

- Uniform deviation**

$$\Delta_n(Z^n) := \|P_n - P\|_{\mathcal{F}} \equiv \sup_{f \in \mathcal{F}} |P_n(f) - P(f)|$$

$\mathbb{E}_{P_n} f(Z) = L_{P_n}(f)$   
 $\uparrow \quad \quad \quad \uparrow \quad \quad \uparrow$

## An abstract framework for ERM

- Set  $\mathcal{Z}$
- A class  $\mathcal{P}$  of probability distributions on  $\mathcal{Z}$
- A class  $\mathcal{F}$  of functions  $f : \mathcal{Z} \rightarrow [0, 1]$  (induced losses)

$\mathcal{F}$ -seminorm  $\|\cdot\|_{\mathcal{F}}$

For  $P, P' \in \mathcal{P}$

$$\|P - P'\|_{\mathcal{F}} \triangleq \sup_{f \in \mathcal{F}} |P(f) - P'(f)|$$

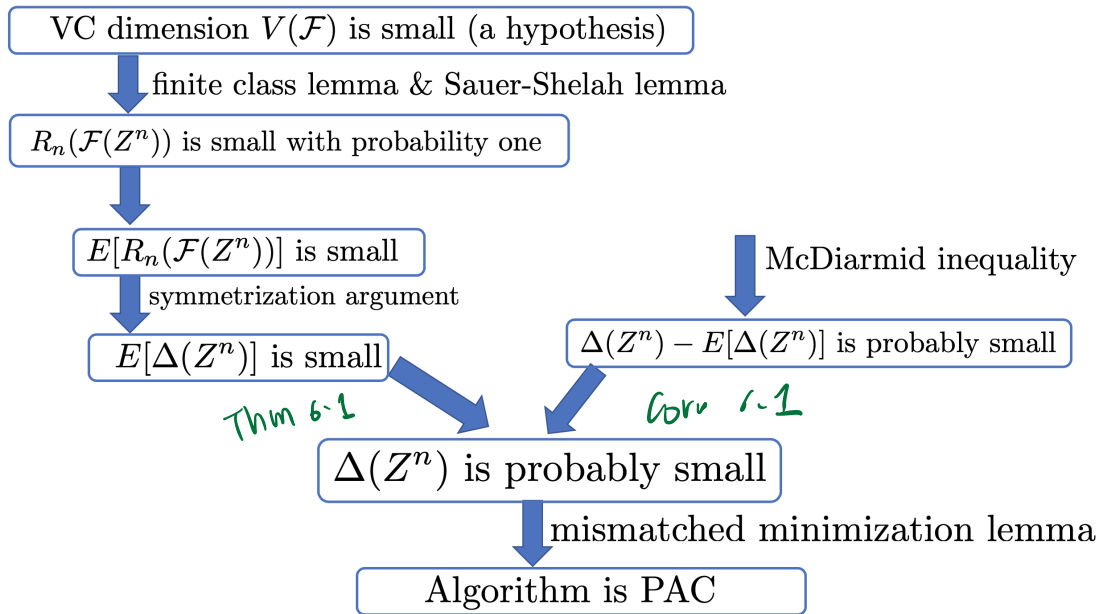
# Why bounding the uniform deviation

## Proposition 6.1

The generalization loss for a learning algorithm satisfies:

$$P\left(\hat{f}_n\right) \leq L^*(\mathcal{F}) + 2\Delta_n\left(Z^n\right) \text{ (if algorithm is ERM)}$$

$$P\left(\hat{f}_n\right) \leq P_n\left(\hat{f}_n\right) + \Delta_n\left(Z^n\right) \text{ (for any algorithm).}$$



