
Rademacher Complexity

Beier Zhu

Definition 1 (Empirical Rademacher Complexity). Let \mathcal{G} be a family of functions mapping from \mathcal{Z} to $[a, b]$ and $S = (z_1, \dots, z_m)$ a fixed sample of size m with elements in \mathcal{Z} . Then, the empirical Rademacher complexity of \mathcal{G} w.r.t the sample S is defined as:

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

where $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_m)^\top$, with σ_i s independent uniform r.v.s taking values $\{-1, +1\}$. The r.v.s σ_i are called Rademacher variables. Let \mathbf{g}_S denote the vector of values taken by function g over the sample S : $\mathbf{g}_S = (g(z_1), \dots, g(z_m))^\top$.

Remark.

- The inner product $\boldsymbol{\sigma}^\top \mathbf{g}_S$ measures the **correlation** of \mathbf{g}_S with the vector of random noise $\boldsymbol{\sigma}$.
- The supremum $\sup_{g \in \mathcal{G}} \frac{\boldsymbol{\sigma}^\top \mathbf{g}_S}{m}$ measures **how well the function class** \mathcal{G} correlates with $\boldsymbol{\sigma}$ over the sample S .
- The expectation on supremum $\mathbb{E}[\sup_{g \in \mathcal{G}} \frac{\boldsymbol{\sigma}^\top \mathbf{g}_S}{m}]$ measures **on average** how well the function class \mathcal{G} correlates with noise on the sample S .

Definition 2 (Rademacher Complexity). Let \mathcal{D} denote the distribution from which samples are drawn. For any $m \geq 1$, the Rademacher complexity of \mathcal{G} is the expectation of the empirical Rademacher complexity over all samples of size m drawn from \mathcal{D} :

$$\mathfrak{R}_m(\mathcal{G}) = \mathbb{E}_{S \sim \mathcal{D}^m} [\widehat{\mathfrak{R}}_S(\mathcal{G})]. \quad (2)$$

Proposition 1. For any $\delta > 0$, with probability at least $1 - \delta$, the following holds:

$$\mathfrak{R}_m(\mathcal{G}) \leq \widehat{\mathfrak{R}}_S(\mathcal{G}) + (b - a) \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (3)$$

Proof. Let S and S' be two samples differing by exactly one point, say z_m in S and z'_m in S' . Then, we have

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^m \sigma_i g(z_i) \right] \quad (4)$$

$$= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^{m-1} \sigma_i g(z_i) + \sigma_m g(z'_m) + \sigma_m g(z_m) - \sigma_m g(z'_m) \right] \quad (5)$$

$$\leq \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^{m-1} \sigma_i g(z_i) + \sigma_m g(z'_m) \right] + \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}_m} \left[\sup_{g \in \mathcal{G}} |\sigma_m(g(z_m) - g(z'_m))| \right] \quad (6)$$

$$\leq \widehat{\mathfrak{R}}_{S'}(\mathcal{G}) + \frac{b - a}{m} \quad (7)$$

Similarly, we obtain $\widehat{\mathfrak{R}}_{S'}(\mathcal{G}) \leq \widehat{\mathfrak{R}}_S(\mathcal{G}) + \frac{b - a}{m}$, thus $|\widehat{\mathfrak{R}}_{S'}(\mathcal{G}) - \widehat{\mathfrak{R}}_S(\mathcal{G})| \leq \frac{b - a}{m}$. Then, we use McDiarmid's inequality to have Eq.(3). \square

Theorem 1. Let \mathcal{G} be a family of functions mapping from \mathcal{Z} to \mathbb{R} . Let $\widehat{\mathbb{E}}_S[g(z)]$ denote the empirical average of g over S : $\widehat{\mathbb{E}}_S[g(z)] = \frac{1}{m} \sum_{i=1}^m g(z_i)$.

$$\mathbb{E}_S \left[\sup_{g \in \mathcal{G}} \left[\widehat{\mathbb{E}}_S[g(z)] - \mathbb{E}[g(z)] \right] \right] \leq 2\mathfrak{R}_m(\mathcal{G}) \quad (8)$$

Proof. Fix $S = [z_1, \dots, z_m]$, the term in the expectation on the LHS of Eq. (8) is

$$\sup_{g \in \mathcal{G}} \left[\widehat{\mathbb{E}}_S[g(z)] - \mathbb{E}[g(z)] \right] = \sup_{g \in \mathcal{G}} \left[\frac{1}{m} \sum_{i=1}^m g(z_i) - \mathbb{E}[g(z)] \right] \quad (9)$$

$$= \sup_{g \in \mathcal{G}} \left[\frac{1}{m} \sum_{i=1}^m g(z_i) - \mathbb{E}_{S'} \left[\frac{1}{m} \sum_{i=1}^m g(z'_i) \right] \right] \quad (\mathbb{E}[g] = \mathbb{E}_{S'}[\widehat{\mathbb{E}}_{S'}[g]])$$

$$= \frac{1}{m} \sup_{g \in \mathcal{G}} \left[\mathbb{E}_{S'} \left[\sum_{i=1}^m g(z_i) - g(z'_i) \right] \right] \quad (10)$$

$$\leq \frac{1}{m} \mathbb{E}_{S'} \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^m g(z_i) - g(z'_i) \right] \quad (\text{sub-add. of the sup.})$$

