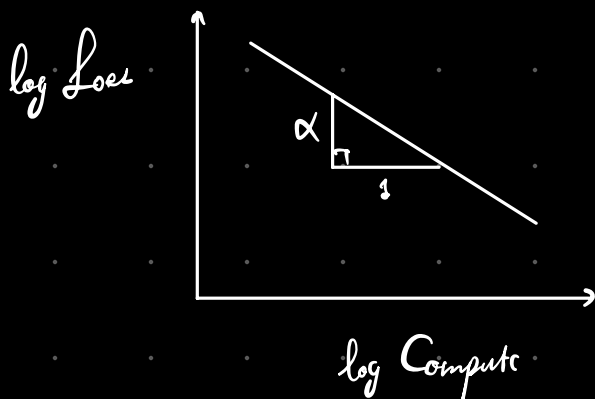


Proportional Scaling Limit

- Scaling law



$$\Rightarrow \text{Loss} = \text{const.} \times \text{Compute}^{-\alpha}$$

- What if we interpret

Compute \rightarrow index $n \in \mathbb{N}$

NN \rightarrow element $x_n \in \mathcal{X}$ in "space of NNs"

Loss \rightarrow test function $L: \mathcal{X} \rightarrow \mathbb{R}$

(morally speaking)

• Then does $x_n \rightarrow x^*$ w/ rate $n^{-\alpha}$

imply $L(x_n) \rightarrow L(x^*)$ w/ rate $n^{-\alpha}$?

Example • Let $\alpha > 0$. $X_i \stackrel{\text{iid}}{\sim} \rho : \mathbb{E} X_i = 0, \mathbb{E} X_i^2 = 1,$

and $\mathbb{E} |X|^{2(1+\alpha)} < \infty$ is the highest moment that exists.

$$\Rightarrow S_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n X_i \xrightarrow{d} \mathcal{N}(0, 1)$$

at rate $n^{-\alpha}$

Remark • Maybe we should study the scaling limit as $\text{Compute} \rightarrow \infty$.

• Roughly speaking

$$\text{Compute} \propto \underbrace{\text{Width}^2}_{\substack{\downarrow \\ \# \text{ param.} \\ \text{per layer}}} \times \underbrace{\text{Depth}}_{\substack{\downarrow \\ \# \text{ layers}}} \times \underbrace{\text{Batch}}_{\substack{\downarrow \\ \# \text{ data} \\ \text{per iter.}}} \times \underbrace{\text{Train Iter.}}_{\substack{\downarrow \\ \# \text{ iter.}}}$$

⇒ want to study a joint limit of width, depth, data, # iter.
 need to specify relationship "easier" "hard"

• In this talk, we will study the limit as

$$\begin{cases} \text{width } n, \text{ depth } d \rightarrow \infty. \\ \text{ratio } \frac{d}{n} \rightarrow \bar{\tau} > 0. \end{cases} \rightarrow \text{"proportional"} \quad (\text{w/o skip connections})$$

Warm Up

• Linear network. one sample $x \in \mathbb{R}^{n_0}$ \rightarrow input dim.

$$h_1 = \frac{1}{\sqrt{n_0}} \overset{n \times n_0}{W_0} \overset{n_0 \times 1}{x} \quad W_{e,ij} \overset{\text{iid}}{\sim} \mathcal{N}(0,1)$$

$$h_{e+1} = \frac{1}{\sqrt{n}} \overset{n \times n}{W_e} \overset{n \times 1}{h_e} \quad \ell \in [d]$$

$$f = \frac{1}{\sqrt{n}} \overset{1 \times n}{W_d} \overset{n \times 1}{h_d}$$

Goal • Analyze $L(h_e)$, $L(f)$ as $d, n \rightarrow \infty$.

