



## Classification: Rademacher Complexity

Machine Learning: Jordan Boyd-Graber  
University of Colorado Boulder

LECTURE 6

Slides adapted from Rob Schapire

## Setup

---

Nothing new ...

- Samples  $S = ((x_1, y_1), \dots, (x_m, y_m))$
- Labels  $y_i = \{-1, +1\}$
- Hypothesis  $h: X \rightarrow \{-1, +1\}$
- Training error:  $\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i]$

## An alternative derivation of training error

---

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

(2)

(3)

(4)

## An alternative derivation of training error

---

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

$$= \frac{1}{m} \sum_i^m \begin{cases} 1 & \text{if } (h(x_i, y_i) == (1, -1) \text{ or } (-1, 1)) \\ 0 & \text{if } (h(x_i, y_i) == (1, 1) \text{ or } (-1, -1)) \end{cases} \quad (2)$$

$$(3)$$

$$(4)$$

## An alternative derivation of training error

---

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

$$= \frac{1}{m} \sum_i^m \begin{cases} 1 & \text{if } (h(x_i), y_i) == (1, -1) \text{ or } (-1, 1) \\ 0 & \text{if } (h(x_i), y_i) == (1, 1) \text{ or } (-1, -1) \end{cases} \quad (2)$$

$$= \frac{1}{m} \sum_i^m \frac{1 - y_i h(x_i)}{2} \quad (3)$$

$$(4)$$

## An alternative derivation of training error

---

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

$$= \frac{1}{m} \sum_i^m \begin{cases} 1 & \text{if } (h(x_i), y_i) == (1, -1) \text{ or } (-1, 1) \\ 0 & \text{if } (h(x_i), y_i) == (1, 1) \text{ or } (-1, -1) \end{cases} \quad (2)$$

$$= \frac{1}{m} \sum_i^m \frac{1 - y_i h(x_i)}{2} \quad (3)$$

$$= \frac{1}{2} - \frac{1}{2m} \sum_i^m y_i h(x_i) \quad (4)$$

## An alternative derivation of training error

---

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

$$= \frac{1}{m} \sum_i^m \begin{cases} 1 & \text{if } (h(x_i), y_i) == (1, -1) \text{ or } (-1, 1) \\ 0 & \text{if } (h(x_i), y_i) == (1, 1) \text{ or } (-1, -1) \end{cases} \quad (2)$$

$$= \frac{1}{m} \sum_i^m \frac{1 - y_i h(x_i)}{2} \quad (3)$$

$$= \frac{1}{2} - \frac{1}{2m} \sum_i^m y_i h(x_i) \quad (4)$$

Correlation between predictions and labels

## An alternative derivation of training error

---

$$\hat{R}(h) = \frac{1}{m} \sum_i^m \mathbb{1}[h(x_i) \neq y_i] \quad (1)$$

$$= \frac{1}{m} \sum_i^m \begin{cases} 1 & \text{if } (h(x_i), y_i) == (1, -1) \text{ or } (-1, 1) \\ 0 & \text{if } (h(x_i), y_i) == (1, 1) \text{ or } (-1, -1) \end{cases} \quad (2)$$

$$= \frac{1}{m} \sum_i^m \frac{1 - y_i h(x_i)}{2} \quad (3)$$

$$= \frac{1}{2} - \frac{1}{2m} \sum_i^m y_i h(x_i) \quad (4)$$

Minimizing training error is thus equivalent to maximizing correlation

$$\arg \max_h \frac{1}{m} \sum_i^m y_i h(x_i) \quad (5)$$



## Playing with Correlation

---

Imagine where we replace true labels with *Rademacher random variables*

$$\sigma_i = \begin{cases} +1 & \text{with prob .5} \\ -1 & \text{with prob .5} \end{cases} \quad (6)$$

## Playing with Correlation

---

Imagine where we replace true labels with *Rademacher random variables*

$$\sigma_i = \begin{cases} +1 & \text{with prob .5} \\ -1 & \text{with prob .5} \end{cases} \quad (6)$$

