#### More Rademacher

CPSC 532S: Modern Statistical Learning Theory 31 January 2022 <u>cs.ubc.ca/~dsuth/532S/22/</u>

### Admin

- Back to normal office hours (Tuesday 10-11, Thursday 4-5)
- A1 grading: hopefully by end of the week
- A2 release: hopefully later this week

### Admin

- Back to normal office hours (Tuesday 10-11, Thursday 4-5)
- A1 grading: hopefully by end of the week
- A2 release: hopefully later this week
- The website for SSBD is currently down (???) If you don't have the pdf saved, I put it on the Canvas Files section

### Admin

- $\frac{1}{4}(\langle w, x \rangle^2 0) + \omega \cos \frac{\omega \cos 2}{20} + \omega \cos 2} + \omega \cos \frac{\omega \cos 2}{20} + \omega \cos 2$ If you don't have the pdf saved, I put it on the Canvas Files section

- Back to normal office hours (Tuesday 10-11, Thursday 4-5) • A1 grading: hopefully by end of the week • A2 release: hopefully later this week The website for SSBD is currently down (???)
- Note on VC dim of homogenous linear classifiers: SSBD are being kinda sloppy
  - Radon's Theorem implies non-homogeneous halfspaces can't shatter size d + 2(Basically the same proof as we talked about)
  - Can do a reduction between homogeneous in  $\mathbb{R}^d$  and non-homog in  $\mathbb{R}^{d-1}$





#### Forever ago: "Fundamental Theorem of Learning"

- These are all equivalent:
- For binary classification with 0-1 loss:
- 1.  $\mathscr{H}$  has the uniform convergence property 2. Any ERM rule agnostically PAC learns  $\mathscr{H}$ Simmediate
- 3. *H* is agnostic PAC learnable
- 4. Any ERM rule PAC learns  $\mathcal{H}$
- 5.  $\mathcal{H}$  is PAC learnable 6.  $VCdim(\mathcal{H}) < \infty$

- If  $VCdim(\mathcal{H}) = d$ :

  - .  ${\mathscr H}$  is agnostic PAC learnable,
  - $\mathcal{H}$  is PAC learnable,

•  $\mathscr{H}$  has uniform convergence property,  $\frac{C_1}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right] \le n_{\mathscr{H}}^{UC} \le \frac{C_2}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right]$  $\frac{C_1}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right] \le n_{\mathcal{H}} \le \frac{C_2}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right]$  $\frac{C_1}{\varepsilon} \left| d + \log \frac{1}{\delta} \right| \le n_{\mathcal{H}} \le \frac{C_2}{\varepsilon} \left| d \log \frac{1}{\varepsilon} + \log \frac{1}{\delta} \right|$ 





#### Last time: Rademacher complexity • The empirical Rademacher complexity of $\mathcal{G}$ on a set $S = (z_1, ..., z_n)$ is

$$\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right]$$
  
$$\sigma \sim \operatorname{Rad}^{n} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right]$$
  
$$\sigma_{i} \sim \operatorname{Rad} \text{ means } \Pr(\sigma_{i} = -1) = \frac{1}{2} = \Pr(\sigma_{i} = 1)$$



#### Last time: Rademacher complexity • The empirical Rademacher complexity of $\mathcal{G}$ on a set $S = (z_1, \dots, z_n)$ is $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right] \qquad \mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ ""bow well confunctions from

 $\sigma_i \sim \text{Rad means } \Pr(\sigma_i = -1) = \frac{1}{2} = \Pr(\sigma_i = 1)$ G has function g: X→ E-1, 13 -() +( -( )-)-( +( +(

"how well can functions from  ${\mathscr G}$ correlate with random noise?"





#### Last time: Rademacher complexity • The empirical Rademacher complexity of $\mathcal{G}$ on a set $S = (z_1, \dots, z_n)$ is

## $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \left| \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right| = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right]$

 $\sigma_i \sim \text{Rad means } \Pr(\sigma_i = -1) = \frac{1}{2} = \Pr(\sigma_i = 1)$ 

$$\mathbf{g}_S = \left(g(z_1), \dots, g(z_n)\right)$$

"how well can functions from  ${\mathscr G}$ correlate with random noise?"

• The ("average-case") Rademacher complexity is just  $\mathfrak{R}_{n}(\mathscr{G}) = \mathbb{E}_{S \sim \mathfrak{N}^{n}}[\widehat{\mathfrak{R}}_{S}(\mathscr{G})]$ 





#### Last time: Rademacher complexity • The empirical Rademacher complexity of $\mathcal{G}$ on a set $S = (z_1, \dots, z_n)$ is

- - $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \left| \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right| = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right]$  $\sigma_i \sim \text{Rad means } \Pr(\sigma_i = -1) = \frac{1}{2} = \Pr(\sigma_i = 1)$
- Take  $\mathfrak{R}_n(\mathscr{G})$  for  $\mathscr{G} = \{z \mapsto \ell(h, z) : h \in \mathscr{G}\}$ 
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ ,

 $\frac{1}{7} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) = \frac{1}{7} \sum_{i=1}^{n} \ell(h_{i})$   $\frac{1}{7} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) = 1$  $= \frac{1}{9} \sum_{i=0}^{n} \frac{1}{2i} \sum_{i=0}^{n} \frac{1}{2$  $\frac{1}{2}L_{S_{+}}[h]$ St= Ezits

$$\mathbf{g}_S = \left(g(z_1), \dots, g(z_n)\right)$$

"how well can functions from  ${\mathscr G}$ correlate with random noise?"

• The ("average-case") Rademacher complexity is just  $\mathfrak{R}_n(\mathscr{G}) = \mathbb{E}_{S \sim \mathcal{D}^n}[\widehat{\mathfrak{R}}_S(\mathscr{G})]$ 

$$= \mathcal{H} \}$$
and  $\widehat{\mathcal{R}}_{S}(\mathcal{G}) = \frac{1}{4} \widehat{\mathcal{R}}_{S|_{x}}(\mathcal{H})$ 

$$= \frac{1}{2} \mathcal{L}_{S|_{x}}(\mathcal{H})$$

$$= \frac{1}{2} \mathcal{L}_{S|_{x}}(\mathcal{H})$$

$$= \frac{1}{2} \mathcal{L}_{S|_{x}}(n) - \mathcal{L}_{S|_{x}}(n)$$







#### Last time: Rader

The empirical Rademacher complexity

#### $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \left| \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right| = \sigma_{i} g(z_{i})$ $\sigma_i \sim \text{Rad means } \Pr(\sigma_i = -1) =$

- Take  $\mathfrak{R}_n(\mathscr{G})$  for  $\mathscr{G} = \{ z \mapsto \ell(h, z) : h \in \mathscr{H} \}$ 
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\Re}_S(\mathscr{G}) = \widehat{\Re}_{SL}(\mathscr{H})$
- Binary clf:  $\Re_n(\mathcal{H}) \le \sqrt{2 \operatorname{VCdim}(\mathcal{H}) \log(n) / n}$   $(n \ge d \ge 3;$  Massart, Sauer-Shelah)

**macher complexit**  
**y** of 
$$\mathscr{G}$$
 on a set  $S = (z_1, ..., z_n)$  is  

$$= \mathbb{E}_{\sigma} \begin{bmatrix} \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_S \\ g \in \mathscr{G} \end{bmatrix} \qquad \mathbf{g}_S = (g(z_1), ..., g(z_n))$$

$$\stackrel{\text{thow well can functions from correlate with random noise}$$

• The ("average-case") Rademacher complexity is just  $\Re_n(\mathscr{G}) = \mathbb{E}_{S \sim \mathcal{D}^n} [\widehat{\Re}_S(\mathscr{G})]$ 







#### Last time: Rademacher complexity • The empirical Rademacher complexity of $\mathcal{G}$ on a set $S = (z_1, \dots, z_n)$ is $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right] \qquad \mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ $\mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ "how well can functions from $\mathscr{G}$

- - $\sigma_i \sim \text{Rad means } \Pr(\sigma_i = -1) = \frac{1}{2} = \Pr(\sigma_i = 1)$
- Take  $\mathfrak{R}_n(\mathscr{G})$  for  $\mathscr{G} = \{z \mapsto \ell(h, z) : h \in \mathscr{G}\}$ 
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ ,
- Binary clf:  $\Re_n(\mathcal{H}) \leq \sqrt{2 \operatorname{VCdim}(\mathcal{H})}$

• Theorem: if  $\mathscr{G}$  maps to [0,1],  $\sup_{g \in \mathscr{G}} |\mathbb{E}[g(g)]|$ 

correlate with random noise?"

