

# Machine Learning Theory 2022

## Lecture 5

Tim van Erven

Download these slides now from [elo.mastermath.nl!](http://elo.mastermath.nl)

Focus on binary classification:

- ▶ Review
- ▶ Remaining proof:  
growth function controls uniform convergence



homework roulette  
in the break

# Uniform Convergence Upper Bound with VC-Dimension

## Theorem

Consider binary classification. Suppose  $\text{VCdim}(\mathcal{H}) \leq v < \infty$ . Then there exists an absolute constant  $C > 0$  such that

$$\sup_{h \in \mathcal{H}} |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon \quad \text{with probability } \geq 1 - \delta,$$

whenever

$$m \geq C \frac{v \ln(1/\epsilon) + \ln(1/\delta) + 1}{\epsilon^2}.$$

## Proof Approach

**Growth function:**  $\tau_{\mathcal{H}}(m) = \max_{|\mathcal{C}|=m} |\mathcal{H}_{\mathcal{C}}|$

- ▶ Interpretation: How many truly different hypotheses are there when we only observe  $m$  inputs  $\mathcal{C} = \{x_1, \dots, x_m\}$ ?

# Proof Approach

**Growth function:**  $\tau_{\mathcal{H}}(m) = \max_{|\mathcal{C}|=m} |\mathcal{H}_{\mathcal{C}}|$

- ▶ Interpretation: How many truly different hypotheses are there when we only observe  $m$  inputs  $\mathcal{C} = \{x_1, \dots, x_m\}$ ?

**Part I: Growth function controls uniform convergence:**

$$\sup_{h \in \mathcal{H}} |L_S(h) - L_{\mathcal{D}}(h)| \leq c \sqrt{\frac{\ln \tau_{\mathcal{H}}(m)}{m}} + c \sqrt{\frac{\ln(4/\delta)}{m}} \quad \text{with probability } \geq 1 - \delta$$

**Part II: VC-dimension controls growth function (Sauer's Lemma):**

$$\ln \tau_{\mathcal{H}}(m) \leq v \ln \left( \frac{em}{v} \right) \quad \text{for } m > v.$$

- ▶ Finish: combine Parts I and II, and find lower bound on  $m$  s.t.  
 $\sup_{h \in \mathcal{H}} |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon.$

**Proof Part I:**  
**Growth Function Controls Uniform  
Convergence**

## Part I: Proof Outline

### Lemma (Two-sided Bound)

Consider binary classification. Then there exists an absolute constant  $c > 0$  such that, for any  $\delta \in (0, 1]$ ,

$$\sup_{h \in \mathcal{H}} |L_S(h) - L_{\mathcal{D}}(h)| \leq c \sqrt{\frac{\ln \tau_{\mathcal{H}}(m)}{m}} + c \sqrt{\frac{\ln(4/\delta)}{m}} \quad \text{w.p.} \geq 1 - \delta.$$

## Part I: Proof Outline

### Lemma (Two-sided Bound)

Consider binary classification. Then there exists an absolute constant  $c > 0$  such that, for any  $\delta \in (0, 1]$ ,

$$\sup_{h \in \mathcal{H}} |L_S(h) - L_{\mathcal{D}}(h)| \leq c \sqrt{\frac{\ln \tau_{\mathcal{H}}(m)}{m}} + c \sqrt{\frac{\ln(4/\delta)}{m}} \quad \text{w.p.} \geq 1 - \delta.$$

Note: could measure loss in binary classification differently.

Sufficient to show:

### Lemma (One-sided Bound)

For any loss function  $\ell(h, \mathbf{X}, Y) = \tilde{\ell}(h(\mathbf{X}), Y)$  with range  $[0, 1]$ :

$$\sup_{h \in \mathcal{H}} \{L_S(h) - L_{\mathcal{D}}(h)\} \leq c \sqrt{\frac{\ln \tau_{\mathcal{H}}(m)}{m}} + c \sqrt{\frac{\ln(2/\delta)}{m}} \quad \text{w.p.} \geq 1 - \delta.$$

