

# Course Note for CPSC532D STAT Learn Theory

L2:

$$L_D(\hat{h}_S) - L_D(h^*) \leq L_D(\hat{h}_S) - L_S(\hat{h}_S) + L_S(h^*) - L_D(h^*)$$

uniform convergence  
estimation error

Hoeffding  
 $\frac{1}{m} \sum_{i=1}^m l(h^*, z_i) - \mathbb{E}_z l(h^*, z)$

$$L_D(\hat{h}_S) - L_{\text{Bayes}} = \underbrace{L_D(\hat{h}_S) - \inf_{h \in H} L_D(h)}_{\text{excess error}} + \underbrace{\inf_{h \in H} L_D(h) - L_{\text{Bayes}}}_{\text{estimation error}} = \inf_{h \in H} L(h, y)$$

$\not\in$  measurable

approximation

increase  $H$  =  $\rightarrow$

$$\left\{ \begin{array}{l} \text{central limit} \\ \text{concentration inequalities} \end{array} \right. \begin{array}{l} = O(m) \\ = \sqrt{O(m)} \end{array}$$

**Hoeffding:** Let  $x_1, \dots, x_m$  be iid with =

- $\mathbb{E} x_i = \mu$
- $\Pr(a \leq x_i \leq b) = 1$

Let  $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$ , Then

$$\Pr(\bar{x} \leq \mu + (b-a)\sqrt{\frac{\log(1/\delta)}{2m}}) \geq 1-\delta \Leftrightarrow \Pr(\bar{x} > \mu + (b-a)\sqrt{\frac{\log(1/\delta)}{2m}}) < \delta$$

$$\therefore Y_i = -x_i. \text{ Then } \Pr(\bar{Y} \leq -\mu + (b-a)\sqrt{\frac{\log(1/\delta)}{2m}}) \geq 1-\delta$$

$$\hookrightarrow \Pr(\bar{x} \geq \mu - (b-a)\sqrt{\frac{\log(1/\delta)}{2m}}) \geq 1-\delta \Leftrightarrow \Pr(\bar{x} < \mu - (b-a)\sqrt{\frac{\log(1/\delta)}{2m}}) < \delta$$

$$\Rightarrow \Pr(|x - \mu| > \mu - (b-a)\sqrt{\frac{\log(1/\delta)}{2m}}) \leq 2\delta$$

$$L_S(h^*) - L_D(h^*)$$

$$\Pr(L_S(h^*) - L_D(h^*) \geq (b-a)\sqrt{\frac{1/\delta}{2m}}) \leq \delta$$

$\nabla$  we can't use Hoeffding on  $L_D(\hat{h}_S) - L_S(\hat{h}_S)$

$l(\hat{h}_S, z_i)$  are not independent of each other

change  $z_i$ , will change  $\hat{h}_S$

$\hookrightarrow$  Uniform convergence !! : For every single  $h \in H$ , gap is not too big

$$\sup_{h \in H} L_D(h) - L_S(h) \leq \varepsilon \Rightarrow (\forall h \in H L_D(h) - L_S(h) \leq \varepsilon)$$

① if  $|H| < \infty$  : union bound

$$q_S(h) = L_D(h) - L_S(h)$$

$$\sup_{h \in H} q_S(h) < \varepsilon$$

Not tight enough because of  $P(A \cup B) \leq P(A) + P(B)$

$$\Pr(\exists h \in H, q_S(h) > \varepsilon) \leq \sum_{h \in H} \Pr(q_S(h) > \varepsilon)$$

now  $h$  can be fixed

Assume bounded

$\{lh, z\} \in [a, b] \forall z, h,$   
 $H$  finite  
 $\hat{h}_S$  is an ERM

$$\begin{aligned} L_D(\hat{h}_S) - \min_{h \in H} L_D(h) &\leq L_D(\hat{h}_S) - L_S(\hat{h}_S) + L_S(\hat{h}^*) - L_D(\hat{h}^*) \\ &\leq \max_{h \in H} [L_D(h) - L_S(h)] + L_S(\hat{h}^*) - L_D(\hat{h}^*) \end{aligned}$$

by showing:

- $\Pr(\forall h \in H, \Pr(L_D(h) - L_S(h) > \varepsilon) \leq \frac{\delta}{|H|+1})$
- $\Pr(L_S(\hat{h}^*) - L_D(\hat{h}^*) > \varepsilon) \leq \frac{\delta}{|H|+1}$

$$\therefore \Pr(L_D(\hat{h}_S) - \min_{h \in H} L_D(h) \leq 2\varepsilon) \geq 1 - \delta$$

$$\varepsilon = (b-a) \sqrt{\frac{1}{2m} \log \frac{|H|+1}{\delta}}$$

$$\Pr(L_D(\hat{h}_S) - \min_{h \in H} L_D(h) \leq \sqrt{\frac{2}{m} \log \frac{|H|+1}{\delta}}) \geq 1 - \delta$$

### L3. MARKOV

$X$  cannot be too big with high prob  $\Pr(X \geq 0) = 1$

- Markov's inequality:  $\Pr(X \geq t) \leq \frac{E[X]}{t}$  for all  $t > 0$  is weak because only assume non-neg

Pf  $\begin{cases} x \geq 0 \\ x \geq t \text{ when } x \geq t \end{cases} \therefore \Pr(X \geq t) \leq \Pr\left(\frac{E[X]}{t} \geq t\right)$   
 $\text{let } \delta = \frac{E[X]}{t} \text{ then } \Pr\left(X \geq \frac{E[X]}{\delta}\right) \leq \delta \Rightarrow \Pr\left(X \leq \frac{E[X]}{\delta}\right) \geq 1 - \delta$

- Chebychev's inequality: for any  $X$ ,  $\Pr(|X - E[X]| \geq \epsilon) \leq \frac{\text{Var}[X]}{\epsilon^2}$

Pf  $\Pr(|X - E[X]| \geq \epsilon) = \Pr((X - E[X])^2 \geq \epsilon^2) \leq \frac{1}{\epsilon^2} E[(X - E[X])^2] = \frac{\text{Var}[X]}{\epsilon^2}$   
Markov inequality  
with probability at least  $1 - \delta$ ,  $|X - E[X]| \leq \sqrt{\frac{\text{Var}[X]}{\delta}} = \frac{\sigma}{\sqrt{\delta}}$

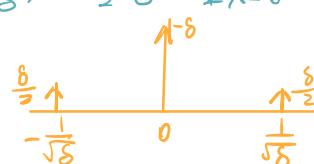
- Consider iid  $X_1, X_2, \dots, X_m$ ,  $E[X_i] = \mu$ ,  $\text{Var}(X_i) = \sigma^2$   
 $\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i$     $E[\bar{X}] = \frac{1}{m} \sum_{i=1}^m E[X_i] = \mu$     $\text{Var}(\bar{X}) = \frac{1}{m^2} \cdot m \sigma^2 = \frac{\sigma^2}{m}$   
 $\Rightarrow$  with  $\Pr \geq 1 - \delta$ ,  $|\bar{X} - \mu| \leq \frac{\sigma}{\sqrt{m\delta}}$ .

Chebychev's inequality may not be good enough!

$\sqrt{m}$  is fine, but the dependence on  $\delta$  is the problem  $\sqrt{\log \frac{1}{\delta}}$  better than  $\sqrt{\delta}$

$$\Pr(X=0) = 1 - \delta, \Pr(X=\frac{1}{\sqrt{\delta}}) = \Pr(X=-\frac{1}{\sqrt{\delta}}) = \frac{1}{2}\delta \quad E[X]=0 \quad \text{Var}[X]=1$$

$$\Pr(|\bar{X}| \leq \frac{1}{\sqrt{m\delta}}) \geq 1 - \delta$$



Chebychev's cares about the worst case

- Chernoff bounds: construct a non-negative random variable

$$\lambda > 0 \text{ arbitrary} \quad Y = e^{\lambda(X-\mu)} \quad \Pr(Y \geq t) \leq \frac{1}{t} E[Y]$$

$$\Pr(e^{\lambda(X-\mu)} \geq e^{\lambda\epsilon}) \leq e^{-\lambda\epsilon} E[e^{\lambda(X-\mu)}]$$

$$= \Pr(\lambda(X-\mu) \geq \lambda\epsilon) \quad \text{centred moment-generating function } M_X(t)$$

$$e^t = 1 + t + \frac{t^2}{2!} + \frac{t^3}{3!} + \dots$$

$$M_X(\lambda) = \mathbb{E} e^{\lambda(X-\mu)} = 1 + \lambda \underbrace{\mathbb{E}(X-\mu)}_0 + \frac{\lambda^2}{2!} \mathbb{E}[(X-\mu)^2] + \frac{\lambda^3}{3!} \mathbb{E}[(X-\mu)^3] + \dots$$

$k$ th derivative of  $M_X(\lambda)$  at  $\lambda=0$  :  $M_X^{(k)}(0) = \mathbb{E}[(X-\mu)^k]$

If  $X \sim \mathcal{N}(\mu, \sigma^2)$ , then  $\mathbb{E} e^{\lambda(X-\mu)} = e^{\frac{1}{2}\lambda^2\sigma^2}$

for  $X \sim \mathcal{N}(0, 1)$

$$\mathbb{E} e^{\lambda X} = \int_{X \sim \mathcal{N}(0, 1)} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} e^{\lambda x} dx = e^{-\frac{1}{2}\lambda^2} \underbrace{\int_{X \sim \mathcal{N}(0, 1)} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\lambda)^2} dx}_{= \mathcal{N}(\lambda, 1)} = e^{-\frac{1}{2}\lambda^2}$$

for  $Y \sim \mathcal{N}(0, \sigma^2)$

$$\mathbb{E} e^{\lambda Y} = \mathbb{E} e^{\lambda \sigma X} = \mathbb{E} e^{\lambda \sigma X} = e^{-\frac{1}{2}\lambda^2\sigma^2}$$

$$\Pr(X \geq \mu + \varepsilon) \leq e^{-\lambda\varepsilon} e^{\frac{1}{2}\lambda^2\sigma^2}$$