• The ("average-case") Rademacher complexity is just  $\mathfrak{R}_n(\mathscr{G}) = \mathbb{E}_{S \sim \mathfrak{N}^n} [\widehat{\mathfrak{R}}_S(\mathscr{G})]$ 

$$= \mathscr{H} \}$$
and  $\widehat{\Re}_{S}(\mathscr{G}) = \widehat{\Re}_{S|_{x}}(\mathscr{H})$ 

$$\overline{\log(n)/n} \quad (n \ge d \ge 3; \text{ Massart, Sauer-Sh}$$

$$(z)] - \frac{1}{n} \sum_{i=1}^{n} g(z_{i}) \bigg| \le 2\Re_{n}(\mathscr{G}) + \sqrt{\frac{1}{2n} \log 2}$$







#### Last time: Rademacher complexity • The empirical Rademacher complexity of $\mathcal{G}$ on a set $S = (z_1, \dots, z_n)$ is $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right] \qquad \mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ $\mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ "how well can functions from $\mathscr{G}$

## $\sigma_i \sim \text{Rad means } \Pr(\sigma_i = -1) = \frac{1}{2} = \Pr(\sigma_i = 1)$

- Take  $\mathfrak{R}_n(\mathscr{G})$  for  $\mathscr{G} = \{z \mapsto \ell(h, z) : h \in \mathscr{G}\}$ 
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ ,
- Binary clf:  $\Re_n(\mathcal{H}) \leq \sqrt{2 \operatorname{VCdim}(\mathcal{H})}$

• Theorem: if  $\mathscr{G}$  maps to [0,1],  $\sup_{g \in \mathscr{G}} | \mathbb{E}[g(g)] \in \mathscr{G}$ 

correlate with random noise?"

• The ("average-case") Rademacher complexity is just  $\mathfrak{R}_n(\mathscr{G}) = \mathbb{E}_{S \sim \mathfrak{N}^n} [\widehat{\mathfrak{R}}_S(\mathscr{G})]$ 

$$\{ \mathcal{H} \}$$

$$\text{and} \quad \widehat{\mathfrak{R}}_{S}(\mathcal{G}) = \underbrace{\widehat{\mathfrak{R}}_{S|_{x}}(\mathcal{H})}_{\text{og}(n)/n} \quad (n \ge d \ge 3; \text{ Massart, Sauer-Sh})$$

$$(z)] - \frac{1}{n} \sum_{i=1}^{n} g(z_{i}) \leq 2\mathfrak{R}_{n}(\mathcal{G}) + \sqrt{\frac{1}{2n} \log 2n}$$







## Last time: Rademacher complexity • The empirical Rademacher complexity of $\mathcal{G}$ on a set $S = (z_1, \dots, z_n)$ is

### $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right] \qquad \mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ $\mathbf{\sigma} \sim \operatorname{Rad}^{n} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right] \qquad \operatorname{constant}^{*} \mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ "how well can functions from $\mathscr{G}$ $\sigma_i \sim \text{Rad means } \Pr(\sigma_i = -1) = \frac{1}{2} = \Pr(\sigma_i = 1)$ correlate with random noise?"

• Take  $\mathfrak{R}_n(\mathscr{G})$  for  $\mathscr{G} = \{ z \mapsto \ell(h, z) : h \in \mathscr{H} \}$ • For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathscr{G}) = \widehat{\mathfrak{R}}_{S|_x}(\mathscr{H})$ • Binary clf:  $\Re_n(\mathcal{H}) \le \sqrt{2 \operatorname{VCdim}(\mathcal{H}) \log(n) / n}$   $(n \ge d \ge 3;$  Massart, Sauer-Shelah) • Theorem: if  $\mathscr{G}$  maps to [0,1],  $\sup_{g \in \mathscr{G}} \left| \mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^{n} g(z_i) \right| \le 2\Re_n(\mathscr{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}} L_{\mathscr{G}}(h_{\mathcal{A}})$ 

• The ("average-case") Rademacher complexity is just  $\mathfrak{R}_n(\mathscr{G}) = \mathbb{E}_{S \sim \mathcal{D}^n}[\widehat{\mathfrak{R}}_S(\mathscr{G})]$ 







#### Last time: Rademacher complexity • The empirical Rademacher complexity of $\mathcal{G}$ on a set $S = (z_1, \dots, z_n)$ is $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right] \qquad \mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ $\mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ "how well can functions from $\mathscr{G}$ correlate with random noise?"

# $\sigma_i \sim \text{Rad means } \Pr(\sigma_i = -1) = \frac{1}{2} = \Pr(\sigma_i = 1)$

- Take  $\mathfrak{R}_n(\mathscr{G})$  for  $\mathscr{G} = \{ z \mapsto \ell(h, z) : h \in \mathscr{H} \}$ 
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathscr{G}) = \widehat{\mathfrak{R}}_{S|_u}(\mathscr{H})$

• The ("average-case") Rademacher complexity is just  $\mathfrak{R}_n(\mathscr{G}) = \mathbb{E}_{S \sim \mathcal{D}^n}[\widehat{\mathfrak{R}}_S(\mathscr{G})]$ 

• Binary clf:  $\Re_n(\mathscr{H}) \le \sqrt{2 \operatorname{VCdim}(\mathscr{H}) \log(n) / n}$   $(n \ge d \ge 3;$  Massart, Sauer-Shelah)

Theorem: if  $\mathscr{G}$  maps to [0,1],  $\sup_{g \in \mathscr{G}} \left| \mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^{n} g(z_i) \right| \leq 2\Re_n(\mathscr{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$  $L_{\mathscr{D}}(h_{4}) \xrightarrow{bounds} \mathbb{E} \sup |L_{\mathscr{D}}(h) - L_{\mathcal{S}}(h)|$ 







#### Last time: Rademacher complexity • The empirical Rademacher complexity of $\mathcal{G}$ on a set $S = (z_1, \dots, z_n)$ is $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right] \qquad \mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ $\mathbf{\sigma} \sim \operatorname{Rad}^{n} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} g(z_{i}) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \frac{1}{n} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right] \qquad \operatorname{constant}^{*} \mathbf{g}_{S} = \left( g(z_{1}), \dots, g(z_{n}) \right)$ "how well can functions from $\mathscr{G}$

## $\sigma_i \sim \text{Rad means } \Pr(\sigma_i = -1) = \frac{1}{2} = \Pr(\sigma_i = 1)$

- Take  $\mathfrak{R}_n(\mathscr{G})$  for  $\mathscr{G} = \{ z \mapsto \ell(h, z) : h \in \mathscr{H} \}$ 
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathscr{G}) = \widehat{\mathfrak{R}}_{S_{|_{\mathcal{I}}}}(\mathscr{H})$

**Theorem:** if  $\mathscr{G}$  maps to [0,1],  $\sup_{g \in \mathscr{G}} \left| \mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^{n} g(z_i) \right|_{L_{\infty}(h)}$ 

correlate with random noise?"