Trick • Condition on the "right" object

• In this case $\mathcal{F}_\ell := \sigma(\{h_k\}_{k \leq \ell}) \rightarrow$ Filtration

"deterministic" \rightarrow all the prev. layers

$$h_{e+1} | \mathcal{F}_e = \frac{1}{\sqrt{n}} \underbrace{W_e}_{\text{id} \sim \mathcal{N}(0,1)} h_e | \mathcal{F}_e$$

$$\sim \mathcal{N}\left(0, \underbrace{\frac{1}{n} \|h_e\|^2}_{=: \Phi_e} \cdot I_n\right)$$

\Rightarrow key quantity

$$\begin{aligned} \Phi_{e+1} &= \frac{1}{n} \|h_{e+1}\|^2 \\ &= \frac{1}{n} \left\| \frac{1}{\sqrt{n}} W_e h_e \right\|^2 \\ &= \frac{1}{n} \left\| \underbrace{W_e \frac{h_e}{\|h_e\|}}_{\mathcal{N}(0, I_n)} \right\|^2 \underbrace{\frac{1}{n} \|h_e\|^2}_{\Phi_e} \\ &= \frac{1}{n} \chi_n^2 \Phi_e \quad \rightarrow \text{note } \Phi_{e+1} | \mathcal{F}_e = \Phi_{e+1} | \sigma(\Phi_e) \\ &= \Phi_e + \underbrace{\Phi_e \left(\frac{1}{n} \chi_n^2 - 1 \right)}_{\Rightarrow \Phi_{e+1} \text{ is a Markov chain!}} \end{aligned}$$

$$\underbrace{\Phi_{e+1}}_{\substack{\mathbb{E} = 0, \text{Var} = 2/n \\ \hookrightarrow \mathbb{E} = 0, \text{Var} = 1}} =: \Phi_e + \sqrt{\frac{2}{n}} \Phi_e \xi_e$$

Euler discretization of $d\Phi_t = \sqrt{2} \Phi_t dB_t$

w/ step size $\frac{1}{n}$.

(Geometric Brownian motion)
 $\sim \exp(\mathcal{N}(-\tau, 2\tau))$.

$l \mapsto l+1$ one step \Rightarrow total d steps (step size $\frac{1}{n}$)

\Rightarrow total depth time $\bar{\tau} := \frac{d}{n}$, current time $\tau = \frac{l}{n}$.

hence depth/width should be proportional.

• Furthermore observe that

$$f | \mathcal{F}_d \sim \mathcal{N}(0, \Phi_d) \Rightarrow \text{Sample } \Phi_d \rightarrow \Phi_{\bar{z}} \rightarrow \text{terminal time}$$

$$h_{t+1} | \mathcal{F}_t \sim \mathcal{N}(0, \Phi_t) \quad \Phi_t \rightarrow \Phi_{\bar{c}}$$

Remark • Can write down a limiting Fokker-Planck equation for $\mathcal{L}(h_{t,i})$ or $\frac{1}{n} \sum_{i=1}^n \delta_{h_{t,i}}$

Multiple Data Points

• Let $x = [x^1 \dots x^m] \in \mathbb{R}^{n \times m} \rightarrow \# \text{ data}$

$$h_{t+1}^\alpha = \frac{1}{\sqrt{n}} W_t h_t^\alpha \quad \alpha \in [m]$$

$$h_t = [h_t^1 \dots h_t^m]$$

• Let $u, v \in \mathbb{R}^n$ be deterministic, then

$$\begin{bmatrix} W_t u \\ W_t v \end{bmatrix} \sim \mathcal{N}\left(0, \begin{bmatrix} \|u\|^2 & \langle u, v \rangle \\ \langle u, v \rangle & \|v\|^2 \end{bmatrix} \otimes I_n\right)$$

where $A \otimes B = [a_{ij} B]$

$$\Rightarrow \text{vec}(h_{t+1}) | \mathcal{F}_t \stackrel{d}{=} \mathcal{N}\left(0, \left[\frac{1}{n} \langle h_t^\alpha, h_t^\beta \rangle \right]_{\alpha, \beta=1}^m \otimes I_n\right) \quad (*)$$

$$\text{vec}: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^{nm}$$

$$\frac{1}{n} h_t^\top h_t =: \Phi_t \in \mathbb{R}^{m \times m}$$

Feature covariance matrix.