Sub-additivity of the supremum functions: $\sup(U + V) \leq \sup(U) + \sup(V)$. Similarly, $\sup(\mathbb{E}_X[f(X)]) \leq \mathbb{E}_X[\sup(f(X))]$. Now, take the expectation over S for both side :

$$\mathbb{E}_S \left[\sup_{g \in \mathcal{G}} \left[\widehat{\mathbb{E}}_S[g(z)] - \mathbb{E}[g(z)] \right] \right] \leq \frac{1}{m} \mathbb{E}_{S,S'} \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^m g(z_i) - g(z'_i) \right] \quad (11)$$

$$= \frac{1}{m} \mathbb{E}_{S,S',\sigma} \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^m \sigma_i (g(z_i) - g(z'_i)) \right] \quad (12)$$

$$\leq \frac{1}{m} \mathbb{E}_{S,\sigma} \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^m \sigma_i g(z_i) \right] + \frac{1}{m} \mathbb{E}_{S',\sigma} \left[\sup_{g \in \mathcal{G}} \sum_{i=1}^m -\sigma_i g(z'_i) \right] \\ (\sigma_i \stackrel{d}{=} -\sigma_i)$$

$$= 2\mathfrak{R}_m(\mathcal{G}) \quad (13)$$

Eq. (12) holds because $\sigma_i(g(z_i) - g(z'_i)) \stackrel{d}{=} g(z_i) - g(z'_i)$, since $g(z_i) - g(z'_i)$ and $g(z'_i) - g(z_i)$ have symmetric distribution. \square

Theorem 2. Let \mathcal{G} be a family of functions mapping from \mathcal{Z} to $[0, 1]$. Then, for any $\delta > 0$, with probability at least $1 - \delta$ over the draw of an i.i.d. sample $S = [z_1, \dots, z_m]$ of size m , each of the following holds for all $g \in \mathcal{G}$:

$$\mathbb{E}[g(z)] \leq \widehat{\mathbb{E}}_S[g(z)] + 2\mathfrak{R}_m(\mathcal{G}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (14)$$

$$\mathbb{E}[g(z)] \leq \widehat{\mathbb{E}}_S[g(z)] + 2\widehat{\mathfrak{R}}_S(\mathcal{G}) + 3\sqrt{\frac{\log \frac{1}{\delta}}{2m}}. \quad (15)$$

Compared to Theorem 2, we can rewrite the above as

$$\sup_{g \in \mathcal{G}} \left[\mathbb{E}[g(z)] - \widehat{\mathbb{E}}_S[g(z)] \right] \leq 2\mathfrak{R}_m(\mathcal{G}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (16)$$

Proof. The proof uses McDiarmid's inequality to function Φ defined for any sample S by

$$\Phi(S) = \sup_{g \in S} \left(\mathbb{E}[g] - \widehat{\mathbb{E}}_S[g] \right). \quad (17)$$

Let S and S' be two samples differing by exactly one point, say z_m in S and z'_m in S' . Then, since the difference of suprema does not exceed the supremum of the difference, we have

$$\Phi(S') - \Phi(S) \leq \sup_{g \in \mathcal{G}} (\widehat{\mathbb{E}}_S[g] - \widehat{\mathbb{E}}_{S'}[g]) = \sup_{g \in \mathcal{G}} \frac{g(z_m) - g(z'_m)}{m} \leq \frac{1}{m}. \quad (18)$$

Similarly, we obtain $\Phi(S) - \Phi(S') \leq \frac{1}{m}$, thus $|\Phi(S) - \Phi(S')| \leq \frac{1}{m}$. Then, by McDiarmid's inequality, for any $\delta > 0$, with probability at least $1 - \delta$, the following holds:

$$\Phi(S) \leq \mathbb{E}_S[\Phi(S)] + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}. \quad (19)$$