• The ("average-case") Rademacher complexity is just  $\mathfrak{R}_n(\mathscr{G}) = \mathbb{E}_{S \sim \mathcal{D}^n}[\widehat{\mathfrak{R}}_S(\mathscr{G})]$ 

• Binary clf:  $\Re_n(\mathcal{H}) \le \sqrt{2 \operatorname{VCdim}(\mathcal{H}) \log(n) / n}$   $(n \ge d \ge 3;$  Massart, Sauer-Shelah)

 $\log \frac{1}{s}$ bounds how much bigger  $\mathbb{E}\sup|L_{\mathcal{D}}(h)-L_{\mathcal{S}}(h)|$ 







### X of binary classifiers vs loss

## $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \Big[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \mathbb{I}(h(x_{i}) \neq y_{i}) \Big]$



### X of binary classifiers vs loss If $h: \mathcal{X} \to \{0,1\}$ , define $\tilde{h}: \mathcal{X} \to \{-1,1\}$ by $\tilde{h}(x) = 2h(x) - 1; \quad \tilde{y}_i = 2y_i - 1$ $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \Big[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \mathbb{I}(h(x_{i}) \neq y_{i}) \Big]$



### X of binary classifiers vs loss If $h: \mathcal{X} \to \{0,1\}$ , define $\tilde{h}: \mathcal{X} \to \{-1,1\}$ by $\tilde{h}(x) = 2h(x) - 1; \quad \tilde{y}_i = 2y_i - 1$ $\widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \Big[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \mathbb{I}(h(x_{i}) \neq y_{i}) \Big] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \frac{1 - \tilde{y}_{i} \tilde{h}(x_{i})}{2} \right]$



$$\begin{array}{l}
\boldsymbol{\mathfrak{R} of binary c} \\
\text{If } h: \mathcal{X} \to \{0,1\}, \text{ define } \tilde{h}: \mathcal{X} \to \{-1\}, \\
\widehat{\boldsymbol{\mathfrak{R}}}_{S}(\mathcal{G}) = \mathbb{E}_{\sigma} \Big[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \mathbb{I}(h(x_{i}) \neq y_{i}) \Big]
\end{array}$$

Adding constants doesn't change  $\Re_S$ :  $\mathbb{E}_{\boldsymbol{\sigma}}\left[\sup_{g\in\mathscr{G}}\boldsymbol{\sigma}^{\mathsf{T}}(\mathbf{g}_{S}+c\mathbf{1})\right] = \mathbb{E}_{\boldsymbol{\sigma}}\left[\left(\sup_{g\in\mathscr{G}}\boldsymbol{\sigma}^{\mathsf{T}}\mathbf{g}_{S}\right)+c\,\boldsymbol{\sigma}^{\mathsf{T}}\mathbf{1}\right] = \mathbb{E}_{\boldsymbol{\sigma}}\left[\sup_{g\in\mathscr{G}}\boldsymbol{\sigma}^{\mathsf{T}}\mathbf{g}_{S}\right]$ 

### lassifiers vs loss 1,1} by $\tilde{h}(x) = 2h(x) - 1;$ $\tilde{y}_i = 2y_i - 1$ $[v_i] = \mathbb{E}_{\sigma} \left| \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right|$



$$\begin{split} & \Re \text{ of binary c} \\ & \text{If } h : \mathscr{X} \to \{0,1\}, \text{ define } \tilde{h} : \mathscr{X} \to \{-1\} \\ & \widehat{\Re}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \Big[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \mathbb{I}(h(x_{i}) \neq y_{i}) \\ & = \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} - \sigma_{i} \tilde{y}_{i} \tilde{h}(x_{i}) \right] \end{split}$$

Adding constants doesn't change 
$$\widehat{\mathfrak{R}}_{S}$$
:  

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \sigma^{\mathsf{T}}(\mathbf{g}_{S} + c\mathbf{1}) \right] = \mathbb{E}_{\sigma} \left[ \left( \sup_{g \in \mathscr{G}} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right) + c \sigma^{\mathsf{T}} \mathbf{1} \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathscr{G}} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right]$$

### lassifiers vs loss 1,1} by $\tilde{h}(x) = 2h(x) - 1;$ $\tilde{y}_i = 2y_i - 1$ $y_{i})] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \frac{1 - \tilde{y}_{i} \tilde{h}(x_{i})}{2} \right]$





$$\begin{split} & \Re \text{ of binary c} \\ & \text{If } h : \mathscr{X} \to \{0,1\}, \text{ define } \tilde{h} : \mathscr{X} \to \{-1\} \\ & \widehat{\Re}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \Big[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \mathbb{I}(h(x_{i}) \neq y_{i}) \\ & = \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} - \sigma_{i} \tilde{y}_{i} \tilde{h}(x_{i}) \right] \end{split}$$

Adding constants doesn't change 
$$\widehat{\mathfrak{R}}_{S}$$
:  
 $\mathbb{E}_{\sigma}\left[\sup_{g\in\mathscr{G}}\sigma^{\mathsf{T}}(\mathbf{g}_{S}+c\mathbf{1})\right] = \mathbb{E}_{\sigma}\left[\left(\sup_{g\in\mathscr{G}}\sigma^{\mathsf{T}}\mathbf{g}_{S}\right)+c\sigma^{\mathsf{T}}\mathbf{1}\right] = \mathbb{E}_{\sigma}\left[\sup_{g\in\mathscr{G}}\sigma^{\mathsf{T}}\mathbf{g}_{S}\right]$ 

### lassifiers vs loss 1,1} by $\tilde{h}(x) = 2h(x) - 1;$ $\tilde{y}_i = 2y_i - 1$ $y_{i})] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \frac{1 - \tilde{y}_{i} \tilde{h}(x_{i})}{2} \right]$

for fixed  $\tilde{y}_i$ :  $-\sigma_i \tilde{y}_i \sim \text{Rad}$ 



$$\begin{split} & \Re \text{ of binary c} \\ & \text{If } h : \mathcal{X} \to \{0,1\}, \text{ define } \tilde{h} : \mathcal{X} \to \{-1\} \\ & \widehat{\Re}_{S}(\mathcal{G}) = \mathbb{E}_{\sigma} \Big[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \mathbb{I}(h(x_{i}) \neq y_{i}) \\ & = \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} - \sigma_{i} \tilde{y}_{i} \tilde{h}(x_{i}) \\ & = \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \tilde{h}(x_{i}) \right] \\ & \text{Adding constants doesn't change } \\ & \widehat{\Re}_{s \in \mathcal{S}} \Big] = \mathbb{E}_{\sigma} \Big[ \left( \sup_{g \in \mathcal{G}} \sigma^{\mathsf{T}} g_{s} \right) + c \sigma^{\mathsf{T}} \Big] = \mathbb{E}_{\sigma} \Big[ \sup_{g \in \mathcal{G}} \sigma^{\mathsf{T}} g_{s} \Big] \end{split}$$

