

STAT 946 - Topics in Probability and Statistics: Mathematical
 Foundations of Deep Learning
 Lecture 15
 Professor Mufan Li

Lucas Noritomi-Hartwig
 University of Waterloo

November 5, 2025 from 16h00 to 17h20 in M3 3103

Recall: From Hanin and Nica (2019) “Products of Many Random Matrices”

Result:

$$\begin{aligned} h_{l+1} &= \frac{1}{\sqrt{n}} W_l h_l \quad W_{l,jk} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1) \\ h_1 &= \frac{1}{\sqrt{n_0}} W_0 x \\ \log \left(\frac{1}{n} |h_d|^2 \right) - \log \left(\frac{1}{n_0} |x|^2 \right) &\xrightarrow{d, n \rightarrow \infty, \frac{d}{n} \rightarrow \bar{\tau}} \mathcal{N}(-\bar{\tau}, 2\bar{\tau}) \end{aligned}$$

This is the law of geometric Brownian motion:

$$dX_\tau = \sqrt{2} X_\tau dB_\tau$$

Li (this Professor Mufan Li), Nica, Roy (2022) “Neural Covariance SDEs . . .”

$$\begin{aligned} \text{Recall: } h_{l+1} \mid \mathcal{F}_l \left(= \sigma \left((h_k)_{k \leq l} \right) \right) &\stackrel{d}{=} \mathcal{N} \left(0, \underbrace{\frac{1}{n} |h_l|^2}_{\Phi_l} \otimes I_n \right) \\ \Phi_{l+1} &= \Phi_l \frac{1}{n} |\xi_l|^2, \quad \xi_l \sim \mathcal{N}(0, I) \\ &= \Phi_l \left(1 + \sqrt{\frac{2}{n}} \xi_l \varepsilon_l \right), \quad \varepsilon_l \sim \mathcal{D} (\mathbb{E} [\varepsilon_l] = 0, \mathbb{E} [\varepsilon_l^2] = 1) \\ \implies \Phi_{l+1} &= \Phi_l + \underbrace{\frac{1}{\sqrt{n}}}_{\text{step } \frac{1}{n}} \underbrace{\sqrt{2} \Phi_l}_{\sigma(\Phi_l)} \underbrace{\varepsilon_l}_{\text{zero mean}} \\ \Phi_\tau^{(n)} &:= \Phi_{\lfloor \tau n \rfloor} \\ &\xrightarrow{d, n \rightarrow \infty, \frac{d}{n} \rightarrow \bar{\tau}} d\Phi_\tau \\ &= \sqrt{2} \Phi_\tau dB_\tau \end{aligned}$$

Geometric Brownian motion!

$$\text{Itô} \quad d \log(\Phi_\tau) = (-1) d\tau + \sqrt{2} dB_\tau$$

Thus,

$$\begin{aligned} \log(\Phi_\tau) - \log(\Phi_0) &= \int_0^\tau (-1) d\tau' + \int_0^\tau \sqrt{2} dB_{\tau'} \\ &= -\tau + \sqrt{2} B_\tau \\ &\sim \mathcal{N}(-\tau, 2\tau) \end{aligned}$$

Mutliple data points: $m \geq 2$.

$$\begin{aligned} h_{l+1}^\alpha &= \frac{1}{\sqrt{n}} W_l h_l^\alpha \\ h_1^\alpha &= \frac{1}{\sqrt{n_0}} W_0 x^\alpha \end{aligned}$$

Recall:

$$\begin{aligned} [h_{l+1}^\alpha]_{\alpha=1}^m \mid \mathcal{F}_l \left(= \sigma((h_k)_{k \leq l}) \right) &\stackrel{\text{d}}{=} \mathcal{N} \left(0, \underbrace{\left[\frac{1}{n} \langle h_l^\alpha, h_l^\beta \rangle \right]_{\alpha, \beta=1}^m}_{\Phi_l = [\Phi_l^{\alpha\beta}]_{\alpha\beta}} \otimes I_n \right) \\ \mathbb{E} [\Phi_{l+1}^{\alpha\beta} \mid \mathcal{F}_l] &= \mathbb{E} \left[\frac{1}{n} \sum_{j=1}^n h_{l+1, j}^\alpha h_{l+1, j}^\beta \mid \mathcal{F}_l \right] \\ &= \mathbb{E} [h_{l+1, j}^\alpha h_{l+1, j}^\beta \mid \mathcal{F}_l], \quad \forall j \in [1 : n] \\ &= \Phi_l^{\alpha\beta} \end{aligned}$$

Subtract the mean:

$$\begin{aligned} \Phi_{l+1}^{\alpha\beta} - \Phi_l \mid \mathcal{F}_l &= \frac{1}{n} \sum_{j=1}^n h_{l+1, j}^\alpha h_{l+1, j}^\beta - \Phi_l^{\alpha\beta} \mid \mathcal{F}_l \\ &= \underbrace{\frac{1}{\sqrt{n}}}_{\text{step size } \frac{1}{n}} \underbrace{\frac{1}{\sqrt{n}} \sum_{j=1}^n \left(h_{l+1, j}^\alpha h_{l+1, j}^\beta - \Phi_l^{\alpha\beta} \right)}_{\text{zero mean}} \mid \mathcal{F}_l \\ &=: \xi_l^{\alpha\beta} \text{ CLT term } = \Theta(1) \end{aligned}$$