# Bounding the uniform deviation: Rademacher averages

## Theorem 6.1

Fix a space  $\mathcal{Z}$  and let  $\mathcal{F}$  be a class of functions  $f : \mathcal{Z} \rightarrow [0, 1]$ . Then for any  $P \in \mathcal{P}(\mathcal{Z})$

$$\mathbf{E}\Delta_n(Z^n) \leq 2\mathbf{E}R_n(\mathcal{F}(Z^n))$$



# Bounding the uniform deviation: Rademacher averages

## Theorem 6.1

Fix a space  $\mathcal{Z}$  and let  $\mathcal{F}$  be a class of functions  $f : \mathcal{Z} \rightarrow [0, 1]$ . Then for any  $P \in \mathcal{P}(\mathcal{Z})$

$$\mathbf{E} \Delta_n(Z^n) \leq 2 \mathbf{E} R_n(\mathcal{F}(Z^n))$$

$\mathcal{F}(Z^n) \triangleq \{ (f(z_1), \dots, f(z_n)) : f \in \mathcal{F} \}$

## Definition 6.1.

Let  $\mathcal{A} \subset \mathbb{R}^n$  with  $\mathcal{A}$  bounded. The Rademacher average of  $\mathcal{A}$ , denoted by  $R_n(\mathcal{A})$ , is defined by

$$R_n(\mathcal{A}) = \mathbf{E} \left[ \sup_{a \in \mathcal{A}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i a_i \right| \right]$$

where  $\varepsilon_1, \dots, \varepsilon_n$  are independent Rademacher (i.e.,  $\pm 1$  with equal probability) random variables.

# Bounding the uniform deviation: Rademacher averages

## Theorem 6.1

Fix a space  $\mathcal{Z}$  and let  $\mathcal{F}$  be a class of functions  $f : \mathcal{Z} \rightarrow [0, 1]$ . Then for any  $P \in \mathcal{P}(\mathcal{Z})$

$$\mathbf{E} \Delta_n(Z^n) \leq 2 \mathbf{E} R_n(\mathcal{F}(Z^n))$$

**Proof:**

$$\mathbf{E} \Delta_n(Z^n) = \mathbf{E} \left[ \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(Z_i) - P(f) \right| \right]$$

Fact 1.  $\{Z_i\}$  i.i.d.  $P$ ,  $\mathbf{E} \left[ \frac{1}{n} \sum_{i=1}^n f(Z_i) \right] = \mathbf{E}_P f(Z) = P(f)$

Fact 2:  $\phi : \left\{ y(f) = \mathbf{E}_y f(Z) : f \in \mathcal{F} \right\} \rightarrow \mathbb{R}$   
 $y \in \mathcal{P}$

Jason 1

$$\eta \mapsto \mathbb{E} \left[ \sup_{t \in \mathcal{T}} \left| \frac{1}{n} \sum_{i=1}^n f(z_i) - \eta(t) \right| \right]$$

$$\Delta_n(z^n) = \phi(\overbrace{p(\mathcal{T})}) = \phi \left( \mathbb{E}_{\bar{z}_i} \left[ \frac{1}{n} \sum_{i=1}^n f(\bar{z}_i) \right] \right)$$

$t \in \mathbb{R}$ .

$$M_1(t) = \mathbb{E} \left[ e^{\underbrace{t(f(z_i) - f(\bar{z}_i))}_{\text{red}}} \right]$$

$$= \mathbb{E} \left[ e^{\underbrace{t f(z_i)}_{\text{red}}} \right] \cdot \mathbb{E} \left[ e^{-\underbrace{t f(\bar{z}_i)}_{\text{red}}} \right]$$

$$M_2(t) = \mathbb{E} \left[ e^{\underbrace{t f(\bar{z}_i)}_{\text{red}}} \right] \cdot \mathbb{E} \left[ e^{-\underbrace{t f(z_i)}_{\text{red}}} \right] = \mathbb{E}_{\bar{z}_i} \left[ \sup_{t \in \mathcal{T}} \left| \frac{1}{n} \sum_{i=1}^n \underbrace{(f(z_i) - f(\bar{z}_i))}_{\text{red}} \right| \right]$$

$$\underbrace{f(z_i) - f(\bar{z}_i)}_{\textcircled{1}} \stackrel{d}{=} \underbrace{f(\bar{z}_i) - f(z_i)}_{\textcircled{2}} \stackrel{d}{=} \underbrace{\varepsilon_i (f(z_i) - f(\bar{z}_i))}_{\text{red}}$$

$$\mathbb{E}_{z^n} \Delta_n(z^n) \leq \mathbb{E}_{z^n, \bar{z}_i, \varepsilon_i, t \in \mathcal{T}} \left[ \sup_{t \in \mathcal{T}} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (f(z_i) - f(\bar{z}_i)) \right| \right]$$

$$\begin{aligned}
 &\leq \mathbb{E} \left[ \sup \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(z_i) \right| \right] + \mathbb{E} \left[ \sup \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(\bar{z}) \right| \right] \\
 &= 2 R_n(\mathcal{F}(z^n))
 \end{aligned}$$

