

# Computational and Statistical Learning Theory

TTIC 31120

**Prof. Nati Srebro**

Lecture 9:

Real-Valued Loss:

Covering Numbers, VC-Subgraph Dimension  
and Rademacher Complexities

# Beyond Binary Classification

$$\min_h L(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\text{loss}(h(x); y)]$$

- Until now:  $h: \mathcal{X} \rightarrow \{\pm 1\}$ ,  $\text{loss}(\hat{y}; y) = \mathbb{1}[\hat{y} \neq y]$

- More generally:

$$\mathcal{D}(\mathcal{X}, \mathcal{Y}), \quad h: \mathcal{X} \rightarrow \hat{\mathcal{Y}}, \quad \text{loss}: \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow \mathbb{R}$$

- Regression:

$$\mathcal{Y} = \hat{\mathcal{Y}} = \mathbb{R}, \quad \text{loss}(\hat{y}, y) = (y - \hat{y})^2 \text{ or } \text{loss}(\hat{y}, y) = |\hat{y} - y|$$

- Multiclass:

$$\mathcal{Y} = \hat{\mathcal{Y}} = \{1, \dots, 10\}, \quad \text{loss}(\hat{y}, y) = \mathbb{1}[\hat{y} \neq y] \text{ or cost-sensitive loss}$$

- More complex label spaces,  $\mathcal{Y}$  = sequences, graphs, images, sentences...

- Even for binary classification  $\mathcal{Y} = \{\pm 1\}$ , might want  $\hat{\mathcal{Y}} = \mathbb{R}$

- Hinge loss:  $\text{loss}(\hat{y}, y) = [1 - \hat{y} \cdot y]_+$

- Logistic loss:  $\text{loss}(\hat{y}, y) = \log(1 + e^{-\hat{y} \cdot y})$

- Exp-loss:  $\text{loss}(\hat{y}, y) = e^{-\hat{y} \cdot y}$

# General Learning

$$\min_{h \in \mathcal{H}} L(h) = \mathbb{E}_{z \sim \mathcal{D}} [\ell(h, z)]$$

for a given  $\ell: \overline{\mathcal{H}} \times \mathcal{Z} \rightarrow \mathbb{R}$  and an unknown  $\mathcal{D}(\mathcal{Z})$ , based on an i.i.d. sample  $z_1, \dots, z_m \sim \mathcal{D}$

- **Supervised learning:**

- $\mathcal{Z} = \mathcal{X} \times \mathcal{Y} = \{z = (x, y) \mid x \in \mathcal{X}, y \in \mathcal{Y}\}$
- $\overline{\mathcal{H}} = \{h: \mathcal{X} \rightarrow \hat{\mathcal{Y}}\}$
- $\ell(h, z) = \text{loss}(h(x); y)$
- Proper Learning: compete with  $\inf_{h \in \overline{\mathcal{H}}} L(h)$  restricted to  $\overline{\mathcal{H}} = \mathcal{H}$
- Improper Learning: compete with  $\inf_{h \in \mathcal{H}} L(h)$  using  $\mathcal{H} = \hat{\mathcal{Y}}^{\mathcal{X}}$

# General Learning—Examples

$$\min_{h \in \mathcal{H}} L(h) = \mathbb{E}_{z \sim \mathcal{D}} [\ell(h, z)] \text{ based on } z_1, \dots, z_m \sim \text{iid } \mathcal{D}$$

- **Supervised learning:**
  - $\mathcal{Z} = \mathcal{X} \times \mathcal{Y} = \{z = (x, y) \mid x \in \mathcal{X}, y \in \mathcal{Y}\}$
  - $h: \mathcal{X} \rightarrow \mathcal{Y}$
  - $\ell(h, z) = \text{loss}(h(x); y)$
  - Improper learning: compete with  $\mathcal{H} \subset \overline{\mathcal{H}} = \mathcal{Y}^{\mathcal{X}}$
  - Proper learning:  $\mathcal{H} = \overline{\mathcal{H}} \subset \mathcal{Y}^{\mathcal{X}}$
- **Unsupervised learning, e.g.  $k$ -means clustering:**
  - $z = x \in \mathbb{R}^d$ ,
  - $h = (\mu[1], \mu[2], \dots, \mu[k]) \in \mathbb{R}^{d \times k}$  specified  $k$  cluster centers
  - $\ell((\mu[1], \mu[2], \dots, \mu[k]), x) = \min_i \|\mu[i] - x\|^2$
- **Density estimation:**
  - $z = x$  in some measurable space  $\mathcal{Z}$  (e.g.  $\mathbb{R}^d$ )
  - $h$  specifies probability density  $p_h(z)$
  - $\ell(h, z) = -\log p_h(z)$
  - Proper learning: fit density model of certain form
  - Improper learning:  $\overline{\mathcal{H}} = \{ \text{all probability densities } p(z) \text{ over } \mathcal{Z} \}$
- **Learning a good route with random traffic:**
  - $z =$  traffic delays on each road segment
  - $h =$  route chosen (indicator over road segments in route)
  - $\ell(h, z) = \langle h, z \rangle =$  total delay along route

