

Towards Understanding the Role of Over-Parametrization in Generalization of Neural Networks

Lecturer: Harrison Zhou

Scribe: Cynthia Yue, Feb 27, 2019

1 Introduction

In their 2018 paper *Towards Understanding the Role of Over-Parametrization in Generalization of Neural Networks*, Neyshabur et al. introduce a way to understand why neural networks generalize better with over-parametrization. They show that the difference between test and training error keeps decreasing as the number of parameters is increased. Specifically, they find tighter upper and lower bounds for the difference between test and training error. To do this, they involve the Rademacher complexity in their upper bound and are able to find a matching lower bound for the Rademacher complexity, therefore showing not only that this complexity decreases to 0 but also the speed at which it decreases.

They prove these bounds for the particular model of a two layer ReLU network, where $V \in \mathbb{R}^{1 \times h}$, $U \in \mathbb{R}^{h \times d}$, and $x \in \mathbb{R}^{d \times 1}$:

$$\begin{aligned} f_{V,U}(x) &= V[Ux]_+ \\ &= (v_1 \dots v_h) \left(\begin{pmatrix} U_1 \\ U_2 \\ \vdots \\ U_h \end{pmatrix} x \right)_+. \end{aligned}$$

They consider using this network for the classification task with data (x_i, y_i) i.i.d. with $i = 1, \dots, m$ where the label $y_i = \pm 1$. For this classification, they use the loss function of ramp loss:

$$l(f(x), y) = \begin{cases} 0 & f(x)y \geq 1 \\ 1 - f(x)y & 0 \leq f(x)y \leq 1 \\ 1 & f(x)y < 0. \end{cases}$$

They have expected loss $L(f(x), y) = \mathbb{E}(l(f(x), y))$ and empirical loss $\hat{L}(f(x), y) = \frac{1}{m} \sum_{i=1}^m l(f(x_i), y_i)$.

When finding bounds for the difference between test and training error, they look at the difference between test error $L(f(x), y)$ and training error $\hat{L}(f(x), y)$: $L(f(x), y) - \hat{L}(f(x), y)$.

2 Upper Bound

For the upper bound, they find the bound with respect to the following parameter space:

$$\mathcal{W} = \{(V, U) : \|v_j\| \leq \alpha_j, \|u_j - u_j^0\|_2 \leq \beta_j\}.$$

The upper bound for this parameter space is found to be:

$$L(f(x), y) - \hat{L}(f(x), y) \lesssim \frac{1}{m} \sum_{j=1}^h \alpha_j (\beta_j \|X\|_F + \|u_j^0 X\|_2) + \sqrt{\frac{\log \frac{2}{s}}{m}}. \quad (1)$$

In order to find this upper bound, there are two parts. We first show that the difference between test and training error is bounded above by the Rademacher complexity. We then show that the Rademacher complexity is bounded above by (1).

2.1 Difference between test and training error bounded above by Rademacher complexity

In order to show $L(f(x), y) - \hat{L}(f(x), y) \lesssim \frac{1}{m} \sum_{j=1}^h \alpha_j (\beta_j \|X\|_F + \|u_j^0 X\|_2) + \sqrt{\frac{\log \frac{2}{\delta}}{m}}$ (1), we first need to show that the difference between test and training error is bounded above by the Rademacher complexity:

$$L(f(x), y) - \hat{L}(f(x), y) \lesssim \frac{1}{m} \mathbb{E} \sup_l \sum_{i=1}^m \xi_i l(f(x_i), y_i) + \sqrt{\frac{\log \frac{2}{\delta}}{m}}, \quad (2)$$

where $\frac{1}{m} \mathbb{E} \sup_l \sum_{i=1}^m \xi_i l(f(x_i), y_i)$ is the Rademacher complexity, denoted \mathcal{R} , with ξ i.i.d. $= \begin{cases} 1 & \text{w.p. } \frac{1}{2} \\ -1 & \text{w.p. } \frac{1}{2} \end{cases}$.

