

# More VC + Rademacher

CPSC 532S: Modern Statistical Learning Theory

26 January 2022

[cs.ubc.ca/~dsuth/532S/22/](https://cs.ubc.ca/~dsuth/532S/22/)

# Admin

- Reminder: no office hours this week (ICML...)
  - But feel free to Piazza / schedule a meeting if needed
- A1 grading: probably late next week

# Admin

- Reminder: no office hours this week (ICML...)
  - But feel free to Piazza / schedule a meeting if needed
- A1 grading: probably late next week
- A2 will be released probably next week, due ~2-3 weeks after release. A weird plan:
  - Groups of up to 3 allowed
  - You can use a different group *per question* if you want
    - (e.g. do one problem alone, one with person A, one with B+C)
  - You'll hand in questions as separate Gradescope assignments
    - (one per group per question, using Gradescope group feature)
  - Drop lowest: still for total assignment grade (or more advantageous, TBD)
  - Trying to encourage **actively participating in each question**
    - Please **don't** just split assignment in thirds
  - Dropping lowest assignment grade is still per student
  - Will try to calibrate difficulty/length a bit, but you'll have groups

# Previously: “Fundamental Theorem of Learning”

For binary classification  
with 0-1 loss:

These are all equivalent:

1.  $\mathcal{H}$  has the uniform convergence property
2. Any ERM rule agnostically PAC learns  $\mathcal{H}$
3.  $\mathcal{H}$  is agnostic PAC learnable
4. Any ERM rule PAC learns  $\mathcal{H}$
5.  $\mathcal{H}$  is PAC learnable
6.  $\text{VCdim}(\mathcal{H}) < \infty$

# Previously: “Fundamental Theorem of Learning”

For binary classification  
with 0-1 loss:

These are all equivalent:

1.  $\mathcal{H}$  has the uniform convergence property
2. Any ERM rule agnostically PAC learns  $\mathcal{H}$
3.  $\mathcal{H}$  is agnostic PAC learnable
4. Any ERM rule PAC learns  $\mathcal{H}$
5.  $\mathcal{H}$  is PAC learnable
6.  $\text{VCdim}(\mathcal{H}) < \infty$

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

- $\mathcal{H}$  has uniform convergence property,  $\frac{C_1}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right] \leq n_{\mathcal{H}}^{UC} \leq \frac{C_2}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right]$
- $\mathcal{H}$  is agnostic PAC learnable,  $\frac{C_1}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right] \leq n_{\mathcal{H}} \leq \frac{C_2}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right]$
- $\mathcal{H}$  is PAC learnable,  $\frac{C_1}{\varepsilon_3} \left[ d + \log \frac{1}{\delta} \right] \leq n_{\mathcal{H}} \leq \frac{C_2}{\varepsilon} \left[ d \log \frac{1}{\varepsilon} + \log \frac{1}{\delta} \right]$

# Previously: “Fundamental Theorem of Learning”

For binary classification  
with 0-1 loss:

These are all equivalent:

1.  $\mathcal{H}$  has the uniform convergence property
2. Any ERM rule agnostically PAC learns  $\mathcal{H}$  *← saw today*
3.  $\mathcal{H}$  is agnostic PAC learnable
4. Any ERM rule PAC learns  $\mathcal{H}$  *← immediate*
5.  $\mathcal{H}$  is PAC learnable
6.  $\text{VCdim}(\mathcal{H}) < \infty$  *← no free lunch*

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

- $\mathcal{H}$  has uniform convergence property,  $\frac{C_1}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right] \leq n_{\mathcal{H}}^{UC} \leq \frac{C_2}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right]$
- $\mathcal{H}$  is agnostic PAC learnable,  $\frac{C_1}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right] \leq n_{\mathcal{H}} \leq \frac{C_2}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right]$
- $\mathcal{H}$  is PAC learnable,  $\frac{C_1}{\varepsilon_3} \left[ d + \log \frac{1}{\delta} \right] \leq n_{\mathcal{H}} \leq \frac{C_2}{\varepsilon} \left[ d \log \frac{1}{\varepsilon} + \log \frac{1}{\delta} \right]$

# Previously: “Fundamental Theorem of Learning”

For binary classification  
with 0-1 loss:

These are all equivalent:

1.  $\mathcal{H}$  has the uniform convergence property
  2. Any ERM rule agnostically PAC learns  $\mathcal{H}$  *← saw today*
  3.  $\mathcal{H}$  is agnostic PAC learnable *← immediate*
  4. Any ERM rule PAC learns  $\mathcal{H}$  *← immediate*
  5.  $\mathcal{H}$  is PAC learnable
  6.  $\text{VCdim}(\mathcal{H}) < \infty$  *← no free lunch*
- to show!* (red arrow from 6 to 1)

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

- $\mathcal{H}$  has uniform convergence property,  $\frac{C_1}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right] \leq n_{\mathcal{H}}^{UC} \leq \frac{C_2}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right]$
- $\mathcal{H}$  is agnostic PAC learnable,  $\frac{C_1}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right] \leq n_{\mathcal{H}} \leq \frac{C_2}{\varepsilon^2} \left[ d + \log \frac{1}{\delta} \right]$
- $\mathcal{H}$  is PAC learnable,  $\frac{C_1}{\varepsilon_3} \left[ d + \log \frac{1}{\delta} \right] \leq n_{\mathcal{H}} \leq \frac{C_2}{\varepsilon} \left[ d \log \frac{1}{\varepsilon} + \log \frac{1}{\delta} \right]$



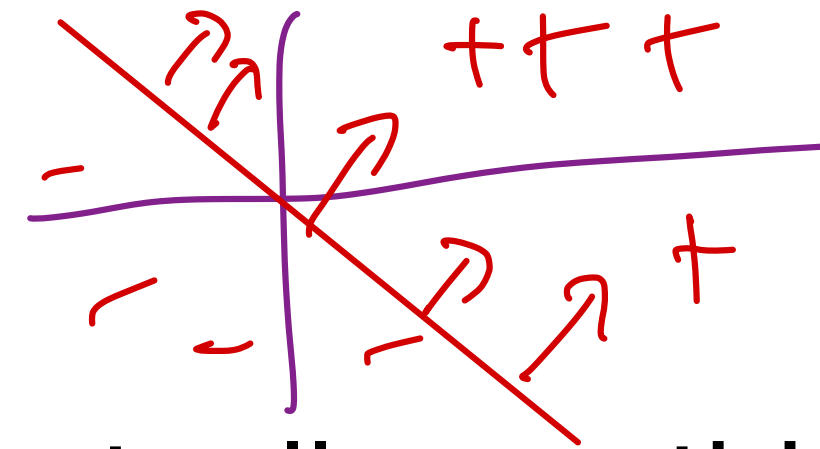
# Last time: Finite VCdim implies uniform convergence

