

# Machine Learning Theory Lecture 3: Rademacher Complexity

Richard Xu

April 30, 2022

## 1 Definition

let each of  $S = \{Z_i\}$  be distributed from a data distribution  $\mathcal{D}$

$$\text{Rad}_n(\mathcal{H}) = \mathbb{E}_S \left[ \mathbb{E}_{\vec{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{\sum_{i=1}^n \sigma_i h(Z_i)}{n} \right] \right] \quad (1)$$

1. In words, We sample  $n$  data  $\{Z_i\}_{i=1}^n$  at random from  $\mathcal{D}$ ; We also sample  $n$  random binary labels from Radmarcher distribution. What is the “average of the best correlations” can hypothesis set  $\mathcal{H}$  achieve? Obviously, the higher the correlations that  $h \in \mathcal{H}$  can achieve between the set  $\{Z_i\}_{i=1}^n$  and the set  $\{\sigma_i\}_{i=1}^n$ , a better performance (or complexity) for  $\mathcal{H}$ .
2. Obviously, the most difficult for computing  $\text{Rad}_n(\mathcal{H})$  is to max over a possibly infinite hypothesis set  $\mathcal{H}$  (for example all the lines in linear classifications). Luckily, we can take advantage of for example:
  - (a) there are infinite  $h$ , but  $h(Z_i)$  has finite outcome.
  - (b) or the algebraic property, for example: when  $h(w^\top \mathbf{x}) = \sup_w (w^\top \mathbf{x})$

### 1.1 alternative definition

however, some text are using definition:

$$\text{Rad}_n(\mathcal{H}) = \mathbb{E}_S \left[ \mathbb{E}_{\vec{\sigma}} \left[ \sup_{h \in \mathcal{H}} \left| \frac{\sum_{i=1}^n \sigma_i h(Z_i)}{n} \right| \right] \right] \quad (2)$$

**QUESTION** is the above definition also valid?

## 1.2 Empirical Rademacher Complexity

$$\widehat{\text{Rad}}_S(\mathcal{H}) = \mathbb{E}_{\vec{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{\sum_{i=1}^n \sigma_i h(Z_i)}{n} \right] \quad (3)$$

which is precisely the stuff inside  $\text{Rad}_n(\mathcal{H})$ , i.e.,

$$\text{Rad}_n(\mathcal{H}) = \mathbb{E}_S \left[ \widehat{\text{Rad}}_S(\mathcal{H}) \right] \quad (4)$$

## 2 Bounds

### 2.1 Expected hypothesis bound

Before we state the expression for the generalization error bound in terms of expected and empirical loss, let's begin with a generic bound in terms of Rademacher complexity:

**Theorem 1** *Let  $Z, Z_1, \dots, Z_n$  be i.i.d random variables sampled from  $\mathcal{D}$ , and consider every hypothesis  $f \in \mathcal{F}$  is bounded by  $[a, b]$ . Then,  $\forall \delta > 0$ , with probability of at least  $1 - \delta$ , and respect to sample  $S$ , we have:*

$$\forall f \in \mathcal{F} : \quad \mathbb{E}_Z[f(Z)] \leq \frac{1}{n} \sum_{i=1}^n f(Z_i) + 2\text{Rad}_n(\mathcal{F}) + (b - a) \sqrt{\frac{\log(1/\delta)}{2n}} \quad (5)$$

$$\forall f \in \mathcal{F} : \quad \mathbb{E}_Z[f(Z)] \leq \frac{1}{n} \sum_{i=1}^n f(Z_i) + 2\widehat{\text{Rad}}_S(\mathcal{F}) + 3(b - a) \sqrt{\frac{\log(2/\delta)}{2n}} \quad (6)$$

Please note the following:

1. some literature write it as:

$$\forall f \in \mathcal{F} : \quad \mathbb{E}_Z[f(Z)] \leq \frac{1}{n} \sum_{i=1}^n f(Z_i) + 2\text{Rad}_n(\mathcal{F}) + (b - a) \sqrt{\frac{\log(1/\delta)}{n}} \quad (7)$$

which is also valid.

2. note how I deliberately write it using  $f \in \mathcal{F}$  instead of  $h \in \mathcal{H}$  because we have:

$$f_h(Z_i) = \ell(h(\mathbf{x}_i), y_i) \quad (8)$$

3. note that Theorem (1) can be applied more generically, and it does not need to apply specifically for Generalization error bound (in terms of expected and empirical loss) as stated in Theorem generalization bound rademacher complexity. In order to reflect this, let's give a symbol for:

$$\begin{aligned} T_n &= \mathbb{E}_Z[f(Z)] \\ \widehat{T}_S &= \frac{1}{n} \sum_{i=1}^n f(Z_i) \end{aligned} \quad (9)$$

4. I will show the proof of Theorem (1) later using McDiarmid inequality. in Section (3)

## 2.2 Generalization error bound

Firstly, Generalization error bound is what we are interested in. The expected and empirical risks are expressed as:

$$\begin{aligned} R(h) &= \mathbb{E}_Z[f_h(Z)] \\ &= \mathbb{E}_Z[\ell(h(\mathbf{x}), y)] \\ \hat{R}_S(h) &= \frac{1}{n} \sum_{i=1}^n f_h(Z_i) \\ &= \frac{1}{n} \sum_{i=1}^n \ell(h(\mathbf{x}_i), y_i) \end{aligned} \tag{10}$$

**Theorem 2** *let  $\mathcal{H}$  be set of hypothesis taking values in  $\{-1, +1\}$  and  $\mathcal{F} = L(\mathcal{H}) = \{f_h(x, y) \rightarrow \mathbb{1}_{\{h(x) \neq y\}} : h \in \mathcal{H}\}$ . For any  $\delta > 0$ , with probability at least  $1 - \delta$  over a sample  $S$  of size  $n$  drawn from  $\mathcal{D}$ :*

$$\forall h \in \mathcal{H} : R(h) \leq \hat{R}_S(h) + \widehat{\text{Rad}}_S(\mathcal{H}) + 3\sqrt{\frac{\log(2/\delta)}{2n}} \tag{11}$$

note that we need to apply Rademacher calculus to  $\mathcal{H}$  taking values in other ranges.