Using a theorem about Rademacher complexity [2, Theorem 3.3]:

To do this, we first need the following theorem:

$$\mathbb{E} g(z) \leq \frac{1}{m} \sum_{i=1}^m g(z_i) + 2\mathcal{R}_m(\mathcal{G}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}$$

$$\text{where } \mathcal{R}_m(\mathcal{G}) = \frac{1}{m} \mathbb{E} \sup_{g \in \mathcal{G}} \sum_{i=1}^m \xi_i g(z_i).$$

Let us prove this. However, the proof involves an inequality called McDiarmid's inequality, so before we prove this, let us briefly discuss McDiarmid's inequality [4].

McDiarmid's inequality is a concentration inequality. Concentration inequalities are used when we want to show that some random quantity is close to its mean with high probability [1]. Some examples of concentration inequalities are Markov's inequality, which states that, for $Z \geq 0$ and $\epsilon > 0$, $\mathbb{P}(Z > \epsilon) \leq \frac{\mathbb{E}(Z)}{\epsilon}$, and Chebyshev's inequality, which states that $\mathbb{P}(|Z - \mathbb{E}Z| > \epsilon) \leq \frac{\text{Var}(Z)}{\epsilon^2}$. McDiarmid's inequality is stronger.

McDiarmid's inequality states: For Z_1, \dots, Z_m independent random variables and $i = 1, \dots, m$, suppose we have a function f that satisfies

$$\sup_{z_1, \dots, z_m, z'_i} \left| f(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_m) - f(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_m) \right| \leq c_i.$$

Then

$$\mathbb{P} \left(\left| f(Z_1, \dots, Z_m) - \mathbb{E}(f(Z_1, \dots, Z_m)) \right| \geq \epsilon \right) \leq 2 \exp \left(- \frac{2\epsilon^2}{\sum_{i=1}^m c_i^2} \right).$$

Now, let us use McDiarmid's inequality to actually prove the theorem about Rademacher complexity:

Proof. For any sample $S = (z_1, \dots, z_m)$ and any $g \in \mathcal{G}$, we first choose function f to be

$$f(S) = \sup_{g \in \mathcal{G}} \left(\mathbb{E}[g] - \hat{\mathbb{E}}_S[g] \right).$$

Now, let us have two samples S and S' , which differ by exactly one point, say z_m in S and z'_m in S' . Then,

$$\begin{aligned} f(S') - f(S) &\leq \sup_{g \in \mathcal{G}} \left(\hat{\mathbb{E}}_S[g] - \hat{\mathbb{E}}_{S'}[g] \right) \\ &= \sup_{g \in \mathcal{G}} \frac{g(z_m) - g(z'_m)}{m} \\ &\leq \frac{1}{m}. \end{aligned}$$

We can similarly show that $f(S) - f(S') \leq \frac{1}{m}$, which indicates that $|f(S) - f(S')| \leq \frac{1}{m}$. As a result, the c_i in $\sup_{z_1, \dots, z_m, z'_i} \left| f(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_m) - f(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_m) \right| \leq c_i$ can be considered to be $\frac{1}{m}$.

We can now apply McDiarmid's inequality to get: for any $\delta > 0$, with probability at least $1 - \delta/2$,

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

From here, let's break down the expectation on the right-hand side:

$$\begin{aligned} \mathbb{E}_S[f(S)] &= \mathbb{E}_S \left[\sup_{g \in \mathcal{G}} (\mathbb{E}[g] - \hat{\mathbb{E}}_S(g)) \right] \\ &= \mathbb{E}_S \left[\sup_{g \in \mathcal{G}} \mathbb{E}_{S'} [\hat{\mathbb{E}}_{S'}(g) - \hat{\mathbb{E}}_S(g)] \right] \\ &\leq \mathbb{E}_{S, S'} \left[\sup_{g \in \mathcal{G}} (\hat{\mathbb{E}}_{S'}(g) - \hat{\mathbb{E}}_S(g)) \right] \\ &= \mathbb{E}_{S, S'} \left[\sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m (g(z'_i) - g(z_i)) \right] \\ &= \mathbb{E}_{\sigma, S, S'} \left[\sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m \sigma_i (g(z'_i) - g(z_i)) \right] \\ &\leq \mathbb{E}_{\sigma, S'} \left[\sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m \sigma_i g(z'_i) \right] + \mathbb{E}_{\sigma, S} \left[\sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m -\sigma_i g(z_i) \right] \\ &= 2 \mathbb{E}_{\sigma, S} \left[\sup_{g \in \mathcal{G}} \frac{1}{m} \sum_{i=1}^m \sigma_i g(z_i) \right] \\ &= 2\mathcal{R}_m(\mathcal{G}). \end{aligned}$$