We next bound the expectation of the right-hand side as follows:

$$\mathbb{E}_S[\Phi(S)] = \mathbb{E}_S \left[\sup_{g \in \mathcal{G}} (\mathbb{E}(g) - \widehat{\mathbb{E}}_S(g)) \right] \quad (20)$$

$$= \mathbb{E}_S \left[\sup_{g \in \mathcal{G}} \mathbb{E}_{S'}[\widehat{\mathbb{E}}_{S'}(g) - \widehat{\mathbb{E}}_S(g)] \right] \quad (\mathbb{E}[g] = \mathbb{E}_{S'}[\widehat{\mathbb{E}}_{S'}[g]])$$

$$\leq \mathbb{E}_{S,S'} \left[\sup_{g \in \mathcal{G}} (\widehat{\mathbb{E}}_{S'}(g) - \widehat{\mathbb{E}}_S(g)) \right] \quad (\text{sub-add. of the sup.})$$

$$= \mathbb{E}_{S,S'} \left[\frac{1}{m} \sum_{i=1}^m \sup_{g \in \mathcal{G}} g(z'_i) - g(z_i) \right] \quad (21)$$

$$= \mathbb{E}_{\sigma, S, S'} \left[\frac{1}{m} \sum_{i=1}^m \sup_{g \in \mathcal{G}} \sigma_i (g(z'_i) - g(z_i)) \right] \quad (22)$$

$$\leq \mathbb{E}_{\sigma, S'} \left[\frac{1}{m} \sum_{i=1}^m \sup_{g \in \mathcal{G}} \sigma_i g(z'_i) \right] + \mathbb{E}_{\sigma, S} \left[\frac{1}{m} \sum_{i=1}^m \sup_{g \in \mathcal{G}} -\sigma_i g(z_i) \right] \quad (\text{sub-add. of the sup.})$$

$$= 2\mathfrak{R}_m(\mathcal{G}) \quad (-\sigma \text{ is Rademacher r.v.)}$$

In Eq.(22), we introduce Rademacher variables σ_i which do not change the expectation appearing in Eq.(21): when $\sigma_i = 1$, the associated summand remains unchanged; when $\sigma_i = -1$, the associated summand flips signs, which equivalent to swapping z_i and z'_i between S and S' . Since (z_i, z'_i) and (z'_i, z_i) have the same joint distribution, this swap does not affect the overall expectation. Substituting it back into Eq. (19) and using the definition of $\Phi(S)$, we obtain

$$\sup_{g \in S} (\mathbb{E}[g] - \widehat{\mathbb{E}}_S[g]) \leq 2\mathfrak{R}_m(\mathcal{G}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (23)$$

Obviously, $\forall g \in \mathcal{G}$, we have:

$$\mathbb{E}[g(z)] \leq \widehat{\mathbb{E}}_S[g(z)] + 2\mathfrak{R}_m(\mathcal{G}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (24)$$

Using Proposition 1, we have $\mathfrak{R}_m(\mathcal{G}) \leq \widehat{\mathfrak{R}}_S(\mathcal{G}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$. Then,

$$\mathbb{E}[g(z)] \leq \widehat{\mathbb{E}}_S[g(z)] + 2\widehat{\mathfrak{R}}_S(\mathcal{G}) + 3\sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (25)$$

□

Lemma 1. Let \mathcal{H} be a family of functions taking values in $\{-1, +1\}$ and let \mathcal{G} be the family of loss functions associated to \mathcal{H} for the zero-one loss: $\mathcal{G} = \{(x, y) \mapsto \mathbb{1}[h(x) \neq y] : h \in \mathcal{H}\}$. For any sample $S = ((x_1, y_1), \dots, (x_m, y_m))$ of elements in $\mathcal{X} \times \{-1, +1\}$, let $S_{\mathcal{X}}$ denote its projection over \mathcal{X} : $S_{\mathcal{X}} = (x_1, \dots, x_m)$. Then, the following relation holds between the empirical Rademacher complexities of \mathcal{G} and \mathcal{H} :

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \frac{1}{2} \widehat{\mathfrak{R}}_{S_{\mathcal{X}}}(\mathcal{H}). \quad (26)$$

Proof.

$$\widehat{\mathfrak{R}}_S(\mathcal{G}) = \frac{1}{m} \mathbb{E} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i \mathbf{1}[h(x_i) \neq y_i] \right] \quad (27)$$

$$= \frac{1}{m} \mathbb{E} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i \frac{1 - h(x_i)y_i}{2} \right] \quad (28)$$

$$= \frac{1}{2} \frac{1}{m} \mathbb{E} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m (-y_i \sigma_i) h(x_i) \right] \quad (29)$$

$$= \frac{1}{2} \widehat{\mathfrak{R}}_{S_X}(\mathcal{H}), \quad (30)$$

where we use the fact that $\mathbf{1}[h(x_i) \neq y_i] = \frac{1-h(x_i)y_i}{2}$ and the fact that for fixed $y_i = \{-1, +1\}$, $(-\sigma_i y_i)$ are also Rademacher r.v.s. \square

Theorem 3. Let \mathcal{H} be a family of functions taking values in $\{-1, +1\}$ and let \mathcal{D} be the distribution over the input space \mathcal{X} . Then, for any $\delta > 0$, with probability at least $1 - \delta$ over a sample S of size m drawn according to \mathcal{D} , each of the following holds for any $h \in \mathcal{H}$:

$$\mathcal{R}(h) \leq \widehat{\mathcal{R}}_S(h) + \mathfrak{R}_m(\mathcal{H}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (31)$$

$$\mathcal{R}(h) \leq \widehat{\mathcal{R}}_S(h) + \widehat{\mathfrak{R}}_S(\mathcal{H}) + 3\sqrt{\frac{\log \frac{1}{\delta}}{2m}} \quad (32)$$

Proof. The result follows immediately by Theorem 2 and Lemma 1. \square

Remark. The second bound is data-dependent: the empirical Rademacher complexity $\widehat{\mathfrak{R}}_S(\mathcal{H})$ is a function of a specific sample S .