One-sided Bound  $\Rightarrow$  Two-sided Bound

$$\text{Let } z = \sup_{h \in \mathcal{H}} \{ L_S(h) - L_D(h) \}$$

$$z' = \sup_{h \in \mathcal{H}} \{ L_D(h) - L_S(h) \}$$

Applying one-sided bound with  $\ell' = 1 - \ell$   
reduces  $z'$ , because

$$\begin{aligned} L'_S(h) - L'_D(h) &= \frac{1}{m} \sum_{i=1}^m (1 - \ell(h, x_i, y_i)) \\ &\quad - \mathbb{E}[1 - \ell(h, x, y)] \\ &= \mathbb{E}[\ell(h, x, y)] - \frac{1}{m} \sum_{i=1}^m \ell(h, x_i, y_i) \\ &= L_D(h) - L_S(h) \end{aligned}$$

Then

$$\sup_{h \in \mathcal{H}} |L_S(h) - L_D(h)| = \max \{z, z'\}$$
$$\leq C \sqrt{\frac{\ln(2\delta/m)}{m}} + C \sqrt{\frac{\ln(4/\delta)}{m}} \text{ w.p. } \geq 1-\delta$$

by one-sided bounds for  $z$  and  $z'$  with  $\delta' = \frac{\delta}{2}$   
+ union bound.

# Approach for One-Sided Bound

## Lemma (One-sided Bound)

For any loss function  $\ell(h, \mathbf{X}, Y) = \tilde{\ell}(h(\mathbf{X}), Y)$  with range  $[0, 1]$ :

$$\sup_{h \in \mathcal{H}} \{L_S(h) - L_D(h)\} \leq c\sqrt{\frac{\ln \tau_{\mathcal{H}}(m)}{m}} + c\sqrt{\frac{\ln(2/\delta)}{m}} \quad \text{w.p. } \geq 1 - \delta.$$

- Remark:**
- ▶ Book first derives suboptimal dependence on  $\delta$  in Chapter 6
  - ▶ I am taking a shortcut through Chapters 6, 26 and 28

# Approach for One-Sided Bound

## Lemma (One-sided Bound)

For any loss function  $\ell(h, \mathbf{X}, Y) = \tilde{\ell}(h(\mathbf{X}), Y)$  with range  $[0, 1]$ :

$$\sup_{h \in \mathcal{H}} \{L_S(h) - L_D(h)\} \leq c \sqrt{\frac{\ln \tau_{\mathcal{H}}(m)}{m}} + c \sqrt{\frac{\ln(2/\delta)}{m}} \quad \text{w.p. } \geq 1 - \delta.$$

Proof consists of 3 steps:

1. **Concentration:** Abbreviate  $Z = \sup_{h \in \mathcal{H}} \{L_S(h) - L_D(h)\}$ . Then, for any loss function  $\ell(h, \mathbf{X}, Y)$  with range  $[0, 1]$ ,

$$Z \leq \mathbb{E}_S[Z] + c \sqrt{\frac{\ln(2/\delta)}{m}} \quad \text{w.p. } \geq 1 - \delta.$$

# Approach for One-Sided Bound

## Lemma (One-sided Bound)

For any loss function  $\ell(h, \mathbf{X}, Y) = \tilde{\ell}(h(\mathbf{X}), Y)$  with range  $[0, 1]$ :

$$\sup_{h \in \mathcal{H}} \{L_S(h) - L_D(h)\} \leq c \sqrt{\frac{\ln \tau_{\mathcal{H}}(m)}{m}} + c \sqrt{\frac{\ln(2/\delta)}{m}} \quad \text{w.p. } \geq 1 - \delta.$$

Proof consists of 3 steps:

1. **Concentration:** Abbreviate  $Z = \sup_{h \in \mathcal{H}} \{L_S(h) - L_D(h)\}$ . Then, for any loss function  $\ell(h, \mathbf{X}, Y)$  with range  $[0, 1]$ ,

$$Z \leq \mathbb{E}_S[Z] + c \sqrt{\frac{\ln(2/\delta)}{m}} \quad \text{w.p. } \geq 1 - \delta.$$

2. **Symmetrization:** For any loss function:

$$\mathbb{E}[Z] \leq 2 \mathbb{E}[\mathcal{R}(\ell, \mathcal{H}, S)].$$

# Approach for One-Sided Bound

## Lemma (One-sided Bound)

For any loss function  $\ell(h, \mathbf{X}, Y) = \tilde{\ell}(h(\mathbf{X}), Y)$  with range  $[0, 1]$ :

$$\sup_{h \in \mathcal{H}} \{L_S(h) - L_{\mathcal{D}}(h)\} \leq c \sqrt{\frac{\ln \tau_{\mathcal{H}}(m)}{m}} + c \sqrt{\frac{\ln(2/\delta)}{m}} \quad \text{w.p. } \geq 1 - \delta.$$