This gives us Rademacher correlation—what's the best that a random classifier could do?

$$\hat{\mathcal{R}}_S(H) \equiv \mathbb{E}_\sigma \left[ \max_{h \in H} \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right] \quad (7)$$

## Playing with Correlation

---

Imagine where we replace true labels with *Rademacher random variables*

$$\sigma_i = \begin{cases} +1 & \text{with prob .5} \\ -1 & \text{with prob .5} \end{cases} \quad (6)$$

This gives us Rademacher correlation—what's the best that a random classifier could do?

$$\hat{\mathcal{R}}_{\mathbf{S}}(H) \equiv \mathbb{E}_{\sigma} \left[ \max_{h \in H} \frac{1}{m} \sum_i \sigma_i h(\mathbf{x}_i) \right] \quad (7)$$

Note: Empirical Rademacher complexity is with respect to a sample.

## Rademacher Extrema

---

- What are the maximum values of Rademacher correlation?

## Rademacher Extrema

---

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$|H| = 2^m$$

## Rademacher Extrema

---

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$\mathbb{E}_{\sigma} \left[ \max_{h \in H} \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right]$$

$$|H| = 2^m$$

## Rademacher Extrema

---

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$\mathbb{E}_{\sigma} \left[ \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right]$$

$$|H| = 2^m$$

## Rademacher Extrema

---

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$\mathbb{E}_{\sigma} \left[ \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right] = 0$$

$$|H| = 2^m$$



## Rademacher Extrema

---

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$\mathbb{E}_{\sigma} \left[ \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right] = 0$$

$$|H| = 2^m$$

$$\mathbb{E}_{\sigma} \left[ \max_{h \in H} \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right]$$

## Rademacher Extrema

---

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$\mathbb{E}_{\sigma} \left[ \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right] = 0$$

$$|H| = 2^m$$

$$\frac{m}{m} = 1$$

## Rademacher Extrema

---

- What are the maximum values of Rademacher correlation?

$$|H| = 1$$

$$\mathbb{E}_{\sigma} \left[ \frac{1}{m} \sum_i^m \sigma_i h(x_i) \right] = 0$$

$$|H| = 2^m$$

$$\frac{m}{m} = 1$$

- Rademacher correlation is larger for more complicated hypothesis space.
- What if you're right for stupid reasons?

## Generalizing Rademacher Complexity

We can generalize Rademacher complexity to consider all sets of a particular size.

$$\mathcal{R}_m(H) = \mathbb{E}_{S \sim D^m} [\hat{\mathcal{R}}_S(H)] \quad (8)$$

### Theorem

**Convergence Bounds** Let  $F$  be a family of functions mapping from  $Z$  to  $[0, 1]$ , and let sample  $S = (z_1, \dots, z_m)$  where  $z_i \sim D$  for some distribution  $D$  over  $Z$ . Define  $\mathbb{E}[f] \equiv \mathbb{E}_{z \sim D}[f(z)]$  and  $\hat{\mathbb{E}}_S[f] \equiv \frac{1}{m} \sum_{i=1}^m f(z_i)$ . With probability greater than  $1 - \delta$  for all  $f \in F$ :

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\mathcal{R}_m(F) + \mathcal{O}\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (9)$$

## Generalizing Rademacher Complexity

We can generalize Rademacher complexity to consider all sets of a particular size.

$$\mathcal{R}_m(H) = \mathbb{E}_{S \sim D^m} [\hat{\mathcal{R}}_S(H)] \quad (8)$$

### Theorem

**Convergence Bounds** Let  $F$  be a family of functions mapping from  $Z$  to  $[0, 1]$ , and let sample  $S = (z_1, \dots, z_m)$  where  $z_i \sim D$  for some distribution  $D$  over  $Z$ . Define  $\mathbb{E}[f] \equiv \mathbb{E}_{z \sim D}[f(z)]$  and  $\hat{\mathbb{E}}_S[f] \equiv \frac{1}{m} \sum_{i=1}^m f(z_i)$ . With probability greater than  $1 - \delta$  for all  $f \in F$ :