Proposition 2 (Talagrand's Lemma). Let Φ_1, \dots, Φ_m be l -Lipschitz functions from \mathbb{R} to \mathbb{R} and $\sigma_1, \dots, \sigma_m$ be Rademacher variable. Then, for any hypothesis set \mathcal{H} of real-valued functions, the following inequality holds:

$$\frac{1}{m} \mathbb{E} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i (\Phi_i \circ h)(x_i) \right] \leq \frac{l}{m} \mathbb{E} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i) \right] = l \widehat{\mathfrak{R}}_S(\mathcal{H}). \quad (33)$$

In particular, if $\Phi_i = \Phi$ for all $i \in [m]$, then the following holds:

$$\widehat{\mathfrak{R}}_S(\Phi \circ \mathcal{H}) \leq l \widehat{\mathfrak{R}}_S(\mathcal{H}). \quad (34)$$

Proof. First we fix a sample $S = (x_1, \dots, x_m)$, then, by definition,

$$\frac{1}{m} \mathbb{E} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i (\Phi_i \circ h)(x_i) \right] = \frac{1}{m} \mathbb{E}_{\sigma_{1:m-1}} \left[\mathbb{E}_{\sigma_m} \left[\sup_{h \in \mathcal{H}} u_{m-1}(h) + \sigma_m (\Phi_m \circ h)(x_m) \right] \right], \quad (35)$$

where $u_{m-1}(h) = \sum_{i=1}^{m-1} \sigma_i (\Phi_i \circ h)(x_i)$. By the definition of the supremum ($s = \sup A$, then for any $\epsilon > 0$, $\exists a_\epsilon \in A$ s.t. $a_\epsilon > s - \epsilon$), for any $\epsilon > 0$, there exist $h_1, h_2 \in \mathcal{H}$ such that

$$u_{m-1}(h_1) + (\Phi_m \circ h_1)(x_m) \geq (1 - \epsilon) \left[\sup_{h \in \mathcal{H}} u_{m-1}(h) + (\Phi_m \circ h)(x_m) \right] \quad (36)$$

$$u_{m-1}(h_2) - (\Phi_m \circ h_2)(x_m) \geq (1 - \epsilon) \left[\sup_{h \in \mathcal{H}} u_{m-1}(h) - (\Phi_m \circ h)(x_m) \right] \quad (37)$$

Thus, for any $\epsilon > 0$, by definition of \mathbb{E}_{σ_m} ,

$$(1 - \epsilon) \mathbb{E}_{\sigma_m} \left[\sup_{h \in \mathcal{H}} u_{m-1}(h) + \sigma_m(\Phi_m \circ h)(x_m) \right] \quad (38)$$

$$= (1 - \epsilon) \left[\frac{1}{2} \sup_{h \in \mathcal{H}} [u_{m-1}(h) + (\Phi_m \circ h)(x_m)] + \frac{1}{2} \sup_{h \in \mathcal{H}} [u_{m-1}(h) - (\Phi_m \circ h)(x_m)] \right] \quad (39)$$

$$\leq \frac{1}{2} [u_{m-1}(h_1) + (\Phi_m \circ h_1)(x_m)] + \frac{1}{2} [u_{m-1}(h_2) - (\Phi_m \circ h_2)(x_m)] \quad (40)$$

Let $s = \text{sign}(h_1(x_m) - h_2(x_m))$. Then, the previous inequality implies:

$$\begin{aligned} &= \frac{1}{2} [u_{m-1}(h_1) + u_{m-1}(h_2) + \Phi_m(h_1(x_m)) - \Phi_m(h_2(x_m))] && \text{(rearranging)} \\ &\leq \frac{1}{2} [u_{m-1}(h_1) + u_{m-1}(h_2) + s(l(h_1(x_m) - h_2(x_m)))] && \text{(l-Lipschitzness)} \\ &= \frac{1}{2} [u_{m-1}(h_1) + slh_1(x_m)] + \frac{1}{2} [u_{m-1}(h_2) - slh_2(x_m)] && \text{(rearranging)} \\ &\leq \frac{1}{2} \sup_{h \in \mathcal{H}} [u_{m-1}(h) + slh(x_m)] + \frac{1}{2} \sup_{h \in \mathcal{H}} [u_{m-1}(h) - slh(x_m)] && \text{(definition of sup)} \\ &= \mathbb{E}_{\sigma_m} \left[\sup_{h \in \mathcal{H}} u_{m-1}(h) + \sigma_m l h(x_m) \right] && \text{(definition of } \mathbb{E}_{\sigma_m} \text{)} \end{aligned}$$

