

# Generalization Bounds

CMPUT 296: Basics of Machine Learning

Textbook Ch.12

# Logistics

- Thought questions #3 marked (*finally...*)
- **Thought questions #4** due **TODAY at 11:59pm** (Mountain time)
  - Anything between logistic regression and generalization bounds is fair game
- **Last class** will be on **Tuesday (Dec 1)**
  - It will be a review class for the final exam
- **Assignment #3** due **Thursday Dec 3**
  - But there is no class on Dec 3
- **FINAL EXAM** will be **Friday Dec 18**

# Outline

1. Recap & Logistics
2. McDiarmid's Inequality
3. Proof of Generalization Bound for Binary Classification

# Recap: Rademacher Complexity

The **empirical Rademacher complexity** of  $\mathcal{F}$  with respect to  $\mathcal{D}$  is

$$\hat{R}_{\mathcal{D}}(\mathcal{F}) = \mathbb{E} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right]$$

where

- $\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid 1 \leq i \leq n\}$  is a dataset
- $\sigma = (\sigma_1, \sigma_2, \dots, \sigma_n)$  is a vector of  $n$  random variables, with  $\sigma_i \stackrel{i.i.d}{\sim} \text{Uniform}\{-1, +1\}$  ("Rademacher variables")
- $\mathcal{F}$  is a hypothesis class

The **Rademacher complexity** of a hypothesis class  $\mathcal{F}$  is the expected **empirical Rademacher complexity** over all datasets  $\mathcal{D}$  (of size  $n$ ):

$$R_n(\mathcal{F}) = \mathbb{E} \left[ \hat{R}_{\mathcal{D}}(\mathcal{F}) \right]$$

# Recap: Uniform Generalization Bound for Binary Classification

**Theorem:**

Let  $\mathcal{F}$  be a family of binary classification functions taking values in  $\{-1, +1\}$ , and let **cost** be the 0-1 classification cost.

Then for every  $f \in \mathcal{F}$ , and every  $\delta > 0$ ,

$$C(f) \leq \hat{C}(f) + R_n(\mathcal{F}) + \sqrt{\frac{\log 1/\delta}{2n}}$$

with probability  $1 - \delta$ .

# Implication: Training Error vs. Test Error

$$C(f) \leq \hat{C}(f) + R_n(\mathcal{F}) + \sqrt{\frac{\log 1/\delta}{2n}}, \quad \forall f \in \mathcal{F}$$

1. The **more data** you have, the closer your training error will be to the true generalization error
2. The **simpler** your hypothesis class, the closer your training error will be to the true error
  - Notice that this bound applies to **all**  $f \in \mathcal{F}$

**Question:** Why does this bound not go to zero for infinite data?

# McDiarmid's Inequality

Let  $S = (X_1, \dots, X_n) \in \mathcal{X}^n$  be a vector of  $n \geq 1$  independent random variables.

If there exist constants  $c_1, \dots, c_n$  such that  $f : \mathcal{X}^n \rightarrow \mathbb{R}$  satisfies

$$\left| f(x_1, \dots, x_k, \dots, x_n) - f(x_1, \dots, x'_k, \dots, x_n) \right| \leq c_k$$

for every  $1 \leq k \leq n$  and any values  $x_1, \dots, x_n, x'_k \in \mathcal{X}$ , then for all  $\epsilon > 0$ ,

1.  $\Pr [f(S) - \mathbb{E}[f(S)] \geq \epsilon] \leq \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^n c_i^2}\right)$ , and

2.  $\Pr [f(S) - \mathbb{E}[f(S)] \leq -\epsilon] \leq \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^n c_i^2}\right)$ .

# Proof: Generalization Bound for Binary Classification

Update our notation for empirical cost to indicate the dataset:

$$\hat{C}_{\mathcal{D}}(f) = \frac{1}{n} \sum_{i=1}^n \text{cost}(f(\mathbf{x}_i), y_i)$$

Define  $\Phi(\mathcal{D})$  as the maximum difference between the true cost of a predictor and the empirical cost on dataset  $\mathcal{D}$ :

$$\Phi(\mathcal{D}) \doteq \max_{f \in \mathcal{F}} C(f) - \hat{C}_{\mathcal{D}}(f)$$

(A bit sloppy; should really use  $\sup_{f \in \mathcal{F}}$ )