- **Uniform convergence:**  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \varepsilon$  with probability at least  $1 - \delta$
- Will prove in terms of the **growth function**:  $\tau_{\mathcal{H}}(n) = \max_{C \subseteq \mathcal{X}: |C|=n} |\mathcal{H}_C|$ 
  - How many *actually different* functions from  $\mathcal{H}$  are there on sets of size  $n$ ?
  - If  $\text{VCdim}(\mathcal{H}) = d$ , then  $\tau_{\mathcal{H}}(n) = 2^n$  for  $n \leq d$
- **Theorem** (SSBD 6.11):  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2n))}}{\delta \sqrt{2n}}$  **This is the part we didn't prove yet**
- **Sauer-Shelah lemma:** when  $n \geq d$ ,  $\tau_{\mathcal{H}}(n) \leq (en/d)^d = \mathcal{O}(n^d)$
- Plugging together: uniform convergence when  $n \geq 4 \frac{2d}{(\delta \varepsilon)^2} \log \left( \frac{2d}{(\delta \varepsilon)^2} \right) + \frac{4d \log(2e/d)}{(\delta \varepsilon)^2}$



# VC dimension of linear classifiers

- $h_w(x) = \mathbb{I}(w^T x \geq 0)$  on  $\mathbb{R}^d$



- Can shatter, e.g.,  $\{e_1, e_2, \dots, e_d\}$  (actually, anything in "general position")

$$\Sigma_{ii}^T = \begin{cases} 1/\epsilon_{ii} & \text{if } \epsilon_{ii} > 0 \\ 0 & \text{if } \epsilon_{ii} = 0 \end{cases}$$

$$(1, 0, \dots, 0) \in \mathbb{R}^d$$

$$w = (2\gamma_1^{-1}, \dots, 2\gamma_d^{-1})$$

$$w^T e_i = 2\gamma_i^{-1}$$

let  $X = U \Sigma V^T$  be its SVD

$$W = \underbrace{V \Sigma^{-1} U^T}_{X^+} Y$$

(pseudo-inverse)

$$(\gamma_1, \dots, \gamma_n)$$

$\gamma_i \in \{-1, 1\}$

$$r \leq \min(n, d)$$

$$\begin{aligned} Xw &= X V \Sigma^{-1} U^T Y \\ &= U \underbrace{\Sigma V^T V \Sigma^{-1}}_{I_r} U^T Y = \underbrace{U U^T}_{n \times 1} Y = Y \quad \text{if } r = n \end{aligned}$$

projection onto rank-r subspace

# VC dimension of linear classifiers

- $h_w(x) = \mathbb{I}(w^\top x \geq 0)$  on  $\mathbb{R}^d$
- Can't shatter anything of size  $d + 1$

$x_1, \dots, x_{d+1}$

not all  $a_i = 0$

$$\sum_{i=1}^{d+1} a_i x_i = 0$$

$$\sum_{i: a_i > 0} a_i x_i = \sum_{i: a_i < 0} (-a_i) x_i$$

$\uparrow w^\top x_i > 0$

$\uparrow w^\top x_i < 0$

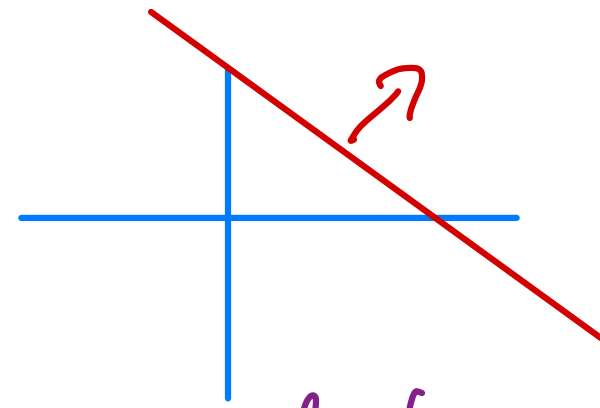
$$0 < \sum_{i: a_i > 0} \underbrace{a_i}_{>0} \underbrace{x_i^\top w}_{>0} = \sum_{i: a_i < 0} \underbrace{(-a_i)}_{>0} \underbrace{x_i^\top w}_{<0} < 0$$

What about all  $a_i \geq 0$ ,  
 $w^\top x_i = 0$

# VC dimension of non-homogenous linear classifiers

- $h_w(x) = \mathbb{I}(w_0 + w^T x \geq 0)$

$\tilde{x} \Rightarrow (1, x) \in \mathbb{R}^{d+1}$   
 $\tilde{w} = (w_0, w) \in \mathbb{R}^{d+1}$



Can't shatter sets of size  $d+2$

Can shatter  $\{e_1, \dots, e_{d+1}\}$

VC dim =  $d+1$

# M-ing the ER for halfspaces

- Just showed that ERM will agnostically PAC-learn a linear classifier (halfspace) with  $\Omega\left(\frac{1}{\varepsilon} \left[d + \log \frac{1}{\delta}\right]\right)$  samples (with 0-1 loss)

# M-ing the ER for halfspaces

- Just showed that ERM will agnostically PAC-learn a linear classifier (halfspace) with  $\Omega\left(\frac{1}{\varepsilon} \left[d + \log \frac{1}{\delta}\right]\right)$  samples (with 0-1 loss)
- But...turns out that finding the optimal halfspace (for 0-1 loss) is NP-hard to even *approximate* (i.e. to get loss  $(1 + \gamma) L^*$ )  
[Ben-David+Simon, NeurIPS 2000]

# M-ing the ER for halfspaces

- Just showed that ERM will agnostically PAC-learn a linear classifier (halfspace) with  $\Omega\left(\frac{1}{\varepsilon} \left[d + \log \frac{1}{\delta}\right]\right)$  samples (with 0-1 loss)
- But...turns out that finding the optimal halfspace (for 0-1 loss) is NP-hard to even *approximate* (i.e. to get loss  $(1 + \gamma) L^*$ )  
[Ben-David+Simon, NeurIPS 2000]
- In the realizable (separable) case ( $L_S(h_S) = 0$ ), easy algorithms in polynomial time
  - Perceptron
  - Linear programming
  - Logistic regression
  - SVMs
  - ...