□

# Bounding the uniform deviation: Rademacher averages

## Corollary 6.1

For any  $P \in \mathcal{P}(\mathcal{Z})$  and any  $n$ , with probability at least  $1 - \delta$

$$\Rightarrow P(\hat{f}_n) \leq L^*(\mathcal{F}) + 4\mathbf{E}R_n(\mathcal{F}(Z^n)) + \sqrt{\frac{2\log(\frac{1}{\delta})}{n}} \quad (\text{if ERM is used})$$

$$P(\hat{f}_n) \leq P_n(\hat{f}_n) + 2\mathbf{E}R_n(\mathcal{F}(Z^n)) + \sqrt{\frac{\log(\frac{1}{\delta})}{2n}} \quad (\text{for any algorithm})$$

**Proof:**  $P(\hat{f}_n) \leq P_n(\hat{f}_n) + \Delta_n(z^n)$

$$\text{Coro 6.1: } P(\hat{f}_n) \leq L^*(\mathcal{F}) + 2\Delta_n(z^n) + 2\underbrace{\mathbb{E}\Delta_n(z^n)} - 2\mathbb{E}\Delta_n(z^n)$$

$$\leq L^*(\mathcal{F}) + 4R_n(\mathcal{F}(z^n)) + 2\underbrace{(\Delta_n(z^n) - \mathbb{E}\Delta_n(z^n))}$$

$$\text{w.p. } 1-\delta \leq \sqrt{\frac{\log(\frac{1}{\delta})}{2n}}$$

McDiarmid inequality :

$X^n = (X_1 \dots X_n)$   $n$ -tuple of independent r.v's

If a func  $g: X^n \rightarrow \mathbb{R}$  has bound difference -i.e.,

$$\sup_{x \in X} g(x_1 \dots x_{i-1}, x, x_{i+1} \dots x_n) - \inf_{x \in X} g(x_1 \dots x_{i-1}, x, x_{i+1} \dots x_n) \leq C_i$$

Then, for all  $t > 0$

$$\mathbb{P}(g(X^n) - \mathbb{E} g(X^n) \geq t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n C_i^2}\right).$$

$$\Delta_n(\mathbf{z}^n) = \sup_{f \in \mathcal{F}} \left| \frac{1}{n} \sum_{i=1}^n f(z_i) - P(f) \right|$$

$$C_1 = \frac{1}{n}$$

$$\text{Let } \varepsilon = \sqrt{\frac{\log(\frac{1}{\delta})}{2n}}$$

$$\mathbb{P}(\Delta(\mathbf{z}^n) - \mathbb{E} \Delta_n(\mathbf{z}^n) \geq \varepsilon) \leq \exp\left(-\frac{2\varepsilon^2}{n \cdot \frac{1}{n^2}}\right) = \delta$$

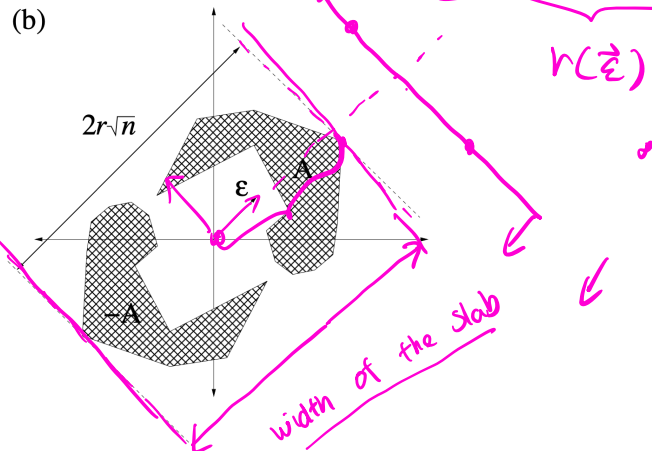
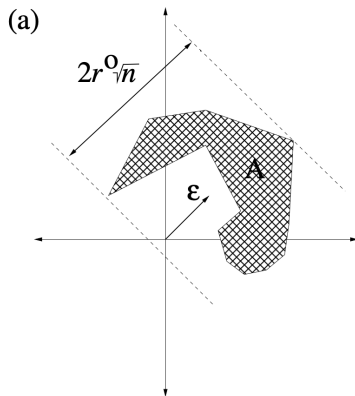
□