## 2.3 first attempt

Now, a straight application of Theorem 1 will give the following expression, let  $a = -1$ ,  $b = 1$ , and let  $f(Z) \equiv f_h(Z)$ :

$$\begin{aligned} \forall h \in \mathcal{H} : \mathbb{E}_Z[f_h(Z)] &\leq \frac{1}{n} \sum_{i=1}^n f_h(Z_i) + 2\widehat{\text{Rad}}_S(\mathcal{F}) + 3\sqrt{\frac{\log(2/\delta)}{2n}} \\ \implies \forall h \in \mathcal{H} : R(h) &\leq \hat{R}_S(h) + 2\widehat{\text{Rad}}_S(\mathcal{F}) + 3\sqrt{\frac{\log(2/\delta)}{2n}} \\ &= \hat{R}_S(h) + 2\widehat{\text{Rad}}_S(L(\mathcal{H})) + 3\sqrt{\frac{\log(2/\delta)}{2n}} \end{aligned} \tag{12}$$

where  $\mathcal{F} = L(\mathcal{H}) = \{f_h(x, y) \rightarrow \mathbb{1}_{\{h(x) \neq y\}} : h \in \mathcal{H}\}$

Note that we do not like to express it in terms of  $\widehat{\text{Rad}}_S(L(\mathcal{H}))$ , we need instead to express it using  $\widehat{\text{Rad}}_S(\mathcal{H})$ . Because  $\widehat{\text{Rad}}_S(\mathcal{H})$  is the actual complexity of the model, not just the output loss function! Therefore, we must find a way to express  $\widehat{\text{Rad}}_S(L(\mathcal{H}))$  in terms of  $\widehat{\text{Rad}}_S(\mathcal{H})$ .

## 2.4 express $\widehat{\text{Rad}}_S(L(\mathcal{H}))$ in terms of $\widehat{\text{Rad}}_S(\mathcal{H})$

Firstly, there is no general expression for all  $\mathcal{F}$  and  $\mathcal{H}$ . However, imagine we have a hypothesis set  $\mathcal{H}$  of binary functions, i.e.,  $\mathcal{H} \subset \{-1, 1\}^{\mathcal{X}}$ . This notation means output of hypothesis set  $\mathcal{H}$  is a sub-set of all possible  $\{-1, 1\}^{\mathcal{X}}$ .

Then, imagine we define a composite function where we add loss to it:  $f_h \in \mathcal{F}$ :

$$\mathcal{F} = L(\mathcal{H}) = \{\mathbb{1}_{\{h(x) \neq y\}} \mid h \in \mathcal{H}\} \tag{13}$$

what is then the relationship between  $\widehat{\text{Rad}}_S(\mathcal{H})$  and  $\widehat{\text{Rad}}_S(L(\mathcal{H}))$ :

**Lemma 3** given  $L(\mathcal{H}) = \{\mathbb{1}_{\{h(x) \neq y\}} \mid h \in \mathcal{H}\}$  and  $h(x) \in \{-1, 1\}$ :

$$\widehat{\text{Rad}}_S(L(\mathcal{H})) = \frac{1}{2} \widehat{\text{Rad}}_S(\mathcal{H}) \quad (14)$$

The results is as expected, as  $f_h \in \mathcal{F}$  outputs  $\{0, 1\}$ , therefore  $\widehat{\text{Rad}}_S(\mathcal{F})$  is much more restrictive than  $\widehat{\text{Rad}}_S(\mathcal{H})$ , hence has less complexity. The proof is easy to understand from the original source. Firstly, notice that conversion between  $h \rightarrow f_h$  is:

$$\begin{aligned} f_h &= \mathbb{1}_{h(X_i) \neq Y_i} & Y_i &\in \{-1, 1\} \\ &= \frac{1 - Y_i h(X_i)}{2} & \text{easy to see: } \frac{1 - Y_i h(X_i)}{2} &= \begin{cases} 0 & Y_i = h(X_i) \\ 1 & Y_i \neq h(X_i) \end{cases} \end{aligned} \quad (15)$$

$$\begin{aligned} \widehat{\text{Rad}}_S(L(\mathcal{H})) &= \widehat{\text{Rad}}_S(\mathcal{F}) \\ &= \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \frac{1 - Y_i h(X_i)}{2} \right] \\ &= \mathbb{E}_\sigma \left[ \frac{1}{2n} \sum_{i=1}^n \sigma_i + \frac{1}{2} \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \times (-Y_i) h(X_i) \right] \\ &= \underbrace{\frac{1}{2n} \mathbb{E}_\sigma \left[ \sum_{i=1}^n \sigma_i \right]}_0 + \frac{1}{2} \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i \times (-Y_i) h(X_i) \right] \\ &= \frac{1}{2} \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(X_i) \right] \quad \because \Pr(\sigma_i(-Y_i)) = \Pr(\sigma_i) \text{ and } Y_i \text{ is symmetric} \\ &= \frac{1}{2} \widehat{\text{Rad}}_S(\mathcal{H}) \end{aligned} \quad (16)$$

However, the result is significant. Using lemma (3), you can then apply Theorem (1), where you can express  $R(h)$  and  $\hat{R}_S(h)$  in terms of the  $\text{Rad}_n(\mathcal{H})$  and  $\widehat{\text{Rad}}_S(\mathcal{H})$  instead of expressing them in terms of  $\widehat{\text{Rad}}_S(L(\mathcal{H}))$ , i.e.,:

$$\begin{aligned} \forall h \in \mathcal{H} : \quad R(h) &\leq \hat{R}_S(h) + 2\widehat{\text{Rad}}_S(L(\mathcal{H})) + 3\sqrt{\frac{\log(2/\delta)}{2n}} \\ \implies \forall h \in \mathcal{H} : \quad R(h) &\leq \hat{R}_S(h) + \widehat{\text{Rad}}_S(\mathcal{H}) + 3\sqrt{\frac{\log(2/\delta)}{2n}} \end{aligned} \quad (17)$$

**Lemma 4** given any function class  $\mathcal{H}$  and constants  $a, b \in \mathbb{R}$ , if we denote a function class,  $\{f_h \mid f(x) = ah(x) + b\}$  by  $a\mathcal{F} + b$ , then:

immediately, you noticed that if one let  $f_h = \frac{1 - Yh(\mathbf{x})}{2}$ , then:

$$\begin{aligned} \widehat{\text{Rad}}_S(a\mathcal{H} + b) &= |a| \widehat{\text{Rad}}_S(\mathcal{H}) \\ \widehat{\text{Rad}}_S(\mathcal{F}) &= \widehat{\text{Rad}}_S\left(\frac{1 - Y\mathcal{H}(\mathbf{x})}{2}\right) \\ &= \left| \frac{-Y}{2} \right| \widehat{\text{Rad}}_S(\mathcal{H}) \\ &= \frac{1}{2} \widehat{\text{Rad}}_S(\mathcal{H}) \end{aligned} \quad (18)$$

### 2.4.1 Proof

the proof is rather straightforward:

$$\begin{aligned}
\widehat{\text{Rad}}_S(a\mathcal{H} + b) &= \mathbb{E}_\sigma \left[ \sup_{f_h \in a\mathcal{H} + b} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i f_h(Z_i) \right) \right] \\
&= \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i a h(Z_i) + b \right) \right] \\
&= \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i a h(Z_i) + \frac{1}{n} \sum_{i=1}^n \sigma_i b \right) \right] \\
&= \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i a h(Z_i) \right) \right] \quad \because b \mathbb{E}_\sigma \left[ \sum_{i=1}^n \sigma_i \right] = 0 \\
&= |a| \mathbb{E}_\sigma \left[ \sup_{h \in \mathcal{H}} \left( \frac{1}{n} \sum_{i=1}^n \sigma_i h(Z_i) \right) \right] \\
&= |a| \widehat{\text{Rad}}_S(\mathcal{H})
\end{aligned} \tag{19}$$

the last two lines can be understood by:

$$\begin{aligned}
\Pr(\sigma_i = 1) &= \Pr(\sigma_i = -1) = \frac{1}{2} \\
\implies \Pr(\sigma_i a = a) &= \Pr(\sigma_i a = -a) = \frac{1}{2}
\end{aligned} \tag{20}$$