**(pause)**



# Generalization bound from growth functions

- **SSBD's theorem 6.11:** 
$$\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2n))}}{\delta \sqrt{2n}}$$

with probability at least  $1 - \delta$  over the choice of  $S \sim \mathcal{D}^n$ , for any  $\mathcal{D}$

# Generalization bound from growth functions

- **SSBD's theorem 6.11:** 
$$\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2n))}}{\delta \sqrt{2n}}$$

with probability at least  $1 - \delta$  over the choice of  $S \sim \mathcal{D}^n$ , for any  $\mathcal{D}$

- Follows from 
$$\mathbb{E} \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2n))}}{\sqrt{2n}}$$

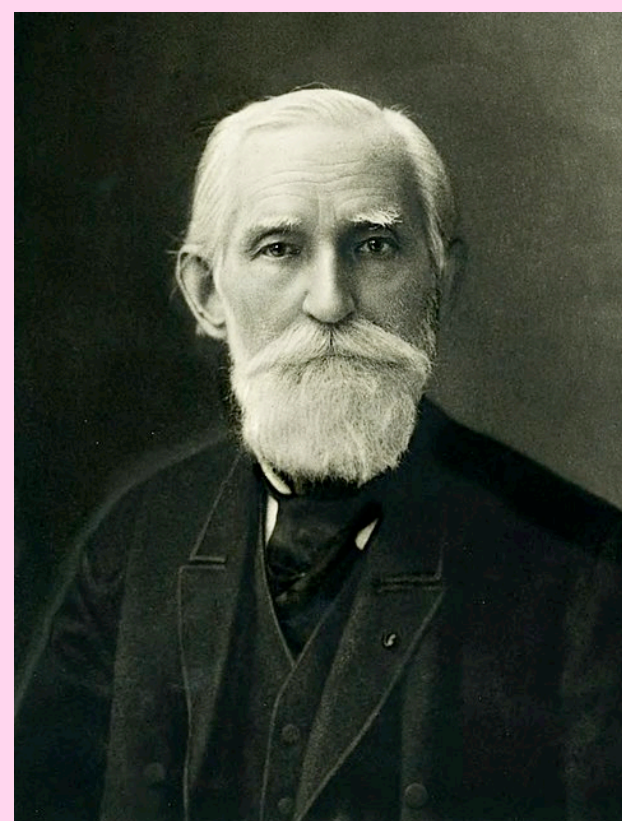
# Generalization bound from growth functions

- **SSBD's theorem 6.11:**  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2n))}}{\delta \sqrt{2n}}$

with probability at least  $1 - \delta$  over the choice of  $S \sim \mathcal{D}^n$ , for any  $\mathcal{D}$

- Follows from  $\mathbb{E} \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2n))}}{\sqrt{2n}}$

**Markov's  
inequality**  
(~1860s)



Andrey Markov Pafnuty Chebyshev



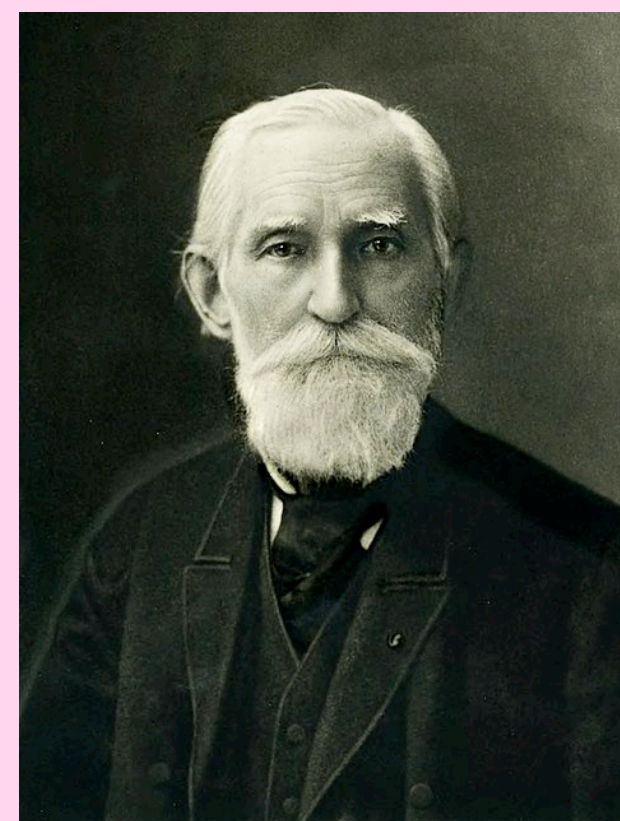
# Generalization bound from growth functions

- **SSBD's theorem 6.11:**  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2n))}}{\delta \sqrt{2n}}$

with probability at least  $1 - \delta$  over the choice of  $S \sim \mathcal{D}^n$ , for any  $\mathcal{D}$

- Follows from  $\mathbb{E} \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2n))}}{\sqrt{2n}}$

**Markov's  
inequality**  
(~1860s)



Andrey Markov Pafnuty Chebyshev

If  $\Pr(X \geq 0) = 1$ , then  $\Pr\left(X > \frac{1}{\delta} \mathbb{E}[X]\right) \leq \delta$



# Generalization bound from growth functions

- **SSBD's theorem 6.11:**  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2n))}}{\delta \sqrt{2n}}$

with probability at least  $1 - \delta$  over the choice of  $S \sim \mathcal{D}^n$ , for any  $\mathcal{D}$

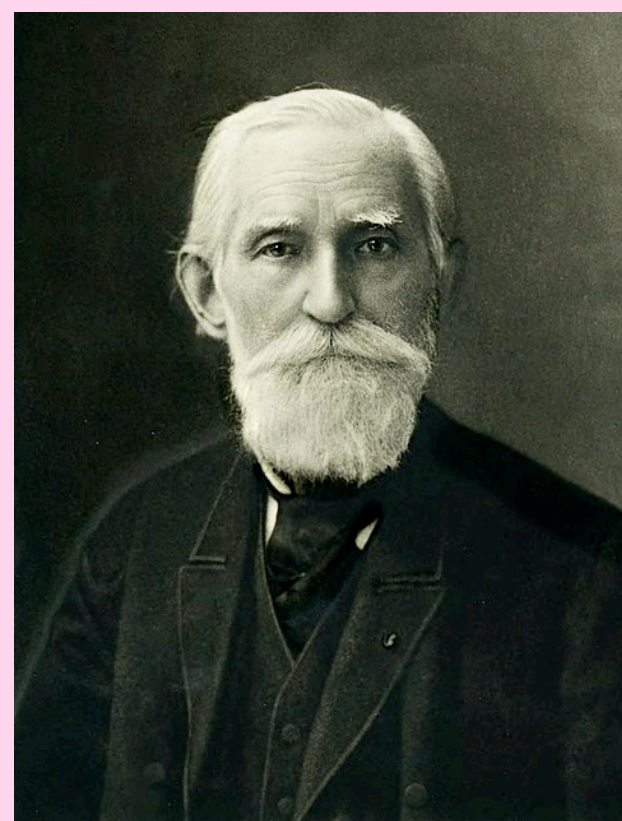
- Follows from  $\mathbb{E} \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)| \leq \frac{4 + \sqrt{\log(\tau_{\mathcal{H}}(2n))}}{\sqrt{2n}}$

$$\begin{aligned} \frac{1}{\delta} &= e^a \\ \log \frac{1}{\delta} &= a \\ \sqrt{\log \frac{1}{\delta}} &= 3 \end{aligned}$$

