

## Lecture 7: Property of Two Layer Neural Network

July 16

Lecturer: Chao Ma

Scribe: Yuncong Peng, Yilei Han, Jingsong Zhang

## 1 Estimation of NN vs. RF at a limited time scale

Consider NN vs. RF

(1)  $m \geq O(n^2)$ , NN is always close to RF

(2)  $m < O(n^2)$ , NN is only close to RF at the beginning

Assume target function as  $f^*$

$$f^* = \int_{S^{d-1}} a^* \sigma(b^T x) d(\pi_0(b))$$

$$\gamma(f^*) = \max\{1, \sup |a^*(b)|\} < \infty$$

- RF can learn  $f^*$  well,
- $NN \approx RF, t \in [0, t_0]$
- NN with early stopping can also learn well.

**Theorem 1.1** (RF).  $\forall \delta$ , with prob  $\geq 1 - \delta$  over the chance of  $\{b_0\}$  there exist  $a^*$ , s.t.

$$R(a^*, B) \leq \frac{\gamma^2(f^*)}{m} (1 + \sqrt{2 \log(\frac{1}{\delta})})^2$$

$$\|a^*\|_1 \leq \frac{\gamma(f^*)}{\sqrt{m}}$$

Notations:

$$f(x; \mathbf{a}, B) = \sum_{k=1}^m a_k \sigma(b_k^T x)$$

$$\mathcal{R}(\mathbf{a}, B) = \|f(x; \mathbf{a}, B) - f^*\|^2$$

$$\hat{\mathcal{R}}(\mathbf{a}, B) = \frac{1}{n} \sum_{i=1}^n (f(x_i; \mathbf{a}, B) - f^*(x_i))^2$$

**Theorem 1.2.** For a two-layer Neural Network, which is defined as  $f(x; \mathbf{a}, B) = \sum_{k=1}^m a_k \sigma(b_k^T x)$ , with the assumption below:

- $b_k(0) \sim \pi_0$  (same as Random Feature)
- $a_k(0) = \pm\beta, \beta = \frac{C}{m}$
- $\|f^*\|_\infty \leq 1$

Then  $\forall \delta > 0$ , with  $p \geq 1 - 4\delta$ , we have

$$\hat{\mathcal{R}}_n(\mathbf{a}_t, B_t) \leq C\left(\frac{1}{m} + \frac{1}{mt} + \frac{1}{\sqrt{n}}\right)$$

**proof:**

Let  $\hat{\rho} = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}$ , then we have

$$\begin{aligned} \hat{\mathcal{R}}_n(\mathbf{a}_t, B_t) &= \|f(\cdot; \mathbf{a}_t, B_t) - f^*(\cdot)\|_{\hat{\rho}}^2 \\ &\leq 3 \left( \|f(\cdot; \mathbf{a}_t, B_t) - f(\cdot; \tilde{\mathbf{a}}_t, B_t)\|_{\hat{\rho}}^2 + \|f(\cdot; \tilde{\mathbf{a}}_t, B_t) - f(\cdot; \tilde{\mathbf{a}}_t, B_0)\|_{\hat{\rho}}^2 + \hat{\mathcal{R}}_n(\tilde{\mathbf{a}}_t, B_0) \right) \end{aligned}$$

By Cauchy-Schwartz, we have

$$\begin{aligned} \|f(\mathbf{x}; \mathbf{a}_t, B_t) - f(\mathbf{x}; \tilde{\mathbf{a}}_t, B_t)\|_{\hat{\rho}}^2 &\leq \|\mathbf{a}_t - \tilde{\mathbf{a}}_t\|^2 \|B_t\|^2 \\ \|f(\mathbf{x}; \tilde{\mathbf{a}}_t, B_t) - f(\mathbf{x}; \tilde{\mathbf{a}}_t, B_0)\|_{\hat{\rho}}^2 &\leq \|\tilde{\mathbf{a}}_t\|^2 \|B_t - B_0\|^2 \end{aligned}$$