$\downarrow$  we can optimize  $\lambda$  to get tighter bound  
 $\sigma^2\lambda = \varepsilon \rightarrow \lambda = \frac{\varepsilon}{\sigma^2} \rightarrow e^{\frac{1}{2}\frac{\varepsilon^2}{\sigma^2} - \frac{\varepsilon^2}{\sigma^2}} = e^{-\frac{\varepsilon^2}{2\sigma^2}}$

$$\Pr(X \geq \mu + \varepsilon) \leq e^{-\frac{\varepsilon^2}{2\sigma^2}} \quad e^{-\frac{\varepsilon^2}{2\sigma^2}} = \delta \quad \varepsilon = \sqrt{2\sigma^2 \log \frac{1}{\delta}}$$

For  $X \sim \mathcal{N}(\mu, \sigma^2)$ :

\* with probability at least  $1 - \delta$   $X \in [\mu - \sqrt{2 \log \frac{1}{\delta}}, \mu + \sqrt{2 \log \frac{1}{\delta}}]$

If  $X$  is s.t.  $\mathbb{E} e^{\lambda(X-\mathbb{E}X)} \leq e^{\frac{1}{2}\lambda^2\sigma^2}$ , then  $\Pr(X \in [\mathbb{E}X - \sqrt{2 \log \frac{1}{\delta}}, \mathbb{E}X + \sqrt{2 \log \frac{1}{\delta}}]) \geq 1 - \delta$

SG( $\sigma$ ) sub-gaussian with  $\sigma$  (no heavier tails than Gaussian)

$\sigma$  does not imply the variance of the sub-gaussian  $X$   
may have relation but not equal!

- sub-Gaussian -  $X \sim SG(\sigma)$

random variable  $X$  with mean  $\mu = \mathbb{E}[X]$  is sub-Gaussian with  $\sigma$   
 if  $\mathbb{E}[e^{\lambda(X-\mu)}]$  exists and for all  $\lambda \in \mathbb{R}$ ,  $\mathbb{E}[e^{\lambda(X-\mu)}] \leq e^{\frac{1}{2}\lambda^2\sigma^2}$

if  $\sigma_1 < \sigma_2$  then  $X \sim SG(\sigma_1) \Rightarrow X \sim SG(\sigma_2)$   $SG(\sigma_1) \subseteq SG(\sigma_2)$

(i) If  $X \in SG(\sigma)$ ,  $aX \in SG(|a|\sigma)$  for  $\forall a \in \mathbb{R}$

$$M_X(\lambda) \leq e^{\frac{1}{2}\lambda^2\sigma^2} \Rightarrow \mathbb{E} e^{\lambda(aX-\mathbb{E}aX)} = \mathbb{E} e^{a\lambda(X-\mathbb{E}X)} \leq e^{\frac{1}{2}(a\lambda)^2\sigma^2} = e^{\frac{1}{2}\lambda^2(a\sigma)^2}$$

(ii) If  $X_1 \in SG(\sigma_1)$ ,  $X_2 \in SG(\sigma_2)$  are independent  $X_1 + X_2 \in SG(\sqrt{\sigma_1^2 + \sigma_2^2})$

$$\mathbb{E} e^{\lambda(X_1+X_2-\mathbb{E}(X_1+X_2))} = \mathbb{E} e^{\lambda(X_1-\mathbb{E}X_1)} \mathbb{E} e^{\lambda(X_2-\mathbb{E}X_2)} \leq e^{\frac{1}{2}\lambda^2\sigma_1^2} \cdot e^{\frac{1}{2}\lambda^2\sigma_2^2}$$

we can use the above properties to construct a set of  $X$

Hoeffding's lemma :

a real-valued random variable bounded in  $[a, b]$  is  $SG(\frac{b-a}{2})$   
i.e.  $\Pr(a \leq X \leq b) = 1, X \in SG(\frac{b-a}{2})$  see pf in notes 3.

- Hoeffding: If  $X_1, X_2, \dots, X_m$  iid with  $\mathbb{E}X_i = \mu$ ,  $\Pr(a \leq X_i \leq b) = 1$   
then  $\Pr\left(\frac{1}{m} \sum_{i=1}^m X_i > \mu + (b-a)\sqrt{\frac{1}{2m} \log \frac{1}{\delta}}\right) \leq \delta$   
 $\frac{1}{m} \sum_{i=1}^m X_i \in SG\left(\frac{1}{m} \sqrt{m} \left(\frac{b-a}{2}\right)^2\right) = SG\left(\frac{b-a}{2\sqrt{m}}\right)$   
with prob at least  $1-\delta$   $\frac{1}{m} \sum_{i=1}^m X_i \leq \mu + \frac{b-a}{2\sqrt{m}} \sqrt{2 \log \frac{1}{\delta}} = \mu + (b-a)\sqrt{\frac{\log \frac{1}{\delta}}{2m}}$

## Lec 4 . PAC Learning (infinite $\mathcal{H}$ )

Defn : an algorithm  $A$  is agnostically PAC learns  $\mathcal{H}$  with loss  $l$   
if there is a function  $m(\epsilon, \delta)$  sample complexity function  
s.t. for any  $D$ , for any  $\epsilon, \delta \in (0, 1)$   
if  $S \sim D^m$  with  $m \geq m(\epsilon, \delta)$   
then  $\Pr_{\text{next}}(L_0(A(S)) \leq \inf_{h \in \mathcal{H}} L_0(h) + \epsilon) \geq 1 - \delta$

solve for  $m$  will get sample complexity

efficient = polynomial run time

Def  $\mathcal{H}$  is agnostically PAC learnable if  $\exists A$  that agnostically PAC learn  $\mathcal{H}$   
agnostic PAC is the worst case,  $\mathcal{H}$  has nothing to do with data distribution.

so  $m(\epsilon, \delta)$  should work for any  $D$

PAC learnable does not show <sup>how</sup> learn quickly in terms of  $m$ .

{ Agnostically PAC learning : work for any distribution ERM infinite  
PAC learning is weaker : work only for realizable distribution

Def (PAC-Learn) A PAC-learns  $\mathcal{H}$  if

$$\exists m: (0, 1)^2 \rightarrow N \text{ s.t.}$$

for any  $(\epsilon, \delta)$ , for any realizable  $\mathcal{D}$ , if  $m \geq m(\epsilon, \delta)$

$$\Pr_{\mathcal{D}}(L_D(A(S)) \leq \epsilon) \geq 1 - \delta$$

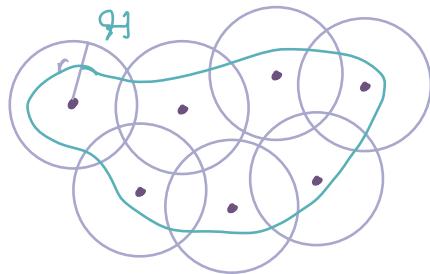
$$\sup_m$$

Def (realizable). for a non-negative loss  $l(h, z) \geq 0$ .

$\mathcal{D}$  is realizable by  $\mathcal{H}$  if there exists an  $h^* \in \mathcal{H}$  s.t.  $L_D(h^*) = 0$

We already know finite case, what about infinite  $\mathcal{H}$ ?

Idea of covering number



still: uniform convergence

$$\begin{aligned} & \sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)| \\ &= \sup_{h \in \mathcal{H}} |L_D(h) - L_D(h_j) + L_D(h_j) - L_S(h_j) + L_S(h_j) - L_S(h)| \\ &\leq \sup_{h \in \mathcal{H}} [L_D(h) - L_D(h_j)] + \max_{h \in T} [L_D(h_j) - L_S(h_j)] + \sup_{h \in \mathcal{H}} [L_S(h_j) - L_S(h)] \end{aligned}$$

Lipschitz union bound on size of  $T$

as  $\rho$  decrease, the set to cover points become tighter

$T$  is a  $\rho$ -cover for  $\mathcal{H}$  of size  $N(\mathcal{H}, \rho)$

specific case: Logistic regression

$$Z = X \times Y \quad X = \mathbb{R}^d \quad Y = \{-1, 1\}$$

$$\mathcal{H} = \{x \mapsto w \cdot x : w \in \mathbb{R}^d, \|w\| \leq B\} \quad \text{equivalent to } L_2 \text{ reg}$$

$$l_{\log}(h, (x, y)) = l(y, h(x)) = \log(1 + \exp(-y h(x)))$$

$$|L_D(h_w) - L_D(h_v)| = \left| \mathbb{E}_{(x, y) \sim D} l_{\log}(h_w, (x, y)) - l_{\log}(h_v, (x, y)) \right|$$

use absolute value to make analysis easier, but results looser

$$\begin{aligned} &\leq \mathbb{E} |l_y(h_w(x)) - l_y(h_v(x))| \quad l_y \text{ is } G\text{-Lipschitz} \\ &\leq G \mathbb{E} |h_w(x) - h_v(x)| \\ &\quad G=1 \text{ for } l_y \end{aligned}$$

input close  $\rightarrow$  output close  
smoothness



$G$ -Lipschitz:  $l_y$  is  $G$ -Lips if

$$|l_y(\hat{y}_1) - l_y(\hat{y}_2)| \leq G |\hat{y}_1 - \hat{y}_2| \quad \forall \hat{y}_1, \hat{y}_2$$

Euclidean distance (can be other distance if particular specified) if  $f$  is differentiable