**Markov's inequality**  
(~1860s)



Andrey Markov



Pafnuty Chebyshev

If  $\Pr(X \geq 0) = 1$ , then  $\Pr\left(X > \frac{1}{\delta} \mathbb{E}[X]\right) \leq \delta$

**Proof:** take  $a = \frac{1}{\delta} \mathbb{E}X$  in:

$$a \mathbb{1}_{[X \geq a]} \leq X, \text{ so } \mathbb{E}[a \mathbb{1}_{[X \geq a]}] = a \Pr(X \geq a) \leq \mathbb{E}X$$

but, actually, we're going to do better than that

but, actually, we're going to do better than that

using (a) Rademacher complexity  
and (b) McDiarmid's inequality rather than Markov's



but, actually, we're going to do better than that

using (a) Rademacher complexity  
and (b) McDiarmid's inequality rather than Markov's

the overall proof uses the same core techniques,  
and the basic version gets *much* closer to the optimal rate,  
using machinery that we were going to do pretty soon anyway

most of this is in SSBD chapter 26,  
but today's presentation will more or less follow MRT chapter 3

# Rademacher complexities

- Measure the complexity of a set of functions  $\mathcal{G}$  from  $\mathcal{Z}$  to  $\mathbb{R}$ 
  - For example, could have  $\mathcal{Z} = \mathcal{X}$  and  $\mathcal{G} = \mathcal{H}$  of binary classifiers
  - But (unlike VC), will be easy to use for much more general settings

# Rademacher complexities

- Measure the complexity of a set of functions  $\mathcal{G}$  from  $\mathcal{Z}$  to  $\mathbb{R}$ 
  - For example, could have  $\mathcal{Z} = \mathcal{X}$  and  $\mathcal{G} = \mathcal{H}$  of binary classifiers
  - But (unlike VC), will be easy to use for much more general settings
- The **empirical Rademacher complexity** of  $\mathcal{G}$  on a set  $S = (z_1, \dots, z_n)$  is

# Rademacher complexities

- Measure the complexity of a set of functions  $\mathcal{G}$  from  $\mathcal{Z}$  to  $\mathbb{R}$ 
  - For example, could have  $\mathcal{Z} = \mathcal{X}$  and  $\mathcal{G} = \mathcal{H}$  of binary classifiers
  - But (unlike VC), will be easy to use for much more general settings
- The **empirical Rademacher complexity** of  $\mathcal{G}$  on a set  $S = (z_1, \dots, z_n)$  is

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(z_i) \right]$$

$\sigma \sim \text{Rad}^n$

# Rademacher complexities

- Measure the complexity of a set of functions  $\mathcal{G}$  from  $\mathcal{Z}$  to  $\mathbb{R}$ 
  - For example, could have  $\mathcal{Z} = \mathcal{X}$  and  $\mathcal{G} = \mathcal{H}$  of binary classifiers
  - But (unlike VC), will be easy to use for much more general settings
- The **empirical Rademacher complexity** of  $\mathcal{G}$  on a set  $S = (z_1, \dots, z_n)$  is

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(z_i) \right]$$

$\sigma \sim \text{Rad}^n$

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

# Rademacher complexities

- Measure the complexity of a set of functions  $\mathcal{G}$  from  $\mathcal{Z}$  to  $\mathbb{R}$ 
  - For example, could have  $\mathcal{Z} = \mathcal{X}$  and  $\mathcal{G} = \mathcal{H}$  of binary classifiers
  - But (unlike VC), will be easy to use for much more general settings
- The **empirical Rademacher complexity** of  $\mathcal{G}$  on a set  $S = (z_1, \dots, z_n)$  is

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(z_i) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \boldsymbol{\sigma}^\top \mathbf{g}_S \right] \quad \mathbf{g}_S = (g(z_1), \dots, g(z_n))$$

$\boldsymbol{\sigma} \sim \text{Rad}^n$

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

# Rademacher complexities

- Measure the complexity of a set of functions  $\mathcal{G}$  from  $\mathcal{Z}$  to  $\mathbb{R}$ 
  - For example, could have  $\mathcal{Z} = \mathcal{X}$  and  $\mathcal{G} = \mathcal{H}$  of binary classifiers
  - But (unlike VC), will be easy to use for much more general settings
- The **empirical Rademacher complexity** of  $\mathcal{G}$  on a set  $S = (z_1, \dots, z_n)$  is

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(z_i) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sigma^\top \mathbf{g}_S \right] \quad \mathbf{g}_S = (g(z_1), \dots, g(z_n))$$

$\sigma \sim \text{Rad}^n$

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

“how well can functions from  $\mathcal{G}$  correlate with random noise?”



# Rademacher complexities

- Measure the complexity of a set of functions  $\mathcal{G}$  from  $\mathcal{Z}$  to  $\mathbb{R}$ 
  - For example, could have  $\mathcal{Z} = \mathcal{X}$  and  $\mathcal{G} = \mathcal{H}$  of binary classifiers
  - But (unlike VC), will be easy to use for much more general settings
- The **empirical Rademacher complexity** of  $\mathcal{G}$  on a set  $S = (z_1, \dots, z_n)$  is

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \sigma_i g(z_i) \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \frac{1}{n} \sigma^\top \mathbf{g}_S \right] \quad \mathbf{g}_S = (g(z_1), \dots, g(z_n))$$

$\sigma \sim \text{Rad}^n$

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

“how well can functions from  $\mathcal{G}$  correlate with random noise?”

- The (“average-case”) **Rademacher complexity** is just  $\mathfrak{R}_n(\mathcal{G}) = \mathbb{E}_{S \sim \mathcal{D}^n} [\widehat{\mathfrak{R}}_S(\mathcal{G})]$ 
  - Distribution-dependent notion of complexity!

# $\mathfrak{R}$ of singleton sets

- If  $\mathcal{H} = \{h\}$

$$\hat{R}_S(\mathcal{H}) = \mathbb{E}_{\sigma} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_i \sigma_i h(z_i)$$

$$= \mathbb{E}_{\sigma} \frac{1}{n} \sum_i \sigma_i h(z_i)$$

$$= \frac{1}{n} \sum_i \underbrace{\mathbb{E}[\sigma_i]}_0 h(z_i)$$

$$= 0$$

# $\mathfrak{R}$ of linear functions

- If  $\mathcal{H} = \{x \mapsto w^\top x : \|w\| \leq B\}$

Turns out  $\hat{R}_S(2t) \leq \frac{2B \sup_{x \in \mathcal{X}} \|x\|}{\sqrt{n}}$



Pascal Massart

# Relating back to VC dim

- **Massart's lemma:** for  $\mathcal{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathcal{A}} \|a\| \leq r$ ,  $\mathbb{E}_{\sigma} \left[ \max_{a \in \mathcal{A}} \frac{1}{n} \sigma^{\top} a \right] \leq \frac{1}{n} r \sqrt{2 \log |\mathcal{A}|}$