For  $\hat{\mathcal{R}}_n(\tilde{\mathbf{a}}_t, B_0)$ , with probability  $1 - \delta$ , there exists  $\mathbf{a}^*$  as a good estimator. Thus we have

$$\begin{aligned} \hat{\mathcal{R}}_n(\tilde{\mathbf{a}}_t, B_0) &= \left( \hat{\mathcal{R}}_n(\tilde{\mathbf{a}}_t, B_0) - \hat{\mathcal{R}}_n(\mathbf{a}^*, B_0) \right) + \left( \hat{\mathcal{R}}_n(\mathbf{a}^*, B_0) - \mathcal{R}(\mathbf{a}^*, B_0) \right) + \mathcal{R}(\mathbf{a}^*, B_0) \\ &=: I_1 + I_2 + I_3 \end{aligned}$$

As the *generalization error*, We can bound  $I_2$  as follows:

$$I_2 \leq \frac{2(2\sqrt{m}\|\mathbf{a}^*\| + 1)^2}{\sqrt{n}} \left( 1 + \sqrt{2 \ln \left( \frac{2}{\delta} \left( \|\mathbf{a}^*\| + \frac{1}{\|\mathbf{a}^*\|} \right) \right)} \right)$$

For  $I_1$ , consider the Lyapunov function

$$J(t) = t \left( \hat{\mathcal{R}}_n(\tilde{\mathbf{a}}_t, B_0) - \hat{\mathcal{R}}_n(\mathbf{a}^*, B_0) \right) + \frac{1}{2} \|\tilde{\mathbf{a}}_t - \mathbf{a}^*\|^2$$

Since  $\hat{\mathcal{R}}_n(\tilde{\mathbf{a}}_t, B_0)$  is convex with respect to  $\tilde{\mathbf{a}}_t$ , we have  $\frac{d}{dt} J(t) \leq 0$ , which implies  $J(t) \leq J(0)$ . Hence we have

$$\hat{\mathcal{R}}_n(\tilde{\mathbf{a}}_t, B_0) \leq \hat{\mathcal{R}}_n(\mathbf{a}^*, B_0) + \frac{\|\mathbf{a}_0 - \mathbf{a}^*\|^2}{2t}$$

Estimate the norm of parameters.

$$\begin{aligned}
& \frac{d}{dt} (a_t - \tilde{a}_t) \\
&= -\frac{1}{n} \sum_{i=1}^n [(f(x_i, a_t, B_t) - y_i) \tau(B_H, x_i) - (f(x_i; \tilde{a}_t, B_0) - y_i) \tau(B(0)x_i)] \\
&= -\frac{1}{n} \sum_{i=1}^n (a_k^T \tau(B_t x_i) \tau(B_t x_i) - \tilde{a}_t^T \tau(B_0 x_i) \tau(B_0 x_i)) + \frac{1}{n} \sum_{i=1}^n y^T (\tau(B_t x_i) - \tau(B_0 x_i)) \\
&= -\frac{1}{n} \sum_{i=1}^n \tau(B_t x_i) \tau(B_t x_i)^T (a_t - \tilde{a}_t) + \frac{1}{n} \sum_{i=1}^n \left( \tau(B_t x_i) \tau(B_t x_i)^T - \tau(B_0 x_i) \tau(B_0 x_i)^T \right) \tilde{a}_t \\
&\quad + \frac{1}{n} \sum_{i=1}^n y^T (\tau(B_t x_i) - \tau(B_0 x_i))
\end{aligned}$$