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\mathcal{R}_m(F) + \mathcal{O}\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (9)$$

## Generalizing Rademacher Complexity

We can generalize Rademacher complexity to consider all sets of a particular size.

$$\mathcal{R}_m(H) = \mathbb{E}_{S \sim D^m} [\hat{\mathcal{R}}_S(H)] \quad (8)$$

### Theorem

**Convergence Bounds** Let  $F$  be a family of functions mapping from  $Z$  to  $[0, 1]$ , and let sample  $S = (z_1, \dots, z_m)$  where  $z_i \sim D$  for some distribution  $D$  over  $Z$ . Define  $\mathbb{E}[f] \equiv \mathbb{E}_{z \sim D}[f(z)]$  and  $\hat{\mathbb{E}}_S[f] \equiv \frac{1}{m} \sum_{i=1}^m f(z_i)$ . With probability greater than  $1 - \delta$  for all  $f \in F$ :

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\mathcal{R}_m(F) + \mathcal{O}\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (9)$$

## Aside: McDiarmid's Inequality

---

If we have a function:

$$|f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x'_i, \dots, x_m)| \leq c_i \quad (10)$$

then:

$$\Pr[f(x_1, \dots, x_m) \geq \mathbb{E}[f(X_1, \dots, X_m)] + \epsilon] \leq \exp \left\{ \frac{-2\epsilon^2}{\sum_i^m c_i^2} \right\} \quad (11)$$

## Aside: McDiarmid's Inequality

---

If we have a function:

$$|f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x'_i, \dots, x_m)| \leq c_i \quad (10)$$

then:

$$\Pr[f(x_1, \dots, x_m) \geq \mathbb{E}[f(X_1, \dots, X_m)] + \epsilon] \leq \exp \left\{ \frac{-2\epsilon^2}{\sum_i^m c_i^2} \right\} \quad (11)$$

Proof in the back of the textbook (requires Martingales).



## Aside: McDiarmid's Inequality

---

If we have a function:

$$|f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x'_i, \dots, x_m)| \leq c_i \quad (10)$$

then:

$$\Pr[f(x_1, \dots, x_m) \geq \mathbb{E}[f(X_1, \dots, X_m)] + \epsilon] \leq \exp \left\{ \frac{-2\epsilon^2}{\sum_i^m c_i^2} \right\} \quad (11)$$

Proof in the back of the textbook (requires Martingales).

What function do we care about for Rademacher complexity? Let's define

$$\Phi(S) = \sup \left( \mathbb{E}[f] - \hat{\mathbb{E}}_S[f] \right) = \sup \left( \mathbb{E}[f] - \frac{1}{m} \sum_i f(z_i) \right) \quad (12)$$

## Step 1: Bounding divergence from true Expectation

---

### Lemma

**Moving to Expectation** *With probability at least  $1 - \delta$ ,*

$$\Phi(S) \leq \mathbb{E}_S[\Phi(S)] + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}}$$

Since  $f(z_1) \in [0, 1]$ , changing any  $z_i$  to  $z'_i$  in the training set will change  $\frac{1}{m} \sum_i f(z_i)$  by at most  $\frac{1}{m}$ , so we can apply McDiarmid's inequality.

## Step 2: Comparing two different empirical expectations

---

Define a ghost sample  $S' = (z'_1, \dots, z'_m) \sim D$ . How much can two samples from the same distribution vary?

### Lemma

#### Two Different Samples

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[ \sup_f (\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (13)$$

(14)

## Step 2: Comparing two different empirical expectations

Define a ghost sample  $S' = (z'_1, \dots, z'_m) \sim D$ . How much can two samples from the same distribution vary?

### Lemma

#### Two Different Samples

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[ \sup_f (\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (13)$$

$$= \mathbb{E}_S \left[ \sup_{f \in F} (\mathbb{E}_{S'} [\hat{\mathbb{E}}_{S'}[f]] - \hat{\mathbb{E}}_S[f]) \right] \quad (14)$$

$$(15)$$

The expectation is equal to the expectation of the empirical expectation of all sets  $S'$

## Step 2: Comparing two different empirical expectations

Define a ghost sample  $S' = (z'_1, \dots, z'_m) \sim D$ . How much can two samples from the same distribution vary?