## 3 Expected hypothesis bound proof

### 3.1 facts about $\sup(\cdot)$

Before we start the proof, we need to use few properties to about  $\sup(\cdot)$  function. There are few properties about  $\sup(\cdot)$  are all obvious, but we try to prove them formally:

1. **fact 1**  $\sup_x f(x) + \sup_x g(x) \geq \sup_x (f(x) + g(x))$

**proof**, even though it's obvious:

$$\begin{aligned}
\sup_{x, y \in A, x=y} (f(x) + g(y)) &\leq \sup_{x, y \in A} (f(x) + g(y)) \quad \text{LHS has smaller domain set} \\
\implies \sup_{x \in A} (f(x) + g(x)) &\leq \sup_{x \in A} f(x) + \sup_{x \in A} g(x)
\end{aligned} \tag{21}$$

2. **fact 2**  $\sup_x f(x) - \sup_x g(x) \leq \sup_x (f(x) - g(x))$

could just use above and write:

$$\begin{aligned}
\sup_x ((f(x) - g(x)) + g(x)) &\leq \sup_x (f(x) - g(x)) + \sup_x (g(x)) \\
\implies \sup_x f(x) &\leq \sup_x (f(x) - g(x)) + \sup_x (g(x)) \\
\implies \sup_x f(x) - \sup_x g(x) &\leq \sup_x (f(x) - g(x))
\end{aligned} \tag{22}$$

3. **fact 3**

proving direction of  $\sup_x f(x) - \sup_x g(x) \leq \sup_x (f(x) - g(x))$  only has one way. However, starting from  $\sup_x (f(x) - g(x))$  has two different routes:

$$\begin{cases} \sup_x (f(x) - g(x)) & \geq \sup_x f(x) - \sup_x g(x) \\ \sup_x (f(x) - g(x)) = \underbrace{\sup_x (f(x) + g(-x))}_{\text{if possible for } g(x) = -g(x)} & \leq \sup_x f(x) + \sup_x g(-x) \end{cases} \quad (23)$$

second line is a lot more useful as we usually derive bound for LHS with a  $\leq$  sign.

4. **fact 4**  $\sup(\cdot)$  is a convex function

**proof** let  $\bar{X} = \{x_1, \dots, x_n\}$  and  $\bar{Y} = \{y_1, \dots, y_n\}$ :

$$\theta \sup(\bar{X}) + (1 - \theta) \sup \bar{Y} \geq \sup(\theta \bar{X} + (1 - \theta) \bar{Y}) \quad (24)$$

to prove this, we can pick arbitrary  $x_i$  and  $y_k$  from the two sets respectively:

$$\begin{aligned} \theta \sup(\bar{X}) + (1 - \theta) \sup \bar{Y} &\geq \theta x_i + (1 - \theta) y_k \quad \forall i, k \\ \implies \theta \sup(\bar{X}) + (1 - \theta) \sup \bar{Y} &\geq \sup(\theta \bar{X} + (1 - \theta) \bar{Y}) \end{aligned} \quad (25)$$

### 3.2 Bound $\sup_{h \in \mathcal{H}} \left( \mathbb{E}_Z[h(Z)] - \frac{1}{n} \sum_{i=1}^n h(Z_i) \right)$ by $\text{Rad}_n(\mathcal{H})$

We do not want to confuse  $\mathbb{E}_Z[h(Z)]$  by  $R(h)$ , therefore in Eq.(9), we defined:

$$T(h) = \mathbb{E}_Z[h(Z)], \quad \hat{T}_S(h) = \frac{1}{n} \sum_{i=1}^n h(Z_i) \quad (26)$$

$$\begin{aligned} \phi(S) &= \sup_{h \in \mathcal{H}} \left( \mathbb{E}_Z[h(Z)] - \frac{1}{n} \sum_{i=1}^n h(Z_i) \right) \quad Z_i \in S \\ &= \sup_{h \in \mathcal{H}} \left( T(h) - \hat{T}_S(h) \right) \end{aligned} \quad (27)$$

in word,  $\phi(S)$  is maximum difference between empirical and expected hypothesis output for any  $h \in \mathcal{H}$ . Note that there is no absolute value here yet.

we then apply this to see the difference between two sets:

$$\begin{aligned} S &= \{Z_1, \dots, Z_{i-1}, Z_i, Z_{i+1}, \dots, Z_n\} \\ S' &= \{Z_1, \dots, Z_{i-1}, Z'_i, Z_{i+1}, \dots, Z_n\} \end{aligned} \quad (28)$$

$$\begin{aligned} \phi(S) - \phi(S') &= \sup_{h \in \mathcal{H}} \underbrace{\left( T[h] - \hat{T}_S(h) \right)}_{f(z)} - \sup_{h \in \mathcal{H}} \underbrace{\left( T[h] - \hat{T}_{S'}(h) \right)}_{g(z)} \quad Z_i \in S, Z'_i \in S' \\ &\leq \sup_{h \in \mathcal{H}} \left( (T(h) - \hat{T}_S(h)) - (T(h) - \hat{T}_{S'}(h)) \right) \quad \text{use fact 1} \\ &= \sup_{h \in \mathcal{H}} \left( \hat{T}_{S'}(h) - \hat{T}_S(h) \right) \quad \text{remove } T(h) \text{ which is hard to deal with} \\ &= \frac{1}{n} \sup_{h \in \mathcal{H}} \left( h(Z'_i) - h(Z_i) \right) \quad \text{use fact that only single } Z_i \neq Z'_i \text{ so, only one term remain in sum} \\ &\leq \frac{b-a}{n} \quad \text{assume a single } h \in [a, b] \end{aligned} \quad (29)$$