$\Rightarrow$  if  $f'$  exist everywhere, and  $\sup |f'(x)| \leq G$ ,  $f$  is  $G$ -Lips the Lips constance is

$$|f(x) - f(x')| = \left| \int_x^{x'} f'(t) dt \right| \leq \int_x^{x'} |f'(t)| dt \leq \int_x^{x'} G dt = G|x' - x|$$

the largest |derivative| (upper bound of derivative |)

for the LR:  $|h_w(x) - h_v(x)| = |w \cdot x - v \cdot x| = |\langle w - v \rangle \cdot x| \leq \|w - v\| \cdot \|x\|$

$$|l_y(\hat{y})| = \log(1 + \exp(-y\hat{y})) \quad l_y'(\hat{y}) = \frac{-y \exp(-y\hat{y})}{1 + \exp(-y\hat{y})} = \frac{-y}{\exp(y\hat{y}) + 1} = \begin{cases} -\frac{1}{e^{\hat{y}} + 1}, & y=1 \\ \frac{1}{e^{-\hat{y}} + 1}, & y=-1 \end{cases}$$

$$|h_w(x) - h_v(x)| = |w \cdot x - v \cdot x| = |\langle w - v \rangle \cdot x| \leq \|w - v\| \cdot \|x\|$$

Then ⑥  $|\mathcal{L}_D(h_w) - \mathcal{L}_D(h_v)| \leq \underbrace{(\mathbb{E} \|x\|)}_{\text{constant of distribution}} \cdot \underbrace{\|w - v\|}_{\text{distance from } h \text{ to } h_j \text{ at most } p}$

Assume  $\|x\| \leq c$  on  $\mathcal{D}$ .

$$\textcircled{2} \quad |\mathcal{L}_S(h_w) - \mathcal{L}_S(h_v)| \leq \frac{1}{m} \sum_{i=1}^m \|x_i\| \|w - v\| \leq CP$$

$$\textcircled{3} \quad |\mathcal{L}_D(h_j) - \mathcal{L}_S(h_j)| \leq \underbrace{(b-a)\sqrt{\frac{1}{2m} \log \frac{N(B, P)}}}_{\text{Ly is not bounded, but we assume } \|w\| \leq B} \quad \text{and } \|x_i\| \leq c \text{ so we can use Hoeffdings}$$

$$b-a \leq \log(1 + \exp(BC)) \leq BC+1 \quad (\text{just simpler})$$

$$\text{So } \sup_{h \in \mathcal{H}} \mathcal{L}_D(h) - \mathcal{L}_S(h) \leq 2(CP + (BC+1)\sqrt{\frac{1}{2m} \log \frac{N(B, P)}{\delta}})$$

what is the covering number?

$$N(B, P) \leq \left(\frac{3B}{P}\right)^d \quad (\text{proof in notes } \Psi)$$

$$\sup_{h \in \mathcal{H}} \mathcal{L}_D(h) - \mathcal{L}_S(h) \leq 2(CP + (BC+1)\sqrt{\frac{1}{2m}(\log \frac{1}{\delta} + d \log \frac{3B}{P})})$$

how to choose  $P$  to minimize the right term?

$$\begin{aligned} \text{let } P = \frac{B}{\sqrt{m}} \quad &\leq \frac{2CB}{\sqrt{m}} + (BC+1) \sqrt{\frac{\log \frac{1}{\delta}}{2m} + \frac{d \log(9m)}{2m}} \\ (\text{roughly optimal}) \quad &\text{Big O notation means some random variable is stochastically bounded. (Hoeffding)} \\ &\text{dominate by } O_p(\sqrt{\frac{\log m}{m}}) \end{aligned}$$

$$\Pr_{h \in \mathcal{H}} (\sup_{h \in \mathcal{H}} L_D(h) - L_S(h) > \frac{2BC}{\sqrt{m}} + (BC)^{\frac{1}{2}} \sqrt{\frac{\log(1/\delta) + \frac{d}{2} \log(9m)}{2m}} ) < \delta$$

- the sample complexity depends on the input distribution  $\mathbb{E}\|X\| \leq c$   
X not allowed in agnostically PAC learning  
 $\text{ERM}^+$   
So this bound doesn't show Linear Regression is agnostically PAC learnable
- we can never achieve 0 loss on any  $D$ , so it's not realizable

## L5 Rademacher complexity distribution-specific complexity

} covering number approach : depend on dimension d.  
need bounded norm. B. scale sensitive

Rademacher complexity : do not depend on dimension d.

here we are going to bound on average instead of high probability.

$$\mathbb{E}_{S \sim D^m} \sup_{h \in \mathcal{H}} L_D(h) - L_S(h) = \mathbb{E}_{S \sim D^m} \left( \sup_{h \in \mathcal{H}} \mathbb{E}_{S' \sim D^m} [L_{S'}(h) - L_S(h)] \right)$$

how much can I overfit to samples

$$\textcircled{1} \sup_x \mathbb{E}_y f_y(x) \leq \mathbb{E}_x \sup_y f_y(x)$$

Pf: for any  $y$ , we have  $f_y(x) \leq \sup_{y'} f_{y'}(x)$

$$\text{then } \mathbb{E}_x f_y(x) \leq \mathbb{E}_x \sup_{y'} f_{y'}(x)$$

$$\text{and then } \sup_x \mathbb{E}_y f_y(x) \leq (\sup_x) \mathbb{E}_y \sup_{y'} f_{y'}(x)$$

$$\stackrel{\textcircled{1}}{\leq} \mathbb{E}_{S, S' \sim D^m} \sup_{h \in \mathcal{H}} [L_{S'}(h) - L_S(h)]$$

$$= \mathbb{E}_{S, S'} \sup_h \frac{1}{m} \sum_{i=1}^m [\ell(h, z_i) - \ell(h, z'_i)]$$

$$\textcircled{2} = \mathbb{E}_{S, S'} \mathbb{E}_{U, U'} \left[ \sup_h \frac{1}{m} \sum_{i=1}^m \sigma_i (\ell(h, u_i) - \ell(h, u'_i)) \mid S, S', \vec{\sigma} \right]$$

not doing anything (deterministic w.r.t  $(S, S')$  condition)

switch order of expectation  $P(S), P(\vec{\sigma}) P(U, U') P(S, S')$

$$= \mathbb{E}_{U, U'} \mathbb{E}_{S, S'} \left[ \sup_h \frac{1}{m} \sum_{i=1}^m \sigma_i (\ell(h, u_i) - \ell(h, u'_i)) \mid U, U', \vec{\sigma} \right]$$

$P(U, U') P(S, S')$

deterministic  
we can drop S

\textcircled{2} let  $\sigma_i \in \{-1, 1\}$  for  $i \in [m]$ ,

$$(u_i, u'_i) = \begin{cases} (z_i, z'_i) & \text{if } \sigma_i = 1 \\ (z'_i, z_i) & \text{if } \sigma_i = -1 \end{cases}$$

for any  $\vec{\sigma} = (\sigma_1, \dots, \sigma_m)$

$$\text{then } \ell(h, z_i) - \ell(h, z'_i) = \sigma_i (\ell(h, z_i) - \ell(h, z'_i))$$

$$\sigma_i \sim \text{Unif}; P(\sigma_i = 1) = P(\sigma_i = -1) = \frac{1}{2}$$

$U, U' | S, S', \vec{\sigma}$  is determinated  
 but  $U, U' \sim D^m$ ,  $U, U', \vec{\sigma}$  are independent  
 $S, S' | U, U', \vec{\sigma}$  is determinated

$$\begin{aligned}
 &= \mathbb{E}_{U, U'} \mathbb{E}_{\vec{\sigma}} \left[ \sup_{h} \frac{1}{m} \sum_{i=1}^m \sigma_i (l(h, u_i) - l(h, u'_i)) \right] \\
 &\quad \text{rename } U \rightarrow S \\
 &= \mathbb{E}_{S, S' \sim D^m} \mathbb{E}_{\vec{\sigma}, h} \left[ \sup_{h} \frac{1}{m} \sum_{i=1}^m \sigma_i (l(h, z_i) - l(h, z'_i)) \right] \\
 &\quad \sup (\sigma_i (l(h, z_i) + (-\sigma_i) l(h, z'_i)))
 \end{aligned}$$

$$\begin{aligned}
 &\leq \mathbb{E}_{S, \sigma} \mathbb{E}_h \sup \frac{1}{m} \sum_{i=1}^m \sigma_i l(h, z_i) + \mathbb{E}_{S, \sigma} \mathbb{E}_h \sup \frac{1}{m} \sum_{i=1}^m (-\sigma_i) l(h, z'_i) \\
 &= 2 \mathbb{E}_{S, \sigma} \sup \frac{1}{m} \sum_{i=1}^m \sigma_i l(h, z_i) \\
 &\stackrel{(3)}{=} 2 \mathbb{E}_S \text{Rad}((l \circ H)/s)
 \end{aligned}$$

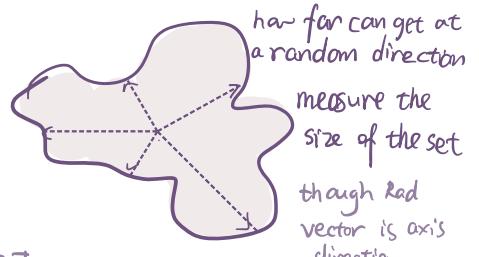
(3)

def  $\ell \circ H = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

$$F|_s = \{(f(z_1), f(z_2), \dots, f(z_m)) : f \in F\} \subseteq \mathbb{R}^m \text{ per-sample risk}$$

$$\text{Rad}(V) = \mathbb{E}_{\vec{\sigma} \sim \text{Rad}^m} \sup_{v \in V} \underbrace{\frac{1}{m} \sum_{i=1}^m \sigma_i v_i}_{\vec{\sigma} \cdot v}$$

the biggest product one can get given a Rademacher vector



(i) suppose one-element hypothesis :

$$\text{Rad}(\{v\}) = \mathbb{E}_{\vec{\sigma}} \sup_{v \in \{v\}} \frac{\vec{\sigma} \cdot v}{m} = \mathbb{E}_{\vec{\sigma}} \frac{\vec{\sigma} \cdot v}{m} = \frac{(\mathbb{E} \vec{\sigma}) \cdot v}{m} = 0$$