### lassifiers vs loss 1,1} by $\tilde{h}(x) = 2h(x) - 1;$ $\tilde{y}_i = 2y_i - 1$ $y_{i})] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \frac{1 - \tilde{y}_{i}\tilde{h}(x_{i})}{2} \right]$

for fixed  $\tilde{y}_i$ :  $-\sigma_i \tilde{y}_i \sim \text{Rad}$ 



$$\begin{split} & \mathbf{\mathfrak{R} of binary c} \\ & \text{If } h: \mathcal{X} \to \{0,1\}, \text{ define } \tilde{h}: \mathcal{X} \to \{-1\} \\ & \widehat{\mathbf{\mathfrak{R}}}_{S}(\mathcal{G}) = \mathbb{E}_{\sigma} \Big[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \mathbb{I}(h(x_{i}) \neq y_{i}) \\ & = \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} - \sigma_{i} \tilde{y}_{i} \tilde{h}(x_{i}) \\ & = \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \tilde{h}(x_{i}) \right] = \\ & \text{Adding constants doesn't change } \\ & \widehat{\mathbf{\mathfrak{R}}}_{g \in \mathcal{G}}^{\mathsf{T}}(\mathbf{g}_{S} + c\mathbf{1}) \Big] = \mathbb{E}_{\sigma} \Big[ \left( \sup_{g \in \mathcal{G}} \sigma^{\mathsf{T}} \mathbf{g}_{S} \right) + c\sigma^{\mathsf{T}} \mathbf{1} \Big] = \mathbb{E}_{\sigma} \Big[ \sup_{g \in \mathcal{G}} \sigma^{\mathsf{T}} \mathbf{g}_{S} \Big] \end{split}$$

### lassifiers vs loss 1,1} by $\tilde{h}(x) = 2h(x) - 1;$ $\tilde{y}_i = 2y_i - 1$ $v_i) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right]$ for fixed $\tilde{y}_i$ : $-\sigma_i \tilde{y}_i \sim \text{Rad}$





5



$$\begin{split} & \mathbf{\mathfrak{R} of binary classifiers vs loss} \\ & \text{If } h: \mathcal{X} \to \{0,1\}, \text{ define } \tilde{h}: \mathcal{X} \to \{-1,1\} \text{ by } \tilde{h}(x) = 2h(x) - 1; \quad \tilde{y}_i = 2y_i - \widehat{\mathfrak{R}}_{\mathcal{S}}(\mathscr{G}) = \mathbb{E}_{\sigma} \Big[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \Big] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right] \\ & = \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^n - \sigma_i \tilde{y}_i \tilde{h}(x_i) \right] \quad \text{for fixed } \tilde{y}_i: \quad -\sigma_i \tilde{y}_i \sim \text{Rad} \\ & = \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{h}(x_i) \right] \quad = \frac{1}{2} \widehat{\mathfrak{R}}_{\mathcal{S}|_{\mathcal{X}}}(\widetilde{\mathscr{H}}) \\ & \text{Adding constants doesn't change } \widehat{\mathfrak{R}}_{\mathcal{S}}: \qquad \text{Scaling by } c \text{ gives } |c| \widehat{\mathfrak{R}}_{\mathcal{S}}: \\ & \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s + c\mathbf{1}) \right] = \mathbb{E}_{\sigma} \left[ \left( \sup_{\substack{x \in \mathcal{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}\mathbf{g}_s \right) + c \sigma^{\mathsf{T}}\mathbf{1} \right] = \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{\substack{x \in \mathscr{Y} \\ x \in \mathscr{Y}} \sigma^{\mathsf{T}}(\mathbf{g}_s) \right] = |c| \mathbb$$





$$\begin{split} & \mathbf{\mathfrak{R} of binary classifiers vs loss} \\ & \text{If } h : \mathcal{X} \to \{0,1\}, \text{ define } \tilde{h} : \mathcal{X} \to \{-1,1\} \text{ by } \tilde{h}(x) = 2h(x) - 1; \qquad \tilde{y}_i = 2y_i - \\ & \widehat{\mathfrak{R}}_{S}(\mathscr{G}) = \mathbb{E}_{\sigma} \Big[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_i \mathbb{I}(h(x_i) \neq y_i) \Big] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right] \\ & = \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} - \sigma_i \tilde{y}_i \tilde{h}(x_i) \right] \quad \text{for fixed } \tilde{y}_i: \quad -\sigma_i \tilde{y}_i \sim \text{Rad} \\ & = \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathscr{H}} \frac{1}{n} \sum_{i=1}^{n} \sigma_i \tilde{h}(x_i) \right] \quad = \frac{1}{2} \widehat{\mathfrak{R}}_{S|_x}(\widetilde{\mathscr{H}}) = \frac{1}{q} \widehat{\mathfrak{R}}_{S|_x}(\mathscr{H}) \\ & \text{Adding constants doesn't change } \widehat{\mathfrak{R}}_S: \qquad \text{Scaling by } c \text{ gives } |c| \widehat{\mathfrak{R}}_S: \\ & \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}(\mathsf{g}_S + c1) \right] = \mathbb{E}_{\sigma} \left[ \left( \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right) + c\sigma^{\mathsf{T}} \right] = \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}(\mathsf{g}_S) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{s \in \mathscr{F}} \sigma^{\mathsf{T}}\mathsf{g}_S \right] = |c|$$





## To start: a "low-probability" bound We'll show $\mathbb{E} \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right| \le 2\Re_{n}(\mathscr{G})$ (for any $\mathscr{G}$ ; don't need bounded)



## **To start: a "low-probability" bound** We'll show $\mathbb{E} \sup_{x \in \mathscr{C}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right| \le 2\Re_{n}(\mathscr{G})$ (for any $\mathscr{G}$ ; don't need bounded)

- Immediately implies by Markov's inequality that  $\sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right| \leq \frac{2}{\delta} \Re_{n}(\mathscr{G})$  with prob at least  $1 - \delta$ 
  - Pretty similar to SSBD 6.11, but more general + better bound

Markov's inequality (~1860s)



Andrey Markov Pafnuty Chebyshev

If Pr( Proo  $a\mathbb{I}_{[X>}$ 

$$\begin{aligned} &(X \ge 0) = 1, \text{ then } \Pr\left(X > \frac{1}{\delta} \mathbb{E}[X]\right) \le \delta \\ & \text{f: take } a = \frac{1}{\delta} \mathbb{E}X \text{ in:} \\ & \underset{e \ge a_{-6}}{\ge} \le X, \text{ so } \mathbb{E}[a\mathbb{I}_{[X \ge a]}] = a \Pr(X \ge a) \le \mathbb{I}. \end{aligned}$$





## To start: a "low-probability" bound We'll show $\mathbb{E} \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right| \le 2\Re_{n}(\mathscr{G})$ (for any $\mathscr{G}$ ; don't need bounded)

- Immediately implies by Markov's inequality that  $\sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right| \leq \frac{2}{\delta} \Re_{n}(\mathscr{G})$  with prob at least  $1 - \delta$ 
  - Pretty similar to SSBD 6.11, but more general + better bound
- Will get a better bound with McDiarmid's inequality after

Markov's inequality (~1860s)





If Pr( Proo  $a\mathbb{I}_{[X>}$ 

**Andrey Markov Pafnuty Chebyshev** 

$$\begin{aligned} &(X \ge 0) = 1, \text{ then } \Pr\left(X > \frac{1}{\delta} \mathbb{E}[X]\right) \le \delta \\ &\text{f: take } a = \frac{1}{\delta} \mathbb{E}X \text{ in:} \\ & \underset{\delta}{\ge a_{6}} \le X, \text{ so } \mathbb{E}[a\mathbb{I}_{[X \ge a]}] = a \Pr(X \ge a) \le \mathbb{I}. \end{aligned}$$