### Lemma

#### Two Different Samples

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[ \sup_f (\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (13)$$

$$= \mathbb{E}_S \left[ \sup_{f \in F} (\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[f]] - \hat{\mathbb{E}}_S[f]) \right] \quad (14)$$

$$= \mathbb{E}_S \left[ \sup_{f \in F} (\mathbb{E}_{S'}[\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f]]) \right] \quad (15)$$

$$(16)$$

$S$  and  $S'$  are distinct random variables, so we can move inside the expectation

## Step 2: Comparing two different empirical expectations

Define a ghost sample  $S' = (z'_1, \dots, z'_m) \sim D$ . How much can two samples from the same distribution vary?

### Lemma

#### Two Different Samples

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[ \sup_f (\mathbb{E}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (13)$$

$$= \mathbb{E}_S \left[ \sup_{f \in F} (\mathbb{E}_{S'} [\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (14)$$

$$\leq \mathbb{E}_{S, S'} \left[ \sup_f (\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f]) \right] \quad (15)$$

The expectation of a max over some function is at least the max of that expectation over that function

### Step 3: Adding in Rademacher Variables

---

From  $S, S'$  we'll create  $T, T'$  by swapping elements between  $S$  and  $S'$  with probability .5. This is still iid from  $D$ . They have the same distribution:

$$\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f] \sim \hat{\mathbb{E}}_{T'}[f] - \hat{\mathbb{E}}_T[f] \quad (16)$$

### Step 3: Adding in Rademacher Variables

---

From  $S, S'$  we'll create  $T, T'$  by swapping elements between  $S$  and  $S'$  with probability .5. This is still iid from  $D$ . They have the same distribution:

$$\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f] \sim \hat{\mathbb{E}}_{T'}[f] - \hat{\mathbb{E}}_T[f] \quad (16)$$

Let's introduce  $\sigma_i$ :

$$\hat{\mathbb{E}}_{T'}[f] - \hat{\mathbb{E}}_T[f] = \frac{1}{m} \begin{cases} f(z_i) - f(z'_i) & \text{with prob .5} \\ f(z'_i) - f(z_i) & \text{with prob .5} \end{cases} \quad (17)$$

$$= \frac{1}{m} \sum_i \sigma_i (f(z'_i) - f(z_i)) \quad (18)$$



### Step 3: Adding in Rademacher Variables

---

From  $S, S'$  we'll create  $T, T'$  by swapping elements between  $S$  and  $S'$  with probability .5. This is still iid from  $D$ . They have the same distribution:

$$\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f] \sim \hat{\mathbb{E}}_{T'}[f] - \hat{\mathbb{E}}_T[f] \quad (16)$$

Let's introduce  $\sigma_i$ :

$$\hat{\mathbb{E}}_{T'}[f] - \hat{\mathbb{E}}_T[f] = \frac{1}{m} \begin{cases} f(z_i) - f(z'_i) & \text{with prob .5} \\ f(z'_i) - f(z_i) & \text{with prob .5} \end{cases} \quad (17)$$

$$= \frac{1}{m} \sum_i \sigma_i (f(z'_i) - f(z_i)) \quad (18)$$

Thus:

$$\mathbb{E}_{S,S'} \left[ \sup_{f \in F} \left( \hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f] \right) \right] = \mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \left( \sum_i \sigma_i (f(z'_i) - f(z_i)) \right) \right].$$

## Step 4: Making These Rademacher Complexities

---

Before, we had  $\mathbb{E}_{S, S', \sigma} \left[ \sup_{f \in F} \sum_i \sigma_i (f(z'_i) - f(z_i)) \right]$

## Step 4: Making These Rademacher Complexities

---

Before, we had  $\mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \sum_i \sigma_i (f(z'_i) - f(z_i)) \right]$

$$\leq \mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \sum_i \sigma_i f(z'_i) + \sup_{f \in F} \sum_i (-\sigma_i) f(z_i) \right] \quad (19)$$

(20)

Taking the sup jointly must be less than or equal the individual sup.