Thus,

$$\Phi_{l+1}^{\alpha\beta} = \Phi_l^{\alpha\beta} + \frac{1}{\sqrt{n}} \xi_l^{\alpha\beta}, \quad \alpha, \beta \in [1 : m]$$

We will flatten the Markov chain:

- Φ is symmetric (PSD)
- $\text{vec} : \text{Sym}(m) \rightarrow \mathbb{R}^{m=m(m+1)/2}$, $\text{vec} : \Phi \mapsto [\Phi^{\alpha\beta}]_{\alpha \leq \beta} \rightarrow$ upper-triangular entries (unique entries)

Markov chain is always in the space $\mathbb{R}^{\bar{m}}$. So,

$$\text{vec}(\Phi_{l+1}) = \text{vec}(\Phi_l) + \frac{\text{vec}(\xi_l)}{\sqrt{n}}, \quad \xi_l \rightarrow [\xi_l^{\alpha\beta}]_{\alpha\beta}$$

Abusing notation:

$$\Phi_{l+1} = \Phi_l + \frac{1}{\sqrt{n}} \xi_l$$

We need the covariance matrix: $\text{Cov}(\xi_l)$:

$$\begin{aligned} \mathbb{E}[\xi_l^{\alpha\beta} \xi_l^{\gamma\delta}] &= \mathbb{E}\left[\left(h_{l+1,j}^\alpha h_{l+1,j}^\beta - \Phi_l^{\alpha\beta}\right)\left(h_{l+1,j}^\gamma h_{l+1,j}^\delta - \Phi_l^{\gamma\delta}\right)\right] \\ &= \mathbb{E}[h^\alpha h^\beta h^\gamma h^\delta] - \Phi^{\alpha\beta} \Phi^{\gamma\delta} \\ &= (\mathbb{E}[h^\alpha h^\beta] \mathbb{E}[h^\gamma h^\delta] + \mathbb{E}[h^\alpha h^\gamma] \mathbb{E}[h^\beta h^\delta] + \mathbb{E}[h^\alpha h^\delta] \mathbb{E}[h^\gamma h^\beta]) - \Phi^{\alpha\beta} \Phi^{\gamma\delta} \quad (\text{Wick's Isserlis}) \\ &= (\Phi^{\alpha\beta} \Phi^{\gamma\delta} + \Phi^{\alpha\delta} \Phi^{\beta\gamma} + \Phi^{\alpha\gamma} \Phi^{\beta\delta}) - \Phi^{\alpha\beta} \Phi^{\gamma\delta} \\ &= \Phi^{\alpha\delta} \Phi^{\beta\gamma} + \Phi^{\alpha\gamma} \Phi^{\beta\delta} \end{aligned}$$

where $\xi_l^{\alpha\beta} \in \mathbb{R}^{\bar{m}}$ and $\alpha \leq \beta, \gamma \leq \delta$.