(ii) suppose  $\mathcal{H}$  is bigger :

$$\text{Rad}(\{-1, 1\}^m) = \mathbb{E}_{\vec{\sigma}} \sup_{v \in \{-1, 1\}^m} \frac{\vec{\sigma} \cdot v}{m} = 1$$

$$\text{Rad}([-1, 1]^m) = 1 \quad \text{if } \mathcal{H} \text{ bigger, but Rademacher is equal}$$

$$\text{Rad}(\{v : \|v\| \leq \sqrt{m}\}) = 1 \quad \text{because Rad vector only looks at } \{\pm 1\}^m \\ \text{Gaussian complexity would look at any direction}$$

property :

$$\begin{aligned}
 (i) \text{Rad}(cV) &= \text{Rad}(\{cv : v \in V\}) = \frac{1}{m} \mathbb{E}_{\vec{\sigma}} \sup_{v \in V} \sigma \cdot (cv) = \frac{|c|}{m} \mathbb{E}_{\vec{\sigma}} \sup_{v \in V} \text{sign}(c) \vec{\sigma} \cdot v \\
 &= |c| \text{Rad}(V)
 \end{aligned}$$

$$\begin{aligned}
 (ii) \text{Rad}(v+w) &= \text{Rad}(\{v+w : v \in V, w \in W\}) = \frac{1}{m} \mathbb{E}_{v, w} \sup_{\vec{\sigma}} \sigma \cdot (v+w) = \frac{1}{m} \mathbb{E}_{\vec{\sigma}} \sup_v \sigma \cdot v + \sup_w \sigma \cdot w \\
 &= \text{Rad}(v) + \text{Rad}(w)
 \end{aligned}$$

If  $w$  is constant vector  $\Rightarrow \text{Rad}(w)=0$

How to compute  $\text{Rad}(\ell \circ H|_S)$  practically?

- Telgrard's contraction lemma :

$l$ -Lipschitz function

Let  $\phi: \mathbb{R}^m \rightarrow \mathbb{R}^m$  be given by  $\phi(t) = (\varphi_1(t_1), \dots, \varphi_m(t_m))$ , where each  $\varphi_i$  is  $l$ -lip then

$$\text{Rad}(\phi \circ V) = \text{Rad}(\{\phi(v) : v \in V\}) \leq l \text{Rad}(V)$$

$\leftarrow l$ -Lipschitz

for linear regression,  $\ell_Y$  is  $l$ -Lipschitz

$$\begin{aligned} \text{Rad}(\ell_Y \circ H|_S) &= \text{Rad}(\{\ell_{Y_i}(h(x_1)), \ell_{Y_2}(h(x_2)), \dots, \ell_{Y_m}(h(x_m))\}) \\ &\leq \text{Rad}(H|_S) \end{aligned}$$

$$\text{Rad}(H_B) = \mathbb{E} \sup_{\sigma} \frac{1}{m} \sum_{i=1}^m \sigma_i \langle w, x_i \rangle$$

$$= \mathbb{E} \sup_{\sigma} \frac{1}{m} \langle w, \sum_{i=1}^m \sigma_i x_i \rangle$$

$$\leq \frac{1}{m} \mathbb{E} \sup_{\sigma} B \|\sum_{i=1}^m \sigma_i x_i\| \quad \text{to get rid of } \sigma$$

$$\leq \frac{1}{m} B \sqrt{\mathbb{E} \|\sum_{i=1}^m \sigma_i x_i\|^2}$$

$$\begin{aligned} \textcircled{1} \quad \|\sum_i a_i x_i\|^2 &= \langle \sum_i a_i, \sum_i a_i \rangle \\ &= \sum_{ij} \langle a_i, a_j \rangle \\ &= \sum_{ij} \langle x_i, x_j \rangle \end{aligned}$$

$$\begin{aligned} &\stackrel{\textcircled{2}}{=} \frac{B}{m} \sqrt{\mathbb{E} \sum_{ij} \langle \sigma_i x_i, \sigma_j x_j \rangle} = \frac{B}{m} \sqrt{\mathbb{E} \sum_{ij} \sum_{ik} \sigma_i \sigma_k x_i x_k} \\ &= \frac{B}{m} \sqrt{\sum_i B \sigma_i^2 \|x_i\|^2 + \sum_{i \neq j} B \sigma_i \sigma_j \langle x_i, x_j \rangle} \end{aligned}$$

$\sum_{i \neq j} \sigma_i \sigma_j \langle x_i, x_j \rangle$  independent  $= \mathbb{E} \sigma_i \mathbb{E} \sigma_j = 0$

$$= \frac{B}{m} \sqrt{\sum_i \|x_i\|^2}$$

$$\mathbb{E}_S \text{Rad}(H|_S) \leq \frac{B}{\sqrt{m}} \mathbb{E}_S \sqrt{\frac{1}{m} \sum_i \|x_i\|^2} \leq \frac{B}{\sqrt{m}} \sqrt{\mathbb{E} \|x\|^2}$$

## Lec 6 Radermacher Complexity (2)

$$\text{Rad}(V) = \mathbb{E}_{\sigma \sim \text{unif}} \frac{\sigma \cdot V}{m} \text{ for } V \in \mathbb{R}^m$$

Lemma: if  $\varphi: \mathbb{R} \rightarrow \mathbb{R}$  is  $l$ -Lipschitz,  $\text{Rad}(\{\varphi(V_1), V_2, \dots, V_m\}) \leq \text{Rad}(V)$

① use Lemma 2 to prove Telegrond's contraction Lemma

Pf:

$$\begin{aligned} & \text{Rad}(\{\underbrace{\frac{1}{l}\varphi_1(V_1)}, V_2, \dots, V_m : V \in V\}) \leq \text{Rad}(V) \\ & \text{Rad}(\{\underbrace{V_2, \frac{1}{l}\varphi_1(V_1)}, \dots, V_m : V \in V\}) \text{ Rotate doesn't change Rad} \\ & \text{Rad}(\{V_2, V_3, \dots, \frac{1}{l}\varphi_1(V_1) : V \in V\}) \\ & \text{the same for } \varphi_1, \varphi_2, \dots \\ & \text{Rad}(\{\underbrace{\frac{1}{l}\varphi_1(V_1)}, \frac{1}{l}\varphi_2(V_2), \dots, V_m : V \in V\}) \leq \text{Rad}(V) \\ & \text{Rad}(\{\frac{1}{l}\varphi_2(V_2), \frac{1}{l}\varphi_1(V_1), \dots, V_m : V \in V\}) \end{aligned}$$

$$\therefore \text{Rad}(\frac{1}{l}\varphi \circ V) \leq \text{Rad}(V)$$

$$\therefore \text{Rad}(\varphi \circ V) \leq l \text{ Rad}(V)$$

② prove this Lemma 2

$$m \cdot \text{Rad}(\{\varphi(V_1), V_2, \dots, V_m : V \in V\})$$

$$\begin{aligned} &= \mathbb{E}_{\vec{\sigma}_2 \sim \text{Rad}^m} \sup_{V \in V} \underbrace{\sigma_1 \varphi(V_1) + \sigma_2 \cdot V_2 + \dots + \sigma_m \cdot V_m}_{\vec{\sigma}_2 \cdot V_2} \\ &= \frac{1}{2} \mathbb{E}_{\vec{\sigma}_2 \sim \text{Rad}^m} \sup_{V \in V} [\varphi(V_1) + \vec{\sigma}_2 \cdot V_2] + \frac{1}{2} \mathbb{E}_{\vec{\sigma}_2 \sim \text{Rad}^m} \sup_{V \in V} [-\varphi(V_1) + \vec{\sigma}_2 \cdot V_2] \xleftarrow{\text{two independent parts}} \\ &= \frac{1}{2} \mathbb{E}_{\vec{\sigma}_2 \sim V, V' \in V} \sup_{V, V' \in V} [|\varphi(V_1) - \varphi(V_1')| + \vec{\sigma}_2 \cdot (V_2 + V_2')] \\ &\leq \frac{1}{2} \mathbb{E}_{\vec{\sigma}_2 \sim V, V' \in V} \sup_{V, V' \in V} [|V_1 - V_1'| + \vec{\sigma}_2 \cdot (V_2 + V_2')] \quad \text{1-Lipschitz} \\ &= \frac{1}{2} \mathbb{E}_{\vec{\sigma}_2 \sim V, V' \in V} [|V_1 - V_1'| + \vec{\sigma}_2 \cdot (V_2 + V_2')] \end{aligned}$$

$$= \frac{1}{2} \mathbb{E}_{\vec{\sigma}_2 \sim V} [V_1 + \vec{\sigma}_2 \cdot V_2] + \sup_{V \in V} [V_1 + \vec{\sigma}_2 \cdot V_2] = \frac{1}{2} \mathbb{E}_{\vec{\sigma}_2 \sim V} \vec{\sigma}_2 \cdot V = m \text{ Rad}(V)$$

Then for typical supervised learning losses

$$\begin{aligned} (\ell \circ H)|_S &= \{(l(h, z_1), \dots, l(h, z_m)) : h \in H\} \\ &= \{(l_{y_1}(h(x_1)), \dots, l_{y_m}(h(x_m))) : h \in H\} \\ &= (l_{S_y} \circ H)|_{S_x} \end{aligned}$$

might depend on  $S_y$

if  $l_{y_i}$  (the loss func of a prediction for  $y_i$ ) are all  $\tilde{p}$ -Lipschitz.  
then Talagrand's lemma gives:

$$S_x = (x_1, \dots, x_m)$$

$$\text{Rad}((\ell \circ H)|_S) \leq \tilde{p} \text{Rad}(H|_{S_x})$$

only holds for real-valued hypothesis ( $H$  don't need to be function)  
(can be param of Gaussian..)

consider logistic regression  $H_B = \{x \mapsto w \cdot x : \|w\| \leq B\}$ .