Equally, by symmetry, we can show  $\phi(S') - \phi(S) \leq \frac{b-a}{n}$ . Therefore we have:

$$|\phi(S') - \phi(S)| \leq \frac{b-a}{n} \quad (30)$$

which satisfies the condition for McDiarmid's inequality, i.e.,

$$\begin{aligned} |f(x_1, \dots, x_i, \dots, x_m) - f(x_1, \dots, x'_i, \dots, x_m)| &\leq c_i \\ |\phi(S') - \phi(S)| &\leq \frac{b-a}{n} \end{aligned} \quad (31)$$

the only complicated part is that  $f(\cdot) \equiv \phi(\cdot)$  involves a  $\sup_h$  term. We can derive:

$$\begin{aligned} \Pr(f(X_1, X_2, \dots, X_n) - \mathbb{E}[f(X_1, X_2, \dots, X_n)] \geq \epsilon) &\leq \exp\left(-\frac{2\epsilon^2}{\sum_{i=1}^n c_i^2}\right) \\ \implies \Pr(\phi(S) - \mathbb{E}_S[\phi(S)] \geq \epsilon) &\leq \exp\left(-\frac{2\epsilon^2}{n \left(\frac{b-a}{n}\right)^2}\right) \\ &\leq \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right) \end{aligned} \quad (32)$$

now, simplify with R.H.S:

$$\begin{aligned} \delta &= \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right) \\ \log(\delta) &= -\frac{2n\epsilon^2}{(b-a)^2} \\ \log(1/\delta) &= \frac{2n\epsilon^2}{(b-a)^2} \\ \implies \epsilon^2 &= \frac{\log(1/\delta)(b-a)^2}{2n} \\ \implies \epsilon &= (b-a)\sqrt{\frac{\log(1/\delta)}{2n}} \end{aligned} \quad (33)$$

therefore, with probability of **at least**  $1 - \delta$ :

$$\begin{aligned} \phi(S) - \mathbb{E}_S[\phi(S)] &\leq (b-a)\sqrt{\frac{\log(1/\delta)}{2n}} \\ \implies \phi(S) &\leq \mathbb{E}_S[\phi(S)] + (b-a)\sqrt{\frac{\log(1/\delta)}{2n}} \end{aligned} \quad (34)$$

however, the above was just bounding  $\phi(S) = \sup_{h \in \mathcal{H}} (T(h) - \hat{T}_S(h))$ , this is independent of Rademacher complexity.

### 3.2.1 introduce $\text{Rad}_n(\mathcal{H})$

looking at Eq.(34), let's see if we can bound the quantity  $\mathbb{E}_S[\phi(S)]$  with  $\text{Rad}_n(\mathcal{H})$ .

Firstly, we realize that:

$$\begin{aligned}
T(h) &= \mathbb{E}_Z[h(Z)] && \text{expectation of single sample} \\
&= \mathbb{E}_S[\widehat{T}_S(h)] && \text{expectation of sample average} \\
&= \mathbb{E}_S\left[\frac{1}{n} \sum_{i=1}^n h(Z_i)\right] && Z_i \in S
\end{aligned} \tag{35}$$

This is analogous to Empirical and Expected loss.

Note that we let  $T(h) = \mathbb{E}_S\left[\frac{1}{n} \sum_{i=1}^n h(Z_i)\right]$  instead of  $T(h) = \mathbb{E}_Z[h(Z)]$ , so that the equations will align with  $\widehat{\text{Rad}}_S(\mathcal{H})$ .

Using the expectation of sample average version from Eq.(35). In here, we need to bring a ghost set  $\tilde{S}$  into the equation, which needs to be **independent** from the  $S$  in the outer expectation.

$$\begin{aligned}
\mathbb{E}_S[\phi(S)] &= \mathbb{E}_S\left[\sup_{h \in \mathcal{H}} \left(T(h) - \widehat{T}_S(h)\right)\right] \\
&= \mathbb{E}_S\left[\sup_{h \in \mathcal{H}} \left(\mathbb{E}_{\tilde{S}}\left[\frac{1}{n} \sum_{i=1}^n h(\tilde{Z}_i)\right] - \frac{1}{n} \sum_{i=1}^n h(Z_i)\right)\right] \quad Z_i \in S
\end{aligned} \tag{36}$$

note that in here, all elements  $S = \{Z_i\}$  differ from  $\tilde{S} = \{\tilde{Z}_i\}$ . Do not confuse this with Eq.(28), in here, **all** elements in sets  $S$  and  $\tilde{S}$  are independent.

$$\begin{aligned}
\mathbb{E}_S[\phi(S)] &= \mathbb{E}_S\left[\sup_{h \in \mathcal{H}} \left(\mathbb{E}_{\tilde{S}}\left[\frac{1}{n} \sum_{i=1}^n h(\tilde{Z}_i)\right] - \frac{1}{n} \sum_{i=1}^n h(Z_i)\right)\right] \quad Z_i \in S, \quad \tilde{Z}_i \in \tilde{S} \\
&= \mathbb{E}_S\left[\underbrace{\sup_{h \in \mathcal{H}} \frac{1}{n} \left(\mathbb{E}_{\tilde{S}}\left[\sum_{i=1}^n (h(\tilde{Z}_i) - h(Z_i))\right]\right)}_{\sup(\mathbb{E}[\cdot])}]\right] \quad Z_i \in S, \quad \tilde{Z}_i \in \tilde{S} \\
&\leq \mathbb{E}_S\left[\underbrace{\mathbb{E}_{\tilde{S}}\left[\sup_{h \in \mathcal{H}} \frac{1}{n} \left(\sum_{i=1}^n (h(\tilde{Z}_i) - h(Z_i))\right)\right]}_{\mathbb{E}[\sup(\cdot)]}\right] \\
&= \mathbb{E}_S\left[\mathbb{E}_{\tilde{S}}\left[\mathbb{E}_{\tilde{\sigma}}\left[\sup_{h \in \mathcal{H}} \frac{1}{n} \left(\sum_{i=1}^n \sigma_i (h(\tilde{Z}_i) - h(Z_i))\right)\right]\right]\right] \quad \because h(\tilde{Z}_i) - h(Z_i) \text{ are symmetric, same trick prove hoeffding lemma} \\
&= \mathbb{E}_{S, \tilde{S}, \tilde{\sigma}}\left[\underbrace{\sup_{h \in \mathcal{H}} \frac{1}{n} \left(\sum_{i=1}^n \sigma_i h(\tilde{Z}_i) + (-\sigma_i h(Z_i))\right)}_{\sup_x (f(x) + (-g(x))) \leq \sup_x f(x) + \sup_x -g(x)}\right] \quad \text{use short notation} \\
&\leq \mathbb{E}_{\tilde{S}, \tilde{\sigma}}\left[\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(\tilde{Z}_i)\right] + \mathbb{E}_{S, \tilde{\sigma}}\left[\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n -\sigma_i h(Z_i)\right] \\
&= \mathbb{E}_{\tilde{S}, \tilde{\sigma}}\left[\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(\tilde{Z}_i)\right] + \mathbb{E}_{S, \tilde{\sigma}}\left[\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(Z_i)\right] \quad \text{both sets are from } \mathcal{D} \\
&= 2\mathbb{E}_{S, \tilde{\sigma}}\left[\sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(Z_i)\right] \\
&= 2\text{Rad}_n(\mathcal{H})
\end{aligned} \tag{37}$$



therefore, with probability of **at least**  $1 - \delta$ :

$$\begin{aligned} \phi(S) &\leq \mathbb{E}_S[\phi(S)] + (b-a)\sqrt{\frac{\log(1/\delta)}{2n}} \\ \sup_{h \in \mathcal{H}} \left( T(h) - \widehat{T}_S(h) \right) &\leq 2\text{Rad}_n(\mathcal{H}) + (b-a)\sqrt{\frac{\log(1/\delta)}{2n}} \end{aligned} \quad (38)$$

**Fatou lemma**

$$\underbrace{\sup_{h \in \mathcal{H}} \frac{1}{n} \left( \mathbb{E}_{S'} \left[ \sum_{i=1}^n (h(Z'_i) - h(Z_i)) \right] \right)}_{\sup(\mathbb{E}[\cdot])} \leq \underbrace{\mathbb{E}_{S'} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \left( \sum_{i=1}^n (h(Z'_i) - h(Z_i)) \right) \right]}_{\mathbb{E}[\sup(\cdot)]} \quad (39)$$

Fatou lemmas says:

$$\limsup_{n \rightarrow \infty} \int_S f_n d\mu \leq \int_S \limsup_{n \rightarrow \infty} f_n d\mu. \quad (40)$$

however, as we haven't chosen  $n \rightarrow \infty$ , it is unclear how Fatou Lemma is applied.