# Proof (2): Bounding $\left| \Phi(\mathcal{D}) - \Phi(\mathcal{D}') \right|$

Let  $\mathcal{D}$  and  $\mathcal{D}'$  be two datasets of  $n$  observations that differ in **exactly one** datapoint:  $(x_k, y_k) \neq (x'_k, y'_k)$ , and  $(x_i, y_i) = (x'_i, y'_i)$  for all  $i \neq k$

$$\begin{aligned} \Phi(\mathcal{D}) - \Phi(\mathcal{D}') &= \left( \max_{f \in \mathcal{F}} C(f) - \hat{C}_{\mathcal{D}}(f) \right) - \left( \max_{g \in \mathcal{F}} C(g) - \hat{C}_{\mathcal{D}'}(g) \right) \\ &\leq \max_{f \in \mathcal{F}} \left( C(f) - \hat{C}_{\mathcal{D}}(f) \right) - \left( C(f) - \hat{C}_{\mathcal{D}'}(f) \right) \\ &= \max_{f \in \mathcal{F}} C(f) - \hat{C}_{\mathcal{D}}(f) - C(f) + \hat{C}_{\mathcal{D}'}(f) \\ &= \max_{f \in \mathcal{F}} \hat{C}_{\mathcal{D}'}(f) - \hat{C}_{\mathcal{D}}(f) \end{aligned}$$

Proof (3):

Bounding  $\left| \Phi(\mathcal{D}) - \Phi(\mathcal{D}') \right|$

$$\begin{aligned} \Phi(\mathcal{D}) - \Phi(\mathcal{D}') &\leq \max_{f \in \mathcal{F}} \hat{C}_{\mathcal{D}'}(f) - \hat{C}_{\mathcal{D}}(f) \\ &= \max_{f \in \mathcal{F}} \sum_{i=1}^n \frac{\text{cost}(f(x_i), y_i) - \text{cost}(f(x'_i), y'_i)}{n} \\ &= \max_{f \in \mathcal{F}} \frac{\text{cost}(f(x_k), y_k) - \text{cost}(f(x'_k), y'_k)}{n} = \frac{1}{n} \end{aligned}$$

By an identical argument,  $\Phi(\mathcal{D}') - \Phi(\mathcal{D}) \leq 1/n$ . So

$$\left| \Phi(\mathcal{D}) - \Phi(\mathcal{D}') \right| \leq \frac{1}{n}$$

# Proof (4): Bounding $\Phi(\mathcal{D})$

$$|\Phi(\mathcal{D}) - \Phi(\mathcal{D}')| \leq \frac{1}{n}$$

Apply McDiarmid's Inequality:

$$\Pr [\Phi(\mathcal{D}) - \mathbb{E}[\Phi(\mathcal{D})] \geq \epsilon] \leq \exp\left(\frac{-2\epsilon^2}{\sum_{i=1}^n 1/n^2}\right)$$

$$= \exp(-2n\epsilon^2) = \delta$$

$$\sum_{i=1}^n \frac{1}{n^2} = n \frac{1}{n^2} = \frac{1}{n}$$

$$\delta = \exp(-2n\epsilon^2)$$

$$\iff \log \delta = -2n\epsilon^2$$

$$\iff \frac{-\log \delta}{2n} = \epsilon^2$$

$$\iff \epsilon = \sqrt{\frac{\log 1/\delta}{2n}}$$

$$\Pr \left[ \Phi(\mathcal{D}) - \mathbb{E}[\Phi(\mathcal{D})] \geq \sqrt{\frac{\log 1/\delta}{2n}} \right] \leq \delta$$

$$\Pr \left[ \Phi(\mathcal{D}) \leq \mathbb{E}[\Phi(\mathcal{D})] + \sqrt{\frac{\log 1/\delta}{2n}} \right] \geq 1 - \delta$$

# Proof (5): Bounding $\mathbb{E}[\Phi(\mathcal{D})]$

With probability  $1 - \delta$ , we have  $\Phi(\mathcal{D}) \leq \mathbb{E}[\Phi(\mathcal{D})] + \sqrt{\frac{\log 1/\delta}{2n}}$ .

$$\mathbb{E}[\Phi(\mathcal{D})] = \mathbb{E}_{\mathcal{D}} \left[ \max_{f \in \mathcal{F}} C(f) - \hat{C}_{\mathcal{D}}(f) \right]$$

$$= \mathbb{E}_{\mathcal{D}} \left[ \max_{f \in \mathcal{F}} \mathbb{E}_{\mathcal{D}'} \left[ \hat{C}_{\mathcal{D}'}(f) - \hat{C}_{\mathcal{D}}(f) \right] \right]$$

$$\leq \mathbb{E}_{\mathcal{D}, \mathcal{D}'} \left[ \max_{f \in \mathcal{F}} \hat{C}_{\mathcal{D}'}(f) - \hat{C}_{\mathcal{D}}(f) \right]$$

$$= \mathbb{E}_{\mathcal{D}, \mathcal{D}'} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (\text{cost}(f(x_i), y_i) - \text{cost}(f(x'_i), y'_i)) \right]$$