$$\text{Rad}((\ell_{\log} \circ H_B)|_S) \leq \text{Rad}(H_B|_{S_x})$$

$$\text{Then } m \text{Rad}(H_B|_{S_x}) = \mathbb{E}_{\substack{\sigma \\ \|\omega\| \leq B}} \sum_i \sigma_i \cdot \langle \omega, x_i \rangle$$

$$= \mathbb{E}_{\substack{\sigma \\ \|\omega\| \leq B}} \langle \omega, \sum_i \sigma_i x_i \rangle \quad \text{Cauchy-Schwarz}$$

$$\leq \mathbb{E}_{\substack{\sigma \\ \|\omega\| \leq B}} \|\omega\| \cdot \|\sum_i \sigma_i x_i\|$$

$$= B \mathbb{E}_{\substack{\sigma \\ \|\omega\| \leq B}} \|\sum_i \sigma_i x_i\| \quad (\mathbb{E}T)^2 \leq \mathbb{E}T^2$$

$$\leq B \sqrt{\mathbb{E}_{\substack{\sigma \\ \|\omega\| \leq B}} \|\sum_i \sigma_i x_i\|^2}$$

$$= B \sqrt{\mathbb{E}_{\substack{\sigma \\ \|\omega\| \leq B}} \sum_{i,j} \sigma_i \sigma_j \langle x_i, x_j \rangle}$$

$$= B \sqrt{\sum_i \mathbb{E}_{\substack{\sigma_i \\ \|\omega\| \leq B}} \sigma_i^2 \|x_i\|^2 + \sum_{i \neq j} \mathbb{E}_{\substack{\sigma_i, \sigma_j \\ \|\omega\| \leq B}} [\sigma_i \sigma_j] \langle x_i, x_j \rangle}$$

$$= B \sqrt{\sum_i \|x_i\|^2}$$

$$\text{Rad}(H_B|_{S_x}) \leq \frac{B}{\sqrt{m}} \sqrt{\sum_i \|x_i\|^2}, \quad \mathbb{E}_{S_x} \text{Rad}(H|_{S_x}) \leq \frac{B}{\sqrt{m}} \sqrt{\mathbb{E} \|x\|^2}$$

$$\text{If } \Pr_{\text{hard bound}}(\|x\| \leq C) = 1, \quad \text{Rad}(H_B|_{S_x}) \leq \frac{BC}{\sqrt{m}}$$

It is fine in some case to only look the average

use Randermacher to get a high probability bound

$\Rightarrow$  McDiarmid's inequality (concentration inequality)

McDiarmid's inequality: Let  $X_1, \dots, X_m$  be independent, let  $f(X_1, \dots, X_m)$  be a real-valued function satisfying bounded differences  
 $\forall i \in [m] \sup_{\substack{x_1, \dots, x_i, x'_i \\ x_{i+1}, \dots, x_m}} |f(x_1, \dots, x_m) - f(x_1, \dots, \tilde{x}_{i-1}, x'_i, \tilde{x}_{i+1}, \dots, x_m)| \leq c_i$

$\checkmark$  change any variable at any one index doesn't change result too much

Then with probability at least  $1 - \delta$ ,

$$f(X_1, \dots, X_m) \leq \mathbb{E}f(X_1, \dots, X_m) + \sqrt{\frac{1}{2} (\sum_{i=1}^m c_i^2) \log \frac{1}{\delta}}$$

(special case:  $f(x) = \frac{1}{m} \sum_i x_i$ , each  $x_i \in [a, b]$ )

$$\text{then } \forall i \in [m] \sup_{\substack{x_1, \dots, x_m, x'_i \\ x_k=x_m \forall k \neq i}} |(x_1 + \dots + x_m) - (x_1 + \dots + \tilde{x}_{i-1} + x'_i + x_{i+1} + \dots + x_m)| = \frac{1}{m} \sup_{\substack{x_k=x'_k \\ k \neq i}} |x_k - x'_k| = \frac{1}{m} (b-a)$$

$$\text{so } c_i = \frac{b-a}{m}, \frac{1}{m} \sum_i x_i - \mathbb{E} x_i \leq \sqrt{\frac{m}{2} (\frac{b-a}{m})^2 \log \frac{1}{\delta}} = (b-a) \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

(McDiarmid's inequality implies Hoeffding)

in case of loss  $l(h, z) \in [a, b]$  for all  $h, z$ , then with prob  $\geq 1 - \delta$

$$\sup_{h \in \mathcal{H}} L_0(h) - L_S(h) \leq \mathbb{E} \sup_{h \in \mathcal{H}} [L_0(h) - L_S(h)] + (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}} \quad (1)$$

if  $\hat{h}_S$  is an ERM, with prob  $\geq 1 - \delta$  that

$$L_0(\hat{h}_S) - L_0(h^*) \leq \mathbb{E} \sup_{h \in \mathcal{H}} [L_0(h) - L_S(h)] + (b-a) \sqrt{\frac{2}{m} \log \frac{2}{\delta}} \quad (2)$$

Pf. prove (1): Let  $S^{(i)} = (z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_m)$ . we have

$$L_0(h) - L_S(h) = L_0(h) - L_{S^{(i)}}(h) + L_{S^{(i)}}(h) - L_S(h)$$

$$\sup_h |L_0(h) - L_S(h)| \leq \sup_h |L_0(h) - L_{S^{(i)}}(h)| + \sup_h |L_{S^{(i)}}(h) - L_S(h)|$$

$$|\sup_h |L_0(h) - L_S(h)| - \sup_h |L_0(h) - L_{S^{(i)}}(h)|| \leq |\sup_h |L_{S^{(i)}}(h) - L_S(h)|| \quad S, S^{(i)} \text{ changeable}$$

$$|\sup_h |L_0(h) - L_S(h)| - \sup_h |L_0(h) - L_{S^{(i)}}(h)|| \leq \sup_h |L_{S^{(i)}}(h) - L_S(h)| = \frac{b-a}{m}$$

apply McDiarmid's inequality  $\Rightarrow$

with prob at least of  $1 - \delta$ :

$$\sup_h |L_0(h) - L_S(h)| \leq \mathbb{E} \sup_h |L_0(h) - L_S(h)| + (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}$$

$\neq \sup_h \mathbb{E} |L_0(h) - L_S(h)|$   $\mathbb{E}$  and max does not commute!

to be continued

Pf for McDiarmid

**very general commonly used.**

$X_{1:k-1} = (X_1, \dots, X_{k-1})$ . Fix some  $k \in [m]$  freeze some arbitrary values  $X_{1:k-1}$

$\mathbb{E} f(X_{1:k-1}, X_k, X_{k+1:m})$  is random depend on  $X_k$

note that  $\sup_{X_k} f(X_{1:m}) - \inf_{X_k} f(X_{1:m}) \leq C_k$  by assumption

$$\Rightarrow \mathbb{E}_{X_{k+1:m}} \sup_{X_k} f(X_{1:k-1}, X_k, X_{k+1:m}) + \sup_{X_k} f(X_{1:k-1}, X_k, X_{k+1:m}) \leq C_k \quad \sup_B \leq \mathbb{E} \sup$$

$$\sup_{X_k} \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m}) + \sup_{X_k} \mathbb{E}_{X_{k+1:m}} (-f(X_{1:k-1}, X_k, X_{k+1:m})) \leq C_k$$

$$\sup_{X_k} \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m}) - \inf_{X_k} \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m}) \leq C_k$$

$X \sim SG(\sigma)$

by Hoeffding's Lemma  $\mathbb{E}_{X_{k+1:m}}$  is  $SG(C_k/2)$   $\mathbb{E}[e^{\lambda(X-\mu)}] \leq e^{\frac{1}{2}\lambda^2\sigma^2}$

$$\text{then } \mathbb{E}_{X_k} \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m})) \leq \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:k-1}, X_k, X_{k+1:m}) + \frac{1}{2}\lambda^2(\frac{C_k}{2})^2)$$

holds for any  $X_{1:k-1}$ , then average no  $X_{1:k-1}$

$$\mathbb{E}_{X_{1:k}} \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:m})) \leq \mathbb{E}_{X_{1:k-1}} \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:m}) + \frac{1}{2}\lambda^2 C_k^2)$$

take  $\log$

$$\underbrace{\log \mathbb{E}_{X_{1:k}} \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:m}))}_{\text{sum over } k=1,\dots,m} \leq \underbrace{\log \mathbb{E}_{X_{1:k-1}} \exp(\lambda \mathbb{E}_{X_{k+1:m}} f(X_{1:m}))}_{a_k} + \frac{1}{2}\lambda^2 C_k^2$$

$$\log \mathbb{E}_{X_{1:m}} \exp(\lambda f(X_{1:m})) \leq \log \exp(\lambda \mathbb{E}_{X_{1:m}} f(X_{1:m})) + \sum_{k=1}^m \frac{1}{2} \lambda^2 C_k^2$$

$$\mathbb{E}_{X_{1:m}} \exp(\lambda f(X_{1:m})) \leq \exp(\lambda \mathbb{E}_{X_{1:m}} f(X_{1:m})) \cdot \exp(\sum_{k=1}^m \frac{1}{2} \lambda^2 C_k^2)$$

$$\underbrace{\mathbb{E}_{X_{1:m}} \exp(\lambda f(X_{1:m}) - \lambda \mathbb{E}_{X_{1:m}} f(X_{1:m}))}_{f(X_{1:m}) \in SG(\frac{1}{2} \sqrt{\sum_{k=1}^m C_k^2})} \leq \exp(\frac{1}{2} \lambda^2 \sum_{k=1}^m C_k^2)$$

with Chernoff bound for sub-Gaussians

with prob at least  $1-\delta$   $f(X_{1:m}) \leq \mathbb{E} f(X_{1:m}) + \frac{1}{2} \sqrt{\sum_{i=1}^m C_i^2} \cdot \sqrt{2 \log \frac{1}{\delta}}$

Rad is scale-sensitive  $\Rightarrow$  depend on the choice of function

## Lec 6 Growth Function and VC Dimension

How to bound Rad for binary classifier?