# Structural results for Rademacher averages

$\vec{\varepsilon} \stackrel{\text{def}}{=} (\varepsilon_1, \dots, \varepsilon_n)$   $\frac{1}{\sqrt{n}} \vec{\varepsilon}$  has length 1.

$$\Rightarrow R_n(\mathcal{A}) := \frac{1}{n} \mathbf{E}_{\varepsilon^n} \left[ \sup_{a \in \mathcal{A}} \left| \sum_{i=1}^n \varepsilon_i a_i \right| \right] = \frac{1}{\sqrt{n}} \mathbb{E}_{\varepsilon^n} \left[ \sup_{\vec{a} \in \mathcal{A}} \left| \left\langle \frac{\vec{\varepsilon}}{\sqrt{n}}, \vec{a} \right\rangle \right| \right]$$

$$R_n^\circ(\mathcal{A}) := \frac{1}{n} \mathbf{E}_{\varepsilon^n} \left[ \sup_{a \in \mathcal{A}} \sum_{i=1}^n \varepsilon_i a_i \right] = \mathbb{E}_{\varepsilon^n} \left[ \sup_{\vec{a} \in \mathcal{A}} \left\langle \frac{\vec{\varepsilon}}{\sqrt{n}}, \vec{a} \right\rangle \right]$$





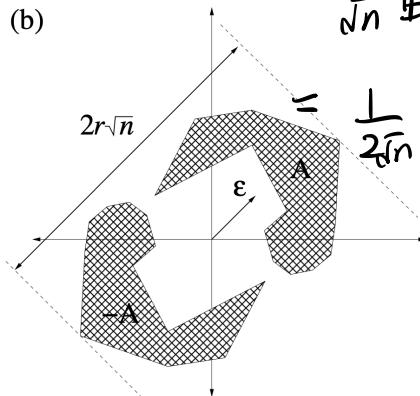
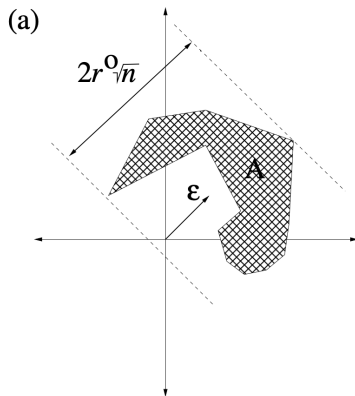
# Structural results for Rademacher averages

$$R_n(\mathcal{A}) := \frac{1}{n} \mathbf{E}_{\varepsilon^n} \left[ \sup_{a \in \mathcal{A}} \left| \sum_{i=1}^n \varepsilon_i a_i \right| \right]$$

$$R_n^\circ(\mathcal{A}) := \frac{1}{n} \mathbf{E}_{\varepsilon^n} \left[ \sup_{a \in \mathcal{A}} \sum_{i=1}^n \varepsilon_i a_i \right] = \frac{1}{\sqrt{n}} \mathbb{E} \left[ \sup_{a \in \mathcal{A}} \left\langle \frac{\vec{\varepsilon}}{\sqrt{n}}, \vec{a} \right\rangle \right]$$

$$= \frac{1}{\sqrt{n}} \mathbb{E} \left[ \inf_{a \in \mathcal{A}} - \left\langle \frac{\vec{\varepsilon}}{\sqrt{n}}, \vec{a} \right\rangle \right]$$

$$= \frac{1}{2\sqrt{n}} \mathbb{E} \left[ \sup_{a \in \mathcal{A}} \left\langle \frac{\vec{\varepsilon}}{\sqrt{n}}, \vec{a} \right\rangle - \inf_{a \in \mathcal{A}} \left\langle \frac{\vec{\varepsilon}}{\sqrt{n}}, \vec{a} \right\rangle \right]$$