### 3.3 bound $\text{Rad}_n(\mathcal{H})$ using $\widehat{\text{Rad}}_S(\mathcal{H})$

The result in Eq.(38) are in terms of  $\text{Rad}_n(\mathcal{H})$ , which may be difficult to evaluate.

We hope to apply this to  $\widehat{\text{Rad}}_S(\mathcal{H})$  making it possible to evaluate. Similar to before, We can apply McDiarmid's inequality. To do so, we introduce a set which differ by a single element  $S'$ :

$$\begin{aligned} S &= \{Z_1, \dots, Z_{i-1}, Z_i, Z_{i+1}, \dots, Z_n\} \\ S' &= \{Z_1, \dots, Z_{i-1}, Z'_i, Z_{i+1}, \dots, Z_n\} \end{aligned} \quad (41)$$

then,

$$\begin{aligned} \widehat{\text{Rad}}_S(\mathcal{H}) - \widehat{\text{Rad}}_{S'}(\mathcal{H}) &= \mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{\sum_{Z_i \in S} \sigma_i h(Z_i)}{n} \right] - \mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{\sum_{Z_i \in S'} \sigma_i h(Z_i)}{n} \right] \\ &= \mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{\sum_{Z_i \in S} \sigma_i h(Z_i)}{n} - \sup_{h \in \mathcal{H}} \frac{\sum_{Z_i \in S'} \sigma_i h(Z_i)}{n} \right] \\ &= \frac{1}{n} \mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \sum_{Z_i \in S} \sigma_i h(Z_i) - \sup_{h \in \mathcal{H}} \sum_{Z_i \in S'} \sigma_i h(Z_i) \right] \\ &\leq \frac{1}{n} \mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \left( \sum_{Z_i \in S} \sigma_i h(Z_i) - \sum_{Z_i \in S'} \sigma_i h(Z_i) \right) \right] \\ &= \frac{1}{n} \mathbb{E}_{\sigma_i} \left[ \sup_{h_i \in \mathcal{H}} \sigma_i (h(Z_i) - \sigma_i h(Z'_i)) \right] \\ &= \frac{1}{n} \sup_{h_i \in \mathcal{H}} (h(Z_i) - h(Z'_i)) \quad \text{symmetry} \\ &\leq \frac{b-a}{n} \end{aligned} \quad (42)$$

then, by proving  $\widehat{\text{Rad}}_{S'}(\mathcal{H}) - \widehat{\text{Rad}}_S(\mathcal{H}) \leq \frac{b-a}{n}$ , let to the condition (not probability) for McDiarmid:

$$|\widehat{\text{Rad}}_S(\mathcal{H}) - \widehat{\text{Rad}}_{S'}(\mathcal{H})| \leq \frac{b-a}{n} \quad (43)$$

Apply McDiarmid's inequality, we can bound the function by the expectation (we need it for  $\text{Rad}_n(\mathcal{H})$ !):

$$\begin{aligned} \widehat{\text{Rad}}_S(\mathcal{H}) - \mathbb{E}_S[\text{Rad}_S(\mathcal{H})] &= \widehat{\text{Rad}}_S(\mathcal{H}) - \text{Rad}_n(\mathcal{H}) \\ &\leq \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right) \\ \implies \Pr\left(\text{Rad}_n(\mathcal{H}) \geq \widehat{\text{Rad}}_S(\mathcal{H}) + (b-a)\sqrt{\frac{\log(1/\delta)}{2n}}\right) &\leq \delta \\ \implies \Pr\left(\text{Rad}_n(\mathcal{H}) \leq \widehat{\text{Rad}}_S(\mathcal{H}) + (b-a)\sqrt{\frac{\log(1/\delta)}{2n}}\right) &\geq 1 - \delta \end{aligned} \quad (44)$$

### 3.4 combine RHS using both $\text{Rad}_n(\mathcal{H})$ and $\widehat{\text{Rad}}_S(\mathcal{H})$

$$\begin{aligned} \Pr\left(\text{Rad}_n(\mathcal{H}) \geq \widehat{\text{Rad}}_S(\mathcal{H}) + (b-a)\sqrt{\frac{\log(1/\delta)}{2n}}\right) &\leq \delta \\ \implies \Pr\left(\text{Rad}_n(\mathcal{H}) \geq \widehat{\text{Rad}}_S(\mathcal{H}) + (b-a)\sqrt{\frac{\log(2/\delta)}{2n}}\right) &\leq \frac{\delta}{2} \\ \implies \Pr\left(\underbrace{2\text{Rad}_n(\mathcal{H}) \geq 2\widehat{\text{Rad}}_S(\mathcal{H}) + 2(b-a)\sqrt{\frac{\log(2/\delta)}{2n}}}\right) &\leq \frac{\delta}{2} \end{aligned} \quad (45)$$

From **theorem** (1):

$$\Pr\left(T(h) \geq \widehat{T}_S(h) + 2\text{Rad}_n(\mathcal{H}) + (b-a)\sqrt{\frac{\log(2/\delta)}{2n}}\right) \leq \frac{\delta}{2} \quad (46)$$

substitute Eq.(45) into Theorem (1):

$$\begin{aligned} \implies \Pr\left(T(h) \geq \widehat{T}_S(h) + \underbrace{2\widehat{\text{Rad}}_S(\mathcal{H}) + 2(b-a)\sqrt{\frac{\log(2/\delta)}{2n}}}_{\text{union bound}} + (b-a)\sqrt{\frac{\log(2/\delta)}{2n}}\right) &\leq \frac{\delta}{2} + \frac{\delta}{2} \\ \implies \Pr\left(T(h) \geq \widehat{T}_S(h) + 2\widehat{\text{Rad}}_S(\mathcal{H}) + 3(b-a)\sqrt{\frac{\log(2/\delta)}{2n}}\right) &\leq \delta \end{aligned} \quad (47)$$

## 4 Finite Class Lemma (Massart)

**Lemma 5** *Let  $\mathcal{A}$  be some finite subset of  $\mathbb{R}^m$  and let:*

$$r = \sup_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a}\| \quad \text{or} \quad \|\mathbf{a}\| \leq r \quad (49)$$

*then:*

$$\mathbb{E}_{\bar{\sigma}} \left[ \sup_{\mathbf{a} \in \mathcal{A}} \frac{1}{n} \sum_{i=1}^n \sigma_i a_i \right] \leq \frac{r \sqrt{2 \log(|\mathcal{A}|)}}{n} \quad (50)$$

It is obvious that we are using it for Rademacher complexity. To make it look like  $\mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(Z_i) \right]$ , by equating:

$$\begin{aligned} a_i &\equiv h(Z_i) \\ \mathbf{a} &= [h(Z_1) \quad h(Z_2) \quad \dots \quad h(Z_n)] \end{aligned} \quad (51)$$

Also  $|\mathcal{A}|$  is bounded, i.e., we have a finite class, meaning the value of  $h(Z)$  is finite.