For binary  $H|_{S_x} = \{h(x_1), h(x_2), \dots, h(x_n) : h \in H\} \subseteq \{0, 1\}^m$

At most  $2^m$  possible vectors even if  $H$  is finite

- $\text{L}_{0-1}(h, (x, y)) = \mathbf{1}(h(x) \neq y)$

$H|_{S_x}$  finite

- If  $|V| < \infty$ ,  $\|v\| \leq B$  for all  $v \in V$

Then  $\text{Rad}(V) \leq \frac{B}{m} \sqrt{2 \log |V|}$

special case: if  $V = H|_{S_x}$ ,  $\|V\| = \sqrt{|H| + \dots + |H|} = \sqrt{m}$

$$\text{Rad}(H|_{S_x}) \leq \sqrt{\frac{2}{m} \log |H|_{S_x}|}$$

Pf  $\text{Rad}(V) = \mathbb{E}_{\sigma} \sup_{v \in V} \sum_{i=1}^m \frac{\sigma_i v_i}{m}$

$\sigma_i \in SG(1)$

$$\frac{\sum_i v_i}{m} \in SG\left(\frac{|V|}{m}\right) \Rightarrow \sum_i \frac{\sigma_i v_i}{m} \in SG\left(\sqrt{\sum_i \frac{v_i^2}{m^2}}\right) \text{ independent}$$

$$= SG\left(\frac{1}{m} \|V\|\right) \leq SG\left(\frac{B}{m}\right)$$

- $T_i \in SG(\sigma)$ ,  $\mathbb{E} T_i = 0$ ,  $T_i$  can be dependent

then  $\mathbb{E}_{i \sim m} T_i \leq \sigma \sqrt{2 \log m}$

Pf. see in A2 Q2.4 Jensen's Inequality:  $\exp(\mathbb{E} Y) \leq \mathbb{E} \exp(Y)$

when taking  $\sup \vec{\sigma} v$  can be dependent, so  $\mathbb{E} \sup_{v \in V} \sum_{i=1}^m \frac{\sigma_i v_i}{m} \leq \frac{B}{m} \sqrt{2 \log |V|}$

$\text{Rad}(H|_{S_x})$  depends on particular distribution

we can use  $|H|_{S_x} < |H|$  but it's very loose!

use growth function to drop dependence on particular  $S_x$

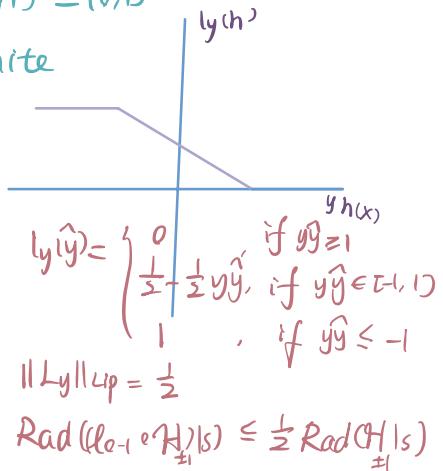
- The growth function of  $H$  is

$$I_H(m) = \sup_{x_1, \dots, x_m \in X} |H|_{(x_1, \dots, x_m)} , \quad |H|_{S_x} \leq I_H(m)$$

- (shatter) = if  $H$  shatters a set  $S$ , means  $H$  can assign any possible label to the set i.e.  $|H|_{S_x} = 2^m$

if the theory could explain any outcome then it's too general

VC is defined on the size of set you can shatter.



- (VC dimension)  $\text{VCdim}(\mathcal{H}) = \max \{ m \geq 0 : P_{\mathcal{H}}(m) = 2^m \}$  worst case  
there is A set  
that can be shattered
- example: threshold  $h_a(x) = \mathbb{1}(x \geq a)$   $\text{VCdim} = 1$



$$\Rightarrow \mathcal{H} = \{x \mapsto \text{sgn}(w \cdot x) : w \in \mathbb{R}^d\} \quad \text{want to show } \text{VCdim}(\mathcal{H}) = d$$

$$\text{sgn}(t) = \begin{cases} 1 & \text{if } t \geq 0 \\ -1 & \text{if } t < 0 \end{cases}$$

①  $\mathcal{H}$  can shatter  $\{e_1, e_2, \dots, e_d\}$  to get  $(y_1, \dots, y_d)$

$$\text{use } w = (y_1, \dots, y_d) \quad w \cdot e_j = y_j$$

(assume) ② Let  $x_1, \dots, x_{d+1}$  be shatterable by  $\mathcal{H}$  (can't be linearly dependent)

$$\therefore \exists \alpha \in \mathbb{R}^{d+1} \setminus \{0\}, \text{ s.t. } \sum_i \alpha_i x_i = 0$$

$$\text{Let } I_+ = \{i \in [d+1] : \alpha_i > 0\}, I_0 = \{i \in [d+1] : \alpha_i = 0\}, I_- = \{i \in [d+1] : \alpha_i < 0\}$$

$$\exists w \text{ s.t. } h_w(x_j) = \begin{cases} 1 & \text{if } j \in I_+ \\ -1 & \text{if } j \in I_- \end{cases}$$

$$0 = w \cdot 0 = w \sum_{i=1}^{d+1} \alpha_i x_i = w \sum_{i \in I_+}^{>0} \alpha_i x_i + w \sum_{i \in I_-}^{<0} \alpha_i x_i + w \sum_{i \in I_0}^{=0} \alpha_i x_i$$

$$= \sum_{i \in I_+} \alpha_i w x_i + \sum_{i \in I_-} \alpha_i w x_i \geq 0 \quad \text{取等 iff } I_- = \emptyset$$

$w x_i = 0$  for  $i \in I_0$

we can also find some  $\tilde{w}$  that  $\tilde{w} \cdot x_i < 0$  for  $\forall i \in I_+$

$$0 = \tilde{w} \cdot 0 = \tilde{w} \cdot \sum_{i \in I_+} \alpha_i x_i = \sum_{i \in I_+} \alpha_i \tilde{w} x_i < 0 \Rightarrow \text{contradiction!}$$

$$\text{VCdim}(\{x \mapsto w \cdot x + b : w \in \mathbb{R}^d, b \in \mathbb{R}\}) = d+1$$

$$\cdot P_{\mathcal{H}}(m) = O(m^{\text{VCdim}(\mathcal{H})})$$

- Binary classifier with  $\text{VCdim}(\mathcal{H}) = d$ , zero-one loss for any  $m > d$ :

$$\mathbb{E}_{h \in \mathcal{H}} [\text{Lo}(h) - \text{L}(h)] \leq \mathbb{E}_{S_x} \sqrt{\frac{2}{m} \log |\mathcal{H}|_{S_x}} \leq \sqrt{\frac{2}{m} \log P_{\mathcal{H}}(m)} \leq \sqrt{\frac{2d}{m} [\log(m+1) + \log d]}$$

Corollary, If  $m \geq d := \text{VCdim}(\mathcal{H})$ , then  $P_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^d$   $\log P_{\mathcal{H}}(m) \leq d \log \frac{m}{d} + 1$

Sauer-Shelah Lemma (1) If  $d = \text{VCdim}(\mathcal{H}) < \infty$ , then  $\Pr_{\mathcal{H}}(m) \leq \sum_{i=0}^d \binom{m}{i} \leq \left(\frac{em}{d}\right)^d$

(we do not prove the lemma)

use S-S Lemma to prove corollary. want to show  $\sum_{i=0}^d \binom{m}{i} \leq \left(\frac{em}{d}\right)^d$  for  $m \geq d$

$$\begin{aligned} \sum_{i=0}^d \binom{m}{i} &\leq \sum_{i=0}^d \binom{m}{i} \left(\frac{m}{d}\right)^{d-i} \\ \text{Pf for corollary} &\leq \sum_{i=0}^m \binom{m}{i} \left(\frac{m}{d}\right)^{d-i} \quad \text{add non-negative} \\ &= \left(\frac{m}{d}\right)^d \sum_{i=0}^m \binom{m}{i} \left(\frac{d}{m}\right)^i \\ &= \left(\frac{m}{d}\right)^d (1 + \frac{d}{m})^m \quad (1 + \frac{x}{n})^n \leq e^x \\ &\leq \left(\frac{m}{d}\right)^d e^d \end{aligned}$$

Lemma: For all finite  $S_x \subseteq X$ ,  $|\mathcal{H}|_{S_x}| \leq |\{T \subseteq S_x : T \text{ is shattered by } \mathcal{H}\}|$

Pf for (1) use (2)  $|\mathcal{H}|_{S_x}|$  is upper bounded by the number of subsets of  $S_x$  of size at most  $d$ .