Proof consists of 3 steps:

1. **Concentration:** Abbreviate  $Z = \sup_{h \in \mathcal{H}} \{L_S(h) - L_{\mathcal{D}}(h)\}$ . Then, for any loss function  $\ell(h, \mathbf{X}, Y)$  with range  $[0, 1]$ ,

$$Z \leq \mathbb{E}_S[Z] + c \sqrt{\frac{\ln(2/\delta)}{m}} \quad \text{w.p. } \geq 1 - \delta.$$

2. **Symmetrization:** For any loss function:

$$\mathbb{E}[Z] \leq 2 \mathbb{E}[\mathcal{R}(\ell, \mathcal{H}, S)].$$

3. For any loss  $\ell(h, \mathbf{X}, Y) = \tilde{\ell}(h(\mathbf{X}), Y)$  with range  $[0, 1]$ :

$$\mathcal{R}(\ell, \mathcal{H}, S) \leq \sqrt{\frac{2 \ln |\mathcal{H}_S|}{m}} \leq \sqrt{\frac{2 \ln \tau_{\mathcal{H}}(m)}{m}} \quad \text{for all } S.$$

### Step I: Concentration

To show: loss  $\ell(h, x_i, y_i) \in [0, 1]$

$$Z = \sup_{h \in H} L_S(h) - L_D(h)$$

$$Z \leq \mathbb{E}[Z] + c \sqrt{\frac{\ln(2/\delta)}{m}} \quad \text{w.p.} \geq 1 - \delta$$

---

Proof:

$Z = f(A_1, \dots, A_m)$  for  $A_i = (x_i, y_i)$

Bounded differences property:

\* If change  $A_i \rightarrow A'_i$ , then

$Z$  changes by at most  $\frac{1}{m}$

(because  $L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, x_i, y_i)$   
changes by at most  $1/m$ )

## McDiarmid's Inequality

Suppose  $A_1, \dots, A_m$  are independent random variables, and  $f: \mathcal{A}^m \rightarrow \mathbb{R}$  satisfies for all  $i$

$$\sup_{\substack{a_1, \dots, a_m \\ a'_i}} |f(a_1, \dots, a_m) - f(a_1, \dots, a_{i-1}, a'_i, a_{i+1}, \dots, a_m)| \leq b.$$

Then, with probability  $\geq 1-\delta$ ,

$$|f(A_1, \dots, A_m) - \mathbb{E}[f(A_1, \dots, A_m)]| \leq b \sqrt{\frac{m}{2} \ln(\frac{2}{\delta})}$$

$Z$  satisfies this with  $b = \epsilon/m$ :

$$|Z - \mathbb{E}[Z]| \leq \sqrt{\frac{\ln(\frac{2}{\delta})}{2m}} \text{ w.p. } \geq 1-\delta$$

$$\leq c \sqrt{\frac{\ln(2/\delta)}{m}} \text{ for } c \geq \frac{1}{\sqrt{2}}$$

□

# Rademacher Complexity

How much can the losses of  $h \in \mathcal{H}$  on  $S$   
**correlate with random errors?**

**Rademacher random variables:** Let  $\sigma = (\sigma_1, \dots, \sigma_m) \in \{-1, +1\}^m$  be i.i.d. with  $\Pr(\sigma_i = -1) = \Pr(\sigma_i = +1) = 1/2$ .

**Rademacher complexity:**

$$\mathcal{R}(\ell, \mathcal{H}, S) = \frac{1}{m} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i \ell(h, \mathbf{X}_i, Y_i) \right]$$

- ▶ Interpret  $\sum_{i=1}^m \sigma_i \ell(h, \mathbf{X}_i, Y_i)$  as correlation of losses with random errors

## Step 2: Symmetrization

$$\mathcal{R}(\ell, \mathcal{H}, S) = \frac{1}{m} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i \ell(h, X_i, Y_i) \right]$$

Lemma

$$\mathbb{E}_S \left[ \sup_{h \in \mathcal{H}} \{L_S(h) - L_{\mathcal{D}}(h)\} \right] \leq 2 \mathbb{E}_S [\mathcal{R}(\ell, \mathcal{H}, S)]$$

## Step 2: Symmetrization

$$\mathcal{R}(\ell, \mathcal{H}, S) = \frac{1}{m} \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i \ell(h, \mathbf{X}_i, Y_i) \right]$$

### Lemma

$$\mathbb{E}_S \left[ \sup_{h \in \mathcal{H}} \{L_S(h) - L_{\mathcal{D}}(h)\} \right] \leq 2 \mathbb{E}_S [\mathcal{R}(\ell, \mathcal{H}, S)]$$

### Amazing because:

- ▶  $\sup_{h \in \mathcal{H}} \{L_S(h) - L_{\mathcal{D}}(h)\}$  may be large for very unlikely  $S$
- ▶ But Rademacher complexity  $\mathcal{R}(\ell, \mathcal{H}, S)$  is small for all  $S$ !