# Learning in the General Setting

$$\min_{h \in \mathcal{H}} L(h) = \mathbb{E}_{z \sim \mathcal{D}} [\ell(h, z)] \text{ based on } z_1, \dots, z_m \sim \text{iid } \mathcal{D}$$

- Generalization of **(agnostic) PAC learning**: Hypothesis class  $\mathcal{H} \subseteq \overline{\mathcal{H}}$  agnostically learnable using learning rule  $A: \mathcal{Z}^* \rightarrow \overline{\mathcal{H}}$  with sample complexity  $m(\epsilon, \delta)$  if for all  $\mathcal{D}$  and all  $\epsilon, \delta > 0$ ,  $\forall_{S \sim \mathcal{D}^{m(\epsilon, \delta)}} L(A(S)) \leq \inf_{h \in \mathcal{H}} L(h) + \epsilon$
- **ERM**:  $ERM_{\mathcal{H}}(S) = \arg \min_{h \in \mathcal{H}} L_S(h)$        $L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i)$
- **Uniform Law of Large Numbers (ULLN)**:  
$$\Pr_{S \sim \mathcal{D}^m} [\forall_{h \in \mathcal{H}} |L(h) - L_S(h)| \leq \epsilon] \geq 1 - \delta$$
- Basic ingredients:
  - LLN:  $\forall_h \Pr_{S \sim \mathcal{D}^m} [ |L(h) - L_S(h)| \leq \epsilon ] \geq 1 - \delta$
  - Bound number of different behaviors on  $S$
  - Union bound over different behaviors
  - Symmetrization

# Law of Large Numbers— Real Valued Random Variables

- **Theorem (Hoeffding Bound):** Let  $V_1, V_2, \dots, V_m$  be independent random variables, with  $\forall_i \Pr[a \leq V_i \leq b] = 1$ , then for  $\bar{V} = \frac{1}{m} \sum_{i=1}^m V_i$ :

$$\Pr \left[ |\bar{V} - \mathbb{E}[\bar{V}]| > (b - a) \sqrt{\frac{\log^2 / \delta}{2m}} \right] < \delta$$

- In our case:  $V_i = \ell(h, z_i)$ ,  $\bar{V} = L_S(h)$ ,  $\mathbb{E}[\bar{V}] = \mathbb{E}[V_i] = L(h)$
- Conclusion:  $\forall_h \forall_{S \sim \mathcal{D}^m}^\delta |L_S(h) - L(h)| \leq a \sqrt{\frac{\log^2 / \delta}{2m}}$
- Applying union bound we can get (same as binary case):
  - Cardinality-based bound for finite classes
  - Description-length bounds
  - PAC-Bayes bounds

# Counting Real-Valued Behaviors

- For a class of functions  $\mathcal{F} = \{f: \mathcal{Z} \rightarrow \mathbb{R}\} \subset \mathbb{R}^{\mathcal{Z}}$ , and a sample  $S \in \mathcal{Z}^m$ , consider all possible behaviors
 
$$\mathcal{F}[S] = \{ (f(z_1), f(z_2), \dots, f(z_m)) \in \mathbb{R}^m \mid f \in \mathcal{F} \}$$
- Maybe infinitely many behaviors, but how many up to some resolution  $\alpha$ ?
- More specifically: how many behaviors  $V \subset \mathbb{R}^m$  can capture all behaviors  $\mathcal{F}[S]$  up to error  $\alpha$ ?