Pascal Massart

# Relating back to VC dim

- **Massart's lemma:** for  $\mathcal{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathcal{A}} \|a\| \leq r$ ,  $\mathbb{E}_{\sigma} \left[ \max_{a \in \mathcal{A}} \frac{1}{n} \sigma^{\top} a \right] \leq \frac{1}{n} r \sqrt{2 \log |\mathcal{A}|}$ 
  - (Proof is a nice + not too complicated result on concentration of max of sums; we'll come back to it)



Pascal Massart

# Relating back to VC dim

- **Massart's lemma:** for  $\mathcal{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathcal{A}} \|a\| \leq r$ ,  $\mathbb{E}_{\sigma} \left[ \max_{a \in \mathcal{A}} \frac{1}{n} \sigma^{\top} a \right] \leq \frac{1}{n} r \sqrt{2 \log |\mathcal{A}|}$ 
  - (Proof is a nice + not too complicated result on concentration of max of sums; we'll come back to it)
- For binary classifiers, with output in  $\{0,1\}$  or  $\{-1,1\}$ :



Pascal Massart

# Relating back to VC dim

- **Massart's lemma:** for  $\mathcal{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathcal{A}} \|a\| \leq r$ ,  $\mathbb{E}_{\sigma} \left[ \max_{a \in \mathcal{A}} \frac{1}{n} \sigma^\top a \right] \leq \frac{1}{n} r \sqrt{2 \log |\mathcal{A}|}$ 
  - (Proof is a nice + not too complicated result on concentration of max of sums; we'll come back to it)
- For binary classifiers, with output in  $\{0, 1\}$  or  $\{-1, 1\}$ :
  - Recall  $\mathcal{H}_S = \left\{ (h(x_1), \dots, h(x_n)) : h \in \mathcal{H} \right\} = \left\{ \mathbf{h}_S : h \in \mathcal{H} \right\}$





Pascal Massart

# Relating back to VC dim

- **Massart's lemma:** for  $\mathcal{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathcal{A}} \|a\| \leq r$ ,  $\mathbb{E}_{\sigma} \left[ \max_{a \in \mathcal{A}} \frac{1}{n} \sigma^{\top} a \right] \leq \frac{1}{n} r \sqrt{2 \log |\mathcal{A}|}$ 
  - (Proof is a nice + not too complicated result on concentration of max of sums; we'll come back to it)
- For binary classifiers, with output in  $\{0, 1\}$  or  $\{-1, 1\}$ :
  - Recall  $\mathcal{H}_S = \left\{ (h(x_1), \dots, h(x_n)) : h \in \mathcal{H} \right\} = \left\{ \mathbf{h}_S : h \in \mathcal{H} \right\}$
  - But  $\widehat{\mathfrak{R}}_S(\mathcal{H}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sigma^{\top} \mathbf{h}_S \right] = \mathbb{E}_{\sigma} \left[ \max_{\mathbf{h}_S \in \mathcal{H}_S} \frac{1}{n} \sigma^{\top} \mathbf{h}_S \right]$



Pascal Massart

# Relating back to VC dim

- **Massart's lemma:** for  $\mathcal{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathcal{A}} \|a\| \leq r$ ,  $\mathbb{E}_{\sigma} \left[ \max_{a \in \mathcal{A}} \frac{1}{n} \sigma^{\top} a \right] \leq \frac{1}{n} r \sqrt{2 \log |\mathcal{A}|}$ 
  - (Proof is a nice + not too complicated result on concentration of max of sums; we'll come back to it)
- For binary classifiers, with output in  $\{0, 1\}$  or  $\{-1, 1\}$ :
  - Recall  $\mathcal{H}_S = \left\{ (h(x_1), \dots, h(x_n)) : h \in \mathcal{H} \right\} = \left\{ \mathbf{h}_S : h \in \mathcal{H} \right\}$
  - But  $\widehat{\mathfrak{R}}_S(\mathcal{H}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sigma^{\top} \mathbf{h}_S \right] = \mathbb{E}_{\sigma} \left[ \max_{\mathbf{h}_S \in \mathcal{H}_S} \frac{1}{n} \sigma^{\top} \mathbf{h}_S \right]$ 
$$\leq \frac{1}{n} \sqrt{n} \sqrt{2 \log |\mathcal{H}_S|}$$

Massart's lemma, using  
 $\|\mathbf{h}_S\| \leq \sqrt{1^2 + \dots + 1^2} = \sqrt{n}$



Pascal Massart

# Relating back to VC dim

- **Massart's lemma:** for  $\mathcal{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathcal{A}} \|a\| \leq r$ ,  $\mathbb{E}_{\sigma} \left[ \max_{a \in \mathcal{A}} \frac{1}{n} \sigma^{\top} a \right] \leq \frac{1}{n} r \sqrt{2 \log |\mathcal{A}|}$ 
  - (Proof is a nice + not too complicated result on concentration of max of sums; we'll come back to it)

- For binary classifiers, with output in  $\{0, 1\}$  or  $\{-1, 1\}$ :

- Recall  $\mathcal{H}_S = \left\{ (h(x_1), \dots, h(x_n)) : h \in \mathcal{H} \right\} = \left\{ \mathbf{h}_S : h \in \mathcal{H} \right\}$

- But  $\widehat{\mathfrak{R}}_S(\mathcal{H}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sigma^{\top} \mathbf{h}_S \right] = \mathbb{E}_{\sigma} \left[ \max_{\mathbf{h}_S \in \mathcal{H}_S} \frac{1}{n} \sigma^{\top} \mathbf{h}_S \right]$

$$\leq \frac{1}{n} \sqrt{n} \sqrt{2 \log |\mathcal{H}_S|} \leq \sqrt{\frac{2}{n} \log \tau_{\mathcal{H}}(n)}$$

Massart's lemma, using

$$\|\mathbf{h}_S\| \leq \sqrt{1^2 + \dots + 1^2} = \sqrt{n}$$

by definition of the growth function

$$\tau_{\mathcal{H}}(n) = \sup_{|S|=n} |\mathcal{H}_S|$$



Pascal Massart

# Relating back to VC dim

- **Massart's lemma:** for  $\mathcal{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathcal{A}} \|a\| \leq r$ ,  $\mathbb{E}_{\sigma} \left[ \max_{a \in \mathcal{A}} \frac{1}{n} \sigma^{\top} a \right] \leq \frac{1}{n} r \sqrt{2 \log |\mathcal{A}|}$ 
  - (Proof is a nice + not too complicated result on concentration of max of sums; we'll come back to it)

- For binary classifiers, with output in  $\{0, 1\}$  or  $\{-1, 1\}$ :

- Recall  $\mathcal{H}_S = \left\{ (h(x_1), \dots, h(x_n)) : h \in \mathcal{H} \right\} = \left\{ \mathbf{h}_S : h \in \mathcal{H} \right\}$

