

STAT 946 - Topics in Probability and Statistics:
Mathematical Foundations of Deep Learning Lecture
13: Feature Learning

1 Feature Learning (Yang & Hu, 2021)

Consider an L -layer fully-connected network parameterized in the *maximal update parametrization* (also called μP). The model is defined as:

$$f(x; \theta) = \left[\frac{1}{n} \right]_{(1 \times n)} W_d \varphi(h_d),$$

where the hidden representations satisfy

$$h_{l+1} = \frac{1}{\sqrt{n}} W_l \varphi(h_l), \quad l = 1, \dots, d-1$$

with input layer

$$h_1 = \frac{1}{\sqrt{n_0}} \frac{x}{(n \times n_0)}.$$

The learning rate scales with width:

$$\boxed{\eta = \eta_0 n}.$$

This scaling is known as the **Maximal Update Parametrization** (μP), which ensures stable feature learning as width $n \rightarrow \infty$.

Definition:

- Δh_l : the change in the hidden representation h_l after one gradient step.
- We say that layer l *learns features* if $\Delta h_l = \Theta(1)$ (i.e., the change is order-one in width).

Remark:

$$\Delta h_l \simeq \frac{1}{\sqrt{n}} \Delta W_l \varphi(h_l) + \frac{1}{\sqrt{n}} W_l \Delta \varphi(h_l).$$

- If we **only train** the input weights W_0 , then $\Delta h_1 = \Theta(1)$.

- If we **freeze** all deeper layers (i.e., $\Delta W_l = 0$ for $l \geq 1$),

$$\implies \Delta h_l = \Theta(1), \quad \forall l.$$

Definition: A parameterization (weight prefactors, learning rates, etc.) is said to satisfy maximal update if

$$\frac{1}{\sqrt{n}} \Delta W_l \varphi(h_l) = \Theta(1), \quad \forall l.$$

Remark: If $n \rightarrow \infty$ is the only limit, then μP is the unique parameterization where all W_l contributes to feature learning

Recall:

$$[\theta_1, \theta_2, \dots, \theta_n] / \stackrel{\text{permutation invariant}}{\sim} \implies \frac{1}{n} \sum_{i=1}^n \delta_{\theta_i} \implies \text{mean-field network.}$$

Consider $d = 2$, linear

$$\begin{aligned} f(x; \theta) &= \frac{1}{n} W_2 \frac{1}{\sqrt{n}} W_1 \frac{1}{\sqrt{n_0}} W_0 x \\ &= \frac{1}{n} \sum_{\boxed{i=1}}^n W_{2,i} \frac{1}{\sqrt{n}} \sum_{\boxed{j=1}}^n \boxed{W_{1,ij}} \frac{1}{\sqrt{n_0}} \langle W_{0,j}, x \rangle. \end{aligned}$$

↓ permutation invariant
permutation invariant

$W_1 \in \mathbb{R}^{n \times n}$ is a separately exchangeable array (when W_1, W_0 are not trained and i.i.d.).

$$[W_{1,ij}(t)]_{i,j} \stackrel{d}{=} [W_{1,\sigma(i)\sigma'(j)}(t)]_{i,j}, \quad \sigma, \sigma' \in S_n.$$

Does not admit an explicit quotient representation.

Aldous–Hoover Representation. There exists a measurable function

$$f : [0, 1]^4 \rightarrow \mathbb{R} \quad \text{such that} \quad [W_{1,ij}]_{i,j} \stackrel{d}{=} [f(U, U_i, V_j, U_{ij})]_{i,j},$$

where

$$U, U_i, V_j, U_{ij} \stackrel{\text{i.i.d.}}{\sim} \text{Unif}([0, 1]).$$

Remark. No mean-field PDEs exist for deep networks under μ_P .

State Space

Mean Field (MF): θ is full state \implies no finite representation,

NTK: f outputs \implies poor feature learning model,

DMFT: $[f^\alpha, h^\alpha, z^\alpha]_{\alpha=1}^m \implies$ depends on $m = \#\text{data}$.