## Want $\mathbb{E} \sup |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \leq 2\Re_{n}(\mathscr{G})$ $\int_{\mathcal{G}} g \in \mathscr{G}_{S} \circ \mathscr{D}^{n}$ Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathscr{D}^{n}}[\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_{1}, ..., z'_{n})$



• Then  $\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| = \mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[g]] - \hat{\mathbb{E}}_{S}[g]|$ 



## • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, ..., z'_n)$



#### • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, ..., z'_n)$ Then $\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| = \mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[g]] - \hat{\mathbb{E}}_{S}[g]|$ $\leq \mathbb{E}_{S} \sup_{g \in \mathscr{G}} \mathbb{E}_{S'} |\hat{\mathbb{E}}_{S'}[g] - \hat{\mathbb{E}}_{S}[g]|$ $g \in \mathcal{G}$

• Then  $\mathbb{E}_S \sup_{G} |\mathbb{E}[g] - \hat{\mathbb{E}}_S[g]| = \mathbb{E}_S \sup_{G} |\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[g]] - \hat{\mathbb{E}}_S[g]|$ f(x) = |x - y| is convex so  $|\mathbb{E}X - y| \leq \mathbb{E}|X - y|$ 



#### • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, ..., z'_n)$ $\leq \mathbb{E}_{S} \sup \mathbb{E}_{S'} \left[ \hat{\mathbb{E}}_{S'}[g] - \hat{\mathbb{E}}_{S}[g] \right]$ $g \in \mathcal{G}$

• Then  $\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| = \mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[g]] - \hat{\mathbb{E}}_{S}[g]|$ 

f(x) = |x - y| is convex so  $|\mathbb{E}X - y| \leq \mathbb{E}|X - y|$ 

Jensen's inequality (1906)



If f is convex,  $f(\mathbb{E}X) \leq \mathbb{E}f(X)$ .

Johan Ludwig William Valdemar Jensen



#### • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, ..., z'_n)$ $g \in \mathcal{G}$ $\leq \mathbb{E}_{S} \sup_{S'} \mathbb{E}_{S'} \left[ \hat{\mathbb{E}}_{S'}[g] - \hat{\mathbb{E}}_{S}[g] \right]$ $g \in \mathcal{G}$

Then  $\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| = \mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[g]] - \hat{\mathbb{E}}_{S}[g]|$ 

f(x) = |x - y| is convex so  $|\mathbb{E}X - y| \leq \mathbb{E}|X - y|$ 

Jensen's inequality (1906)



If f is convex,  $f(\mathbb{E}X) \leq \mathbb{E}f(X)$ .

Johan Ludwig William Valdemar Jensen



#### • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, ..., z'_n)$ $g \in \mathscr{G}$ $\leq \mathbb{E}_{S} \sup \mathbb{E}_{S'} \left| \hat{\mathbb{E}}_{S'}[g] - \hat{\mathbb{E}}_{S}[g] \right|$ $g \in \mathcal{G}$



Then  $\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| = \mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[g]] - \hat{\mathbb{E}}_{S}[g]|$ 

f(x) = |x - y| is convex so  $|\mathbb{E}X - y| \leq \mathbb{E}|X - y|$ 

Jensen's inequality (1906)



If f is convex,  $f(\mathbb{E}X) \leq \mathbb{E}f(X)$ .

It is easily seen that the subdifferential is linear in  $y^{[citation needed]}$  (that is false and the assertion requires Hahn-Banach theorem to be proved)

Johan Ludwig William Valdemar Jensen



#### • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, ..., z'_n)$ $g \in \mathscr{G}$ $\leq \mathbb{E}_{S} \sup \mathbb{E}_{S'} \left| \hat{\mathbb{E}}_{S'}[g] - \hat{\mathbb{E}}_{S}[g] \right|$ $g \in \mathcal{G}$



7
#### • Want $\mathbb{E} \sup |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \le 2\mathfrak{R}_{n}(\mathscr{G})$ $g \in \mathcal{G}$

Then  $\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| = \mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[g]] - \hat{\mathbb{E}}_{S}[g]|$ 



### • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, ..., z'_n)$ $\leq \mathbb{E}_{S} \sup \mathbb{E}_{S'} \left| \hat{\mathbb{E}}_{S'}[g] - \hat{\mathbb{E}}_{S}[g] \right|$ $g \in \mathcal{G}$

#### • Want $\mathbb{E} \sup |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \le 2\mathfrak{R}_{n}(\mathscr{G})$ $g \in \mathcal{G}$

Always have  $\sup_{y} \mathbb{E}_{X}[f_{y}(X)] \leq \mathbb{E}_{X}\left[\sup_{y} f_{y}(X)\right] : g \in \mathcal{G} \left[\hat{\mathbb{E}}_{S'}[g] - \hat{\mathbb{E}}_{S}(g)\right]$ 

 $f_y(X) \leq \sup f_y(X)$  for each y so  $\mathbb{E}_X f_v(X) \leq \mathbb{E}_X \sup f_v(X)$  for each y so can take sup of LHS and still true



## • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} \left| \hat{\mathbb{E}}_{S'}[g] \right|$ with $S' = (z'_1, \dots, z'_n)$ Then $\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| = \mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[g]] - \hat{\mathbb{E}}_{S}[g]|$ $\leq \mathbb{E}_{S} \sup_{g \in \mathscr{G}} \mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[g] - \hat{\mathbb{E}}_{S}[g]|$ $g \in \mathcal{G}$

$$\mathbb{E}_{S^{'}} \sup |\mathbb{E}_{S^{'}}[g] - \mathbb{E}_{S}(g)$$

$$g \in \mathcal{G}$$

#### • Want $\mathbb{E} \sup |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \le 2\Re_{n}(\mathscr{G})$ $g \in \mathcal{G}$

Always have  $\sup_{y} \mathbb{E}_{X}[f_{y}(X)] \leq \mathbb{E}_{X} \begin{bmatrix} \sup_{y} f_{y}(X) \\ y \end{bmatrix} = \mathbb{E}_{X} \begin{bmatrix} \sup_{y} f_{y}(X) \\ y \end{bmatrix} =$ 

 $f_y(X) \le \sup f_y(X)$  for each y so  $\mathbb{E}_X f_v(X) \leq \mathbb{E}_X \sup f_v(X)$  for each y so can take sup of LHS and still true



# • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, \dots, z'_n)$ Then $\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| = \mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[g]] - \hat{\mathbb{E}}_{S}[g]|$ $\leq \mathbb{E}_{S} \sup_{g \in \mathscr{G}} \mathbb{E}_{S'} |\hat{\mathbb{E}}_{S'}[g] - \hat{\mathbb{E}}_{S}[g]|$ $g \in \mathcal{G}$

$$= \mathbb{E}_{S,S'} \sup_{g \in \mathcal{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \left[ g(z_i') - g(z_i) \right] \right|$$

#### • Want $\mathbb{E} \sup |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \le 2\mathfrak{R}_{n}(\mathscr{G})$ $g \in \mathcal{G}$

# • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, ..., z'_n)$ Then $\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \le \mathbb{E}_{S,S'} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \left[ g(z_{i}') - g(z_{i}) \right] \right|$



#### Want $\mathbb{E} \sup |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \le 2\Re_{n}(\mathscr{G})$ $g \in \mathcal{G}$