## Jensen's inequality

For any **convex** function

$$f: \mathcal{X} \rightarrow \mathbb{R},$$

$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$$

(max is a convex function)

# Proof (6): Bounding $\mathbb{E}[\Phi(\mathcal{D})]$

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[\Phi(\mathcal{D})] &\leq \mathbb{E}_{\mathcal{D}, \mathcal{D}'} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n (\text{cost}(f(x_i), y_i) - \text{cost}(f(x'_i), y'_i)) \right] \\ &= \mathbb{E}_{\mathcal{D}, \mathcal{D}', \sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (\text{cost}(f(x_i), y_i) - \text{cost}(f(x'_i), y'_i)) \right] \\ &\leq \mathbb{E}_{\mathcal{D}, \mathcal{D}', \sigma} \left[ \left( \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \text{cost}(f(x_i), y_i) \right) + \left( \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n -\sigma_i \text{cost}(f(x'_i), y'_i) \right) \right] \\ &= \mathbb{E}_{\mathcal{D}, \sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \text{cost}(f(x_i), y_i) \right] + \mathbb{E}_{\mathcal{D}', \sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n -\sigma_i \text{cost}(f(x'_i), y'_i) \right]\end{aligned}$$

# Proof (7): Bounding $\mathbb{E}[\Phi(\mathcal{D})]$

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[\Phi(\mathcal{D})] &\leq \mathbb{E}_{\mathcal{D},\sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \text{cost}(f(x_i), y_i) \right] + \mathbb{E}_{\mathcal{D}',\sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n -\sigma_i \text{cost}(f(x'_i), y'_i) \right] \\ &= 2\mathbb{E}_{\mathcal{D},\sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \text{cost}(f(x_i), y_i) \right] \\ &= 2\mathbb{E}_{\mathcal{D},\sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i (1 - I_{f(x_i)=y_i}) \right] \\ &= 2\mathbb{E}_{\mathcal{D},\sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - y_i f(x_i)}{2} \right] \\ &= 2\frac{1}{2}\mathbb{E}_{\mathcal{D},\sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i \right] + \frac{1}{2}2\mathbb{E}_{\mathcal{D},\sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i - y_i f(x_i) \right] \\ &= \mathbb{E}_{\mathcal{D},\sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right] = \mathbb{E}_{\mathcal{D}} \left[ \mathbb{E}_{\sigma} \left[ \max_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \sigma_i f(x_i) \right] \right] = \mathbb{E}_{\mathcal{D}} \left[ \hat{R}_{\mathcal{D}}(\mathcal{F}) \right] = \boxed{R_n(\mathcal{F})}\end{aligned}$$

# Proof: Putting it all together

1. With probability  $1 - \delta$ , we have  $\Phi(\mathcal{D}) \leq \mathbb{E}[\Phi(\mathcal{D})] + \sqrt{\frac{\log 1/\delta}{2n}}$
2.  $\mathbb{E}_{\mathcal{D}}[\Phi(\mathcal{D})] \leq R_n(\mathcal{F})$
3. With probability  $1 - \delta$ ,

$$\Phi(\mathcal{D}) \leq R_n(\mathcal{F}) + \sqrt{\frac{\log 1/\delta}{2n}}$$

$$\max_{f \in \mathcal{F}} C(f) - \hat{C}_{\mathcal{D}}(f) \leq R_n(\mathcal{F}) + \sqrt{\frac{\log 1/\delta}{2n}}$$

$$C(f) \leq \hat{C}_{\mathcal{D}}(f) + R_n(\mathcal{F}) + \sqrt{\frac{\log 1/\delta}{2n}}, \quad \forall f \in \mathcal{F}. \blacksquare$$

# Summary

- **McDiarmid's Inequality** is a generalization of Hoeffding's Inequality for "stable enough" functions
- We can use McDiarmid's Inequality to prove upper bounds on generalization performance such as (for binary classification)

$$C(f) \leq \hat{C}_{\mathcal{D}}(f) + R_n(\mathcal{F}) + \sqrt{\frac{\log 1/\delta}{2n}}, \quad \forall f \in \mathcal{F}$$

- Bounds of this kind shed light on the relationship between **training error** and **generalization error**
  - The more flexible your hypothesis class (larger  $R_n(\mathcal{F})$ ), the bigger the difference can be
  - This difference is the origin of **overfitting**