## Step 4: Making These Rademacher Complexities

---

Before, we had  $\mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \sum_i \sigma_i (f(z'_i) - f(z_i)) \right]$

$$\leq \mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \sum_i \sigma_i f(z'_i) + \sup_{f \in F} \sum_i (-\sigma_i) f(z_i) \right] \quad (19)$$

$$\leq \mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \sum_i \sigma_i f(z'_i) \right] + \mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \sum_i (-\sigma_i) f(z_i) \right] \quad (20)$$

$$(21)$$

Linearity

## Step 4: Making These Rademacher Complexities

---

Before, we had  $\mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \sum_i \sigma_i (f(z'_i) - f(z_i)) \right]$

$$\leq \mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \sum_i \sigma_i f(z'_i) + \sup_{f \in F} \sum_i (-\sigma_i) f(z_i) \right] \quad (19)$$

$$\leq \mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \sum_i \sigma_i f(z'_i) \right] + \mathbb{E}_{S,S',\sigma} \left[ \sup_{f \in F} \sum_i (-\sigma_i) f(z_i) \right] \quad (20)$$

$$= \mathcal{R}_m(F) + \mathcal{R}_m(F) \quad (21)$$

Definition

## Putting the Pieces Together

---

With probability  $\geq 1 - \delta$ :

$$\Phi(S) \leq \mathbb{E}_S[\Phi(S)] + \sqrt{\frac{\ln \frac{q}{\delta}}{2m}} \quad (22)$$

Step 1

## Putting the Pieces Together

---

With probability  $\geq 1 - \delta$ :

$$\sup_f \left( \mathbb{E}[f] - \hat{\mathbb{E}}_S[h] \right) \leq \mathbb{E}_S[\Phi(S)] + \sqrt{\frac{\ln \frac{q}{\delta}}{2m}} \quad (22)$$

Definition of  $\Phi$

## Putting the Pieces Together

---

With probability  $\geq 1 - \delta$ :

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[h] \leq \mathbb{E}_S[\Phi(S)] + \sqrt{\frac{\ln \frac{q}{\delta}}{2m}} \quad (22)$$

Drop the sup, still true



## Putting the Pieces Together

---

With probability  $\geq 1 - \delta$ :

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[h] \leq \mathbb{E}_{S,S'} \left[ \sup_f (\hat{\mathbb{E}}_{S'}[f] - \hat{\mathbb{E}}_S[f]) \right] + \sqrt{\frac{\ln \frac{q}{\delta}}{2m}} \quad (22)$$

Step 2

## Putting the Pieces Together

---

With probability  $\geq 1 - \delta$ :

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[h] \leq \mathbb{E}_{S, S', \sigma} \left[ \sup_{f \in F} \left( \sum_i \sigma_i (f(z'_i) - f(z_i)) \right) \right] + \sqrt{\frac{\ln \frac{q}{\delta}}{2m}} \quad (22)$$

Step 3

## Putting the Pieces Together

---

With probability  $\geq 1 - \delta$ :

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[h] \leq 2\mathcal{R}_m(F) + \sqrt{\frac{\ln \frac{q}{\delta}}{2m}} \quad (22)$$

Step 4

## Putting the Pieces Together

---

With probability  $\geq 1 - \delta$ :

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[h] \leq 2\mathcal{R}_m(F) + \sqrt{\frac{\ln \frac{q}{\delta}}{2m}} \quad (22)$$

Recall that  $\hat{\mathcal{R}}_S(F) \equiv \mathbb{E}_\sigma \left[ \sup_f \frac{1}{m} \sum_i \sigma_i f(z_i) \right]$ , so we apply McDiarmid's inequality again (because  $f \in [0, 1]$ ):

$$\hat{\mathcal{R}}_S(F) \leq \mathcal{R}_m(F) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (23)$$