Since the inequality holds for all $\epsilon > 0$, we have:

$$\mathbb{E}_{\sigma_m} \left[\sup_{h \in \mathcal{H}} u_{m-1}(h) + \sigma_m(\Phi_m \circ h)(x_m) \right] \leq \mathbb{E}_{\sigma_m} \left[\sup_{h \in \mathcal{H}} u_{m-1}(h) + \sigma_m l h(x_m) \right] \quad (41)$$

Proceeding in the same way for all other $\sigma_i (i \neq m)$ proves the lemma. \square

Proposition 3 (Extending Talagrand's Lemma to Vector Valued Functions). Let \mathcal{H} be a hypothesis set of functions mapping \mathcal{X} to \mathbb{R}^c . Assume that for all $i = 1, \dots, m$, $\Psi_i : \mathbb{R}^c \rightarrow \mathbb{R}$ is μ_i -Lipschitz for \mathbb{R}^c equipped with the 2-norm. That is:

$$|\Psi_i(\mathbf{x}') - \Psi_i(\mathbf{x})| \leq \|\mathbf{x}' - \mathbf{x}\|_2, \quad (42)$$

for all $(\mathbf{x}, \mathbf{x}') \in (\mathbb{R}^c, \mathbb{R}^c)$. Then, for any sample S of m points $x_1, \dots, x_m \in \mathcal{X}$, the following inequality holds

$$\frac{1}{m} \mathbb{E} \left[\sup_{\mathbf{h} \in \mathcal{H}} \sum_{i=1}^m \sigma_i \Psi_i(\mathbf{h}(x_i)) \right] \leq \frac{\sqrt{2}}{m} \mathbb{E} \left[\sup_{\mathbf{h} \in \mathcal{H}} \sum_{i=1}^m \sum_{j=1}^c \epsilon_{ij} \mu_i h_j(x_i) \right], \quad (43)$$

where $\boldsymbol{\epsilon} = (\epsilon_{ij})_{i,j}$ and ϵ_{ij} are independent Rademacher variables uniformly distributed over $\{1, -1\}$. In particular, if $\Psi_i = \Psi$ for all $i \in [m]$, then the following holds:

$$\widehat{\mathfrak{R}}_S(\Psi \circ \mathcal{H}) \leq \frac{\sqrt{2}}{m} \mu \widehat{\mathfrak{R}}_S(\mathcal{H}), \quad (44)$$

Proof. First we fix a sample $S = (x_1, \dots, x_m)$, then, by definition,

$$\frac{1}{m} \mathbb{E} \left[\sup_{\mathbf{h} \in \mathcal{H}} \sum_{i=1}^m \sigma_i \Psi_i(\mathbf{h}(x_i)) \right] = \frac{1}{m} \mathbb{E}_{\sigma_1, \dots, \sigma_{m-1}} \left[\mathbb{E}_{\sigma_m} \left[\sup_{\mathbf{h} \in \mathcal{H}} U_{m-1}(\mathbf{h}) + \sigma_m \Psi_m(\mathbf{h}(x_m)) \right] \right], \quad (45)$$

where $U_{m-1}(\mathbf{h}) = \sum_{i=1}^{m-1} \sigma_i \Psi_i(\mathbf{h}(x_i))$. Assume that the suprema can be attained and let $\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}$ be the hypotheses satisfying

$$U_{m-1}(\mathbf{h}_1) + \Psi_m(\mathbf{h}_1(x_m)) = \sup_{\mathbf{h} \in \mathcal{H}} U_{m-1}(\mathbf{h}) + \Psi_m(\mathbf{h}(x_m)) \quad (46)$$

$$U_{m-1}(\mathbf{h}_2) - \Psi_m(\mathbf{h}_2(x_m)) = \sup_{\mathbf{h} \in \mathcal{H}} U_{m-1}(\mathbf{h}) - \Psi_m(\mathbf{h}(x_m)) \quad (47)$$