### Consequence:

- ▶ Can measure complexity of  $\mathcal{H}$  conditional on  $S$
- ▶ So only restriction of  $\mathcal{H}$  to inputs  $\mathbf{X}_1, \dots, \mathbf{X}_m$  in  $S$  matters!

Step 2 : "Symmetrization"

$$\text{To show: } R(\ell, \mathcal{H}, S) = \frac{1}{m} \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^m G_i \ell(h, x_i, y_i) \right]$$

$$\mathbb{E}_S \left[ \sup_{h \in \mathcal{H}} L_S(h) - L_D(h) \right] \leq 2 \mathbb{E}_S [R(\ell, \mathcal{H}, S)]$$

Proof: Let  $S' = (x'_1, y'_1), \dots, (x'_m, y'_m)$  be independent sample.

$$\mathbb{E}_S \left[ \sup_h L_S(h) - L_D(h) \right] = \mathbb{E}_S \left[ \sup_h L_S(h) - \mathbb{E}_{S'} [L_{S'}(h)] \right]$$

$$= \mathbb{E}_S \left[ \sup_h \mathbb{E}_{S'} [L_S(h) - L_{S'}(h)] \right] \leq \mathbb{E}_{S, S'} \left[ \sup_h L_S(h) - L_{S'}(h) \right]$$

$$= \frac{1}{m} \mathbb{E}_{S, S'} \left[ \sup_h \sum_{i=1}^m \{ \ell(h, x_i, y_i) - \ell(h, x'_i, y'_i) \} \right]$$

Homogenize the two samples

N.B. If we swap any  $(\underline{x}_i)$  and  $(\underline{x}'_i)$  between  $S$  and  $S'$ , then distribution does not change.

Hence, for any  $\sigma_i \in \{-1, +1\}$ ,

$$\frac{1}{m} \mathbb{E}_{S, S'} \left[ \sup_h \sum_{i=1}^m \{\ell(h, \underline{x}_i, y_i) - \ell(h, \underline{x}'_i, y'_i)\} \right]$$

$$= \frac{1}{m} \mathbb{E}_{S, S'} \left[ \sup_h \sum_{i=1}^m \sigma_i \{\ell(h, \underline{x}_i, y_i) - \ell(h, \underline{x}'_i, y'_i)\} \right]$$

$$= \frac{1}{m} \mathbb{E}_{G, S, S'} \left[ \sup_h \sum_{i=1}^m \sigma_i \{\ell(h, \underline{x}_i, y_i) - \ell(h, \underline{x}'_i, y'_i)\} \right]$$

$$\leq \frac{1}{m} \mathbb{E}_{G, S, S'} \left[ \sup_h \sum_i \sigma_i \ell(h, \underline{x}_i, y_i) + \sup_h \sum_i -\sigma_i \ell(h, \underline{x}'_i, y'_i) \right]$$

(using that  $-\sigma$  has same distribution as  $\sigma$ )

$$= \frac{2}{m} \mathbb{E}_G \left[ \sup_h \sum_i \sigma_i \ell(h, \underline{x}_i, y_i) \right] = 2 \mathbb{E}_S \{ R(\ell, H, S) \} \quad \square$$

## Step 3: Bound the Rademacher Complexity

### Lemma

For any loss function  $\ell(h, \mathbf{X}, Y) = \tilde{\ell}(h(\mathbf{X}), Y)$  with range  $[0, 1]$  and any sample  $S$ :

$$\mathcal{R}(\ell, \mathcal{H}, S) \leq \sqrt{\frac{2 \ln |\mathcal{H}_S|}{m}} \leq \sqrt{\frac{2 \ln \tau_{\mathcal{H}}(m)}{m}}.$$

