

## Lecture 14 Mean Field Analysis of Neural Networks

Instructor: Yingyu Liang

Date:

Scriber: Yiyou Sun

## 1 Continuous Setting

Consider the traditional classification task where  $x \in \mathbb{R}^d, y \in \mathbb{R}$ . The goal is to find a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , such that:

$$\min_f Q(f) = L(f) + R(f), L(f) = E_{x,y}[l(f(x)), y],$$

where  $l(\cdot)$  is defined to be the loss function and  $R$  is a regularization function. Similar to Kernel methods, consider the two-level network given below to represent  $f$ :

$$f(\omega, \rho, x) = \int_{\mathbb{R}^d} \sigma(\theta, x) \omega(\theta) \rho(\theta) d\theta \quad (1)$$

where  $\sigma(\theta, x) : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a known real-valued function,  $\omega(\theta) : \mathbb{R}^d \rightarrow \mathbb{R}$  is a real value function of  $\theta$ , and  $\rho(\theta)$  is a probability density over  $\theta$ . For regularizer, we use

$$R(\omega, \rho) = \lambda_1 R_1(\omega, \rho) + \lambda_2 R_2(\rho)$$

where

$$R_1(\omega, \rho) = \int r_1(\omega(\theta)) \rho(\theta) d\theta, r_1(\omega) = |\omega|^2$$

$$R_2(\rho) = \int r_2(\theta) \rho(\theta) d\theta, r_2(\theta) = \|\theta\|^2$$

Next, we show a discrete NN approximates the continuous one when hidden nodes go to infinity and then drive the evolution rule of  $\rho(\theta)$  and  $\omega(\theta)$  from the (noisy) GD algorithm when the step size becomes small.

## 2 Discrete Setting

Consider a finite NN with the following form to approximate  $f(\omega, \rho, x)$ :

$$\hat{f}(\mu, \Theta, x) = \frac{1}{m} \sum_{j=1}^m \mu^j \sigma(\theta^j, x) \quad (2)$$

where  $\Theta = \{\theta^j\}_{j=1}^m$ .

The regularization terms are:

$$\hat{R}_1(\mu, \Theta) = \frac{1}{m} \sum_{j=1}^m r_1(\mu^j), \hat{R}_2(\Theta) = \frac{1}{m} \sum_{j=1}^m r_2(\theta^j), \quad (3)$$

and the training objective is:

$$\widehat{Q}(\mu, \Theta) = \mathbb{E}_{x,y} l(\widehat{f}(\mu, \Theta, x), y) + \lambda_1 \widehat{R}_1(\mu, \Theta) + \lambda_2 \widehat{R}_2(\Theta). \quad (4)$$

We can solve it through the standard (noisy) GD, the algorithm is given by:

**Step 0.** Initialize  $\mu_0 \sim P_{\mu,0}(\mu), \theta_0^j \sim P_{\theta,0}(\theta)$

**Step 1.** Update  $\theta^j$  by

$$\theta_{t+\Delta t}^j = \theta_t^j - \Delta t \nabla_{\theta^j} \left[ \widehat{Q}(\mu_t, \Theta_t) \right] - \sqrt{\lambda_3} \zeta_{t+1}^j,$$

where  $\Delta t$  is the step size and  $\zeta_{t+1}^j \sim N(0, \sqrt{2\Delta t} I_d)$ .

**Step 2.** Update  $\mu^j$  by

$$\mu_{t+\Delta t}^j = \mu_t^j - \Delta t \nabla_{\mu^j} \left[ \widehat{Q}(\mu_t, \Theta_t) \right] - \sqrt{\lambda_3} \zeta_{t+1}^j,$$

where  $\zeta_{t+1}^j \sim N(0, \sqrt{2\Delta t})$ .

## 2.1 Plain GD

We first consider the unnoisy setting where  $\lambda_3 = 0$ . We have the following Lemma.

**Lemma 1.** Suppose  $\lambda_3 = 0$ . Suppose at time  $t \geq 0$ , we have  $\theta_t^j \sim \rho_t$ , and suppose  $\mu_t^j = \omega_t(\theta_t^j)$ . Assume  $l'$  is continuous and  $\sigma$  is twice differentiable. For all  $x$ , we have:

$$\lim_{m \rightarrow \infty} \widehat{f}(\mu_t, \Theta_t, x) = f(\omega_t, \rho_t, x) \quad (5)$$

Furthermore, when  $\Delta t \rightarrow 0, m \rightarrow \infty$ , we can derive,

$$\begin{aligned} \frac{d\rho_t(\theta)}{dt} &= -\nabla_{\theta} \cdot [\rho_t(\theta) g_2(t, \theta, \omega_t(\theta))] \\ \frac{d\omega_t(\theta)}{dt} &= g_1(t, \theta, \omega_t(\theta)) - \nabla_{\theta} [\omega_t(\theta)] g_2(t, \theta, \omega_t(\theta)), \end{aligned}$$

where  $\nabla_{\theta} \cdot$  means the divergence,  $g_1$  and  $g_2$  satisfy:

$$\begin{aligned} g_1(t, \theta, v) &= -\mathbb{E}_{x,y} [l'(f(\omega_t, \rho_t, x), y) \sigma(\theta, x)] - \lambda_1 \nabla_v [r_1(v)] \\ g_2(t, \theta, v) &= -\mathbb{E}_{x,y} [l'(f(\omega_t, \rho_t, x), y) v \nabla_{\theta} \sigma(\theta, x)] - \lambda_2 \nabla_{\theta} [r_2(\theta)]. \end{aligned}$$

To prove the lemma, we utilize the tool with Fokker-Planck Equation to compute the evolution.