Now, replacing $\mathbb{E}_S[f(S)]$ with $2\mathcal{R}_m(\mathcal{G})$, we have:

$$f(S) \leq 2\mathcal{R}_m(\mathcal{G}) + \sqrt{\frac{\log \frac{2}{\delta}}{2m}}.$$

We chose f to be $\sup_{g \in \mathcal{G}} (\mathbb{E}[g] - \hat{\mathbb{E}}_S[g])$, so:

$$\sup_{g \in \mathcal{G}} (\mathbb{E}[g] - \hat{\mathbb{E}}_S[g]) \leq 2\mathcal{R}_m(\mathcal{G}) + \sqrt{\frac{\log \frac{2}{\delta}}{2m}}.$$

$\hat{\mathbb{E}}_S[g] = \frac{1}{m} \sum_{i=1}^m g(z_i)$, so:

$$\mathbb{E}g(z) \leq \frac{1}{m} \sum_{i=1}^m g(z_i) + 2\mathcal{R}_m(\mathcal{G}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}}.$$

□

Proving difference between test and training error bounded above by Rademacher complexity:

We now have that

$$\begin{aligned} \mathbb{E}g(z) &\leq \frac{1}{m} \sum_{i=1}^m g(z_i) + 2\mathcal{R}_m(\mathcal{G}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \\ \text{where } \mathcal{R}_m(\mathcal{G}) &= \frac{1}{m} \mathbb{E} \sup_{g \in \mathcal{G}} \sum_{i=1}^m \xi_i g(z_i). \end{aligned}$$

Let us now use the above in order to prove that the difference between test and training error is bounded above by Rademacher complexity

$$L(f(x), y) - \hat{L}(f(x), y) \lesssim \frac{1}{m} \mathbb{E} \sup_l \sum_{i=1}^m \xi_i l(f(x_i), y_i) + \sqrt{\frac{\log \frac{2}{s}}{m}}.$$

Proof. Choose $g = l(f(x), y)$.

$$\begin{aligned} \mathbb{E}l(f(x), y) &= L(f(x), y) \\ L(f(x), y) &\leq \frac{1}{m} \sum_{i=1}^m l(f(x_i), y_i) + 2\mathcal{R}_m(\mathcal{G}) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \\ L(f(x), y) &\leq \frac{1}{m} \sum_{i=1}^m l(f(x_i), y_i) + 2 * \frac{1}{m} \mathbb{E} \sup_{g \in \mathcal{G}} \sum_{i=1}^m \xi_i l(f(x_i), y_i) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \\ \hat{L}(f(x), y) &= \frac{1}{m} \sum_{i=1}^m l(f(x_i), y_i) \\ L(f(x), y) - \hat{L}(f(x), y) &\leq 2 * \frac{1}{m} \mathbb{E} \sup_{g \in \mathcal{G}} \sum_{i=1}^m \xi_i l(f(x_i), y_i) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} \\ &\lesssim \frac{1}{m} \mathbb{E} \sup_l \sum_{i=1}^m \xi_i l(f(x_i), y_i) + \sqrt{\frac{\log \frac{2}{s}}{m}}. \end{aligned}$$

□

2.2 Rademacher complexity bounded above by (1)

Now that we know that the difference between test and training error is bounded above by the Rademacher complexity, we can find an upper bound for the Rademacher complexity in order to get the proposed bound

$$L(f(x), y) - \hat{L}(f(x), y) \lesssim \frac{1}{m} \sum_{j=1}^h \alpha_j (\beta_j \|X\|_F + \|u_j^0 X\|_2) + \sqrt{\frac{\log \frac{2}{\delta}}{m}}.$$

In other words, if we assume $U_0 = 0$, we essentially need to show that

$$\mathbb{E} \sup_l \sum_{i=1}^m \xi_i l(f(x_i), y_i) \leq \sum_{j=1}^h \alpha_j (\beta_j \|X\|_F). \quad (3)$$

To do so, we will need to use the fact

$$\mathcal{R}(\phi \circ f) \leq a \mathcal{R}(f) \quad \text{if} \quad \text{Lip}(\phi) \leq a.$$