$$\Rightarrow \Phi_{t+1}^{\alpha\beta} | \mathcal{F}_t = \frac{1}{n} \langle h_{t+1}^\alpha, h_{t+1}^\beta \rangle | \mathcal{F}_t$$

note $\mathbb{E} \left[\frac{1}{n} \langle h_{t+1}^\alpha, h_{t+1}^\beta \rangle \mid \mathcal{F}_t \right] = \mathbb{E} \left[h_{t+1,i}^\alpha h_{t+1,i}^\beta \mid \mathcal{F}_t \right]$
 $= \Phi_t^{\alpha\beta}$... by (*)

Then $\Phi_{t+1}^{\alpha\beta} = \Phi_t + \frac{1}{\sqrt{n}} \underbrace{\frac{1}{\sqrt{n}} \sum_{i=1}^n (h_{t+1,i}^\alpha h_{t+1,i}^\beta - \Phi_t^{\alpha\beta})}_{\mathbb{E}[\cdot \mid \mathcal{F}_t] = 0}$

$\mathbb{E}[\cdot \mid \mathcal{F}_t] = 0$
 $\text{Var}[\cdot \mid \mathcal{F}_t] = \text{Var}(h_{t+1,i}^\alpha h_{t+1,i}^\beta \mid \mathcal{F}_t)$
 $= \text{Id}(1)$

$\Phi_{t+1}^{\sigma\delta} = \dots \dots \dots (h_{t+1,i}^\sigma h_{t+1,i}^\delta - \Phi_t^{\sigma\delta})$

$\Rightarrow \text{Cov}(h^\alpha h^\beta, h^\sigma h^\delta)$ where $h \sim \mathcal{N}(0, \Phi_t)$

$= \mathbb{E} h^\alpha h^\beta h^\sigma h^\delta - \mathbb{E} h^\alpha h^\beta \mathbb{E} h^\sigma h^\delta$

\downarrow Isserlis/Wick
 $\mathbb{E} h^\alpha h^\beta \mathbb{E} h^\sigma h^\delta + \Phi^{\alpha\sigma} \Phi^{\beta\delta} + \Phi^{\alpha\delta} \Phi^{\beta\sigma}$

$= \Phi^{\alpha\sigma} \Phi^{\beta\delta} + \Phi^{\alpha\delta} \Phi^{\beta\sigma} = \sum (\Phi)^{\alpha\beta, \sigma\delta}$
 $m(m+1)/2 \times m(m+1)/2$

$\Rightarrow \Phi_{t+1} \mid \mathcal{F}_t \stackrel{d}{=} \Phi_t + \frac{1}{\sqrt{n}} \sum (\Phi_t)^{1/2} \xi_t$
 \downarrow flatten to \mathbb{R} $m(m+1)/2$ $\mathbb{E} = 0$ $\text{Cov} = \text{Id}_{m(m+1)/2}$

$\Rightarrow \{\Phi_{\lfloor n\tau \rfloor}\} \xrightarrow{d, n \rightarrow \infty} \begin{cases} d\Phi_\tau = \sum (\Phi_\tau)^{1/2} dB_\tau, & \tau \in [0, \bar{c}] \\ \Phi_0 = \frac{1}{n_0} x^\top x. \end{cases}$
 \downarrow Q.U. on $\mathbb{R}^{m(m+1)/2}$

Remark • There is a matrix form

$$\text{Let } \Sigma_\tau := \frac{1}{\sigma^2} (\tau \Sigma_\tau + \tau \Sigma_\tau^T)$$

B.M on $\text{Sym}(m)$

↪ Brownian motion on $\mathbb{R}^{m \times m}$

$$\Rightarrow d\Phi_\tau = \Phi_\tau^{1/2} d\Sigma_\tau \Phi_\tau^{1/2}$$

\downarrow
 $\mathbb{R}^{m \times m}$

(Optional)

• $\Sigma(\Phi)$ is the affine-invariant metric

i.e. let $\mathcal{M} = \text{SPD}(m)$, $T_p \mathcal{M} = \text{Sym}(m)$

and Riemannian metric

$$g_\Phi(A, B) = \frac{1}{2} \text{Tr}(A \Phi^{-1} B \Phi^{-1})$$

\downarrow
 $\text{Sym}(m)$

$\Rightarrow \Phi_\tau$ is a Brownian motion on (\mathcal{M}, g)
(with a specific dual connection ∇^*).