## Putting the Pieces Together

---

With probability  $\geq 1 - \delta$ :

$$\mathbb{E}[f] - \hat{\mathbb{E}}_S[h] \leq 2\mathcal{R}_m(F) + \sqrt{\frac{\ln \frac{q}{\delta}}{2m}} \quad (22)$$

Recall that  $\hat{\mathcal{R}}_S(F) \equiv \mathbb{E}_\sigma \left[ \sup_f \frac{1}{m} \sum_i \sigma_i f(z_i) \right]$ , so we apply McDiarmid's inequality again (because  $f \in [0, 1]$ ):

$$\hat{\mathcal{R}}_S(F) \leq \mathcal{R}_m(F) + \sqrt{\frac{\ln \frac{1}{\delta}}{2m}} \quad (23)$$

Putting the two together:

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\mathcal{R}_m(F) + \mathcal{O} \left( \sqrt{\frac{\ln \frac{1}{\delta}}{m}} \right) \quad (24)$$

## What about hypothesis classes?

---

Define:

$$Z \equiv X \times \{-1, +1\} \quad (25)$$

$$f_h(x, y) \equiv \mathbb{1}[h(x) \neq y] \quad (26)$$

$$F_H \equiv \{f_h : h \in H\} \quad (27)$$

## What about hypothesis classes?

---

Define:

$$Z \equiv X \times \{-1, +1\} \quad (25)$$

$$f_h(x, y) \equiv \mathbb{1} [h(x) \neq y] \quad (26)$$

$$F_H \equiv \{f_h : h \in H\} \quad (27)$$

We can use this to create expressions for generalization and empirical error:

$$R(h) = \mathbb{E}_{(x,y) \sim D} [\mathbb{1} [h(x) \neq y]] = \mathbb{E} [f_h] \quad (28)$$

$$\hat{R}(h) = \frac{1}{m} \sum_i \mathbb{1} [h(x_i) \neq y] = \hat{\mathbb{E}}_S [f_h] \quad (29)$$

## What about hypothesis classes?

---

Define:

$$Z \equiv X \times \{-1, +1\} \quad (25)$$

$$f_h(x, y) \equiv \mathbb{1} [h(x) \neq y] \quad (26)$$

$$F_H \equiv \{f_h : h \in H\} \quad (27)$$

We can use this to create expressions for generalization and empirical error:

$$R(h) = \mathbb{E}_{(x,y) \sim D} [\mathbb{1} [h(x) \neq y]] = \mathbb{E} [f_h] \quad (28)$$

$$\hat{R}(h) = \frac{1}{m} \sum_i \mathbb{1} [h(x_i) \neq y] = \hat{\mathbb{E}}_S [f_h] \quad (29)$$

We can plug this into our theorem!



## Generalization bounds

---

- We started with expectations

$$\mathbb{E}[f] \leq \hat{\mathbb{E}}_S[f] + 2\hat{\mathcal{R}}_S(F) + \mathcal{O}\left(\sqrt{\frac{\ln \frac{1}{\delta}}{m}}\right) \quad (30)$$

- We also had our definition of the generalization and empirical error:

$$R(h) = \mathbb{E}_{(x,y) \sim D} [\mathbb{1}[h(x) \neq y]] = \mathbb{E}[f_h] \quad \hat{R}(h) = \frac{1}{m} \sum_i \mathbb{1}[h(x_i) \neq y] = \hat{\mathbb{E}}_S[f_h]$$

- Combined with the previous result:

$$\hat{\mathcal{R}}_S(F_H) = \frac{1}{2} \hat{\mathcal{R}}_S(H) \quad (31)$$

- All together:

$$R(h) \leq \hat{R}(h) + \mathcal{R}_m(H) + \mathcal{O}\left(\sqrt{\frac{\log \frac{1}{\delta}}{m}}\right) \quad (32)$$

## Wrapup

---

- Interaction of data, complexity, and accuracy
- Still very theoretical
- Next time: How to evaluate generalizability of specific hypothesis classes