Proof. Using this fact, we have:

$$\begin{aligned} \mathbb{E} \sup_l \sum_{i=1}^m \xi_i l(f(x_i), y_i) &\leq \mathbb{E} \sup_{V, U} \sum_{i=1}^m \xi_i \sum_{j=1}^h v_j (u_j^T x_i)_+ \\ &= \mathbb{E} \sup_{V, U} \sum_{i=1}^m \xi_i \sum_{j=1}^h v_j \|u_j\|_2 \frac{(u_j^T x_i)}{\|u_j\|_2} + \\ &= \mathbb{E} \sup_{V, U} \sum_{j=1}^h v_j \|u_j\|_2 \sum_{i=1}^m \xi_i \left(\frac{u_j^T}{\|u_j\|_2} x_i \right)_+ \\ &\leq \mathbb{E} \sum_{j=1}^h \alpha_j \beta_j \sup_{u_j} \left| \sum_{i=1}^m \xi_i \left(\frac{u_j^T}{\|u_j\|_2} x_i \right)_+ \right| \\ &\leq \mathbb{E} \sum_{j=1}^h \alpha_j \beta_j \sup_{\|q_j\|_2 \leq 1} \left| \sum_{i=1}^m \xi_i (q_j^T x_i)_+ \right| \quad \text{where } q_j^T = \frac{u_j^T}{\|u_j\|_2} \\ &= \mathbb{E} \sum_{j=1}^h \alpha_j \beta_j \left[\sup_{\|q_j\|_2 \leq 1} \left(\sum_{i=1}^m \xi_i (q_j^T x_i)_+ \right) - \inf_{\|q_j\|_2 \leq 1} \left(\sum_{i=1}^m \xi_i (q_j^T x_i)_+ \right) \right] \\ &= \mathbb{E} \sum_{j=1}^h \alpha_j \beta_j \left[\sup_{\|q_j\|_2 \leq 1} \left(\sum_{i=1}^m \xi_i (q_j^T x_i)_+ \right) + \sup_{\|q_j\|_2 \leq 1} \left(\sum_{i=1}^m (-\xi_i) (q_j^T x_i)_+ \right) \right] \\ &= \mathbb{E} \sum_{j=1}^h \alpha_j \beta_j \sup_{\|q_j\|_2 \leq 1} \left(\sum_{i=1}^m \xi_i (q_j^T x_i)_+ \right) + \mathbb{E} \sum_{j=1}^h \alpha_j \beta_j \sup_{\|q_j\|_2 \leq 1} \left(\sum_{i=1}^m (-\xi_i) (q_j^T x_i)_+ \right) \\ &= 2 \mathbb{E} \sum_{j=1}^h \alpha_j \beta_j \sup_{\|q_j\|_2 \leq 1} \sum_{i=1}^m \xi_i (q_j^T x_i)_+. \end{aligned}$$

From here, we can write:

$$\begin{aligned} 2\mathbb{E} \sum_{j=1}^h \alpha_j \beta_j \sup_{\|q_j\|_2 \leq 1} \sum_{i=1}^m \xi_i(q_j^T x_i)_+ &\leq \sum_{j=1}^h \alpha_j \beta_j \mathbb{E} \sup_{\|q_j\|_2 \leq 1} \sum_{i=1}^m \xi_i(q_j^T x_i) \\ &= \sum_{j=1}^h \alpha_j \beta_j \mathbb{E} \sup_{\|q_j\|_2 \leq 1} q_j^T \sum_{i=1}^m \xi_i x_i. \end{aligned}$$

Now, we can look at the expectation on the right-hand side and show that it is the Frobenius norm:

$$\begin{aligned} \mathbb{E} \sup_{\|q_j\|_2 \leq 1} q_j^T \sum_{i=1}^m \xi_i x_i &= \mathbb{E} \left\| \sum_{i=1}^m \xi_i x_i \right\|_2 \\ &\leq \sqrt{\mathbb{E} \left\| \sum_{i=1}^m \xi_i x_i \right\|_2^2} \\ &= \|X\|_F. \end{aligned}$$

As a result, we have

$$\mathbb{E} \sup_l \sum_{i=1}^m \xi_i l(f(x_i), y_i) \leq \sum_{j=1}^h \alpha_j \beta_j \|X\|_F.$$

□

Hence, we find the overall upper bound to be

$$L(f(x), y) - \hat{L}(f(x), y) \lesssim \frac{1}{m} \sum_{j=1}^h \alpha_j (\beta_j \|X\|_F + \|u_j^0 X\|_2) + \sqrt{\frac{\log \frac{2}{s}}{m}}.$$

3 Lower Bound

For the lower bound, the paper finds a lower bound for the Rademacher complexity.