### 4.1 proof

we can just prove **lemma 5** using  $\mathbb{E}_{\bar{\sigma}} \left[ \sup_{a \in \mathcal{A}} \sum_{i=1}^n \sigma_i a_i \right] \leq r \sqrt{2 \log(|\mathcal{A}|)}$ :

$$\begin{aligned}
\exp\left(\lambda \mathbb{E}_{\bar{\sigma}}\left[\sup_{\mathbf{a} \in \mathcal{A}} \sum_{i=1}^n \sigma_i a_i\right]\right) &\leq \mathbb{E}_{\bar{\sigma}}\left[\exp\left(\lambda \sup_{\mathbf{a} \in \mathcal{A}} \sum_{i=1}^n \sigma_i a_i\right)\right] \quad \lambda > 0 : \text{ Jensen} \\
&= \mathbb{E}_{\bar{\sigma}}\left[\sup_{\mathbf{a} \in \mathcal{A}} \left(\exp\left(\lambda \sum_{i=1}^n \sigma_i a_i\right)\right)\right] \quad \exp(\lambda(\cdot)) \text{ monotone} \\
&\leq \mathbb{E}_{\bar{\sigma}}\left[\sum_{\mathbf{a} \in \mathcal{A}} \exp\left(\lambda \sum_{i=1}^n \sigma_i a_i\right)\right] \quad \text{replace sup with } \sum \\
&= \sum_{\mathbf{a} \in \mathcal{A}} \mathbb{E}_{\bar{\sigma}}\left[\exp\left(\lambda \sum_{i=1}^n \sigma_i a_i\right)\right] \\
&= \sum_{\mathbf{a} \in \mathcal{A}} \mathbb{E}_{\bar{\sigma}}\left[\prod_{i=1}^n \exp^{\lambda \sigma_i a_i}\right] \\
&= \sum_{\mathbf{a} \in \mathcal{A}} \prod_{i=1}^n \mathbb{E}_{\sigma_i}\left[\exp^{\lambda \sigma_i a_i}\right] \quad \because \Pr(\bar{\sigma}) = \prod_i \Pr(\sigma_i) \\
&= \sum_{\mathbf{a} \in \mathcal{A}} \prod_{i=1}^n \left(\frac{1}{2} \exp^{\lambda(+1)a_i} + \frac{1}{2} \exp^{\lambda(-1)a_i}\right) \quad \text{write expectation explicitly} \\
&= \sum_{\mathbf{a} \in \mathcal{A}} \prod_{i=1}^n \frac{\exp^{\lambda a_i} + \exp^{-\lambda a_i}}{2} \\
&\leq \sum_{\mathbf{a} \in \mathcal{A}} \prod_{i=1}^n \exp\left(\frac{\lambda^2 a_i^2}{2}\right) \quad \because \frac{\exp^x + \exp^{-x}}{2} \leq \exp^{\frac{x^2}{2}} \\
&\quad \text{remember previous class: } \underbrace{\frac{\exp^{\lambda} + \exp^{-\lambda}}{2}}_{\text{MGF}_{\sigma \sim \text{Rad}}(\lambda)} \leq \underbrace{\exp^{\frac{\lambda^2}{2}}}_{\text{Hoeffding lemma bound}} \\
&= \sum_{\mathbf{a} \in \mathcal{A}} \exp\left(\sum_{i=1}^n \frac{\lambda^2 a_i^2}{2}\right) = \sum_{\mathbf{a} \in \mathcal{A}} \exp\left(\frac{\lambda^2}{2} \sum_{i=1}^n a_i^2\right) \\
&= \sum_{\mathbf{a} \in \mathcal{A}} \exp\left(\frac{\lambda^2}{2} \|\mathbf{a}\|^2\right) \\
&\leq |\mathcal{A}| \exp^{\frac{\lambda^2 r^2}{2}} \quad \text{union bound over } \mathbf{finite} \text{ class set } \mathcal{A}
\end{aligned} \tag{52}$$

$$\text{let } \mu = \mathbb{E}_{\bar{\sigma}}\left[\sup_{\mathbf{a} \in \mathcal{A}} \sum_{i=1}^n \sigma_i a_i\right]:$$

$$\begin{aligned}
\exp\left(\lambda \mathbb{E}_{\bar{\sigma}}\left[\sup_{\mathbf{a} \in \mathcal{A}} \sum_{i=1}^n \sigma_i a_i\right]\right) &\equiv \exp(\lambda \mu) \leq |\mathcal{A}| \exp^{\frac{\lambda^2 r^2}{2}} \\
\lambda \mu &\leq \log(|\mathcal{A}|) + \frac{\lambda^2 r^2}{2} \\
\Rightarrow \mu &\leq \frac{\log(|\mathcal{A}|)}{\lambda} + \frac{\lambda r^2}{2}
\end{aligned} \tag{53}$$

$$\begin{aligned}
-\log(|\mathcal{A}|)\lambda^{-2} + \frac{r^2}{2} &= 0 \\
\implies \frac{r^2}{2\log(|\mathcal{A}|)} &= \lambda^{-2} \\
\implies \lambda^2 &= \frac{2\log(|\mathcal{A}|)}{r^2} \\
\implies \lambda &= \frac{\sqrt{2\log(|\mathcal{A}|)}}{r}
\end{aligned} \tag{54}$$

by substitution, we have:

$$\begin{aligned}
\mu &\leq \log(|\mathcal{A}|) \frac{r}{\sqrt{2\log(|\mathcal{A}|)}} + \frac{\sqrt{2\log(|\mathcal{A}|)}r}{2} \\
&= \frac{\sqrt{2\log(|\mathcal{A}|)}r}{2} + \frac{\sqrt{2\log(|\mathcal{A}|)}r}{2} \\
&= r\sqrt{2\log(|\mathcal{A}|)}
\end{aligned} \tag{55}$$