Pf for (2) inductive

These are equivalent for binary class, 0-1 loss:

- $\Pr_{h \in \mathcal{H}} (\sup_{x \in S_x} L_0(h)) - L_0(h) \leq \varepsilon \geq 1 - \delta \quad \text{for } m \geq M(\varepsilon, \delta)$
- ERM agnostically PAC-learns  $\mathcal{H}$
- $\mathcal{H}$  is ag PAC-learnable
- ERM PAC-learns  $\mathcal{H}$
- $\mathcal{H}$  is PAC-learnable
- $\text{VCdim}(\mathcal{H}) < \infty \quad \text{implies uniform convergence.}$

$$|x|=2m$$

Theorem.  $\mathcal{H}$  is a binary classifier on  $X$ ,  $d = \text{VCdim}(\mathcal{H})$

Assume  $m \leq \frac{d}{2}$ . then  $\inf_A \sup_{\substack{D \\ \text{realizable} \\ \text{by } \mathcal{H}}} \Pr_{\substack{S_x \sim D \\ \text{realizable} \\ \text{by } \mathcal{H}}} (L_0(A(S_x)) \geq \frac{1}{8}) \geq \frac{1}{7}$

for a best-case learning algorithm there is a worst-case distribution

$$\mathbb{E}_{f \sim \text{Unif}(\hat{f}+y)} \mathbb{E}_{S_x \sim D^m} L_0(f, \hat{f}(S_x)) = \mathbb{E}_{f \sim S_x} \mathbb{E}_{x \sim D_x} \mathbb{1}(\hat{f}(S_x) \neq f(x)) \quad \text{not in the training set}$$

$$\begin{aligned} \mathbb{E}_{f \sim S_x} \mathbb{1}(\hat{f}(S_x) \neq f(S_x)) &= \mathbb{E}_{f \sim S_x} [\Pr_{x \notin S_x} \mathbb{E}_{x \sim D_x} [\mathbb{1}(\hat{f}(S_x) \neq f(x)) | x \notin S_x]] \\ &\quad + [\Pr_{x \in S_x} \mathbb{E}_{x \sim D_x} [\mathbb{1}(\hat{f}(S_x) \neq f(x)) | x \in S_x]] \end{aligned}$$

$$P(x \notin S_x) = \frac{|\hat{X} \setminus S_x|}{|\hat{X}|} \geq \frac{m}{2m} = \frac{1}{2} \quad \checkmark \quad S_x = ((x_1, f(x_1)), \dots, (x_m, f(x_m)))$$

$$\geq \frac{1}{2} \mathbb{E}_{f, S_x, x} [\mathbb{I}(\hat{h}_S(x) \neq h(x)) | x \notin S_x]$$

$$= \frac{1}{4}$$

average over  $f \geq \frac{1}{4}$ , so  $\exists g \in f . \mathbb{E}_{S \sim D_g} L_{0-1}(g) \geq \frac{1}{4}$

Want to show  $\Pr_{S \sim D^m} (L_{0-1}(\hat{h}_S) < \frac{1}{8}) \leq \frac{6}{7}$

$$\text{Markov : } \Pr(1 - L_{0-1}(g) \geq \frac{7}{8}) \leq \frac{\mathbb{E}(1 - L_{0-1}(g))}{\frac{7}{8}} \leq \frac{1 - \frac{1}{4}}{\frac{7}{8}} = \frac{6}{7}$$

with no prior, all algorithms perform the same on average.

Let  $H$  be a set of binary classifier over  $X$ . For any  $m \geq 1$

$$\inf_{A \text{ realizable by } H} \sup_{S \sim D^m} \Pr_{\text{0-1 loss}} (L_0(A(S)) > \frac{\text{VC dim}(H)-1}{32m}) \geq \frac{1}{100}$$

Let  $d = \text{VC dim}(H)$

- ①  $d=1$ , holds almost trivially
- ②  $d \geq 2m$ , this is a corollary of the above theorem
- ③  $2 \leq d < 2m$

What do we do for approximation error?

Lec Structural Risk Minimization not uniform over all  $\mathcal{H}$

- $\mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \dots$   weight: how much do I like  $H_1, H_2 \dots$   
 $\sum_{k>0} w_k \leq 1$  preference on each hypothesis  
 $w_k > 0$   $m \rightarrow \infty$

$\forall K \forall D$ .  $\Pr_{S \sim D^m} (\sup_{h \in \mathcal{H}_K} L_D(h) - L_S(h) \leq \varepsilon_K(m, \delta)) \geq 1 - \delta$  with  $\varepsilon_K(m, \delta) \rightarrow 0$

$$\therefore \Pr_{S \sim D^m} (\forall h \in \mathcal{H}, L_D(h) \leq L_S(h) + \varepsilon_{Kh}(m, \delta_{Kh})) \geq 1 - \delta$$

$h_k := \arg \min_{\substack{k \geq 1 \\ h \in \mathcal{H}_k}} \varepsilon_k(m, \delta_{Kh})$  index that minimize  $\varepsilon_k$  corresponding to  $h$   
 $h$  lies in some of the hypo classes  
at least one  $H_k$

① Def:  $\text{SRM}_{\mathcal{H}}(S) \in \arg \min_{h \in \mathcal{H}} [L_S(h) + \varepsilon_{Kh}(m, \delta_{Kh})]$  commit to a  $\delta$  beforehand

higher  $w_k$  for simpler hypothesis, tighter for complex hypo

- SRM algorithm

We can implement this minimization by a finite number of calls to an "ERM oracle", as long as our loss is lower-bounded by  $a \leq \ell(h, z)$  (typically  $a = 0$ ):

```
function SRM $_{\mathcal{H}}(S)$ 
    best  $\leftarrow \infty$ 
    for  $k = 1, 2, \dots$  do
         $h_k \leftarrow \text{ERM}_{\mathcal{H}_k}(S)$ 
        cand_loss  $\leftarrow L_S(h_k) + \varepsilon_k(m, w_k \delta)$ 
        if cand < best then
             $\hat{h} \leftarrow h_k$ 
            best  $\leftarrow$  cand
        if  $\min_{k' > k} a + \varepsilon_{k'}(m, w_{k'} \delta) > \text{best}$  then
            break
    return  $\hat{h}$ 
```

wired because depend on  $\delta$

Note that if we "decompose" as  $\mathcal{H}_1 = \mathcal{H}$ , then SRM becomes just  $\text{ERM}_{\mathcal{H}}$ .

- general bound of SRM non-uniform learnability because
 $\hat{h}_S := \text{SRM}_{\mathcal{H}}(S) \in \arg \min_{h \in \mathcal{H}} L_S(h) + \varepsilon_{Kh}(m, \delta_{Kh})$   $m$  function depend on  $h$

$$L_D(\hat{h}_S) \leq L_S(\hat{h}_S) + \varepsilon_{Kh_S}(m, \delta_{Kh_S}) \quad (\text{def. prob} \geq 1 - \delta)$$

$$\leq L_S(h^*) + \varepsilon_{Kh^*}(m, \delta_{Kh^*}) \quad \text{some } h \text{ can be learned with less data}$$

$$= L_D(h^*) + L_S(h^*) - L_D(h^*) + \varepsilon_{Kh^*}(m, \delta_{Kh^*})$$

prob  $\geq 1 - \delta$   $\downarrow$

$$\leq L_D(h^*) + (b-a)\sqrt{\frac{1}{2m} \log \frac{1}{\delta}} + \varepsilon_{Kh^*}(m, \delta_{Kh^*})$$

$$\therefore \Pr(L_D(\hat{h}_S) - L_D(h^*) \leq \varepsilon_{Kh^*}(m, \delta_{Kh^*}) + (b-a)\sqrt{\frac{1}{2m} \log \frac{2}{\delta}}) \geq 1 - \delta$$

$h^* \in \mathcal{H}$  is any fixed hypothesis

- Specific bound use Randomchar

$$R_k = \underset{s \in D^m}{\mathbb{E}} \text{Rad}((\ell_0 \circ h_k) | s)$$

$$L_D(h) \leq L_S(h) + 2R_{kh} + (b-a)\sqrt{\frac{1}{2m} \log \frac{1}{\delta}}$$

$$\begin{aligned} w_k &= \frac{6}{\pi^2 k^2} \log \frac{1}{w_{kh} \delta} = \log(K_h^2 \cdot \frac{\pi^2}{6} \cdot \frac{1}{\delta}) = 2 \log K_h + \log \frac{\pi^2}{6} + \log \frac{1}{\delta} \\ &< 2 \log K_h + \frac{1}{2} + \log \frac{1}{\delta} \\ &= 2 \log K_h + \log \frac{\sqrt{e}}{\delta} \end{aligned}$$

$$\sqrt{\log \frac{1}{w_{kh} \delta}} < \sqrt{2 \log K_h} + \sqrt{\log \frac{\sqrt{e}}{\delta}}$$

$$\Pr(\forall h \in H, L_D(h) \leq L_S(h) + 2R_{kh} + (b-a)\sqrt{\frac{1}{2m} \log \frac{\sqrt{e}}{\delta}}) \geq 1-\delta$$

- ②  $\hat{h}_S \in \arg \min_{\hat{h}_S} L_S(h) + 2R_{k\hat{h}_S} + (b-a)\sqrt{\frac{1}{2m} \log K_{\hat{h}_S}}$  does not commit to a  $\delta$   
 holds for  $\hat{h}_S$  does not depend on  $\delta$ , better in practice  
 $L_D(\hat{h}_S) \leq L_S(\hat{h}_S) + 2R_{k\hat{h}_S} + (b-a)\sqrt{\frac{1}{2m} \log K_{\hat{h}_S}} + (b-a)\sqrt{\frac{1}{2m} \log \frac{\sqrt{e}}{\delta}}$  with  $p = 1-\delta$

then

$$L_D(\hat{h}_S) \leq L_S(\hat{h}_S) + 2R_{k\hat{h}_S} + (b-a)\sqrt{\frac{1}{2m} \log K_{\hat{h}_S}} + (b-a)\sqrt{\frac{1}{2m} \log \frac{3}{\delta}} \quad \text{in order to make } \frac{\sqrt{e}}{3} \delta$$

(holds with probability  $1 - \frac{\sqrt{e}}{3} \delta$ )

$$\stackrel{\text{union}}{\leq} L_S(h^*) + 2R_{kh^*} + (b-a)\sqrt{\frac{1}{2m} \log K_{h^*}} + (b-a)\sqrt{\frac{1}{2m} \log \frac{3}{\delta}} \quad (1 - \frac{\sqrt{e}}{3} \delta)$$

$$\begin{aligned} \text{Hoeffding} \quad L_S(h^*) &\leq L_D(h^*) + (b-a)\sqrt{\frac{1}{2m} \log \frac{3}{\delta}} \quad (1 - \frac{\delta}{3}) \\ &\leq L_D(h^*) + 2R_{kh^*} + (b-a)\sqrt{\frac{1}{2m} \log K_{h^*}} + (b-a)\sqrt{\frac{2}{m} \log \frac{3}{\delta}} \quad (1 - \frac{\sqrt{e+1}}{3} \delta) \geq 1-\delta \end{aligned}$$