The lower bound is found to be:

$$\frac{1}{m} \mathbb{E} \sup_f \sum_{i=1}^m \xi_i f(x_i) \gtrsim \sum_{j=1}^h \alpha_j \beta_j \|X\|_F. \quad (4)$$

Recall that we have

$$f_{V,U}(x) = V[UX]_+.$$

This lower bound is obtained for a smaller parameter space with the specific case when $m = h = d = 2^k$ and x_i is the standard coordinate vector e_i for $i = 1, \dots, 2^k$. We have

$$V = (\alpha_1, \dots, \alpha_{2^k})$$

$$U = \begin{pmatrix} \beta_1 & & \\ & \ddots & \\ & & \beta_{2^k} \end{pmatrix} A_{2^k \times 2^k}$$

where $A_{2^k \times 2^k}$ is the Hadamard matrix.

The Hadamard matrix is a matrix, the elements of which are all $\pm 2^{-k/2}$. All of the columns of the Hadamard matrix are orthogonal and have norm 1, so they form an orthonormal basis of \mathbb{R}^{2^k} . Let us denote the j th column of $A_{2^k \times 2^k}$ as a_j . From this, we see that, as U is given by the columns a_1, \dots, a_{2^k} , we can say f is also given by the columns a_1, \dots, a_{2^k} .

Proof. To prove the lower bound, with this setup, we have:

$$\begin{aligned} f_{V,U}(x) &= V[UX]_+ \\ &= (\alpha_1, \dots, \alpha_{2^k}) \left(\begin{pmatrix} \beta_1 & & \\ & \ddots & \\ & & \beta_{2^k} \end{pmatrix} (a_1, \dots, a_{2^k}) e_i \right)_+ \end{aligned}$$

If we let $i = 1$, this becomes

$$\begin{aligned} &= (\alpha_1, \dots, \alpha_{2^k}) \left(\begin{pmatrix} \beta_1 & & \\ & \ddots & \\ & & \beta_{2^k} \end{pmatrix} a_1 \right)_+ \\ &= \sum_{j=1}^h \alpha_j \beta_j (a_{1j})_+ \end{aligned}$$

We can now look at the following:

$$\begin{aligned} &(\alpha_1, \dots, \alpha_{2^k}) \left(\begin{pmatrix} \beta_1 & & \\ & \ddots & \\ & & \beta_{2^k} \end{pmatrix} a_1 \right)_+ + (\alpha_1, \dots, \alpha_{2^k}) \left(\begin{pmatrix} \beta_1 & & \\ & \ddots & \\ & & \beta_{2^k} \end{pmatrix} -a_1 \right)_+ \\ &= \sum_{j=1}^h \alpha_j \beta_j \left[(a_{1j})_+ + (-a_{1j})_+ \right] \\ &= 2^{-k/2} \sum_{j=1}^h \alpha_j \beta_j. \end{aligned}$$

From this, we can see that either $\sum_{j=1}^h \alpha_j \beta_j (a_{1j})_+$ or $\sum_{j=1}^h \alpha_j \beta_j (-a_{1j})_+$ is greater than $2^{-k/2-1} \sum_{j=1}^h \alpha_j \beta_j$.

Now, let us pick a new $A^* = (a_1^*, \dots, a_{2^k}^*)$ for which

$$a_i^* = \begin{cases} a_i & \text{if } \sum_{j=1}^h \alpha_j \beta_j (a_{1j})_+ \geq 2^{-k/2-1} \sum_{j=1}^h \alpha_j \beta_j \\ -a_i & \text{if } \sum_{j=1}^h \alpha_j \beta_j (-a_{1j})_+ \geq 2^{-k/2-1} \sum_{j=1}^h \alpha_j \beta_j. \end{cases}$$