**Corollary 5.1** *Let  $\mathcal{H}$  be finite set of functions such that  $|h(z)| \leq 1$ :*

$$\text{Rad}_n(\mathcal{H}) \leq \sqrt{\frac{2\log(|\mathcal{H}|)}{n}} \tag{56}$$

Given  $\mathcal{H}$  and  $S = \{z_1, \dots, z_n\}$  and let  $\{\mathbf{h} \equiv (h(z_1), \dots, h(z_n)) : h \in \mathcal{H}\}$ , then for every  $h(z) \in \mathcal{H}$ , we have:

$$\|\mathbf{a}\| \leq \sqrt{\sum_{i=1}^n h(z_i)} \leq \sqrt{n} \tag{57}$$

apply **Massart Lemma**, we have  $r \equiv \sqrt{n}$  and  $a_i \equiv h(z_i)$ :

$$\begin{aligned}
\mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(z_i) \right] &\leq \sqrt{n} \frac{\sqrt{2\log(|\mathcal{H}|)}}{n} \\
&= \sqrt{\frac{2\log(|\mathcal{H}|)}{n}}
\end{aligned} \tag{58}$$

## 5 Rademacher Calculus

$$\mathbf{5.1} \quad c \cdot \mathcal{H} + b = \{c \cdot h + b : h \in \mathcal{H}\} \implies \mathbf{Rad}_n(c \cdot \mathcal{H} + b) = |c| \mathbf{Rad}_n(\mathcal{H})$$

this is already proved in lemma 4.

$$\mathbf{5.2} \quad \mathbf{Rad}_n(\text{conv}(\mathcal{H})) = \mathbf{Rad}_n(\mathcal{H})$$

**Lemma 6** *Let  $\text{conv}\mathcal{H} = \{\sum \theta_i h_i : \{h_i\} \subseteq \mathcal{H}, \theta_i \geq 0, \sum \theta_i = 1\}$   
Then:*

$$\text{Rad}_n(\text{conv}(\mathcal{H})) = \text{Rad}_n(\mathcal{H}) \tag{59}$$

Rademacher complexity of all convex combination of functions does not change the original Rademacher complexity value. This is not so surprising, as convex combination does not change the overall complexity.

### 5.2.1 proof

$$\begin{aligned}
\text{Rad}_n(\text{conv}(\mathcal{H})) &= \mathbb{E} \left[ \sup_{h_j \in \mathcal{H}, \|\theta\|_1=1} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{j=1}^k \theta_j h_j(X_i) \right] \\
&= \mathbb{E} \left[ \sup_{h_j \in \mathcal{H}, \|\theta\|_1=1} \frac{1}{n} \sum_{i=1}^n \sigma_i \sum_{j=1}^k \theta_j h_j(X_i) \right] \\
&= \mathbb{E} \left[ \sup_{h_j \in \mathcal{H}} \sup_{\|\theta\|_1=1} \frac{1}{n} \sum_{j=1}^k \theta_j \left( \sum_{i=1}^n \sigma_i h_j(X_i) \right) \right] \\
&= \mathbb{E} \left[ \sup_{h_j \in \mathcal{H}} \max_j \frac{1}{n} \sum_{i=1}^n \sigma_i h_j(X_i) \right] \quad \text{given } \{h_j\} \text{ in outer sup, the inner } \theta \text{ becomes one hot} \\
&= \mathbb{E} \left[ \sup_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sigma_i h(X_i) \right] \\
&= \text{Rad}_n(\mathcal{H})
\end{aligned} \tag{60}$$

meant learning over distributions of hypotheses is statistically not any harder than learning over individual hypotheses.

### 5.2.2 application of $\text{Rad}_n(\text{conv}(\mathcal{H})) = \text{Rad}_n(\mathcal{H})$

**Lemma 7**

$$\text{conv}(S_{n,s}) \subset K_{n,s} \subset 2\text{conv}(S_{n,s}) \tag{61}$$

where  $\text{conv}(\cdot)$  is the convex hull that contains the set  $S_{n,s}$

$$\begin{aligned}
K_{n,s} &:= \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \leq 1, \|\mathbf{x}\|_1 \leq \sqrt{s}\} = B_2^n \cap \sqrt{s}B_1^n \\
S_{n,s} &:= \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_1 \leq 1, \|\mathbf{x}\|_0 \leq s\}
\end{aligned} \tag{62}$$

1.  $\text{conv}(S_{n,s}) \subset K_{n,s}$

$S_{n,s}$  is **not** a convex set, but  $K_{n,s}$  is the convex set. We need to prove that,  $\forall \mathbf{x} \in S_{n,s} \implies \mathbf{x} \in K_{n,s}$ , i.e.,

$$\|\mathbf{x}\|_1 \leq 1 \cap \|\mathbf{x}\|_0 \leq s \implies \|\mathbf{x}\|_2 \leq 1 \cap \|\mathbf{x}\|_0 \leq \sqrt{s} \tag{63}$$

firstly,  $\|\mathbf{x}\| \leq 1 \implies \|\mathbf{x}\|_2 \leq 1$ . think about the geometric region. Then:

$$\begin{aligned}
\|\mathbf{x}\|_1 &= |\mathbf{x}|^\top \mathbf{1}_s \quad \text{where } |\mathbf{x}| = [x_1 \quad x_2 \quad \dots \quad x_n] \\
&\quad \mathbf{1}_s \text{ is a vector with at most } s \text{ one, and rest are zero} \\
&\leq \|\mathbf{x}\|_2 \|\mathbf{1}_s\|_2 \quad \text{cauchy schwarz inequality} \\
&\leq \|\mathbf{x}\|_2 \sqrt{s} \quad \because \|\mathbf{x}\|_0 \leq s \\
&\leq \sqrt{s} \quad \because \|\mathbf{x}\|_1 \leq 1 \implies \|\mathbf{x}\|_2 \leq 1
\end{aligned} \tag{64}$$

The convex hull  $\text{conv}(S)$  is the **smallest** convex set containing  $S$ . Formally, the convex hull of the set  $S$  of points in  $n$  dimension is the intersection of all convex sets containing  $S$ . For  $N$  points  $\mathbf{x}_1, \dots, \mathbf{x}_N$ , the convex hull  $C$  is:

$$C = \left\{ \sum_{j=1}^N \lambda_j \mathbf{x}_j : \lambda_j \geq 0 \quad \forall j \quad \sum_{j=1}^N \lambda_j = 1 \right\} \quad (65)$$

since  $S_{n,s}$  is the convex hull and since  $K_{n,s}$  is a convex set, then it must imply:

$$\text{conv}(S_{n,s}) \subset K_{n,s} \quad (66)$$

2.  $K_{n,s} \subset 2\text{conv}(S_{n,s})$

will prove it for the next session

here comes the real example:

In sparse linear regression, we want to find a linear predictor with at most  $s$ -nonzero coordinates. Loss function is:

$$\arg \min_w \frac{1}{2n} \sum_{i=1}^n (\langle w, x_i \rangle - y_i)^2 \quad \text{s.t. } \mathbf{function\ class} \quad (67)$$

where **function class** can be defined in terms of a ball:

1.  $\mathcal{W}_0$ :

$$\begin{aligned} \mathcal{W}_0 &= \{w \in \mathbb{R}^d : \|w\|_2^2 \leq 1, \|w\|_0 \leq s\} \\ \arg \min_w \frac{1}{2n} \sum_{i=1}^n (\langle w, x_i \rangle - y_i)^2 \quad \text{s.t. } w \in \mathcal{W}_0 \end{aligned} \quad (68)$$

where  $\mathcal{H}_0 = \{f_w(x) = \langle w, x \rangle, w \in \mathcal{W}_0\}$  is NP-hard.