**Background with Fokker-Planck Equation** Suppose the movement of a particle in  $m$ -dimensional space can be characterized by the stochastic differential equation given below:

$$dx_t = g(x_t, t) dt + \sqrt{2\beta^{-1}}\Sigma dB_t$$

Let  $x_t \sim p(x, t)$ , the evolution of  $p(x, t)$  is given by:

$$\frac{\partial p(x, t)}{\partial t} = \frac{\Sigma\Sigma^\top}{\beta}\nabla^2 p(x, t) - \nabla \cdot [p(x, t)g(x_t, t)]$$

**Proof of Lemma 1.** Let the  $p_t(\theta, v)$  as the joint distribution for  $(\theta, v)$ :

$$(\theta_t^j, \mu_t^j) \sim p_t(\theta, v) = \rho_t \delta(v = \omega_t(\theta))$$

We can rewrite  $f(\omega_t, \rho_t, x)$  as:

$$f(\omega_t, \rho_t, x) = \int_{\mathbb{R}^{d+1}} \sigma(\theta, x) p_t(\theta, v) d\theta dv.$$

By the Law of the Large number, when  $m \rightarrow \infty$ ,

$$\widehat{f}(\mu_t, \Theta_t, x) \rightarrow f(\omega_t, \rho_t, x).$$

Now we denote

$$\widehat{g}_2(t, \theta, v) = -\mathbb{E}_{x,y} \left[ l' \left( \widehat{f}(\mu_t, \Theta_t, x), y \right) v \nabla_\theta \sigma(\theta, x) \right] - \lambda_2 \nabla_\theta [r_2(\theta)]$$

From the update rule of GD, we have  $\theta_{t+1}^j = \theta_t^j + \Delta t \widehat{g}_2(t, \theta_t^j, \mu_t^j)$ . Let  $\Delta t \rightarrow 0$ , using  $\mu_t^j = \omega_t(\theta_t^j)$ , we have

$$\frac{d\theta_t^j}{dt} = \widehat{g}_2(t, \theta_t^j, \omega_t(\theta_t^j))$$

By applying Fokker-Planck equation,

$$\frac{d\rho_t(\theta)}{dt} = -\nabla_\theta \cdot [\rho_t(\theta) \widehat{g}_2(t, \theta, \omega_t(\theta))]$$

As  $m \rightarrow \infty$ , and because  $l'$  is continuous,  $\sigma(\theta, x)$  and  $\rho_t$  are also second-order smooth, we obtain:

$$\nabla_\theta \cdot [\rho_t(\theta) \widehat{g}_2(t, \theta, \omega_t(\theta))] - \nabla_\theta \cdot [\rho_t(\theta) g_2(t, \theta, \omega_t(\theta))] \xrightarrow{\text{a.s.}} 0$$

To prove the evolution form for  $\omega_t(\theta)$ , we let:

$$\widehat{g}_1(t, \theta, v) = -\mathbb{E}_{x,y} \left[ l' \left( \widehat{f}(\mu_t, \Theta_t, x), y \right) \sigma(\theta, x) \right] - \lambda_1 \nabla_v r_1(v).$$

Then, (ignoring the superscript  $j$  since all  $j$  have the same calculation)

$$\begin{aligned}
& \omega_{t+\Delta t}(\theta_{t+\Delta t}) \\
&= \omega_t(\theta_{t+\Delta t}) + \frac{d\omega_t(\theta_{t+\Delta t})}{dt} \Delta t + o(\Delta t) \\
&= \omega_t(\theta_t + \widehat{g}_2(t, \theta_t, \omega_t(\theta_t)) \Delta t + o(\Delta t)) + \frac{d\omega_t(\theta_{t+\Delta t})}{dt} \Delta t + o(\Delta t) \\
&= \omega_t(\theta_t) + [\nabla_{\theta} \omega_t(\theta_t)] \cdot \widehat{g}_2(t, \theta_t, \omega_t(\theta_t)) \Delta t + \frac{d\omega_t(\theta_{t+\Delta t})}{dt} \Delta t + o(\Delta t)
\end{aligned}$$

By the update rule  $\omega_{t+\Delta t}(\theta_{t+\Delta t}) = \omega_t(\theta_t) + \widehat{g}_1(t, \theta_t, \omega_t(\theta)) \Delta t$ , we have:

$$\lim_{\Delta t \rightarrow 0} \frac{d\omega_t(\theta_{t+\Delta t})}{dt} = \widehat{g}_1(t, \theta_t, \omega_t(\theta)).$$

The proof is finished by Let  $\Delta t \rightarrow 0$ , and let  $m \rightarrow \infty$ .

## References