2 History and Refs

Statistical Physics - Martin, Siggia, Rose (1973)
 Recurrent NN - Sompolinsky, Crisanti, Sommers (1988) } Not rigorous

Spin Glass - Ben-Arous, Guionnet (1995)

High-Dimensional SGD – Celentano, Cheng, Montanari (2021)

NN ($n \rightarrow \infty$):

- Yang and Hu (2021)
- Bordelon and Pehlevan (2022) \rightarrow Not rigorous

3 Dynamical MF Theory for Training Dynamics

Recall:

$$\boxed{h_{l+1}^\alpha} = \frac{1}{\sqrt{n}} W_l \varphi(h_l^\alpha), \quad (\text{post-activation})$$

$$g_l^\alpha = \sqrt{n} \frac{\partial f}{\partial h_l^\alpha},$$

$$\boxed{z_l^\alpha} = \frac{1}{\sqrt{n}} W_l^\top g_{l+1}^\alpha = \frac{1}{\sqrt{n}} W_l^\top \text{diag}(\varphi'(h_{l+1}^\alpha)) z_{l+1}^\alpha, \quad (\text{pre-activation})$$

$h_{l+1}^\alpha, z_l^\alpha$ can be viewed as states

Markov View.

$$\theta(k) \xrightarrow{\text{GD}} \theta(k+1)$$

$$h(k) \xrightarrow{\text{GD}} h(k+1)$$

$$z(k) \xrightarrow{\text{GD}} z(k+1)$$

However, note that the weights $\theta(k+1)$ rely on information from $\theta(k)$ (excluding h, z).

$$\theta(k) \xrightarrow{n \rightarrow \infty} \text{depends only on } \{h^\alpha(t), z^\alpha(t)\}_{t \leq k, \alpha \in [m]}.$$

Good: Only neuron-wise states needed as $n \rightarrow \infty$.

Bad: Requires full history ($t \leq k$) \Rightarrow non-Markovian.

Full DMFT Equations

Time derivative:

$$\delta_t \theta(t) = -\eta \nabla_\theta L(\theta(t)), \quad \eta = n.$$

The equations:

$$h_l^\alpha(t) = u_l^\alpha(t) + \int_0^t ds \sum_{\beta=1}^n \left[A_l^{\alpha\beta}(t, s) + \Delta^\beta(s) \Phi_{l-1}^{\alpha\beta}(t, s) \right] z_l^\alpha(s) \varphi'(h_l^\alpha(s)).$$

where:

- h_l^α is a single neuron state (scalar, i -th index),
- $u_l^\alpha(t) \sim \text{GP}\left(0, [\Phi_{l-1}^{\alpha\beta}(t, s)]_{\alpha\beta}\right)$,
- $A_l^{\alpha\beta}(t, s) = \frac{1}{n} \sum_{j=1}^n \frac{\partial \varphi(h_l^\alpha(t))_j}{\partial r_l^\beta(s)_j}$,
- $\Delta^\beta(s) = -\frac{\partial L}{\partial f^\beta} = y^\beta - f^\beta$ (for square loss),
- $\Phi_{l-1}^{\alpha\beta}(t, s) = \lim_{n \rightarrow \infty} \frac{1}{n} \left\langle \varphi(h_{l-1}^\alpha(t)), \varphi(h_{l-1}^\beta(s)) \right\rangle,$

- $z_l^\alpha(s)$ is a single-neuron (scalar, i -th index) adjoint / backward state.

$$z_l^\alpha(t) = r_l^\alpha(t) + \int_0^t ds \sum_{\beta=1}^n \left[B_l^{\alpha\beta}(t, s) + \Delta^\beta(s) G_{l+1}^{\alpha\beta}(t, s) \right] \varphi(h_l^\alpha(s)).$$

where:

- z_l^α is a single-neuron (scalar, i -th index) backward state,
- $r_l^\alpha(t) \sim \text{GP}\left(0, [G_{l-1}^{\alpha\beta}(t, s)]_{\alpha\beta}\right)$,
- $B_l^{\alpha\beta}(t, s) = \frac{1}{n} \sum_{i=1}^n \frac{\partial g_{l+1}^\alpha(t)_i}{\partial u_l^\beta(s)_i}$,
- $G_{l+1}^{\alpha\beta}(t, s) = \lim_{n \rightarrow \infty} \frac{1}{n} \langle g_{l+1}^\alpha(t), g_{l+1}^\beta(s) \rangle$, (average over neuron index)
 $\Leftrightarrow \int \cdots d\rho, \quad d\rho = \frac{1}{n} \sum_{i=1}^n \delta_{(h_{l,i}, z_{l,i})}$.