Then 
$$\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \leq \mathbb{E}_{S}$$
  
arbitrary  $\sigma_{i} \in \{-1,1\}$   
 $= \mathbb{E}_{S,S'} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right|$   
Trick called symmetrization:  
it doesn't matter if we swap  
 $z_{i}$  and  $z_{i}'$ , since everything is iid  
and only looked at once

Trick



# • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathscr{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, \dots, z'_n)$ $E_{S,S'} \sup_{g \in \mathcal{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \left[ g(z_i') - g(z_i) \right] \right|$

#### Want $\mathbb{E} \sup |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \le 2\Re_{n}(\mathscr{G})$ $g \in \mathcal{G}$

Trick called symmetrization: it doesn't matter if we swap  $z_i$  and  $z'_i$ , since everything is iid and only looked at once



# • Rewrite using a ghost sample: $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$ with $S' = (z'_1, \dots, z'_n)$ $\begin{array}{c|c} \text{Then } \mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \leq \mathbb{E}_{S,S'} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \left[ g(z_{i}') - g(z_{i}) \right] \right| \\ \stackrel{\text{arbitrary } \sigma_{i} \in \{-1,1\}}{= \mathbb{E}_{S,S'} \sup_{g \in \mathscr{G}} \frac{1}{n}} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right| \\ = \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right| \\ \stackrel{\text{called symmetrization:}}{= \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right| \\ \stackrel{\text{called symmetrization:}}{= \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right| \\ \stackrel{\text{called symmetrization:}}{= \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right| \\ \stackrel{\text{called symmetrization:}}{= \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right| \\ \stackrel{\text{called symmetrization:}}{= \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right| \\ \stackrel{\text{called symmetrization:}}{= \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right| \\ \stackrel{\text{called symmetrization:}}{= \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right| \\ \stackrel{\text{called symmetrization:}}{= \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right|$ (Rademacher distribution)



#### Want $\mathbb{E} \sup |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \le 2\mathfrak{R}_{n}(\mathscr{G})$ $g \in \mathcal{G}$

Then 
$$\mathbb{E}_{S} \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \leq \mathbb{E}_{S}$$
  
arbitrary  $\sigma_{i} \in \{-1,1\}$   
 $= \mathbb{E}_{S,S'} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z_{i}') - g(z_{i}) \right] \right|$   
Trick called symmetrization:  
it doesn't matter if we swap  
 $z_{i}$  and  $z_{i}'$ , since everything is iid  
and only looked at once  
 $\leq \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} g(z_{i}') \right| + \mathbb{E}_{S,S'} \mathbb{E}_{\sigma}$ 



• Rewrite using a ghost sample:  $\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^n} |\hat{\mathbb{E}}_{S'}[g]|$  with  $S' = (z'_1, \dots, z'_n)$  $\sum_{\substack{S,S' \\ g \in \mathcal{G}}} \sup \frac{1}{n} \left| \sum_{i=1}^{n} \left[ g(z_i') - g(z_i) \right] \right|$  $= \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathcal{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_i \left[ g(z_i') - g(z_i) \right] \right|$  $\sigma_i$  now iid Uniform({-(Rademacher distribution)  $b = \frac{1}{n} \mathcal{E} \sigma_i g(z_i) - \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right) + \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \left( \frac{1}{n} \mathcal{E} \sigma_i g(z_i) \right)$  $\sup_{g \in \mathcal{G}} \frac{1}{n} \left| \sum_{i=1}^{j} -\sigma_i g(z_i) \right|$ 



#### • Want $\mathbb{E} \sup_{s \to 0} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \le 2\Re_{n}(\mathscr{G})$ $g \in \mathcal{G}$

• Rewrite using a ghost sample: 
$$\mathbb{E}[g] = \mathbb{E}_{S' \sim \mathscr{D}^n} [\hat{\mathbb{E}}_{S'}[g]]$$
 with  $S' = (z'_1, \dots, z'_n)$   
• Then  $\mathbb{E}_S \sup_{g \in \mathscr{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_S[g]| \le \mathbb{E}_{S,S'} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^n \left[ g(z'_i) - g(z_i) \right] \right|$   
 $= \mathbb{E}_{S,S'} \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i \left[ g(z'_i) - g(z_i) \right] \right| = \mathbb{E}_{S,S'} \mathbb{E}_\sigma \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i \left[ g(z'_i) - g(z_i) \right] \right|$   
Trick called symmetrization:  
it doesn't matter if we swap  
 $z_i$  and  $z'_i$ , since everything is iid  
and only looked at once  
 $\le \mathbb{E}_{S,S'} \mathbb{E}_\sigma \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^n \sigma_i g(z'_i) \right| + \mathbb{E}_{S,S'} \mathbb{E}_\sigma \sup_{g \in \mathscr{G}} \frac{1}{n} \left| \sum_{i=1}^n -\sigma_i g(z_i) \right| = 2\mathfrak{R}_n (\mathscr{G})$   
if  $\mathscr{G}$  is symmetric:  
 $g \in \mathscr{G}$  implies  $-g$ 









 $\leq$ 

$$\begin{aligned} & \text{Want } \mathbb{E}\sup_{g \in \mathcal{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \leq 2\Re_{n}(\mathcal{G}) & \text{always } \mathbb{E}\sup_{g \in \mathcal{G}} \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \leq 2\Re_{n}(\mathcal{G}) \\ & \text{Rewrite using a ghost sample: } \mathbb{E}[g] = \mathbb{E}_{S' \sim \mathcal{D}^{n}} [\hat{\mathbb{E}}_{S}[g]] \text{ with } S' = (z'_{1}, \dots, z'_{n}) \\ & \text{Then } \mathbb{E}_{S}\sup_{g \in \mathcal{G}} |\mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g]| \leq \mathbb{E}_{S,S'} \sup_{g \in \mathcal{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \left[ g(z'_{i}) - g(z_{i}) \right] \right| \\ & = \mathbb{E}_{S,S'} \sup_{g \in \mathcal{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z'_{i}) - g(z_{i}) \right] \right| \\ & = \mathbb{E}_{s,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathcal{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} \left[ g(z'_{i}) - g(z_{i}) \right] \right| \\ & \text{Trick called symmetrization:} \\ & \text{it doesn't matter if we swap} \\ z_{i} \text{ and } z'_{i}, \text{ since everything is iid} \\ & \text{and only looked at once} \\ & \text{triangle inequality } \sup_{g \in \mathcal{G}} \frac{1}{n} \left| \sum_{i=1}^{n} \sigma_{i} g(z'_{i}) \right| + \mathbb{E}_{S,S'} \mathbb{E}_{\sigma} \sup_{g \in \mathcal{G}} \frac{1}{n} \left| \sum_{i=1}^{n} -\sigma_{i} g(z_{i}) \right| \\ & = 2\Re_{n}(\mathcal{G}) \\ & \text{if } \mathcal{G} \text{ is symmetric:} \\ & g \in \mathcal{G} \text{ implies } -g \in \mathcal{G} \end{aligned}$$



)]|





## So what does that mean?

• For  $\{0,1\}$  classifiers in symmetric  $\mathcal{H}$  with 0-1 loss,  $n \ge d = \operatorname{VCdim}(\mathcal{H}) \ge 3$ ,

$$\begin{split} \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_{S}(h)| &\leq \frac{2}{\delta} \Re_{n}(\mathcal{H}) \leq \frac{2}{\delta} \sqrt{\frac{2d \log n}{n}} = \mathcal{E} \\ \mathcal{H} & \searrow \quad \left(\frac{\delta \mathcal{E}}{2}\right)^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \geq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{E}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H} & (\frac{\delta \mathcal{H}}{2})^{2} \leq \frac{2d \log n}{n} \\ \mathcal{H}$$