When the suprema are not reached, a similar argument to what follows can be given by considering instead hypotheses that are ϵ -close to the suprema for any $\epsilon > 0$. By definition of expectation, since σ_m is uniformly distributed over $\{1, -1\}$, we can write

$$\mathbb{E}_{\sigma_m} \left[\sup_{\mathbf{h} \in \mathcal{H}} U_{m-1}(\mathbf{h}) + \sigma_m \Psi_m(\mathbf{h}(x_m)) \right] \quad (48)$$

$$= \frac{1}{2} \sup_{\mathbf{h} \in \mathcal{H}} U_{m-1}(\mathbf{h}) + \Psi_m(\mathbf{h}(x_m)) + \frac{1}{2} \sup_{\mathbf{h} \in \mathcal{H}} U_{m-1}(\mathbf{h}) - \Psi_m(\mathbf{h}(x_m)) \quad (49)$$

$$= \frac{1}{2} [U_{m-1}(\mathbf{h}_1) + \Psi_m(\mathbf{h}_1(x_m))] + \frac{1}{2} [U_{m-1}(\mathbf{h}_2) - \Psi_m(\mathbf{h}_2(x_m))] \quad (50)$$

$$= \frac{1}{2} [U_{m-1}(\mathbf{h}_1) + U_{m-1}(\mathbf{h}_2) + \Psi_m(\mathbf{h}_1(x_m)) - \Psi_m(\mathbf{h}_2(x_m))] \quad (51)$$

$$\leq \frac{1}{2} [U_{m-1}(\mathbf{h}_1) + U_{m-1}(\mathbf{h}_2) + \mu_m \|\mathbf{h}_1(x_m) - \mathbf{h}_2(x_m)\|_2] \quad (52)$$

$$\leq \frac{1}{2} \left[U_{m-1}(\mathbf{h}_1) + U_{m-1}(\mathbf{h}_2) + \mu_m \sqrt{2} \mathbb{E}_{\epsilon_{m1}, \dots, \epsilon_{mc}} \left[\left| \sum_{j=1}^c \epsilon_{mj} (h_{1j}(x_m) - h_{2j}(x_m)) \right| \right] \right], \quad (53)$$

where we use the μ_m -Lipschitzness of Ψ_m and the Khintchine-Kahane inequality. Let $\epsilon_m = (\epsilon_{m1}, \dots, \epsilon_{mc})$ and $s(\epsilon_m) \in \{1, -1\}$ denote the sign of $\sum_{j=1}^c \epsilon_{mj} (h_{1j}(x_m) - h_{2j}(x_m))$. Then, the following holds:

$$\mathbb{E}_{\sigma_m} \left[\sup_{\mathbf{h} \in \mathcal{H}} U_{m-1}(\mathbf{h}) + \sigma_m \Psi_m(\mathbf{h}(x_m)) \right] \quad (54)$$

$$\leq \frac{1}{2} \mathbb{E}_{\epsilon_m} \left[U_{m-1}(\mathbf{h}_1) + U_{m-1}(\mathbf{h}_2) + \mu_m \sqrt{2} \left| \sum_{j=1}^c \epsilon_{mj} (h_{1j}(x_m) - h_{2j}(x_m)) \right| \right] \quad (55)$$

$$= \frac{1}{2} \mathbb{E}_{\epsilon_m} \left[U_{m-1}(\mathbf{h}_1) + \mu_m \sqrt{2} s(\epsilon_m) \sum_{j=1}^c \epsilon_{mj} h_{1j}(x_m) + U_{m-1}(\mathbf{h}_2) - \mu_m \sqrt{2} s(\epsilon_m) \sum_{j=1}^c \epsilon_{mj} h_{2j}(x_m) \right] \quad (56)$$

$$\leq \frac{1}{2} \mathbb{E}_{\epsilon_m} \left[\sup_{\mathbf{h} \in \mathcal{H}} \left(U_{m-1}(\mathbf{h}) + \mu_m \sqrt{2} s(\epsilon_m) \sum_{j=1}^c \epsilon_{mj} h_j(x_m) \right) + \sup_{\mathbf{h} \in \mathcal{H}} \left(U_{m-1}(\mathbf{h}) - \mu_m \sqrt{2} s(\epsilon_m) \sum_{j=1}^c \epsilon_{mj} h_j(x_m) \right) \right] \quad (57)$$

$$= \mathbb{E}_{\epsilon_m} \left[\mathbb{E}_{\sigma_m} \left[\sup_{\mathbf{h} \in \mathcal{H}} U_{m-1}(\mathbf{h}) + \mu_m \sqrt{2} \sigma_m \sum_{j=1}^c \epsilon_{mj} h_j(x_m) \right] \right] \quad (58)$$

$$= \mathbb{E}_{\epsilon_m} \left[\sup_{\mathbf{h} \in \mathcal{H}} U_{m-1}(\mathbf{h}) + \mu_m \sqrt{2} \sum_{j=1}^c \epsilon_{mj} h_j(x_m) \right] \quad (59)$$