# Structural results for Rademacher averages

## Basic properties of Rademacher averages

- (1)  $R_n^\circ(\mathcal{A}) \leq R_n(\mathcal{A}) = R_n^\circ(\mathcal{A} \cup -\mathcal{A})$
- (2)  $R_n^\circ(\mathcal{A}) = R_n(\mathcal{A})$  if  $\mathcal{A} = -\mathcal{A}$
- (3)  $R_n^\circ(\mathcal{A} + v) = R_n^\circ(\mathcal{A})$  for any  $v \in \mathbb{R}^n$
- (4)  $R_n(\mathcal{A} \cup \mathcal{B}) \leq R_n(\mathcal{A}) + R_n(\mathcal{B})$
- (5)  $R_n^\circ(\mathcal{A} + \mathcal{B}) = R_n^\circ(\mathcal{A}) + R_n^\circ(\mathcal{B})$
- (6)  $R_n(c\mathcal{A}) = |c|R_n(\mathcal{A})$
- (7)  $R_n(\mathcal{A}) = R_n(\text{conv}(\mathcal{A}))$
- (8)  $R_n(\mathcal{A}) = R_n(\text{absconv}(\mathcal{A}))$

$$\mathcal{A} + v \triangleq \{ \vec{a} + v : a \in \mathcal{A} \}$$

$$\sum c_i \vec{a}_i$$

$$\sum |c_i| = 1$$

# Structural results for Rademacher averages

## Lemma 6.1 (Finite class lemma).

If  $\mathcal{A} = \{a^{(1)}, \dots, a^{(N)}\} \subset \mathbb{R}^n$  is a finite set with  $\|a^{(j)}\| \leq L$  for all  $j = 1, \dots, N$  and  $N \geq 2$ , then

$$\mathcal{A} = \mathcal{F}(Z^n) = \left\{ \left( f(z_1) \cdots f(z_n) \right) : f \in \mathcal{F} \right\}$$

$$R_n(\mathcal{A}) \leq \frac{2L\sqrt{\log N}}{n}$$

**Proof:**

def 2.1 A r.v.  $X$  is subgaussian with scale parameter  $v$

if  $X$  has finite mean &  $\mathbb{E}[e^{s(X - \mathbb{E}X)}] \leq e^{\frac{s^2 v^2}{2}}$

for  $s \in \mathbb{R}$ .

$N(0, v^2)$

Hoeffding's lemma (Lemm 2.1)

Let  $X$  r.v.  $[b, c]$ , then  $\mathbb{E}[e^{s(X - \mathbb{E}X)}] \leq e^{s^2 \frac{(c-b)^2}{8}}$

• Any  $X \in [b, c]$ ,  $X$  is subgaussian with scale parameter

$$v = \frac{c-b}{2}.$$

• If  $S = \sum_{i=1}^n X_i$  where  $X_i$  are independent subgaussian with scale parameter  $v_i$ .

Then  $S$  is subgaussian,  $v^2 = \sum_{i=1}^n v_i^2$

$$R_n(A) = \mathbb{E} \sup_{a \in A} \frac{1}{n} \sum_{i=1}^n |\varepsilon_i a_i| \quad \left[ -\frac{a_i}{n}, \frac{a_i}{n} \right]$$

$$= \mathbb{E} \sup_{a \in A \cup -A} \frac{1}{n} \sum_{i=1}^n \varepsilon_i a_i$$

$$\gamma_k = \frac{1}{n} \sum_{i=1}^n \varepsilon_i a_i^{(k)} \quad \downarrow$$

$$v^2 = \frac{1}{n} \sum_{i=1}^n \frac{a_i^2}{n^2} \leq \frac{L^2}{n^2}$$

$$= \mathbb{E} \max \{ \gamma_1, -\gamma_1, \dots, \gamma_N, -\gamma_N \}.$$

want to show  $\leq \frac{2L \sqrt{\log N}}{n}$

Maximal lemma for subgaussian r.v (Lemma 2.3)

Suppose  $X_1 \dots X_n$ ,  $\mathbb{E}X_i = 0$  with scale parameter  $V$ .

$$\text{Then } \mathbb{E} \left[ \max_{i \in [n]} x_i \right] \leq V \sqrt{2 \log n}$$

$$V = \frac{L}{n}$$

$$R(A) \leq \frac{L}{n} \cdot \sqrt{2 \log 2N} \leq \frac{2L}{n} \sqrt{\log N}.$$

□.