2.  $\mathcal{W}_1$ :

However, if we work over  $L_1$  ball:

$$\begin{aligned} \mathcal{W}_1 &= \{w \in \mathbb{R}^d : \|w\|_2^2 \leq 1, \|w\|_1 \leq \sqrt{s}\} \\ \arg \min_w \frac{1}{2n} \sum_{i=1}^n (\langle w, x_i \rangle - y_i)^2 \quad \text{s.t. } w \in \mathcal{W}_1 \end{aligned} \quad (69)$$

where  $\mathcal{H}_1 = \{f_w(x) = \langle w, x \rangle, w \in \mathcal{W}_1\}$ :

then optimization problem is convex, and by applying lemma 7, then we have:

$$\widehat{Rad}_n(\mathcal{W}_1) \leq 2\widehat{Rad}_n(\mathcal{W}_0) \quad (70)$$

Statistically, working with  $L_1$  ball does not substantially increase the complexity.

### 5.3 $L$ -Lipschitz functions

**Lemma 8** Let  $\phi_z$  be  $L$ -Lipschitz functions for every  $z \in \mathcal{X}$ , i.e.  $|\phi_z(a) - \phi_z(b)| \leq L|a - b|$ , or in our context,  $|\phi(f(z)) - \phi(f'(z))| \leq L|f(z) - f'(z)|$ . Denote:

$$\phi \circ \mathcal{H} = \{\phi_z(h(z)) : h \in \mathcal{H}\} \quad (71)$$

Then:

$$\text{Rad}_n(\phi \circ \mathcal{H}) = L \text{Rad}_n(\mathcal{H}) \quad (72)$$

it can be written as to proving:

$$\mathbb{E}_{\vec{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{\sum_{i=1}^n \sigma_i \phi(h(Z_i))}{n} \right] \leq L \mathbb{E}_{\vec{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{\sum_{i=1}^n \sigma_i h(Z_i)}{n} \right] \quad (73)$$

let's remove  $n$  on both sides, and write out expectation of  $h(Z_i)$  explicitly and w.l.o.g, we work with just one element  $\sigma_1$ , which can then be replicated to other elements throughout:

$$\begin{aligned} & \mathbb{E}_{\vec{\sigma}} \left[ \sup_{h \in \mathcal{H}} \left( \sigma_1 \phi(h(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] \\ &= \mathbb{E}_{\sigma_1} \left[ \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( \sigma_1 \phi(h(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] \right] \\ &= \underbrace{\frac{1}{2}}_{\Pr(\sigma_1=+1)} \left[ \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( \underbrace{(+1)}_{\sigma_1=1} \phi(h(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] \right] + \underbrace{\frac{1}{2}}_{\Pr(\sigma_1=-1)} \left[ \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( \underbrace{(-1)}_{\sigma_1=-1} \phi(h(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] \right] \\ &= \frac{1}{2} \left[ \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( \phi(h(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] \right] + \frac{1}{2} \left[ \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( -\phi(h(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] \right] \\ &= \frac{1}{2} \left( \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( \phi(h(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] + \sup_{h \in \mathcal{H}} \left( -\phi(h(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] \\ &= \frac{1}{2} \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h, h' \in \mathcal{H}} \left( \phi(h(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) - \phi(h'(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h'(Z_i)) \right) \right] \quad \text{still two separate sup}(\cdot) \\ &\leq \frac{1}{2} \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h, h' \in \mathcal{H}} \left( L|h(Z_1) - h'(Z_1)| + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) + \sum_{i=2}^n \sigma_i \phi(h'(Z_i)) \right) \right] \quad \because |\phi(h(z)) - \phi(h'(z))| \leq L|h(z) - h'(z)| \\ &= \frac{1}{2} \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h, h' \in \mathcal{H}} \left( L(h(Z_1) - h'(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) + \sum_{i=2}^n \sigma_i \phi(h'(Z_i)) \right) \right] \end{aligned} \quad (74)$$

the reason you can remove absolute operator in  $|h(Z_1) - h'(Z_1)| \rightarrow (h(Z_1) - h'(Z_1))$ , because if we swap  $h$  with  $h'$  it won't affect the sum  $\sum_{i=2}^n \sigma_i h(Z_i) + \sum_{i=2}^n \sigma_i h'(Z_i)$ , so sup can swap  $h$  with  $h'$  to make  $h(Z_1) - h'(Z_1) \geq 0$



$$\begin{aligned}
& \mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \left( \sigma_1 \phi(h(Z_1)) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] \\
&= \frac{1}{2} \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( Lh(Z_1) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) + \sup_{\sup_{h \in \mathcal{H}}} \left( -Lh(Z_1) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] \\
&= \underbrace{\frac{1}{2}}_{\Pr(\sigma_1=1)} \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( \underbrace{+1}_{\sigma_1=1} Lh(Z_1) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] + \underbrace{\frac{1}{2}}_{\Pr(\sigma_1=-1)} \mathbb{E}_{\sigma_2, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( \underbrace{-1}_{\sigma_1=-1} Lh(Z_1) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right] \\
&= \mathbb{E}_{\sigma_1, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( L\sigma_1 h(Z_1) + \sum_{i=2}^n \sigma_i \phi(h(Z_i)) \right) \right]
\end{aligned} \tag{75}$$

we can show the rest by computing the expectation using just  $\sigma_2$  for the following:

$$\begin{aligned}
& \mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \left( L\sigma_1 h(Z_1) + \sigma_2 \phi(h(Z_2)) + \sum_{i=3}^n \sigma_i \phi(h(Z_i)) \right) \right] \\
&= \mathbb{E}_{\sigma_2} \left[ \mathbb{E}_{\sigma_1, \sigma_3, \dots, \sigma_n} \left[ \sup_{h \in \mathcal{H}} \left( L\sigma_1 h(Z_1) + \sigma_2 \phi(h(Z_2)) + \sum_{i=3}^n \sigma_i \phi(h(Z_i)) \right) \right] \right]
\end{aligned} \tag{76}$$

repeat the process we will obtain:

$$\mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{\sum_{i=1}^n \sigma_i \phi(h(Z_i))}{n} \right] \leq L \mathbb{E}_{\bar{\sigma}} \left[ \sup_{h \in \mathcal{H}} \frac{\sum_{i=1}^n \sigma_i h(Z_i)}{n} \right] \tag{77}$$