Then

$$\begin{aligned}
& \frac{d}{dt} \|a_t - \tilde{a}_t\|^2 \\
&= 2(a_t - \tilde{a}_t)^T \frac{d}{dt} (a_t - \tilde{a}_t) \\
&\leq - (a_t - \tilde{a}_t)^T \frac{1}{n} \sum_{i=1}^n \tau(B_t x_i) \tau(B_t x_i)^T (a_t - \tilde{a}_t) \\
&\quad + (a_t - \tilde{a}_t)^T \left[ \frac{2}{n} \sum_{i=1}^n \left( \tau(B_t x_i) \tau(B_t x_i)^T - \tau(B_0 x_i) \tau(B_0 x_i)^T \right) \right] \tilde{a}_t \\
&\quad + (a_t - \tilde{a}_t)^T \left[ \frac{2}{n} \sum_{i=1}^n y^T (\tau(B_t x_i) - \tau(B_0 x_i)) \right] \\
&\leq (a_t - \tilde{a}_t)^T \left[ \frac{2}{n} \sum_{i=1}^n \left( \tau(B_t x_i) \tau(B_t x_i)^T - \tau(B_0 x_i) \tau(B_0 x_i)^T \right) \right] \tilde{a}_t \\
&\quad + (a_t - \tilde{a}_t)^T \left[ \frac{2}{n} \sum_{i=1}^n y^T (\tau(B_t x_i) - \tau(B_0 x_i)) \right] \\
&\leq 2 \|B_t - B_0\| (\|B_t\| \|\tilde{a}_t\| + \|B_0\| \|\tilde{a}_t\| + 1) \|a_t - \tilde{a}_t\|
\end{aligned}$$

Then estimate the bound of  $\|a_t\|$ ,  $\|B_t - B_0\|$ ,  $\|\tilde{a}_t\|$ ,  $\|a_t - \tilde{a}_t\|$

**Lemma 1.3.** *Let  $T$  is Constant,  $\exists C_T$  s.t.  $\forall 0 \leq t \leq T$ ,*

$$\begin{aligned}\|a_t\| &\leq C_T \left( \frac{C}{\sqrt{m}} + \sqrt{mt} \right) \\ \|B_t\| &\leq C_T \left( \frac{ct}{\sqrt{m}} + \sqrt{m} \right) \\ \|B_t - B_0\| &\leq C_T(c+1) \left( \frac{ct}{\sqrt{m}} + \frac{\sqrt{m}}{2}t^2 \right)\end{aligned}$$

**Lemma 1.4.** *Assume that  $m \geq r^2, \forall \epsilon > 0$ , with probability  $\geq 1 - 4\epsilon$ ,  $\forall t \in [0, T]$*

$$\|\tilde{a}_t\| \leq \tilde{C}_T \left( \frac{1}{\sqrt{m}} + \frac{t}{\sqrt{m}} + \frac{\sqrt{t}}{n^{\frac{1}{4}}} \right)$$

**Lemma 1.5.** *With assumptions above,  $\forall t \in [0, T]$*

$$\|a_t - \tilde{a}_t\| \leq C_T \frac{t^2}{m} (1 + mt)(t + m) \left( \frac{1}{\sqrt{m}} + \frac{\sqrt{t}}{\sqrt{m}} + \frac{\sqrt{t}}{n^{\frac{1}{4}}} \right)$$

With all above, we have

$$\begin{aligned}\hat{\mathcal{R}}_n(\mathbf{a}_t, B_t) &\leq C \left( \frac{1}{m} + \frac{1}{mt} + \frac{1}{\sqrt{n}} \left( 1 + \sqrt{t} + \frac{\sqrt{mt}}{n^{1/4}} \right)^2 \right. \\ &\quad \left. + \frac{t^2}{m^2} (1 + mt)^2 \left( 1 + \frac{t^2}{m^2} (t + m)^4 \right) \left( 1 + \sqrt{t} + \frac{\sqrt{mt}}{n^{1/4}} \right)^2 \right)\end{aligned}$$

for  $t \in [0, T]$ , and some constant  $C$ .

If we assume  $m \geq n$ , and take  $t \in \left[0, \frac{\sqrt{n}}{m}\right]$ , then we can take  $T = 1$  and obtain

$$\hat{\mathcal{R}}(a_t, B_t) \leq C \left( \frac{1}{m} + \frac{1}{mt} + \frac{1}{\sqrt{n}} \right) \quad \forall 0 \leq t \leq 1$$

One can refer to *A Comparative Analysis of the Optimization and Generalization Property of Two-layer Neural Network and Random Feature Models Under Gradient Descent Dynamics* (arXiv:1904.04326v1) for more details.