Remarks: Mean Field (particle) equations

- RHS (dynamics) only depends on (h, z) and history \Rightarrow the system is closed
- u, r, A, B (depend on history) contributed by weights

Gaussian conditioning:

$$\begin{aligned} W \mid (W\varphi, g^\top W) &\stackrel{d}{=} (*) \\ (*) &= \underbrace{P_g W + W P_\varphi - P_g W P_\varphi}_{\substack{g, \varphi \text{ contain history} \\ \implies A, B \text{ kernels}}} + \underbrace{P_g^\perp \tilde{W} P_\varphi^\perp}_{\substack{\text{Gaussian process (GP)}}}. \end{aligned}$$

where:

- P_g^\perp corresponds to the backward kernel G ,
- \tilde{W} is an independent copy of W ,

- P_φ^\perp corresponds to the forward kernel Φ .

$$h_l(t) = \text{GP} + \underbrace{\int_0^t ds \Theta(1)}_{\substack{\text{feature learning / maximal update} \\ \Delta W_{l-1} \varphi(h_{l-1})}}$$

4 Heuristic Derivation for Scaling (Not DMFT)

- For simplicity: $m = 1$ (single data point), $\varphi(x) = x$ (linear).
- $f = \frac{1}{\gamma\sqrt{n}} W_d h_d$
- $h_{l+1} = \frac{1}{\sqrt{n}} W_l h_l$
- $h_1 = W_0 X$
- $\partial_t \theta = -\eta \nabla_\theta L(\theta)$

Goal: Set γ, η as functions of n .

Recall: $a_n \sim b_n$ if $\frac{a_n}{b_n} \rightarrow \text{const.}$

i.e. $a_n = \Theta(n^p) \iff b_n = \Theta(n^p)$.

$$W_l(t) = W_l(0) + \frac{\eta}{\gamma n} \int_0^t ds \Delta(s) z_{l+1}(s) h_l(s)^\top.$$

$$\begin{aligned} h_{l+1}(t) &= \underbrace{\frac{1}{\sqrt{n}} W_l(0) h_l(t)}_{\substack{\text{CLT scaling} \\ (\text{GP term})}} \\ &+ \frac{1}{\sqrt{n}} \cdot \frac{\eta}{\gamma n} \int_0^t ds \Delta(s) z_{l+1}(s) \underbrace{h_l(s)^\top h_l(s)}_{\substack{\frac{1}{n} h^\top h = \Phi \\ (\text{kernel term})}}. \\ &\sim \text{GP} + \underbrace{\frac{\eta}{\gamma\sqrt{n}}}_{\sim 1} \int_0^t ds. \end{aligned}$$

$$h_{d+1} = \frac{1}{\sqrt{n}}\, W_d\, h_d, \qquad f = \frac{1}{\gamma}\, h_{d+1}.$$

$$K~=~\left<\nabla_\theta h_{d+1},\,\nabla_\theta h_{d+1}\right>~\sim~1.$$

$$\left<\nabla_\theta f,\,\nabla_\theta f\right> ~=~ \frac{1}{\gamma^2}{\,}K ~\sim~ \frac{1}{\gamma^2}.$$

$$\partial_t f = -\eta \left<\nabla_\theta f,\,\nabla_\theta f\right> \Delta = \underbrace{\frac{\eta}{\gamma^2}}_{\sim 1} {\,}\underbrace{K\Delta}_{\sim 1}.$$

$$\implies \eta \sim \gamma^2$$

$$\frac{n}{\gamma} ~\sim~ \gamma \qquad \implies \qquad \boxed{\gamma ~\sim~ \sqrt{n} \quad \implies \quad \eta ~\sim~ n} \quad \Rightarrow \quad \mu\mathrm{P}.$$