- But  $\widehat{\mathfrak{R}}_S(\mathcal{H}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sigma^{\top} \mathbf{h}_S \right] = \mathbb{E}_{\sigma} \left[ \max_{\mathbf{h}_S \in \mathcal{H}_S} \frac{1}{n} \sigma^{\top} \mathbf{h}_S \right]$ 

where  $d$  is the VC dimension of  $\mathcal{H}$   
and  $n \geq d$  (Sauer-Shelah)

$$\leq \frac{1}{n} \sqrt{n} \sqrt{2 \log |\mathcal{H}_S|} \leq \sqrt{\frac{2}{n} \log \tau_{\mathcal{H}}(n)} \leq \sqrt{\frac{2}{n} d [1 + \log n - \log d]}$$

Massart's lemma, using  
 $\|\mathbf{h}_S\| \leq \sqrt{1^2 + \dots + 1^2} = \sqrt{n}$

by definition of the growth function  
 $\tau_{\mathcal{H}}(n) = \sup_{|S|=n} |\mathcal{H}_S|$





Pascal Massart

# Relating back to VC dim

- Massart's lemma:** for  $\mathcal{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathcal{A}} \|a\| \leq r$ ,  $\mathbb{E}_{\sigma} \left[ \max_{a \in \mathcal{A}} \frac{1}{n} \sigma^{\top} a \right] \leq \frac{1}{n} r \sqrt{2 \log |\mathcal{A}|}$ 
    - (Proof is a nice + not too complicated result on concentration of max of sums; we'll come back to it)
  - For binary classifiers, with output in  $\{0, 1\}$  or  $\{-1, 1\}$ :
    - Recall  $\mathcal{H}_S = \left\{ (h(x_1), \dots, h(x_n)) : h \in \mathcal{H} \right\} = \left\{ \mathbf{h}_S : h \in \mathcal{H} \right\}$
    - But  $\widehat{\mathfrak{R}}_S(\mathcal{H}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sigma^{\top} \mathbf{h}_S \right] = \mathbb{E}_{\sigma} \left[ \max_{\mathbf{h}_S \in \mathcal{H}_S} \frac{1}{n} \sigma^{\top} \mathbf{h}_S \right]$ 

$$\leq \frac{1}{n} \sqrt{n} \sqrt{2 \log |\mathcal{H}_S|} \leq \sqrt{\frac{2}{n} \log \tau_{\mathcal{H}}(n)} \leq \sqrt{\frac{2}{n} d [1 + \log n - \log d]}$$

Massart's lemma, using  
 $\|\mathbf{h}_S\| \leq \sqrt{1^2 + \dots + 1^2} = \sqrt{n}$

by definition of the growth function  
 $\tau_{\mathcal{H}}(n) = \sup_{|S|=n} |\mathcal{H}_S|$

where  $d$  is the VC dimension of  $\mathcal{H}$   
 and  $n \geq d$  (Sauer-Shelah)
- $$\leq \sqrt{\frac{2d \log n}{n}} \quad \text{if } d \geq 3$$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension
  - We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{ VCdim}(\mathcal{H}) \log(n) / n}$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension
  - We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$
- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$



- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension
  - We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$
- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$ 
  - $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension
  - We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$
- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$ 
  - $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|_x}(\mathcal{H})$   
 or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension
  - We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$
- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$ 
  - $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|_x}(\mathcal{H})$  (proof is next)  
 or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension
  - We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$
- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$ 
  - $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|x}(\mathcal{H})$  (proof is next)
- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$ 

or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension
  - We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$
- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$ 
  - $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|x}(\mathcal{H})$  (proof is next)
- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$ 

or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

  - $\mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \leq 2\mathfrak{R}_n(\mathcal{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}}$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension
  - We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$
- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$ 
  - $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”
  - For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|_x}(\mathcal{H})$  (proof is next)
- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$ 

$$\mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \leq 2\mathfrak{R}_n(\mathcal{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}} \leq \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{n}} \left( \sqrt{\text{VCdim}} + \sqrt{\log \frac{1}{\delta}} \right) \right)$$

or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$



- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension

- We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$

- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

- $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”

- For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|x}(\mathcal{H})$  (proof is next)

- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$  or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- $\mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \leq 2\mathfrak{R}_n(\mathcal{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}} \leq \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{n}} \left( \sqrt{\text{VCdim}} + \sqrt{\log \frac{1}{\delta}} \right) \right)$

$L_{\mathcal{D}}(h)$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension

- We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$

- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

- $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”

- For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|x}(\mathcal{H})$  (proof is next)

- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$  or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- $$\underbrace{\mathbb{E}[g(z)]}_{L_{\mathcal{D}}(h)} - \underbrace{\frac{1}{n} \sum_{i=1}^n g(z_i)}_{L_S(h)} \leq 2\mathfrak{R}_n(\mathcal{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}} \leq \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{n}} \left( \sqrt{\text{VCdim}} + \sqrt{\log \frac{1}{\delta}} \right) \right)$$



- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension

- We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$

- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

- $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”

- For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|_x}(\mathcal{H})$  (proof is next)

- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$  or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- $$\underbrace{\mathbb{E}[g(z)]}_{L_{\mathcal{D}}(h)} - \underbrace{\frac{1}{n} \sum_{i=1}^n g(z_i)}_{L_S(h)} \leq \underbrace{2\mathfrak{R}_n(\mathcal{G})}_{\mathbb{E} \sup L_{\mathcal{D}}(h) - L_S(h)} + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}} \leq \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{n}} \left( \sqrt{\text{VCdim}} + \sqrt{\log \frac{1}{\delta}} \right) \right)$$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension

- We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$

- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

- $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”

- For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|x}(\mathcal{H})$  (proof is next)

- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$  or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- $$\underbrace{\mathbb{E}[g(z)]}_{L_{\mathcal{D}}(h)} - \underbrace{\frac{1}{n} \sum_{i=1}^n g(z_i)}_{L_S(h)} \leq \underbrace{2\mathfrak{R}_n(\mathcal{G})}_{\mathbb{E} \sup L_{\mathcal{D}}(h) - L_S(h)} + \underbrace{\sqrt{\frac{1}{2n} \log \frac{1}{\delta}}}_{\text{how much bigger than its mean is likely}} \leq \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{n}} \left( \sqrt{\text{VCdim}} + \sqrt{\log \frac{1}{\delta}} \right) \right)$$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension

- We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$

- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

- $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”

- For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|_x}(\mathcal{H})$  (proof is next)

- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$  or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- $$\left| \underbrace{\mathbb{E}[g(z)]}_{L_{\mathcal{D}}(h)} - \frac{1}{n} \sum_{i=1}^n \underbrace{g(z_i)}_{L_S(h)} \right| \leq 2 \underbrace{\mathfrak{R}_n(\mathcal{G})}_{\text{bounds}} + \underbrace{\sqrt{\frac{1}{2n} \log \frac{1}{\delta}}}_{\text{how much bigger than its mean is likely}} \leq \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{n}} \left( \sqrt{\text{VCdim}} + \sqrt{\log \frac{1}{\delta}} \right) \right)$$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension

