

STAT 946 - Topics in Probability and Statistics: Mathematical
Foundations of Deep Learning
Lecture 12
Professor Mufan Li

Lucas Noritomi-Hartwig
University of Waterloo

Ocotber 20, 2025 from 16h00 to 17h20 in M3 3103

Probability Seminar Mondays 11h30

SOLT November 10: 09h00 - 16h30 Lecture November 10: Starts 16h30

- October 27th: Gustavo Ruihan
- November 17th: Kevin Edward

Project presentations: November 24, Nov 26, Dec 1

Feature learning

- Non-kernel (data independent)
- Non-linear with respect to parameters

Recall: 1-layer network (NTK scaling)

$$f(x; \theta) = \frac{1}{\sqrt{n}} W_1 \varphi(W_0 x), \quad W_1 \in \mathbb{R}^{1 \times n}, W_0 \in \mathbb{R}^{n \times n_0}, x \in \mathbb{R}^{n_0 \times 1}$$

Data: $((x^\alpha, y^\alpha))_{\alpha=1}^m$

MSE:

$$\mathcal{L}(\theta) = \frac{1}{2m} \sum_{\alpha=1}^m (f(x^\alpha; \theta) - y^\alpha)^2$$

Training: $\theta_{k+1} = \theta_k - \eta \nabla \mathcal{L}(\theta_k)$, $\eta > 0$ is constant.

Hidden: $h^\alpha(\theta) = W_0 x^\alpha$.

$$\begin{aligned} \Delta h^\alpha &= h^\alpha(\theta_1) - h^\alpha(\theta_0) \\ &= -\frac{\eta}{m\sqrt{n}} \sum_{\beta=1}^m (f(x^\beta; \theta_0) - y^\beta) \langle x^\alpha, x^\beta \rangle W_1^\top \odot \varphi'(h^\beta) \end{aligned}$$

As $n \rightarrow \infty$,

$$\begin{aligned}\Delta h^\alpha &= \Theta\left(\frac{1}{\sqrt{n}}\right) \quad \text{v.s.} \quad h^\alpha(\theta_0) = \Theta(1) \\ \Delta f^\alpha &\approx -\frac{\eta}{m} K^\alpha(f - y) = \Theta(1)\end{aligned}$$

NTK:

$$\begin{aligned}K^{\alpha\beta} &= \langle \nabla_\theta f^\alpha, \nabla_\theta f^\beta \rangle \\ &= \langle \nabla_{W_1} f^\alpha, \nabla_{W_1} f^\beta \rangle + \langle \nabla_{W_0} f^\alpha, \nabla_{W_0} f^\beta \rangle \\ &= \frac{1}{n} \left\langle \varphi(h^\alpha), \varphi\left(\underbrace{h^\beta}_{\Phi_1^{\alpha\beta}}\right) \right\rangle + \underbrace{\langle x^\alpha, x^\beta \rangle}_{\Phi_0^{\alpha\beta}} \underbrace{\frac{1}{n} \sum_{j=1}^n \varphi'(h_j^\alpha) \varphi'(h_j^\beta)}_{\text{at init. } G_1^{\alpha\beta}}\end{aligned}$$

Need $\Delta h^\alpha = \Theta(1)$ for NTK to evolve \rightarrow feature learning.

Naïve method: Increase learning rate $\eta = \eta_0 \sqrt{n}$, η_0 is constant. Non-NTK.

$$\begin{aligned}\Rightarrow \Delta h^\alpha &= \Theta(1) \quad \text{works!} \\ \Delta f^\alpha &= \Theta(\sqrt{n}) \quad \text{diverges - does not work!}\end{aligned}$$

Compensate: Scale down output.

$$\begin{aligned}f(x; \theta) &= \frac{1}{n} W_1 \varphi(W_0 x) \\ \Rightarrow \Delta h^\alpha &= \Theta\left(\frac{1}{\sqrt{n}}\right) \\ \Delta f^\alpha &= \Theta\left(\frac{1}{\sqrt{n}}\right) \rightarrow 2 \times \text{effect} \\ K^{\alpha\beta} &= \left\langle \underbrace{\nabla_\theta f^\alpha}_{\frac{1}{\sqrt{n}}}, \underbrace{\nabla_\theta f^\beta}_{\frac{1}{\sqrt{n}}} \right\rangle\end{aligned}$$

Mean Field Parameterization

Prefactor $\frac{1}{n}$, learning rate: $\eta = \eta_0 n$

$$\begin{aligned}\Rightarrow \Delta h^\alpha &= \Theta(1) \quad \text{works!} \quad \text{accelerated} \\ \Delta f^\alpha &= \Theta(1) \quad \text{works!}\end{aligned}$$

Effectively three choices:

- NTK:
 - Prefactor: $\frac{1}{\sqrt{n}}$ (or 1)
 - Learning rate: 1 (or $\frac{1}{\sqrt{n}}$)
 - Initial standard deviation: 1 (or $\frac{1}{\sqrt{n}}$)

- Mean field:
 - Prefactor: $\frac{1}{n}$
 - Learning rate: n
 - Initial standard deviation: 1

This kind of equivalence is called ABC-reparameterization (Yang and Hu 2020).

Recall in the NTK: $\Delta f \simeq -\frac{\eta}{m} K(f - y)$

This is the same in the mean field, however, K is now evolving.