- **Empirical covering numbers:**

$$\mathcal{N}_p(\mathcal{F}, \alpha, S) = \min_{V \subseteq \mathbb{R}^m} |V| \text{ s.t. } \forall f \in \mathcal{F} \exists v \in V \left( \underbrace{\frac{1}{m} \sum_{i=1}^m |f(z_i) - v_i|^p}_{\text{For } p = \infty: \sup_i |f(z_i) - v_i|} \right)^{1/p} \leq \alpha$$

$$\mathcal{N}_p(\mathcal{F}, \alpha, m) = \sup_{S \in \mathcal{Z}^m} \mathcal{N}_p(\mathcal{F}, \alpha, S)$$

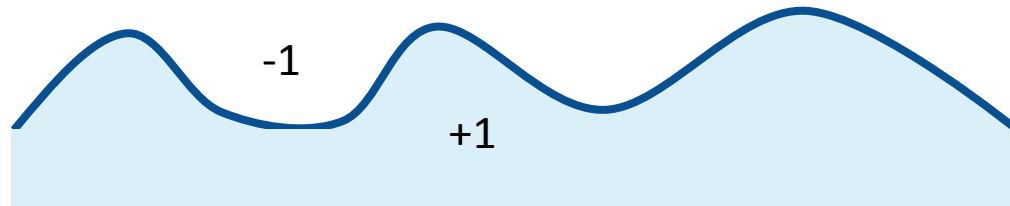
- At scale  $\alpha = 0$ , this is the growth function:  $\mathcal{N}_p(\mathcal{F}, 0, S) = \Gamma_{\mathcal{F}}(S)$
- But even if  $\mathcal{N}_p(\mathcal{F}, 0, S) = \infty$ , might be finite for  $\alpha > 0$

$$\mathcal{N}_p(\mathcal{F}, \alpha', S) \leq \mathcal{N}_p(\mathcal{F}, \alpha, S) \text{ for } \alpha' > \alpha$$

- Also:  $\mathcal{N}_{p'}(\mathcal{F}, \alpha, S) \geq \mathcal{N}_p(\mathcal{F}, \alpha, S)$  for  $p' > p$ , and in particular:
 
$$\mathcal{N}_1(\mathcal{F}, \alpha, S) \leq \mathcal{N}_2(\mathcal{F}, \alpha, S) \leq \mathcal{N}_{\infty}(\mathcal{F}, \alpha, S)$$

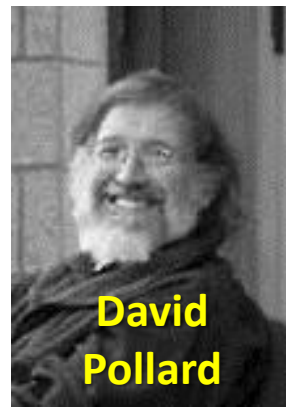
# Controlling Empirical Covering Numbers

- Definition:**  $\mathcal{F} \subset \mathbb{R}^{\mathcal{Z}}$  shatters  $S = \{z_1, \dots, z_m\}$  if  $\exists \theta_1, \theta_2, \dots, \theta_m \in \mathbb{R}$  s.t.  $\forall y_1, y_2, \dots, y_m \in \pm 1 \exists f \in \mathcal{F}$  s.t.  $\forall i$ :
 
$$y_i = +1 \Leftrightarrow f(z_i) > \theta_i$$
- Definition:** The **VC-subgraph dimension**  $\text{VCdim}(\mathcal{F})$  (aka **Pollard's pseudo-dimension**) of  $\mathcal{F}$  is the largest  $m$ , s.t. there exists  $S \in \mathcal{Z}^m$  that it shattered by  $\mathcal{F}$
- Equivalently:  $\text{VCdim}(\mathcal{F}) = \text{VCdim}(\{(z, \theta) \mapsto [f(z) \leq \theta] \mid f \in \mathcal{F}\})$



- Theorem:** For  $\mathcal{F} = \{f: \mathcal{Z} \rightarrow [-a, a]\}$  with  $\text{VCdim}(\mathcal{F}) \leq D$ :

$$\mathcal{N}_p(\mathcal{F}, \alpha, m) \leq \mathcal{N}_\infty(\mathcal{F}, \alpha, m) \leq \sum_{k=1}^D \binom{m}{k} \left(\frac{a}{\alpha}\right)^k \leq \left(\frac{em a}{D \alpha}\right)^D$$