We have the last equality as the product of two independent Rademacher variables ($\sigma_m \epsilon_{mj}$) is still Rademacher variable. Note that $\mathbb{E}[\epsilon_i \epsilon_j a] = 0$ for fixed a , but $\mathbb{E}[\sup_{a \in \mathcal{A}} \epsilon_i \epsilon_j a] \neq 0$. For example, if $\mathcal{A} = \{1, -1\}$, $\mathbb{E}[\sup_{a \in \mathcal{A}} \epsilon_i \epsilon_j a] = 1$. Proceeding in the same way for all other σ_i 's ($i < m$) completes the proof. \square

Rademacher Identities

Fix $m \geq 1$, for any $\alpha \in \mathbb{R}$ and any two hypothesis sets \mathcal{H} and \mathcal{H}' of functions mapping from \mathcal{X} to \mathbb{R} , we have:

- (a) $\mathfrak{R}_m(\alpha \mathcal{H}) = |\alpha| \mathfrak{R}_m(\mathcal{H})$,
where $\alpha \mathcal{H} = \{\alpha h(x) | h \in \mathcal{H}\}$.

- (b) $\mathfrak{R}_m(\alpha + \mathcal{H}) = \mathfrak{R}_m(\mathcal{H})$,
where $\alpha + \mathcal{H} = \{\alpha + h(x) | h \in \mathcal{H}\}$.
- (c) $\mathfrak{R}_m(\mathcal{H} + \mathcal{H}') = \mathfrak{R}_m(\mathcal{H}) + \mathfrak{R}_m(\mathcal{H}')$,
where $\mathcal{H} + \mathcal{H}' = \{h(x) + h'(x) | h \in \mathcal{H}, h' \in \mathcal{H}'\}$.
- (d) $\mathfrak{R}_m(\{\max(h, h') | h \in \mathcal{H}, h' \in \mathcal{H}'\}) \leq \mathfrak{R}_m(\mathcal{H}) + \mathfrak{R}_m(\mathcal{H}')$.
 $\max(h, h') : x \mapsto \max(h(x), h'(x))$.

Fix $\mathbf{x} \in \mathbb{R}^m$ and let $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_m)^\top$ with σ_i s be Rademacher variables, Then:

$$(e) \|\mathbf{x}\|_2 = [\mathbb{E}_{\boldsymbol{\sigma}}(\boldsymbol{\sigma}^\top \mathbf{x})^2]^{\frac{1}{2}}.$$

(f) Khintchine inequality.

$$A_p \|\mathbf{x}\|_2 \leq [\mathbb{E}_{\boldsymbol{\sigma}} |\boldsymbol{\sigma}^\top \mathbf{x}|^p]^{1/p} \leq B_p \|\mathbf{x}\|_2$$

where

$$A_p = \begin{cases} 2^{1/2-1/p} & \text{if } 0 < p \leq p_0, \\ 2^{1/2} (\Gamma(\frac{p+1}{2}) / \sqrt{\pi})^{1/p} & \text{if } p_0 < p < 2, \\ 1 & \text{if } 2 \leq p < \infty, \end{cases} \quad (60)$$

$$\text{and } B_p = \begin{cases} 1 & \text{if } 0 < p \leq 2, \\ 2^{1/2} (\Gamma(\frac{p+1}{2}) / \sqrt{\pi})^{1/p} & \text{if } 2 < p < \infty. \end{cases} \quad (61)$$

$p_0 \approx 1.847$ and Γ is the Gamma function.

Let \mathcal{H}_1 and \mathcal{H}_2 be two families of functions mapping \mathcal{X} to $\{0, 1\}$ and let $\mathcal{H} = \{h_1 h_2 : h_1 \in \mathcal{H}_1, h_2 \in \mathcal{H}_2\}$, Then:

$$(g) \widehat{\mathfrak{R}}_S(\mathcal{H}) \leq \widehat{\mathfrak{R}}_S(\mathcal{H}_1) + \widehat{\mathfrak{R}}_S(\mathcal{H}_2)$$

Proof. We proof the equalities for empirical Rademacher complexity over a sample S , then taking the expectation yields the claimed results.

(a) For fixed sample S , we have

$$\widehat{\mathfrak{R}}_S(\alpha \mathcal{H}) = \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i \alpha h(x_i) \right] = |\alpha| \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i \text{sign}(\alpha) h(x_i) \right] = |\alpha| \widehat{\mathfrak{R}}_S(\mathcal{H}). \quad (62)$$

Note that $\sigma_i \text{sign}(\alpha)$ s are still Rademacher variables and $\sup cA = c \sup A$ if $c \geq 0$.