Nonlinear Network

$$\begin{aligned} \bullet h_{t+1}^\alpha &= \sqrt{\frac{c}{n}} \sum_i \tau_{ic} \varphi(h_{t,i}^\alpha), & \varphi: \mathbb{R} &\rightarrow \mathbb{R} \\ c^{-1} &:= \mathbb{E} \varphi(w)^2, & w &\sim \mathcal{N}(0, 1) \end{aligned}$$

(Optional)

• Why c ? Let $\varphi(x) = \max(x, 0)$ (ReLU), $c=1$

then $h_{t+1}^\alpha \mid \mathcal{F}_t \sim \mathcal{N}(0, \frac{1}{n} \|\varphi(h_t^\alpha)\|^2 \otimes I_n)$

$$\mathbb{E} \left[\frac{1}{n} \|\varphi(h_{t+1}^\alpha)\|^2 \mid \mathcal{F}_t \right] = \mathbb{E} \left[\varphi(h_{t+1}^\alpha)^2 \mid \mathcal{F}_t \right]$$

$$\downarrow \Phi_{t+1}^\alpha = \frac{1}{2} \mathbb{E} \left[(h_{t+1}^\alpha)^2 \mid \mathcal{F}_t \right]$$

$$= \frac{1}{2} \Phi_l^{\alpha\alpha}$$

$$\Rightarrow \mathbb{E}[\Phi_{l+1}^{\alpha\alpha} | \mathcal{F}_l] = \frac{1}{2} \Phi_l^{\alpha\alpha}$$

$$\searrow 0 \text{ as } l \rightarrow \infty.$$

• Let $\varphi_l^\alpha := \varphi(h_l^\alpha)$

then $\text{vec}(h_{l+1}) | \mathcal{F}_l \sim \mathcal{N}(0, \underbrace{\left[\frac{c}{n} \langle \varphi_l^\alpha, \varphi_l^\beta \rangle \right]_{\alpha\beta}}_{=: \Phi_l} \otimes I_n)$.

Proposition • Let $\rho_l^{\alpha\beta} = \frac{\Phi_l^{\alpha\beta}}{(\Phi_l^{\alpha\alpha} \Phi_l^{\beta\beta})^{1/2}}$,

$$\text{then } \mathbb{E}[\rho_{l+1}^{\alpha\beta} | \mathcal{F}_l] = \frac{1}{\pi} \left[\sqrt{1 - \rho_l^{\alpha\beta 2}} + \rho_l^{\alpha\beta} (\pi - \arccos(\rho_l^{\alpha\beta})) \right]$$

$$=: \mathcal{T}(\rho_l^{\alpha\beta})$$

↓
Cho and Saul kernel

and \mathcal{T} has a stable fixed point at 1.

• In other words, $\rho_l^{\alpha\beta} \rightarrow 1$ as $l \rightarrow \infty$.

Remark • This is sometimes known as "rank collapse".

(since $\text{rank}(h_l) \rightarrow 1$)

and is known to cause vanishing gradients.

Solution • Shape the activation function $\varphi(x) \rightarrow x$ as $d \rightarrow \infty$.