David  
Pollard

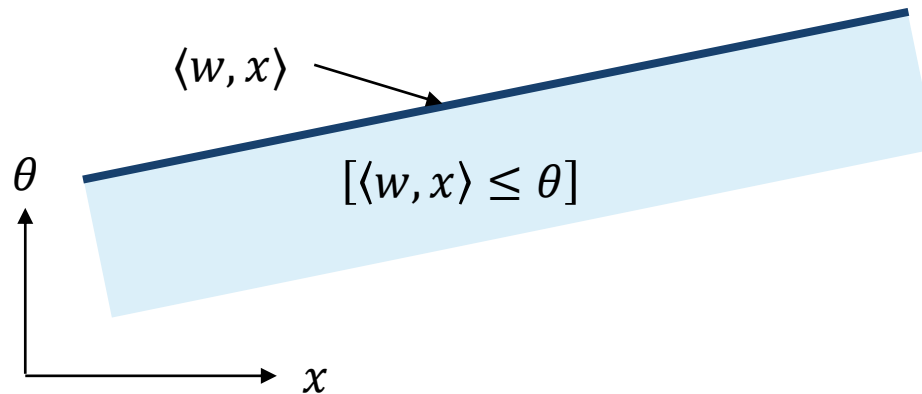


# VC Subgraph Dimension: Example

$$\mathcal{H} = \{x \mapsto \langle w, x \rangle \mid w \in \mathbb{R}^d\}$$

- Can shatter the  $d$  standard basis vectors, with thresholds  $\theta_i = 0$
- Why can't shatter more?
- Consider subgraph class:

$$\{(x, \theta) \mapsto [\langle w, x \rangle \leq \theta] \mid w \in \mathbb{R}^d\} \subset \text{half-spaces in } \mathbb{R}^{d+1}$$



- Conclusion:  $\text{VCdim}(\mathcal{H}) \leq d + 1$
- In fact,  $\text{VCdim}(\mathcal{H}) = d$

# What Function Class Should We Be Covering?

$$\min_{h \in \mathcal{H}} L(h) = \mathbb{E}_{z \sim \mathcal{D}} [\ell(h, z)]$$

Supervised learning:  $\ell(h, z = (x, y)) = \text{loss}(h(x); y)$ ,  $h: \mathcal{X} \rightarrow \mathcal{Y}$

- $\mathcal{N}_p(\mathcal{H}, \alpha, S)$  ?
  - $h: \mathcal{X} \rightarrow \mathcal{Y}$  might not be real valued (e.g.  $\mathcal{Y}$  =sentences)
  - In fact, for general learning problem,  $h$  might not be a function at all
  - Even if  $h: \mathcal{X} \rightarrow \mathbb{R}$ , we care about scale of loss, which might be different from scale of predictions

- Loss class:

$$\mathcal{F} = \{z \mapsto \ell(h, z) \mid h \in \mathcal{H}\}$$

Supervised learning:  $\mathcal{F} = \{(x, y) \mapsto \text{loss}(h(x); y) \mid h \in \mathcal{H}\}$

- Claim: For supervised learning with  $\mathcal{H} = \{h: \mathcal{X} \rightarrow \mathbb{R}\}$ ,
  - If  $\text{loss}(\hat{y}, y)$  is monotone in  $\hat{y}$  (for each  $y$ ):  $\text{VCdim}(\mathcal{F}) \leq \text{VCdim}(\mathcal{H})$
  - If  $\text{loss}(\hat{y}, y)$  is unimodal in  $\hat{y}$  (for each  $y$ ):  $\text{VCdim}(\mathcal{F}) \leq 2\text{VCdim}(\mathcal{H})$

# Back to the Plan

$$\text{ULLN: } \forall_S^\delta \forall_{h \in \mathcal{H}} |L(h) - L_S(h)| \leq \epsilon$$

↓

$$\text{Learning: } \forall_S^\delta L(\text{ERM}_{\mathcal{H}}(S)) \leq \inf_{h \in \mathcal{H}} L(h) + 2\epsilon$$

- Plan for establishing ULLN:

- LLN:  $\forall_h \forall_S^\delta |L(h) - L_S(h)| \leq \epsilon$  ✓

- Counting behaviors ✓

- Union bound

- Symmetrization

- Old-School approach:

- Use  $\mathcal{N}_\infty(\mathcal{F}, \alpha = \epsilon/4, 2m)$

- Approximate  $\ell(h, z_i)$  with  $v$  in cover, incurring  $\leq \epsilon/2$  error