# Structural results for Rademacher averages

## Proposition 6.2 (Contraction principles for Rademacher averages).

If  $\mathcal{A}$  is a bounded subset of  $\mathbb{R}^n$  and for  $i \in [n]$ ,  $\varphi_i : \mathbb{R} \rightarrow \mathbb{R}$  is an  $M$ -Lipschitz continuous function, then  $R_n^\circ(\varphi \circ \mathcal{A}) \leq M R_n^\circ(\mathcal{A})$ . Furthermore, if  $\varphi_i(0) = 0$  for all  $i$  (i.e.,  $\varphi(\mathbf{0}) = \mathbf{0}$ ) then  $R_n(\varphi \circ \mathcal{A}) \leq 2M R_n(\mathcal{A})$

**Proof:**

def:  $\varphi \circ \mathcal{A}$

Let  $\varphi_i : \mathbb{R} \rightarrow \mathbb{R}$ ,  $\varphi \circ \vec{v} \equiv (\varphi_1(v_1), \varphi_2(v_2), \dots, \varphi_n(v_n))$ .

$\varphi \circ \mathcal{A} = \{ \varphi \circ \vec{v} : \vec{v} \in \mathcal{A} \}$ .

def:  $\varphi$  is  $M$ -Lipschitz continuous if for all  $x_1, x_2$

$$|\varphi(x_1) - \varphi(x_2)| \leq M |x_1 - x_2|.$$

$$\varphi \circ \tau = \varphi_1 \circ \varphi_2 \cdots \varphi_n \circ \tau$$

$\uparrow$

WTS:  $R_n^\circ(\varphi_i \circ A) \leq R_n^\circ(A)$        $\varphi_i = \varphi_1$

$$R_n^\circ(A) = \frac{1}{n} \mathbb{E} \left[ \sup_{a \in A} \sum_{i=1}^n \varepsilon_i a_i \right]$$

$$= \frac{1}{n} \mathbb{E} \left[ \sup_{a \in A} \left( a_1 + \sum_{i=2}^n \varepsilon_i a_i \right) \frac{1}{2} \right]$$

$$+ \sup_{a' \in A} \left( -a'_1 + \sum_{i=2}^n \varepsilon_i a'_i \right) \frac{1}{2} \Bigg]$$

$$= \frac{1}{2n} \mathbb{E} \left[ \sup_{a, a' \in A} \underbrace{|a_1 - a'_1|}_{\leq 1} + \sum_{i=2}^n \varepsilon_i a_i + \sum_{i=2}^n \varepsilon_i a'_i \right]$$



$$R_n^\circ(\varphi_1 \circ A) = \frac{1}{2n} \mathbb{E} \left[ \sup_{a, a' \in A} \underbrace{|\varphi(a_1) - \varphi(a'_1)|}_{\leq 1} + \sum_{i=2}^n \varepsilon_i a_i + \sum_{i=2}^n \varepsilon_i a'_i \right]$$

$$R_n^\circ(\varphi_1 \circ A) \leq R_n^\circ(\mathcal{A})$$


---

$$R_n(\varphi \circ A) = R_n^\circ((\varphi \circ A) \cup (-\varphi \circ A))$$

$$= R_n^\circ((\varphi \circ A) \cup (-\varphi \circ A) \cup \{\vec{0}\})$$

$$\begin{array}{l} R_n^\circ(A \cup B) \\ \text{if } \vec{0} \in A, B. \end{array} \leftarrow \leq R_n^\circ((\varphi \circ A) \cup \{\vec{0}\}) + R_n^\circ(-\varphi \circ A \cup \{\vec{0}\})$$

$$\begin{array}{l} \text{then } R_n^\circ(A) + R_n^\circ(B) \\ \geq R_n^\circ(A \cup B) \end{array}$$

$$= R_n^\circ(\varphi \circ (A \cup \{\vec{0}\})) + R_n^\circ(-\varphi \circ (A \cup \{\vec{0}\}))$$

$$\leq R_n^0(A \cup \{\vec{0}\}) + R_n^0(-A \cup \{\vec{0}\}).$$

$$R_n^0(aA) = |a| R_n^0(A) \rightarrow = 2 R_n^0(\underline{A \cup \{\vec{0}\}})$$

$$\leq 2 R_n^0(\underline{A \cup -A})$$

$$= 2 R_n(A).$$