ガモン



# So what does that mean?

• For  $\{0,1\}$  classifiers in symmetric  $\mathscr{H}$  with 0-1 loss,  $n \ge d = \operatorname{VCdim}(\mathscr{H}) \ge 3$ ,

$$\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_{S}(h)| \leq \frac{2}{\delta} \Re_{n}(\mathcal{H}) \leq \frac{2}{\delta} \sqrt{\frac{2d \log n}{n}}$$

- Showed  $d < \infty$  implies  $\mathscr{H}$  has the uniform convergence property
- But optimal rate is  $n_{\mathscr{H}}^{UC} \leq \frac{C_2}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right]$ : much better in  $\delta$

$$n = \mathcal{N}\left(\frac{d}{\varepsilon^2 S^2}\right)$$

• and so finally proved the Fundamental Theorem of Learning! (up to Massart's lemma)





• For  $\{0,1\}$  classifiers,  $\mathcal{H}$  not necessarily symmetric, with 0-1 loss,  $n \ge d = \operatorname{VCdim}(\mathscr{H}) \ge 3$ ,

 $\sup_{h \in \mathcal{H}} L_{\mathcal{D}}(h) - L_{S}(h) \leq \frac{2}{\delta} \Re_{n}(\mathcal{H}) \leq \frac{2}{\delta} \sqrt{\frac{2d \log n}{n}}$ 





• For  $\{0,1\}$  classifiers,  $\mathcal{H}$  not necessarily symmetric, with 0-1 loss,  $n \ge d = \operatorname{VCdim}(\mathscr{H}) \ge 3$ ,

$$\sup_{h \in \mathcal{H}} L_{\mathcal{D}}(h) - L_{S}(h) \leq \frac{2}{\delta} \Re_{n}(\mathcal{H}) \leq \frac{2}{\delta} \sqrt{\frac{2d \log n}{n}}$$

#### • Is what we really care about: for any $h \in \mathcal{H}$ , have $L_{\mathcal{D}}(h) \leq L_{S}(h) + \frac{2}{s} \Re_{n}(\mathcal{H})$

 $V_{S'} = V_{D}$  $E_{V_{S'}} = V_{D}$ 



• For  $\{0,1\}$  classifiers,  $\mathcal{H}$  not necessarily symmetric, with 0-1 loss,  $n \ge d = \operatorname{VCdim}(\mathscr{H}) \ge 3$ ,

$$\sup_{h \in \mathcal{H}} L_{\mathcal{D}}(h) - L_{S}(h) \leq \frac{2}{\delta} \Re_{n}(\mathcal{H}) \leq \frac{2}{\delta} \sqrt{\frac{2d \log n}{n}}$$

- Is what we really care about: for any h
- Rademacher results usually framed this way to avoid the symmetric condition
  - Some other bounds work this way too (e.g. PAC-Bayes)
  - Difference between two-sided and one-sided will be important later!

$$\in \mathscr{H}$$
, have  $L_{\mathscr{D}}(h) \leq L_{S}(h) + \frac{2}{\delta} \Re_{n}(\mathscr{H})$ 





(pause)

 Hoeffding's inequality: "the mean of independent, bounded

"the mean of independent, bounded RVs concentrates around its expectation"

- Hoeffding's inequality: "the mean of independent, bounded RVs concentrates around its expectation"
- McDiarmid's inequality lets us take more general functions of independent RVs:

7

s:

- Hoeffding's inequality: "the mean of independent, bounded RVs concentrates around its expectation"
- McDiarmid's inequality lets us take more general functions of independent RVs:
  Say *f* has bounded differences if for all 1 ≤ *i* ≤ *n*,
  - Say f has bounded differences if for all  $1 \le i \le n$ ,  $\begin{vmatrix} f(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n) \end{vmatrix} \le c_i$   $f(x_1, \dots, x_n) \le \frac{1}{n} \le x_i$   $f(x_1, \dots, x_n) \le \frac{1}{n} \le x_i$   $f(x_1, \dots, x_n) \le \frac{1}{n} \le x_i$   $f(x_1, \dots, x_n) \le \frac{1}{n} \le x_i$

7

s:

- Hoeffding's inequality: "the mean of independent, bounded RVs concentrates around its expectation"
- McDiarmid's inequality lets us take more general functions of independent RVs: • Say *f* has **bounded differences** if for all  $1 \le i \le n$ ,
  - $f(x_1, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$

and same for Pr

**McDiarmid's** inequality (1989)



**Colin McDiarmid** 

$$f(x_1, \dots, x_{i-1}, x'_i, x_{i+1}, \dots, x_n) \le c_i$$

If  $X_1, \ldots, X_n \in \mathbb{R}$  independent, *f* has bounded diffs with  $c_i$ ,

$$\mathbb{E}_{S'}f(S') + \varepsilon \leq \mathbb{P} \exp\left(\frac{-2\varepsilon^2}{\sum_i c_i^2}\right)$$
$$\left(f(S) < \mathbb{E}_{S'}f(S') - \varepsilon\right)$$



More general than Hoeffding  $Pr\left(\left|\frac{1}{n} \xi \times i - \frac{1}{2} \xi \times i\right| \xi = \xi \right) = 2e_{x}p\left(-\frac{2}{2} \xi \times i\right) = 2e_{x}p\left(-\frac{2n \xi^{2}}{2}\right) = 2e_{x}p\left(-\frac{2n \xi^{2}}{2}\right)$  $f(s) = f(X_1, \dots, X_n) = \prod_{n \in X_i} X_i$ 



**Colin McDiarmid** 

# If $X_1, \ldots, X_n \in \mathbb{R}$ independent, *f* has bounded diffs with $c_i$ , McDiarmid's inequality (1989) If $f(S) = \mathbb{E}_{S'}f(S') + \varepsilon \leq \mathbb{E}_{S'}f(S') + \varepsilon \leq \mathbb{E}_{S'}f(S') + \varepsilon$ and same for $\Pr\left(f(S) < \mathbb{E}_{S'}f(S') - \varepsilon\right)$



• Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$ 

- Let S' be S with  $z_i$  replaced by  $z'_i$  (for any one i)



- Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] \hat{\mathbb{E}}_{S}[g] \right|$
- Let S' be S with  $z_i$  replaced by  $z'_i$  (for any one i)  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$



- Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] \hat{\mathbb{E}}_{S}[g] \right|$
- Let S' be S with  $z_i$  replaced by  $z'_i$  (for any one i)  $\Phi(S) = \sup_{\substack{g \in \mathscr{G} \\ g \in \mathscr{G}}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$  $= \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S'}g + \frac{1}{n}g(z'_{i}) - \frac{1}{n}g(z'_{i}) \right|$

$$z_i') - \frac{1}{n}g(z_i)$$



• Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$ 

• Let S' be S with  $z_i$  replaced by  $z'_i$  (for any one i)

$$\Phi(S) = \sup_{\substack{g \in \mathscr{G} \\ g \in \mathscr{G}}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$$
$$= \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S'}g + \frac{1}{n}g(z_{i}') - \frac{1}{n}g(z_{i}) \right|$$
$$\leq \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S'}g \right| + \frac{1}{n}\sup_{g \in \mathscr{G}} |g(z_{i}') - g(z_{i}')|$$