(b) For fixed sample S , we have

$$\widehat{\mathfrak{R}}_S(\alpha + \mathcal{H}) = \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i [\alpha + h(x_i)] \right] \quad (63)$$

$$= \mathbb{E}_{\boldsymbol{\sigma}} \left[\frac{1}{m} \alpha \sum_{i=1}^m \sigma_i \right] + \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] \quad (64)$$

$$= \widehat{\mathfrak{R}}_S(\mathcal{H}) \quad (\mathbb{E}_{\boldsymbol{\sigma}}[\sum_i \sigma_i] = 0)$$

(c) For fixed sample S , we have

$$\widehat{\mathfrak{R}}_S(\mathcal{H} + \mathcal{H}') = \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \frac{1}{m} \sum_{i=1}^m \sigma_i (h(x_i) + h'(x_i)) \right] \quad (65)$$

$$= \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i (h(x_i)) + \sup_{h' \in \mathcal{H}'} \frac{1}{m} \sum_{i=1}^m \sigma_i (h'(x_i)) \right] \quad (66)$$

$$= \widehat{\mathfrak{R}}_S(\mathcal{H}) + \widehat{\mathfrak{R}}_S(\mathcal{H}') \quad (67)$$

(d) Fix sample S and use the identity $\max(a, b) = \frac{1}{2}[a + b + |a - b|]$, we have:

$$\widehat{\mathfrak{R}}_S(\max(\mathcal{H}, \mathcal{H}')) \quad (68)$$

$$= \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \frac{1}{m} \sum_{i=1}^m \sigma_i \max(h(x_i), h'(x_i)) \right] \quad (69)$$

$$= \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \frac{1}{m} \sum_{i=1}^m \sigma_i \frac{1}{2} (h(x_i) + h'(x_i) + |h(x_i) - h'(x_i)|) \right] \quad (70)$$

$$= \frac{1}{2} [\widehat{\mathfrak{R}}_S(\mathcal{H}) + \widehat{\mathfrak{R}}_S(\mathcal{H}')] + \frac{1}{2} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h, h'} \frac{1}{m} \sum_{i=1}^m \sigma_i |h(x_i) - h'(x_i)| \right] \quad (\text{using Eq.(67)})$$

$$\leq \frac{1}{2} [\widehat{\mathfrak{R}}_S(\mathcal{H}) + \widehat{\mathfrak{R}}_S(\mathcal{H}')] + \frac{1}{2} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h, h'} \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) - \sigma_i h'(x_i) \right] \quad (\phi(t) = |t| \text{ is 1-Lipschitz})$$

$$= \frac{1}{2} [\widehat{\mathfrak{R}}_S(\mathcal{H}) + \widehat{\mathfrak{R}}_S(\mathcal{H}')] + \frac{1}{2} \mathbb{E}_h \left[\sup \frac{1}{m} \sum_{i=1}^m \sigma_i h(x_i) \right] + \frac{1}{2} \mathbb{E}_{h'} \left[\sup \frac{1}{m} \sum_{i=1}^m (-\sigma_i) h'(x_i) \right] \quad (71)$$

$$= \widehat{\mathfrak{R}}_S(\mathcal{H}) + \widehat{\mathfrak{R}}_S(\mathcal{H}') \quad (-\sigma_i \text{s are Rademacher r.v.s})$$

(e) Recall that $\mathbb{E}[\sigma_i \sigma_j] = 0$ if $i \neq j$ and $\mathbb{E}[\sigma_i^2] = 1$, we have

$$[\mathbb{E}_{\boldsymbol{\sigma}}(\boldsymbol{\sigma}^\top \mathbf{x})^2]^{\frac{1}{2}} = \left[\mathbb{E}_{\boldsymbol{\sigma}} \left(\sum_i \sigma_i x_i \right)^2 \right]^{\frac{1}{2}} = \left[\sum_i \mathbb{E}[\sigma_i^2] x_i^2 + 2 \sum_{i \neq j} \mathbb{E}[\sigma_i \sigma_j] x_i x_j \right]^{\frac{1}{2}} = \|\mathbf{x}\|_2. \quad (72)$$

(g) Note that for $a, b \in \{0, 1\}^2$, we have $ab = \max(0, a + b - 1)$. Then,

$$\widehat{\mathfrak{R}}_S(\mathcal{H}) = \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h_1 \in \mathcal{H}_1, h_2 \in \mathcal{H}_2} \sum_{i=1}^m \sigma_i h_1(x_i) h_2(x_i) \right] \quad (73)$$

$$= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h_1 \in \mathcal{H}_1, h_2 \in \mathcal{H}_2} \sum_{i=1}^m \sigma_i \max(0, h_1(x_i) + h_2(x_i) - 1) \right] \quad (74)$$

Let $g(t) = \max(0, t - 1)$ which is 1-Lipschitz. Using Talagrand's Lemma, we have:

$$\widehat{\mathfrak{R}}_S(\mathcal{H}) \leq \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h_1 \in \mathcal{H}_1, h_2 \in \mathcal{H}_2} \sum_{i=1}^m \sigma_i (h_1(x_i) + h_2(x_i)) \right] \quad (75)$$

$$= \widehat{\mathfrak{R}}_S(\mathcal{H}_1) + \widehat{\mathfrak{R}}_S(\mathcal{H}_2) \quad (76)$$

□