- We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$

- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

- $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”

- For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|x}(\mathcal{H})$  (proof is next)

- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$  or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- $\left| \mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \right| \leq 2\mathfrak{R}_n(\mathcal{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}} \leq \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{n}} \left( \sqrt{\text{VCdim}} + \sqrt{\log \frac{1}{\delta}} \right) \right)$

$L_{\mathcal{D}}(h)$       $L_S(h)$       $\mathbb{E} \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$

bounds     how much bigger than its mean is likely

- $\mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \leq 2\widehat{\mathfrak{R}}_S(\mathcal{G}) + 3\sqrt{\frac{1}{2n} \log \frac{2}{\delta}}$



- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension

- We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$

- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

- $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”

- For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|x}(\mathcal{H})$  (proof is next)

- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$  or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- $\left| \mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \right| \leq 2\mathfrak{R}_n(\mathcal{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}} \leq \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{n}} \left( \sqrt{\text{VCdim}} + \sqrt{\log \frac{1}{\delta}} \right) \right)$

$L_{\mathcal{D}}(h)$      $L_S(h)$      $\mathbb{E} \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$

bounds    how much bigger than its mean is likely

- $\left| \mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \right| \leq 2\widehat{\mathfrak{R}}_S(\mathcal{G}) + 3\sqrt{\frac{1}{2n} \log \frac{2}{\delta}}$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension

- We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$

- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

- $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”

- For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|x}(\mathcal{H})$  (proof is next)

- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$  or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- $\left| \mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \right| \leq 2\mathfrak{R}_n(\mathcal{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}} \leq \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{n}} \left( \sqrt{\text{VCdim}} + \sqrt{\log \frac{1}{\delta}} \right) \right)$

$L_{\mathcal{D}}(h)$   $L_S(h)$   $\mathbb{E} \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  bounds how much bigger than its mean is likely

- $\left| \mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \right| \leq 2\widehat{\mathfrak{R}}_S(\mathcal{G}) + 3\sqrt{\frac{1}{2n} \log \frac{2}{\delta}}$   
because  $\widehat{\mathfrak{R}}_S$  is close to  $\mathfrak{R}_n$

- Our goal was to upper bound  $\sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$  in terms of VC dimension

- We've shown  $\mathfrak{R}_n(\mathcal{H}) \leq \sqrt{2 \text{VCdim}(\mathcal{H}) \log(n) / n}$

- One more step: consider  $\mathfrak{R}_n(\mathcal{G})$  for  $\mathcal{G} = \{z \mapsto \ell(h, z) : h \in \mathcal{H}\}$

- $g_h : \mathcal{Z} \rightarrow \mathbb{R}$  says, “what would the loss of predictor  $h$  be for a given  $(x, y)$ ?”

- For 0-1 loss,  $g_h((x, y)) = \mathbb{I}(h(x) \neq y)$ , and  $\widehat{\mathfrak{R}}_S(\mathcal{G}) = \widehat{\mathfrak{R}}_{S|x}(\mathcal{H})$  (proof is next)

- **Theorem:** for  $\mathcal{G}$  mapping to  $[0, 1]$ , prob is  $\geq 1 - \delta$  over  $S \sim \mathcal{D}^n$  that for **all**  $g \in \mathcal{G}$  or  $\frac{1}{2}$  that, if  $\mathcal{H}$  maps to  $\pm 1$

- $\left| \mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \right| \leq 2\mathfrak{R}_n(\mathcal{G}) + \sqrt{\frac{1}{2n} \log \frac{1}{\delta}} \leq \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{n}} \left( \sqrt{\text{VCdim}} + \sqrt{\log \frac{1}{\delta}} \right) \right)$

$L_{\mathcal{D}}(h)$      $L_S(h)$      $\mathbb{E} \sup_{h \in \mathcal{H}} |L_{\mathcal{D}}(h) - L_S(h)|$

bounds    how much bigger than its mean is likely

- $\left| \mathbb{E}[g(z)] - \frac{1}{n} \sum_{i=1}^n g(z_i) \right| \leq 2\widehat{\mathfrak{R}}_S(\mathcal{G}) + 3\sqrt{\frac{1}{2n} \log \frac{2}{\delta}}$  (proof on Monday)

because  $\widehat{\mathfrak{R}}_S$  is close to  $\mathfrak{R}_n$



# $\mathcal{R}$ of binary classifiers vs loss

# $\mathfrak{R}$ of binary classifiers vs loss

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \right]$$

# $\mathfrak{R}$ of binary classifiers vs loss

If  $h : \mathcal{X} \rightarrow \{0,1\}$ , define  $\tilde{h} : \mathcal{X} \rightarrow \{-1,1\}$  by  $\tilde{h}(x) = 2h(x) - 1$ ;  $\tilde{y}_i = 2y_i - 1$

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \right]$$

# $\mathfrak{R}$ of binary classifiers vs loss

If  $h : \mathcal{X} \rightarrow \{0,1\}$ , define  $\tilde{h} : \mathcal{X} \rightarrow \{-1,1\}$  by  $\tilde{h}(x) = 2h(x) - 1$ ;  $\tilde{y}_i = 2y_i - 1$

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \right] = \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right]$$

# $\mathfrak{R}$ of binary classifiers vs loss

If  $h : \mathcal{X} \rightarrow \{0,1\}$ , define  $\tilde{h} : \mathcal{X} \rightarrow \{-1,1\}$  by  $\tilde{h}(x) = 2h(x) - 1$ ;  $\tilde{y}_i = 2y_i - 1$

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \right] = \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right]$$

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

$$\mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{G}} \boldsymbol{\sigma}^\top (\mathbf{g}_S + c\mathbf{1}) \right] = \mathbb{E}_\sigma \left[ \left( \sup_{g \in \mathcal{G}} \boldsymbol{\sigma}^\top \mathbf{g}_S \right) + c \boldsymbol{\sigma}^\top \mathbf{1} \right] = \mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{G}} \boldsymbol{\sigma}^\top \mathbf{g}_S \right]$$

# $\mathfrak{R}$ of binary classifiers vs loss

If  $h : \mathcal{X} \rightarrow \{0,1\}$ , define  $\tilde{h} : \mathcal{X} \rightarrow \{-1,1\}$  by  $\tilde{h}(x) = 2h(x) - 1$ ;  $\tilde{y}_i = 2y_i - 1$

$$\begin{aligned}\widehat{\mathfrak{R}}_S(\mathcal{G}) &= \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \right] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right] \\ &= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n -\sigma_i \tilde{y}_i \tilde{h}(x_i) \right]\end{aligned}$$