- Union bound over  $v$  in cover

- Symmetrization as in binary case

- Contemporary approach: use Rademacher Complexity

# Rademacher Complexity

- Recall loss class  $\mathcal{F} = \{z \mapsto \ell(h, z) \mid h \in \mathcal{H}\}$
- For  $f \in \mathcal{F}$ , denote:  $\mathbb{E}_{\mathcal{D}}[f] = \mathbb{E}_{z \sim \mathcal{D}}[f(z)]$  and  $\mathbb{E}_S[f] = \frac{1}{m} \sum_{i=1}^m f(z_i)$
- Symmetrization: as a surrogate for  $\mathbb{E}_{\mathcal{D}}[f] - \mathbb{E}_S[f]$ , use  $\mathbb{E}_{S_1}[f] - \mathbb{E}_{S_2}[f]$
- Starting with  $S = \{z_1, \dots, z_m\}$ , use  $\xi_1, \dots, \xi_m \in \pm 1$  to define  $S_1 = \{z_i \mid \xi_i = +1\}$  and  $S_2 = \{z_i \mid \xi_i = -1\}$ , then if  $|S_1| = |S_2|$ :

$$\mathbb{E}_{S_1}[f] - \mathbb{E}_{S_2}[f] = \frac{2}{m} \sum_{i=1}^m \xi_i f(z_i)$$

- **Empirical Rademacher Complexity of  $\mathcal{F}$ :**

$$\mathcal{R}_S(\mathcal{F}) = \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_m \sim \text{iid unif } \pm 1} \left[ \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \xi_i f(z_i) \right]$$

- **Theorem:**

$$\mathbb{E}_{S \sim \mathcal{D}^m} \left[ \sup_{f \in \mathcal{F}} (\mathbb{E}_{\mathcal{D}}[f] - \mathbb{E}_S[S]) \right] \leq 2 \mathbb{E}_{S \sim \mathcal{D}^m} [\mathcal{R}_S(\mathcal{F})]$$



Dmitry  
Panchenko



Valdimir  
Koltchinskii



Sahar  
Mendelson



Peter  
Bartlett

# Rademacher Complexity $\rightarrow$ ULLN

$$\mathcal{R}_S(\mathcal{F}) = \mathbb{E}_{\xi_1, \xi_2, \dots, \xi_m \sim \text{iid unif } \pm 1} \left[ \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \xi_i f(z_i) \right]$$
$$\mathcal{R}_m(\mathcal{F}) = \sup_{S \in \mathcal{Z}^m} \mathcal{R}_S(\mathcal{F}) \quad \mathcal{R}_{\mathcal{D}^m}(\mathcal{F}) = \mathbb{E}_{S \sim \mathcal{D}^m} [\mathcal{R}_S(\mathcal{F})]$$

- Theorem:**  $\mathbb{E}_{S \sim \mathcal{D}^m} \left[ \sup_{f \in \mathcal{F}} (\mathbb{E}_{\mathcal{D}}[f] - \mathbb{E}_S[f]) \right] \leq 2\mathcal{R}_{\mathcal{D}^m}(\mathcal{F})$

Proof:  $\mathbb{E}_{S \sim \mathcal{D}^m} \left[ \sup_f (\mathbb{E}_{\mathcal{D}}[f] - \mathbb{E}_S[f]) \right] = \mathbb{E}_{S \sim \mathcal{D}^m} \left[ \sup_f \left( \mathbb{E}_{S' \sim \mathcal{D}^m} \left[ \frac{1}{m} \sum_i f(z'_i) \right] - \frac{1}{m} \sum_i f(z_i) \right) \right]$

$$= \mathbb{E}_{S \sim \mathcal{D}^m} \left[ \sup_f \mathbb{E}_{S' \sim \mathcal{D}^m} \left[ \frac{1}{m} \sum_i (f(z'_i) - f(z_i)) \right] \right]$$
$$\leq \mathbb{E}_{S, S' \sim \mathcal{D}^m} \left[ \sup_f \frac{1}{m} \sum_i (f(z'_i) - f(z_i)) \right] = \mathbb{E}_{S, S' \sim \mathcal{D}^m, \xi} \left[ \sup_f \frac{1}{m} \sum_i \xi_i (f(z'_i) - f(z_i)) \right]$$
$$\leq \mathbb{E}_{S' \sim \mathcal{D}^m, \xi} \left[ \sup_f \frac{1}{m} \sum_i \xi_i f(z'_i) \right] + \mathbb{E}_{S \sim \mathcal{D}^m, \xi} \left[ \sup_f \frac{1}{m} \sum_i (-\xi_i) f(z_i) \right] = 2\mathcal{R}_{\mathcal{D}^m}(\mathcal{F})$$