Intuition: weaken a fixed point iteration $\mathcal{T} \circ \mathcal{T} \circ \mathcal{T} \dots$

to an Euler discretization $\mathcal{T}(\rho) \approx \rho + \dots$

e.g. for ReLU, let $\psi_S(x) = S_+ \max(x, 0) + S_- \min(x, 0)$
 where $S_{\pm} = 1 \pm \frac{c_{\pm}}{n^p}$ for some $p > 0$ TBD.

$$\Rightarrow \Phi_{t+1} | \mathcal{F}_t = \Phi_t + \frac{1}{n^p} \left(\overline{\Phi_t^{\alpha\alpha} \Phi_t^{\beta\beta}} \nu(\rho_t^{\alpha\beta}) \right)_{\alpha\beta} \\
 + \frac{1}{\sqrt{n}} \sum (\Phi_t)^{1/2} \sum_{\alpha\beta} \rightarrow b(\Phi_t) \\
 + \text{higher order terms}$$

where $\nu(\rho) := \frac{(c_+ - c_-)^2}{2\pi} \left(\sqrt{1 - \rho^2} - \rho \arccos(\rho) \right)$

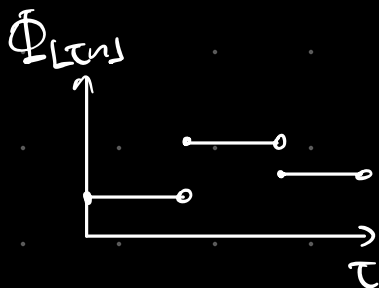
$\Rightarrow \rho = \frac{1}{2}$ matches the step size $\frac{1}{n}$

Theorem (LNR22)

- As $d, n \rightarrow \infty$, $\frac{d}{n} \rightarrow \bar{c}$, $\rho = \frac{1}{2}$
- $\{\Phi_{L(n)}\}_{\tau \in [\bar{c}]}$ converges weakly in the Skorohod topology of $\mathbb{R}^{m(m+1)/2}$ to the solution of

$$\begin{cases} d\Phi_{\tau} = b(\Phi_{\tau}) d\tau + \sum (\Phi_{\tau})^{1/2} dB_{\tau}, \\ \Phi_0 = \frac{1}{n_0} x^T x. \end{cases}$$

Remark



- Skorohod handles the discontinuities in the process.

• Can extend to other architectures e.g. Shaped Transformers.

$$h_{t+1} = \frac{1}{\sqrt{d}} W_e^V h_e A_e \quad [NLL+23]$$

$$A_e = \text{Softmax} \left(\frac{1}{s} (W_e^Q h_e)^T (W_e^K h_e) \right) - \frac{1}{m} \mathbb{1}\mathbb{1}^T + I_m$$

$\hookrightarrow s \sim nm_k$

$$\Rightarrow h_{t+1} \sim \text{linear} + \mathcal{O}\left(\frac{1}{n}\right)$$

• Can extend to training (my research talk, arXiv ~ 1-2 months)

$$W \in \mathbb{R}^{n \times n} \text{ iid } \mathcal{N}(0,1), \quad \varphi, g \in \mathbb{R}^{n \times n} \text{ det.}$$

$$\Rightarrow W \left(\sigma(W\varphi, g^T W) \right)$$

$$\stackrel{\text{d}}{=} W P_\varphi + P_g W - P_g W P_\varphi + P_g^\perp \tilde{W} P_\varphi^\perp$$

proj. onto $\text{col}(\varphi)$

$I - P_g^\perp$

indep. copy of W .

• Isolates the dependence structure in a clean way.

• Can calculate the spectrum of Φ in the linear case. (arXiv this week)

$$\text{specifically let } \rho := \frac{1}{m} \sum_i \delta_{\lambda_i(\Phi_d)},$$

$$\text{then as } d, n, m \rightarrow \infty, \quad \frac{dm}{n} \rightarrow \bar{c}$$

$$\rho \xrightarrow{d} \exp_{\#} \left(\mathcal{N}_{\bar{c}}^{\text{sc}} \boxplus \text{Unif}_{[-\bar{c}, 0]} \right) \boxtimes \rho_0$$

push forward

additive free conv.

multiplicative free conv.

Semicircle law

free log-normal