Now, with Rademacher complexity, when ξ_i is fixed, we can choose a new matrix \tilde{A} with columns $\tilde{a}_1, \dots, \tilde{a}_{2^k}$ to be used in a new function

$$\tilde{f}_{V,U}(x) = V[\tilde{U}x]_+$$

where

$$\tilde{U} = \begin{pmatrix} \beta_1 & & \\ & \ddots & \\ & & \beta_{2^k} \end{pmatrix} \tilde{A}_{2^k \times 2^k}$$

such that

$$\tilde{a}_i = \begin{cases} a_i^* & \text{if } \xi_i > 0 \\ 0 & \text{if } \xi_i \leq 0. \end{cases}$$

From here, we have that

$$\begin{aligned} \xi_i \tilde{f}(x_i) &= \xi_i \tilde{f}(e_i) \\ &= \xi_i(\alpha_1, \dots, \alpha_{2^k}) \left(\begin{pmatrix} \beta_1 & & \\ & \ddots & \\ & & \beta_{2^k} \end{pmatrix} \tilde{A}_{2^k \times 2^k} e_i \right)_+ \\ &= \xi_i \sum_{j=1}^h \alpha_j \beta_j (\tilde{a}_{ij})_+ \\ &= \begin{cases} \sum_{j=1}^h \alpha_j \beta_j (a_{ij}^*)_+ & \text{when } \xi_i > 0, \text{ which has probability } \frac{1}{2} \text{ by Rademacher} \\ 0 & \text{when } \xi_i \leq 0, \text{ which has probability } \frac{1}{2} \text{ by Rademacher.} \end{cases} \end{aligned}$$

Now, we can show:

$$\begin{aligned} \mathbb{E} \sup_f \sum_{i=1}^m \xi_i f(x_i) &\geq \mathbb{E} \sum_{i=1}^m \xi_i \tilde{f}(x_i) \\ &= \sum_{i=1}^m \mathbb{E} \xi_i \tilde{f}(x_i) \\ &= \sum_{i=1}^m \left(\frac{1}{2} * 0 + \frac{1}{2} * \sum_{j=1}^h \alpha_j \beta_j (a_{ij}^*)_+ \right) \\ &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^h \alpha_j \beta_j (a_{ij}^*)_+ \\ &\geq \frac{1}{2} \sum_{i=1}^m 2^{-k/2-1} \sum_{j=1}^h \alpha_j \beta_j \\ &= 2^{-k/2-2} m \sum_{j=1}^h \alpha_j \beta_j \\ &= 2^{-2} m \sum_{j=1}^h \alpha_j \beta_j 2^{-k/2} \quad (\text{Note that } 2^{-k/2} = \frac{\|\mathbf{X}\|_F}{m}.) \\ &= \frac{1}{4} m \sum_{j=1}^h \alpha_j \beta_j \frac{\|\mathbf{X}\|_F}{m} \\ &= \frac{1}{4} \sum_{j=1}^h \alpha_j \beta_j \|\mathbf{X}\|_F. \end{aligned}$$

As a result, we have

$$\frac{1}{m} \mathbb{E} \sup_f \sum_{i=1}^m \xi_i f(x_i) \gtrsim \sum_{j=1}^h \alpha_j \beta_j \|X\|_F.$$

□

Overall, we have an upper bound

$$L(f(x), y) - \hat{L}(f(x), y) \lesssim \frac{1}{m} \sum_{j=1}^h \alpha_j (\beta_j \|X\|_F + \|u_j^0 X\|_2) + \sqrt{\frac{\log \frac{2}{s}}{m}}$$

and a lower bound (for the Rademacher complexity)

$$\frac{1}{m} \mathbb{E} \sup_f \sum_{i=1}^m \xi_i f(x_i) \gtrsim \sum_{j=1}^h \alpha_j \beta_j \|X\|_F.$$

We can also see that we have matching upper and lower bounds for the Rademacher complexity.

References

- [1] John Lafferty, Han Liu, and Larry Wasserman. *Concentration of Measure*, chapter 7. Carnegie Mellon University, 2008-2010.
- [2] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2012.
- [3] Behnam Neyshabur, Zhiyuan Li, Srinadh Bhojanapalli, Yann LeCun, and Nathan Srebro. Towards understanding the role of over-parametrization in generalization of neural networks, 2018.
- [4] Galen Reeves. Concentration inequalities: Hoeffding and mcdiarmid, 2008.