- Corollary:**  $\mathbb{E}_{S \sim \mathcal{D}^m} \left[ \sup_{f \in \mathcal{F}} (\mathbb{E}_S[f] - \mathbb{E}_{\mathcal{D}}[f]) \right] \leq 2\mathcal{R}_{\mathcal{D}^m}(\mathcal{F})$

Proof: consider  $(-\mathcal{F})$  and note that  $\mathcal{R}(-\mathcal{F}) = \mathcal{R}(\mathcal{F})$

- What about high probability guarantees?**

# McDiarmid Inequality

- Definition:  $g: \mathcal{Z}^m \rightarrow \mathbb{R}$  satisfies the **bounded difference property** with  $c_1, \dots, c_m \in \mathbb{R}$  if

$$\forall_i \forall_{z_1, z_2, \dots, z_m, z'_i} |g(z_1, \dots, z_m) - g(z_1, \dots, z'_i, \dots, z_m)| \leq c_i$$

- **Theorem (McDiarmid)**: Let  $Z_1, \dots, Z_m$  be independent random variables over  $\mathcal{Z}$ . If  $g$  satisfies the bounded difference property then  $\forall_{z_1, \dots, z_m}^\delta$

$$|g(z_1, \dots, z_m) - \mathbb{E}[g(Z_1, \dots, Z_m)]| \leq \sqrt{\frac{1}{2} \sum_{i=1}^m c_i^2 \log^2 / \delta}$$

- Generalizes Hoeffding, where  $g(z_1, \dots, z_m) = \frac{1}{m} \sum_i z_i$ 
  - If  $a \leq z_i \leq b$  then  $g$  satisfies bounded difference with  $c_i = \frac{1}{m} (b - a)$

$$\sqrt{\frac{1}{2} \sum_{i=1}^m c_i^2 \log^2 / \delta} = \sqrt{\frac{1}{2} m \left(\frac{b-a}{m}\right)^2 \log^2 / \delta} = (b - a) \sqrt{\frac{\log^2 / \delta}{2m}}$$

# Rademacher $\rightarrow$ High Probability ULLN

- Use  $g(S) = \sup_{f \in \mathcal{F}} (\mathbb{E}_{\mathcal{D}}[f] - \mathbb{E}_S[f])$
- If  $0 \leq f(z) = \ell(h, z) \leq a$  then  $g(S)$  satisfies bounded diff with  $c_i = \frac{a}{m}$
- Recall  $\mathbb{E}_{S \sim \mathcal{D}^m} \left[ \sup_{f \in \mathcal{F}} (\mathbb{E}_{\mathcal{D}}[f] - \mathbb{E}_S[f]) \right] \leq 2\mathcal{R}_{\mathcal{D}^m}(\mathcal{F})$
- Conclusion:  $\forall_S^\delta \sup_{f \in \mathcal{F}} (\mathbb{E}_{\mathcal{D}}[f] - \mathbb{E}_S[f]) \leq 2\mathcal{R}_{\mathcal{D}^m}(\mathcal{F}) + a \sqrt{\frac{\log^2/\delta}{2m}}$
- And similarly for  $\mathbb{E}_S[f] - \mathbb{E}_{\mathcal{D}}[f]$ , yielding ULLN:

$$\forall_{S \sim \mathcal{D}^m}^\delta \sup_{f \in \mathcal{F}} |\mathbb{E}_{\mathcal{D}}[f] - \mathbb{E}_S[f]| \leq 2\mathcal{R}_{\mathcal{D}^m}(\mathcal{F}) + a \sqrt{\frac{\log^4/\delta}{2m}}$$

# Bounding the Rademacher Complexity

- How do we bound  $\mathcal{R}_S(\mathcal{F}) = \mathbb{E}_\xi \left[ \sup_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \xi_i f(z_i) \right]$  ?
- For a finite  $\mathcal{F}$ , apply Hoeffding to the tail of each  $\sum_{i=1}^m \xi_i f(z_i)$  and take a union bound, yielding (Massart's Finite Class Lemma):

$$\mathcal{R}_m(\mathcal{F}) \leq a \sqrt{\frac{\log |\mathcal{F}|}{2m}}$$

- Since  $\mathcal{R}_S(\mathcal{F})$  only depends on behavior inside  $S$ , can take union bound over behaviors:

$$\mathcal{R}_m(\mathcal{F}) \leq a \sqrt{\frac{\log \Gamma_{\mathcal{F}}(m)}{2m}} = a \sqrt{\frac{\log \mathcal{N}_p(\mathcal{F}, \alpha = 0, m)}{2m}}$$