 $g(z_i)$ 



• Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$  and assume  $g(z) \in [0,B]$ 

• Let S' be S with  $z_i$  replaced by  $z'_i$  (for any one i)

$$\Phi(S) = \sup_{\substack{g \in \mathscr{G} \\ g \in \mathscr{G}}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$$
$$= \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S'}g + \frac{1}{n}g(z_{i}') - \frac{1}{n}g(z_{i}) \right|$$
$$\leq \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S'}g \right| + \frac{1}{n}\sup_{g \in \mathscr{G}} |g(z_{i}') - g(z_{i}')|$$

 $g(z_i)$ 



• Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$  and assume  $g(z) \in [0,B]$ 

• Let S' be S with  $z_i$  replaced by  $z'_i$  (for any one i)

 $\Phi(S) = \sup_{\substack{g \in \mathscr{G} \\ g \in \mathscr{G}}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$  $= \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S'}g + \frac{1}{n}g(z'_{i}) \right|$  $\leq \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S'}g \right| + \frac{1}{n} \sup_{g \in \mathscr{G}} \left| g(z'_i) - g(z_i) \right|$ 

 $\leq \Phi(S') + \frac{B}{m}$ 

$$z_i') - \frac{1}{n}g(z_i)$$



• Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$  and assume  $g(z) \in [0,B]$ 

• Let S' be S with  $z_i$  replaced by  $z'_i$  (for any one i)

$$\Phi(S) = \sup_{\substack{g \in \mathscr{G} \\ g \in \mathscr{G}}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$$
$$= \sup_{\substack{g \in \mathscr{G} \\ g \in \mathscr{G}}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S'}g + \frac{1}{n}g(z'_{i}) - \frac{1}{n}g(z_{i}) \right|$$
$$\leq \sup_{\substack{g \in \mathscr{G} \\ g \in \mathscr{G}}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S'}g \right| + \frac{1}{n}\sup_{\substack{g \in \mathscr{G} \\ g \in \mathscr{G}}} |g(z'_{i}) - g(z'_{i})|$$

 $\leq \Phi(S') + \frac{B}{n}$  so  $|\Phi(S) -$ 

 $g(z_i)$ 

$$-\Phi(S')| \leq \frac{B}{n}$$



# A better generalization bound • Consider $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] - \hat{\mathbb{E}}_{S}[g] \right|$ and assume $g(z) \in [0,B]$

- Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] \hat{\mathbb{E}}_{S}[g] \right|$  and assume  $g(z) \in [0,B]$
- $\Phi$  satisfies bounded differences with  $c_i = \frac{B}{-}$

- Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] \hat{\mathbb{E}}_{S}[g] \right|$  and assume  $g(z) \in [0,B]$
- $\Phi$  satisfies bounded differences with  $c_i = \frac{B}{-}$
- Already saw  $\mathbb{E}_{S} \Phi(S) \leq 2\mathfrak{R}_{n}(\mathscr{G})$

- Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] \hat{\mathbb{E}}_{S}[g] \right|$  and assume  $g(z) \in [0,B]$
- $\Phi$  satisfies bounded differences with  $c_i = \frac{D}{m}$
- Already saw  $\mathbb{E}_{S}\Phi(S) \leq 2\mathfrak{R}_{n}(\mathscr{G})$
- So we get

 $\Pr\left(\Phi(S) > 2\Re_n(\mathscr{G}) + \varepsilon\right) \le \Pr\left(\Phi(S) > \mathbb{E}_{S'}\Phi(S') + \varepsilon\right) \le \exp\left(\frac{-2\varepsilon^2}{n \cdot (B/n)^2}\right)$ 



- Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] \hat{\mathbb{E}}_{S}[g] \right|$  and assume  $g(z) \in [0,B]$
- $\Phi$  satisfies bounded differences with  $c_i = \frac{D}{m}$
- Already saw  $\mathbb{E}_{S}\Phi(S) \leq 2\mathfrak{R}_{n}(\mathscr{G})$
- So we get

 $\Pr\left(\Phi(S) > 2\Re_n(\mathscr{G}) + \varepsilon\right) \le \Pr\left(\Phi(S) > \mathbb{E}_{S'}\Phi(S') + \varepsilon\right) \le \exp\left(\frac{-2\varepsilon^2}{n \cdot (B/n)^2}\right)$ 

 $\leq \exp\left(-2n\varepsilon^2/B^2\right)$ 



- Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] \hat{\mathbb{E}}_{S}[g] \right|$  and assume  $g(z) \in [0,B]$
- $\Phi$  satisfies bounded differences with  $c_i = \frac{D}{T}$
- Already saw  $\mathbb{E}_{S}\Phi(S) \leq 2\mathfrak{R}_{n}(\mathscr{G})$
- So we get

 $\Pr\left(\Phi(S) > 2\Re_n(\mathscr{G}) + \varepsilon\right) \le \Pr\left(\Phi(S) > \mathbb{E}_{S'}\Phi(S') + \varepsilon\right) \le \exp\left(\frac{-2\varepsilon^2}{n \cdot (B/n)^2}\right)$ Solving for  $\varepsilon$ , get  $\Phi(S) \le 2\Re_n(\mathscr{G}) + \sqrt{\frac{B^2}{2n}\log\frac{1}{\delta}}$  $\leq \exp\left(-2n\varepsilon^2/B^2\right)$ 



- Consider  $\Phi(S) = \sup_{g \in \mathscr{G}} \left| \mathbb{E}[g] \hat{\mathbb{E}}_{S}[g] \right|$  and assume  $g(z) \in [0,B]$
- $\Phi$  satisfies bounded differences with  $c_i = -$
- Already saw  $\mathbb{E}_{S} \Phi(S) \leq 2\mathfrak{R}_{n}(\mathscr{G})$
- So we get

 $\Pr\left(\Phi(S) > 2\Re_n(\mathscr{G}) + \varepsilon\right) \le \Pr\left(\Phi(S) + \varepsilon\right)$ 

Solving for  $\varepsilon$ , get  $\Phi(S) \leq 2\Re_n(\mathscr{G}) + 1$ 

n Still has pesky  $\log n$  inside the  $\Re_n$ ; can get optimal rate with chaining

$$S) > \mathbb{E}_{S'} \Phi(S') + \varepsilon \le \exp\left(\frac{-2\varepsilon^2}{n \cdot (B/n)^2}\right)$$
$$\sqrt{\frac{B^2}{2n} \log \frac{1}{\delta}} \le \exp\left(-2n\varepsilon^2/B^2\right)$$




- $\Re_S$  also satisfies bounded differences for bounded g
- So  $|\hat{\Re}_{S} \hat{\Re}_{n}|$  concentrates, by McDiarmid
- Can split  $\delta$  failure prob into  $\delta/2$  for this,  $\delta/2$  for previous theorem and bound  $\Phi(S)$  in terms of  $\hat{\Re}_{S}$



If there's time left, let's prove: (If not: we'll come back to it an **Massart's lemma: for**  $\mathscr{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathscr{A}} a \in \mathscr{A}$ 

nother time!)  
$$|a|| \le r$$
,  $\mathbb{E}_{\sigma}\left[\max_{a \in \mathscr{A}} \frac{1}{n} \sigma^{\mathsf{T}} a\right] \le \frac{1}{n} r \sqrt{2\log}$ 



## Summary

- We *finally* proved the fundamental theorem:
  Finite VCdim characterizes PAC learnabili
- Finite VCdim characterizes PAC learnability (realizable or agnostic) + ERM
  Proved a generalization bound via Rademacher complexity
  - Works for any bounded loss function