What is the state space of neural network training?

- Full state space: $\theta_{k+1} = \theta_k - \eta \nabla \mathcal{L}(\theta_k)$
- NTK: f
- Neurons: h, g (in between full and NTK)

We will focus on the full state space.

$\theta = (W_1, W_0) \in \mathbb{R}^{1 \times n} \times \mathbb{R}^{n \times n_0}$. As $n \rightarrow \infty$, θ becomes infinitely dimensional.

The way to write this is:

$$\begin{aligned} f(x; \theta) &= \frac{1}{n} \sum_{j=1}^n W_{1,j} \varphi(\langle W_{0,j}, x \rangle) \\ &= \int w_1 \varphi(\langle w_0, x \rangle) d\rho^{(n)}(w_1, w_0) \end{aligned}$$

where

$$\rho^{(n)} = \frac{1}{n} \sum_{j=1}^n \delta_{(W_{1,j}, W_{0,j})}$$

is the empirical measure.

Recall that $\int q(w) d\delta_{w_0}(w_0) = q(w_0)$.

Back to gradient flow:

$$\partial_t \theta(t) = -\eta \nabla \mathcal{L}(\theta(t))$$

where $\eta = \eta_0 n$.

$$\partial_t W_{1,j} = -\frac{\eta_0}{m} \sum_{\beta=1}^m \left(f_{(t)}^\beta - y^\beta \right) \varphi \left(\left\langle W_{0,j}^{(t)}, x^\beta \right\rangle \right)$$

where

$$\begin{aligned} f_{(t)}^\beta &= \int w_1 \varphi(\langle w_0, x^\beta \rangle) d\rho_t^{(n)} \\ \partial_t W_{0,j} &= -\frac{\eta_0}{m} \sum_{\beta=1}^m \left(f_{(t)}^\beta - y^\beta \right) \varphi \left(\left\langle W_{0,j}^{(t)}, x^\beta \right\rangle \right) x^\beta \end{aligned}$$

$$\begin{aligned}\theta_j &= (W_{1,j}, W_{0,j}) \\ \partial_t \theta_j(t) &= b\left(\theta_j(t), \rho_t^{(n)}\right)\end{aligned}$$

Recall: propagation of chaos.

As $n \rightarrow \infty$, $\rho_t^{(n)} \rightarrow \rho_t \leftarrow \underbrace{\mathcal{L}(\theta_j(t))}_{n \text{ finite}}$ where $\rho_t^{(n)}$ is all particles, and $\theta_j(t)$ is a single particle.

ρ_t is called the mean field measure/distribution, $\rho_t^{(n)}$ is called the mean field ODE (McKean-Vlasov).

Stronger notion of propagation of chaos:

$$\mathcal{L}((\theta_{j_1}, \dots, \theta_{j_k})) \rightarrow \rho_t^{\otimes k}$$

k -particles are independent!

Bound $W_2(\mathcal{L}((\theta_{j_1}, \dots, \theta_{j_k})), \rho_t^{\otimes k}) \lesssim \frac{1}{n} e^{c(t, n_0)}$.

Mean Field PDE

- Cannot differentiate $\rho_t^{(n)}$ or δ_{θ_j} .
- Trick: use test function $q \in C_c^\infty(\mathbb{R}^{n_0+1})$.

$$\int q(\theta) \underbrace{\text{“}\nabla_{\theta} \rho_t^{(n)}(\theta)\text{”}}_{\text{“weak derivative”}} d\theta \stackrel{\text{IBP}}{=} - \int \nabla_{\theta} q(\theta) d\rho_t^{(n)}(\theta)$$

$$\begin{aligned}\partial_t \int q(\theta) d\rho_t^{(n)}(\theta) &= \partial_t \frac{1}{n} \sum_{j=1}^n q(\theta_j(t)) \\ &= \frac{1}{n} \sum_{j=1}^n \langle \nabla_{\theta} q(\theta_j(t)), \partial_t \theta_j(t) \rangle = \frac{1}{n} \sum_{j=1}^n \langle \nabla_{\theta} q(\theta_j(t)), b(\theta_j(t), \rho_t^{(n)}) \rangle \\ &= \int \nabla_{\theta} q(\theta) b(\theta, \rho_t^{(n)}) d\rho_t^{(n)} \\ \underbrace{\int q(\theta) \partial_t \rho_t^{(n)} d\theta}_{\text{not specified}} &\stackrel{\text{IBP}}{=} \int q(\theta) \left(-\operatorname{div}_{\theta} \left(b(\theta, \rho_t^{(n)}) \rho_t^{(n)} \right) \right) d\theta\end{aligned}$$

We say that $\rho_t^{(n)}$ is a weak solution of

$$\begin{cases} \partial_t \rho_t = -\operatorname{div}(b(\theta, \rho_t) \rho_t) & \text{first order non-linear “Transport equation”} \\ \rho_0 = \rho_0^{(n)} \end{cases}$$

Mean field PDE.

References:

- Mei, Montanari, Nguyen (2018)
- Chizat, Bach (2018)
- Nitanda, Suzuki (2017) Japan
- Rotskoff, Vanden-Eijnden (2018) New York
- Sirigano, Spiliopoulos (2018)