- Or over behaviors up to resolution  $\alpha$ :

$$\mathcal{R}_m(\mathcal{F}) \leq \inf_{\alpha} \left( \alpha + a \sqrt{\frac{\log \mathcal{N}_1(\mathcal{F}, \alpha, m)}{2m}} \right)$$



# Bounding the Rademacher Complexity

- Theorem: For  $0 \leq f \leq a$ :

$$\mathcal{R}_S(\mathcal{F}) \leq \inf_{\alpha} \left( \alpha + a \sqrt{\frac{\log \mathcal{N}_1(\mathcal{F}, \alpha, S)}{2m}} \right)$$

Proof: Consider a cover  $V$ ,  $|V| = \log \mathcal{N}_1(\mathcal{F}, \alpha, S)$ . For every  $f$ , let  $v^f \in V$  be s.t.  $\frac{1}{m} \sum_i |f(z_i) - v_i| \leq \alpha$

$$\begin{aligned} \mathcal{R}_S(\mathcal{F}) &= \mathbb{E}_{\xi} \left[ \sup_f \frac{1}{m} \sum_i \xi_i f(z_i) \right] = \mathbb{E}_{\xi} \left[ \sup_f \frac{1}{m} \sum_i \xi_i (v_i^f + f(z_i) - v_i^f) \right] \\ &\leq \mathbb{E}_{\xi} \left[ \sup_f \frac{1}{m} \sum_i \xi_i v_i^f \right] + \frac{1}{m} \sum_i |f(z_i) - v_i^f| \leq a \sqrt{\frac{\log |V|}{2m}} + \alpha \end{aligned}$$

- Conclusion:

$$\mathcal{R}_m(\mathcal{F}) \leq \inf_{\alpha} \left( \alpha + a \sqrt{\frac{\text{VC}(\mathcal{F}) \log \left( \frac{em}{\text{VC}(\mathcal{F})} \frac{a}{\alpha} \right)}{2m}} \right) \leq a \sqrt{\frac{2\text{VC}(\mathcal{F}) \log \left( \frac{em}{\text{VC}(\mathcal{F})} \right)}{m}}$$

$\alpha = a \sqrt{\frac{\text{VCd}(\mathcal{F})}{2m}}$

# Dudley's Integral

- Can improve over

$$\mathcal{R}_S(\mathcal{F}) \leq \inf_{\alpha} \left( \alpha + a \sqrt{\frac{\log \mathcal{N}_1(\mathcal{F}, \alpha, S)}{2m}} \right)$$

- Instead of using cover at a single scale, integrate over all scales:

$$\mathcal{R}_S(\mathcal{F}) \leq \inf_{\alpha_0} \left( 4\alpha_0 + 10 \int_{\alpha_0}^a \sqrt{\frac{\log \mathcal{N}_2(\mathcal{F}, \alpha, S)}{m}} d\alpha \right)$$

- Allows avoiding log-factor and obtaining:

$$\mathcal{R}_m(\mathcal{F}) = O \left( a \sqrt{\frac{\text{VCdim}(\mathcal{F})}{m}} \right)$$

- Even more important in infinite-dim classes



# Putting It All Together

- For supervised learning,  $\ell(h, (x, y)) = \text{loss}(h(x); y)$ , when
  - $0 \leq \text{loss}(\hat{y}; y) \leq a$
  - $\text{loss}(\hat{y}; y)$  is monotone or unimodal in  $\hat{y}$

then with probability  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^m$ :

$$L_{\mathcal{D}}(\text{ERM}_{\mathcal{H}}(S)) \leq \inf_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + O\left(a \sqrt{\frac{\text{VCdim}(\mathcal{H})}{m}}\right)$$

- **Need loss to be bounded**
  - OK for non-convex loss, e.g. 0/1, or truncated squared loss
  - Can bound targets  $y$  and “response”  $h(x)$
  - Only needed for high-probability bound (unavoidable!)---can get small error in expectation even without bounded loss.