### Step 3

To show:  $\ell(h, x, y) = \tilde{\ell}(h(x), y) \in [0, 1]$

$$\text{For any } S: \mathcal{R}(\ell, \mathcal{H}, S) \leq \sqrt{\frac{2 \ln |\mathcal{H}_S|}{m}}$$

$$\begin{aligned}\text{Proof: } m\mathcal{R}(\ell, \mathcal{H}, S) &= \mathbb{E}_{\sigma} \left[ \sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i \cdot \tilde{\ell}(h(x_i), y_i) \right] \\ &= \mathbb{E}_{\sigma} \left[ \max_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i \cdot \hat{\ell}(h(x_i), y_i) \right]\end{aligned}$$

Let  $z_i(h) = \sigma_i \cdot \hat{\ell}(h(x_i), y_i) \in [-1, +1]$ . Then  $\mathbb{E}_{\sigma} [z_i(h)] = 0$

Hoeffding's Lemma (B.2 in Shai's e):

Suppose  $Z$  takes values in  $[a, b]$  and  $\mathbb{E}[Z] = 0$ . Then

$$\mathbb{E}[e^{\lambda Z}] \leq e^{\lambda^2(b-a)^2/8} \quad \text{for any } \lambda > 0.$$

$$\begin{aligned}
 m \cdot R(\ell, H, S) &= \mathbb{E}_{\sigma} \left[ \max_{h \in H_S} \sum_{i=1}^m z_i(h) \right] \\
 &= \frac{1}{\lambda} \mathbb{E}_{\sigma} \left[ \ln \max_{h \in H_S} e^{\sum_{i=1}^m \lambda z_i(h)} \right] \quad \text{for any } \lambda > 0 \\
 &\leq \frac{1}{\lambda} \mathbb{E}_{\sigma} \left[ \ln \sum_{h \in H_S} e^{\sum_{i=1}^m \lambda z_i(h)} \right]
 \end{aligned}$$

(Jensen's inequality)

$$\begin{aligned}
 &\leq \frac{1}{\lambda} \ln \left( \mathbb{E} \left[ \sum_{h \in H_S} e^{\sum_{i=1}^m \lambda z_i(h)} \right] \right) \\
 &= \frac{1}{\lambda} \ln \left( \sum_{h \in H_S} \prod_{i=1}^m \mathbb{E} [e^{\lambda z_i(h)}] \right) \\
 (\text{Hoeffding's Lemma}) \quad &\leq \frac{1}{\lambda} \ln \left( \sum_{h \in H_S} \prod_{i=1}^m e^{\lambda^2/2} \right) = \frac{1}{\lambda} \ln |H_S| + \lambda \frac{m}{2}
 \end{aligned}$$

$$\text{Take } \lambda = \sqrt{\frac{\ln |H_S|}{m/2}} : \quad = \sqrt{2m \ln |H_S|}$$

$$R(\ell, H, S) \leq \sqrt{\frac{2 \ln |H_S|}{m}} \quad \square$$

# Back to the Big Picture

**Part I: Growth function controls uniform convergence:**

$$\sup_{h \in \mathcal{H}} |L_S(h) - L_D(h)| \leq c \sqrt{\frac{\ln \tau_{\mathcal{H}}(m)}{m}} + c \sqrt{\frac{\ln(4/\delta)}{m}} \quad \text{with probability } \geq 1 - \delta$$

**Part II: VC-dimension controls growth function (Sauer's Lemma):**

$$\ln \tau_{\mathcal{H}}(m) \leq v \ln \left( \frac{em}{v} \right) \quad \text{for } m > v.$$

# Back to the Big Picture

**Part I: Growth function controls uniform convergence:**

$$\sup_{h \in \mathcal{H}} |L_S(h) - L_D(h)| \leq c \sqrt{\frac{\ln \tau_{\mathcal{H}}(m)}{m}} + c \sqrt{\frac{\ln(4/\delta)}{m}} \quad \text{with probability } \geq 1 - \delta$$

**Part II: VC-dimension controls growth function (Sauer's Lemma):**

$$\ln \tau_{\mathcal{H}}(m) \leq v \ln \left( \frac{em}{v} \right) \quad \text{for } m > v.$$

---

For  $m > v$ :

$$\sup_{h \in \mathcal{H}} |L_S(h) - L_D(h)| \leq c \sqrt{\frac{v \ln \left( \frac{em}{v} \right)}{m}} + c \sqrt{\frac{\ln(4/\delta)}{m}} \quad \text{with probability } \geq 1 - \delta$$

- ▶ Remaining: find lower bound on  $m$  s.t. bound is at most  $\epsilon$ .