Comparison: if we know the correct  $K_{h^*}$  from the start and run ERM

$$\text{ERM} : L_D(\text{ERM}_{K_{h^*}}) \leq L_D(h^*) + 2R_{kh^*} + (b-a)\sqrt{\frac{2}{m} \log \frac{1}{\delta}}$$

- varying bounded losses

e.g. logistic regression  $\forall h_K = \{x \mapsto w \cdot x = \|w\| \leq B_K\} \quad \Pr(\|x\| \leq C) = 1$   
 $h^*(x) = w^* \cdot x$

if choose  $B_K = 2^K$  then  $2^K < \|w\| \leq 2^K \Rightarrow K_h < \log_2(2\|w\|) \Rightarrow \sqrt{\log K_h} < \sqrt{\log_2(2\|w\|)}$   
 $(b-a)K_h < (2C\|w\| + 1)$ , thus (not commit to a  $\delta$ )

$$\hat{h}_S = h_{\hat{w}_S}; \hat{w}_S \in \arg \min_{w \in \mathbb{R}^d} L_S(h_w) + \frac{4C\|w\|}{\sqrt{m}} + (2C\|w\| + 1)\sqrt{\frac{1}{m} \log(\log_2(2\|w\|))}$$

$$L_D(\hat{h}_S) \leq L_D(h^*) + \frac{4C\|w^*\|}{\sqrt{m}} + (2C\|w^*\| + 1)\sqrt{\frac{1}{m} \log(\log_2(2\|w^*\|))} + (2C\|w^*\| + 2C\|\hat{w}_S\| + 2)\sqrt{\frac{1}{m} \log \frac{3}{\delta}}$$

$\approx$

- relationship to regularization

$$\hat{w} \in \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} L_D(h^*) + \frac{\lambda}{\sqrt{m}} \|w^*\| \quad \text{SRM} \sim \text{regularization}$$

- non-uniform learnability =

we have shown a bound

$$\Pr_{S \sim D^m} (L_D(SRM_{\mathcal{H}, S}) \leq L_D(h^*) + \epsilon_{k \times m, w_{k \times d}} + (b-a) \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}) \geq 1 - \delta$$

the bound does not show PAC learnability because  $\epsilon$  depends on  $h^*$  !!

Def (competes) :  $A(S)(\epsilon, \delta)$  - competes with  $h \in \mathcal{H}$  on  $D$

with  $m$  samples,  $\Pr_{S \sim D^m} (L_D(A(S)) \leq L_D(h) + \epsilon) \geq 1 - \delta$

Def (non-uniform learning) :  $A$  nonuniformly learns  $\mathcal{H}$  if

$\exists$  finite  $m(\epsilon, \delta, h)$ , s.t.  $\forall \epsilon, \delta \in (0, 1)$ ,  $\forall h \in \mathcal{H}$ ,  $\forall D$ ,

$\forall m \geq m(\epsilon, \delta, h)$ ,  $A(S)(\epsilon, \delta)$  - competes with  $h$  on  $D$

this is looser than PAC because  $m$  depend on  $h$

PAC:  $\exists m \geq m(\epsilon, \delta) \quad \Pr_{S \sim D^m} (L_D(A(S)) \leq \inf_h L_D(h) + \epsilon) \geq 1 - \delta \quad \text{for } \forall \epsilon, \delta \in (0, 1)$

If  $\mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \dots$  with  $\text{VCdim}(\mathcal{H}_k) < \infty$

then SRM non-uniformly learns  $\mathcal{H}$  (in 0-1 binary classification)

Pf Let  $\mathcal{H}_K = \{h \in \mathcal{H} : m(\frac{1}{\delta}, \frac{1}{\epsilon}, h) \leq K\}$

$\therefore \mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \dots$

consider  $D$  realizable by  $\mathcal{H}_K$

then  $\Pr_{S \sim D^K} (L_D(A(S)) \leq \frac{\epsilon}{2}) \geq \frac{1}{2}$  using  $h^* \in \mathcal{H}_K$  with  $L_D(h^*) = 0$

polynomial example

$$\mathcal{H} = \{x \mapsto \operatorname{sgn}(w \cdot x) = x \in \mathbb{R}^d\}$$

$$\text{VCdim}(\mathcal{H}) = d$$

$$\sup_{h \in \mathcal{H}} |L_D(h) - L_S(h)| \leq \sqrt{\frac{2d}{m} (\log(m-1) + \log d)} + \sqrt{\frac{1}{2m} \log \frac{1}{\delta}}$$

?

Assume  $\mathcal{H} = \{h_1, h_2, \dots\}$  divide singleton classes

Use  $\mathcal{H}_k = \{h_k\}$

$$\text{Then } \sum_{h \in \mathcal{H}} m_h w_h \leq (b-a) \sqrt{\frac{1}{m} \log \frac{1}{w_h \delta}}$$

$$\leq (b-a) \sqrt{\frac{1}{m} \log \frac{1}{w_h}} + (b-a) \sqrt{\frac{1}{m} \log \frac{1}{\delta}}$$

how to assign weights?

how to decide the order?

Minimum Description Length prefix-free

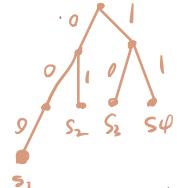
choose a weight according to a binary language for determining  $\mathcal{H}$

Kraft's inequality:

If  $S \subseteq \{0, 1\}^*$  (a set of binary strings) is prefix-free

(if '00' is valid  $s$ , we don't have anything else that starts with '00')  
(or read the data before '0' in (language))

$$\text{then } \sum_{s \in S} 2^{-|s|} \leq 1 \quad (\text{probability distribution})$$



Then, we can choose a representation for  $\mathcal{H}$  and assign  $w_h = 2^{-|h|}$

$$MDL_s \in \operatorname{argmin}_{h \in \mathcal{H}} L_s(h) + (b-a) \sqrt{\frac{1}{m} \log \frac{1}{2^{-|h|}}}$$

$$\text{or } MDL_s \in \operatorname{argmin}_{h \in \mathcal{H}} L_s(h) + \sqrt{\frac{\log 2}{2} \frac{|h|}{m}}$$

$$\therefore L_0(MDL_s) \leq L_0(h^*) + (b-a) \left( \sqrt{\frac{\log 2}{2} \cdot \frac{|h^*|}{m}} + \sqrt{\frac{2}{m} \log \frac{2}{\delta}} \right)$$

Occam's razor: if there are multiple explanations of the data ( $L_s(h_1) = L_s(h_2)$ ) prefer the simpler one (shortest explanation)

If we choose  $|h|$  to be the length of the shortest possible implementation of  $h$  in some programming language, is known Kolmogorov Complexity  
MAP inference with a Kolmogorov prior

# Margins Theory

We do ERM with 0-1 loss

VC theory

$$\mathcal{H} = \{x \mapsto \text{sgn}(w \cdot x) : x \in \mathbb{R}^d\}$$

VCdim(\mathcal{H}) = d

known from previous lec

$$\sup_{h \in \mathcal{H}} L_0(h) - L_S(h) \leq \sqrt{\frac{2d}{m} (\log m + 1 - \log d)} + \sqrt{\frac{2}{m} \log \frac{1}{\delta}} \quad (\text{prob1-}\delta)$$

$$L_0(h_S) - \inf_{h \in \mathcal{H}} L_0(h) \leq \sqrt{\frac{2d}{m} (\log m + 1 - \log d)} + \sqrt{\frac{2}{m} \log \frac{2}{\delta}} \quad (\text{prob2-}\delta)$$

two problems:

① ERM with 0-1 loss is NP hard if D is not realizable by H

② if d is really big, the bound doesn't tell anything until  $\frac{m}{\log m} > 2d$

like in kernel methods, d is sometimes infinite!

we like a better bound when in high dimension (big d) (hope it is not d-dependent)

Rademacher

$$\text{If } \mathcal{H}_B = \{x \mapsto w \cdot x : \|w\| \leq B\} \text{ then } \underset{s}{\mathbb{E}} \text{Rad}(\mathcal{H}_B|s) \leq \frac{B \sqrt{B \|w\|^2}}{\sqrt{m}}$$

but  $\mathbb{P}$  is continuous, logistic loss

$$\text{Rad}(\text{logistic} \circ \mathcal{H}_B|s) \leq 1 \cdot \text{Rad}(\mathcal{H}_B|s)$$

$$\therefore \underset{h \in \mathcal{H}_B}{\mathbb{E}} \left( \underset{h \in \mathcal{H}_B}{\text{argmin}} \underset{h \in \mathcal{H}_B}{L_D^{\text{logistic}}(h)} \right) \leq \frac{2B \sqrt{B \|w\|^2}}{\sqrt{m}} + (B \sqrt{B \|w\|^2} + 1) \sqrt{\frac{2}{m} \log \frac{2}{\delta}}$$

what if we want to bound on accuracy

$$\text{Rad}(\text{0-1} \circ \text{sgn} \circ \mathcal{H}_B)$$

$$\text{let } L_{\text{surr}}(h, z) = \frac{1}{\log 2} L_{\text{logistic}}(h, z)$$

$$\forall h, z \ L_{\text{surr}}(h, z) \geq L_{0-1}(h, z)$$

$$L_D^{\text{surr}}(h) \geq L_D^{0-1}(h)$$

$$L_D^{0-1}(\text{sgn} \circ h) \leq L_D^{\text{surr}}(h) \leq L_S^{\text{surr}}(h) + 2\text{Rad}(L_{\text{surr}} \circ \mathcal{H}|s)$$

$$+ (b-a)_{\text{surr}} \sqrt{\frac{2}{m} \log \frac{2}{\delta}}$$