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

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} (\mathbf{g}_S + c \mathbf{1}) \right] = \mathbb{E}_{\sigma} \left[ \left( \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right) + c \sigma^{\top} \mathbf{1} \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right]$$

# $\mathfrak{R}$ of binary classifiers vs loss

If  $h : \mathcal{X} \rightarrow \{0,1\}$ , define  $\tilde{h} : \mathcal{X} \rightarrow \{-1,1\}$  by  $\tilde{h}(x) = 2h(x) - 1$ ;  $\tilde{y}_i = 2y_i - 1$

$$\begin{aligned}\widehat{\mathfrak{R}}_S(\mathcal{G}) &= \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \right] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right] \\ &= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n -\sigma_i \tilde{y}_i \tilde{h}(x_i) \right] \quad \text{for fixed } \tilde{y}_i, -\sigma_i \tilde{y}_i \sim \text{Rad}\end{aligned}$$

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

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} (\mathbf{g}_S + c \mathbf{1}) \right] = \mathbb{E}_{\sigma} \left[ \left( \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right) + c \sigma^{\top} \mathbf{1} \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right]$$



# $\mathfrak{R}$ of binary classifiers vs loss

If  $h : \mathcal{X} \rightarrow \{0,1\}$ , define  $\tilde{h} : \mathcal{X} \rightarrow \{-1,1\}$  by  $\tilde{h}(x) = 2h(x) - 1$ ;  $\tilde{y}_i = 2y_i - 1$

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \right] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right]$$

$$= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n -\sigma_i \tilde{y}_i \tilde{h}(x_i) \right] \quad \text{for fixed } \tilde{y}_i, -\sigma_i \tilde{y}_i \sim \text{Rad}$$

$$= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{h}(x_i) \right]$$

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

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} (\mathbf{g}_S + c \mathbf{1}) \right] = \mathbb{E}_{\sigma} \left[ \left( \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right) + c \sigma^{\top} \mathbf{1} \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right]$$

# $\mathfrak{R}$ of binary classifiers vs loss

If  $h : \mathcal{X} \rightarrow \{0,1\}$ , define  $\tilde{h} : \mathcal{X} \rightarrow \{-1,1\}$  by  $\tilde{h}(x) = 2h(x) - 1$ ;  $\tilde{y}_i = 2y_i - 1$

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \right] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right]$$

$$= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n -\sigma_i \tilde{y}_i \tilde{h}(x_i) \right] \quad \text{for fixed } \tilde{y}_i, -\sigma_i \tilde{y}_i \sim \text{Rad}$$

$$= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{h}(x_i) \right] = \frac{1}{2} \widehat{\mathfrak{R}}_{S|_x}(\tilde{\mathcal{H}})$$

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

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} (\mathbf{g}_S + c \mathbf{1}) \right] = \mathbb{E}_{\sigma} \left[ \left( \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right) + c \sigma^{\top} \mathbf{1} \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right]$$

# $\mathfrak{R}$ of binary classifiers vs loss

If  $h : \mathcal{X} \rightarrow \{0,1\}$ , define  $\tilde{h} : \mathcal{X} \rightarrow \{-1,1\}$  by  $\tilde{h}(x) = 2h(x) - 1$ ;  $\tilde{y}_i = 2y_i - 1$

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \right] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right]$$

$$= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n -\sigma_i \tilde{y}_i \tilde{h}(x_i) \right] \quad \text{for fixed } \tilde{y}_i, -\sigma_i \tilde{y}_i \sim \text{Rad}$$

$$= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{h}(x_i) \right] = \frac{1}{2} \widehat{\mathfrak{R}}_{S|_x}(\tilde{\mathcal{H}})$$

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

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} (\mathbf{g}_S + c \mathbf{1}) \right] = \mathbb{E}_{\sigma} \left[ \left( \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right) + c \sigma^{\top} \mathbf{1} \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right]$$

Scaling by  $c$  gives  $|c| \widehat{\mathfrak{R}}_S$ :

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} (c \mathbf{g}_S) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \text{sign}(c) \sigma^{\top} \mathbf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right]$$

# $\mathfrak{R}$ of binary classifiers vs loss

If  $h : \mathcal{X} \rightarrow \{0,1\}$ , define  $\tilde{h} : \mathcal{X} \rightarrow \{-1,1\}$  by  $\tilde{h}(x) = 2h(x) - 1$ ;  $\tilde{y}_i = 2y_i - 1$

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \mathbb{I}(h(x_i) \neq y_i) \right] = \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - \tilde{y}_i \tilde{h}(x_i)}{2} \right]$$

$$= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n -\sigma_i \tilde{y}_i \tilde{h}(x_i) \right] \quad \text{for fixed } \tilde{y}_i, -\sigma_i \tilde{y}_i \sim \text{Rad}$$

$$= \frac{1}{2} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \tilde{h}(x_i) \right] = \frac{1}{2} \widehat{\mathfrak{R}}_{S|_x}(\tilde{\mathcal{H}}) = \widehat{\mathfrak{R}}_{S|_x}(\mathcal{H})$$

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

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} (\mathbf{g}_S + c \mathbf{1}) \right] = \mathbb{E}_{\sigma} \left[ \left( \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right) + c \sigma^{\top} \mathbf{1} \right] = \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right]$$

Scaling by  $c$  gives  $|c| \widehat{\mathfrak{R}}_S$ :

$$\mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} (c \mathbf{g}_S) \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \text{sign}(c) \sigma^{\top} \mathbf{g}_S \right] = |c| \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \sigma^{\top} \mathbf{g}_S \right]$$

If there's time left, let's prove: (If not: we'll come back to it another time!)

**Massart's lemma:** for  $\mathcal{A} \subset \mathbb{R}^n$ , if  $\max_{a \in \mathcal{A}} \|a\| \leq r$ ,

$$\mathbb{E}_{\boldsymbol{\sigma}} \left[ \max_{a \in \mathcal{A}} \frac{1}{n} \boldsymbol{\sigma}^\top a \right] \leq \frac{1}{n} r \sqrt{2 \log |\mathcal{A}|}$$

# Summary

- VC dimension for linear classifiers:
  - $d$  for homogenous,  $d + 1$  for not (same as # of params)
  - So ERM works with 0-1 loss...except we can't do ERM in non-separable case!
- We *still* haven't quite proved the fundamental theorem
  - But we'll prove a much better bound when we do, on Monday
- **Empirical Rademacher complexity**: how much can  $\mathcal{G}$  correlate with  $\pm 1$  noise on  $S$ ?
- **Rademacher complexity**: its expectation over a random  $S \sim \mathcal{D}^n$
- Upper bounded in terms of VC dimension
- Stated a generalization bound in terms of Rademacher complexity
  - Will imply (almost) the optimal sample complexity for agnostic case
- Will be easy(ish) to extend to things beyond 0-1 loss binary classification