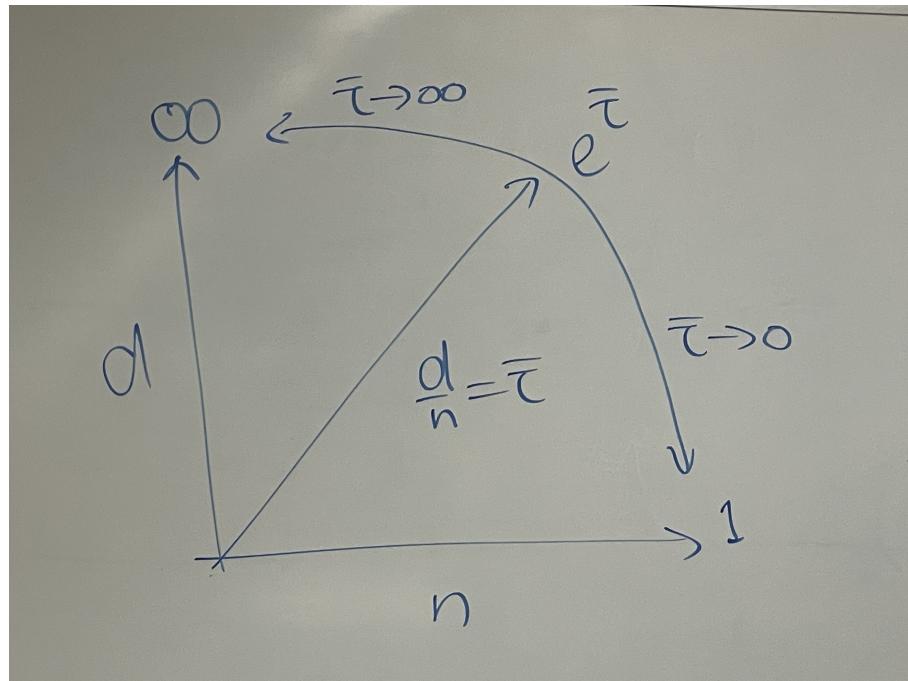
Thus,

$$\begin{aligned} \text{vec}(\Phi_{l+1}) &= \text{vec}(\Phi_l) + \frac{1}{\sqrt{n}} \underbrace{\Sigma(\Phi_l)^{1/2}}_{\mathbb{R}^{\bar{m} \times \bar{m}}} \varepsilon_l, \quad \varepsilon_l \in \mathbb{R}^{\bar{m}}, \varepsilon_l \sim \mathcal{D}(\mathbb{E}[\varepsilon_l] = 0, \mathbb{E}[\varepsilon_l^2] = 1) \\ \Sigma(\Phi)^{\alpha\beta, \gamma\delta} &= \Phi^{\alpha\delta} \Phi^{\beta\gamma} + \Phi^{\alpha\gamma} \Phi^{\beta\delta} \\ \xrightarrow[d, n \rightarrow \infty, \frac{d}{n} \rightarrow \bar{\tau}]{} &\begin{cases} d \underbrace{\Phi_\tau}_{\mathbb{R}^{\bar{m}}} = \Sigma(\Phi_\tau)^{1/2} d \underbrace{B_\tau}_{BM \text{ on } \mathbb{R}^{\bar{m}}} \\ \Phi_0 = \left[\frac{1}{n_0} \langle x^\alpha, x^\beta \rangle \right]_{\alpha, \beta=1}^m \\ \Phi_{\bar{\tau}} \leftarrow \left[\frac{1}{n_0} \langle h_d^\alpha, h_d^\beta \rangle \right]_{\alpha, \beta=1}^m \end{cases} \end{aligned}$$

$$d\Phi_\tau = \underbrace{\Sigma(\Phi_\tau)^{1/2}}_{\text{inverse Riemannian metric (later)}} dB_\tau$$

Intuition:

$$\begin{aligned} \left(\underbrace{1}_{\text{each layer } \approx 1} + \frac{1}{n} \right)^d &\rightarrow \begin{cases} 1, & n \rightarrow \infty \\ \infty, & d \rightarrow \infty \\ e^{\bar{\tau}}, & n \rightarrow \infty, d \rightarrow \infty, \frac{d}{n} \rightarrow \bar{\tau} \end{cases} \\ &\prod_{l=1}^{d-1} \left(1 + \sqrt{\frac{2}{n}} \varepsilon_l \right)^d \end{aligned}$$



Nonlinear Networks

$(m = 1)$

Most naïve thing one can do:

$$\begin{aligned}
 h_{l+1} &= \frac{1}{\sqrt{n}} W_l \varphi(h_l) \\
 h_{l+1} \mid \mathcal{F}_l &\sim \mathcal{N} \left(0, \underbrace{\frac{1}{n} |\varphi(h_l)|^2}_{\Phi_l} \otimes I_n \right) \\
 \mathbb{E} [\Phi_{l+1} \mid \mathcal{F}_l] &= \mathbb{E} \left[\frac{1}{n} \sum_{j=1}^n \varphi(h_{l+1,j})^2 \mid \mathcal{F}_l \right] \\
 &= \mathbb{E} \left[\underbrace{\varphi(h_{l+1,j})^2}_{h_{l+1,j}^2 \mathbb{1}_{\{h_{l+1,j} > 0\}}} \right]
 \end{aligned}$$

$$\begin{aligned}
\mathbb{E} [w^2 \mathbb{1}_{\{w>0\}}] &= \int_0^\infty x^2 d\rho(x), \quad \rho(x) \sim \mathcal{N}(0, \Phi) \\
&= \frac{1}{2} \underbrace{\int_{-\infty}^\infty x^2 d\rho(x)}_{\mathbb{E}[w^2]=\Phi} \\
\implies \mathbb{E} [\Phi_{l+1} | \mathcal{F}_l] &= \frac{1}{2} \Phi_l \\
\Phi_{l+1} &= \underbrace{\frac{\Phi}{2}}_{\text{unstable! } \Phi_l \rightarrow 0 \text{ as } l \rightarrow \infty} + \frac{1}{\sqrt{n}} \xi_l
\end{aligned}$$

(Kaiming)He - Init.

$$\begin{aligned}
\frac{1}{c} &= \mathbb{E} [\varphi(w)^2], \quad w \sim \mathcal{I}, \infty \\
\varphi(|c|x) &= |c|\varphi(x) \\
\implies h_{l+1} &= \sqrt{\frac{c}{n}} W_l \varphi(h_l) \\
h_{l+1} | \mathcal{F}_l &\stackrel{d}{=} \mathcal{N} \left(0, \underbrace{\frac{c}{n} |\varphi(h_l)|^2}_{\Phi_l \text{ Def'n changed!}} \otimes I_n \right)
\end{aligned}$$

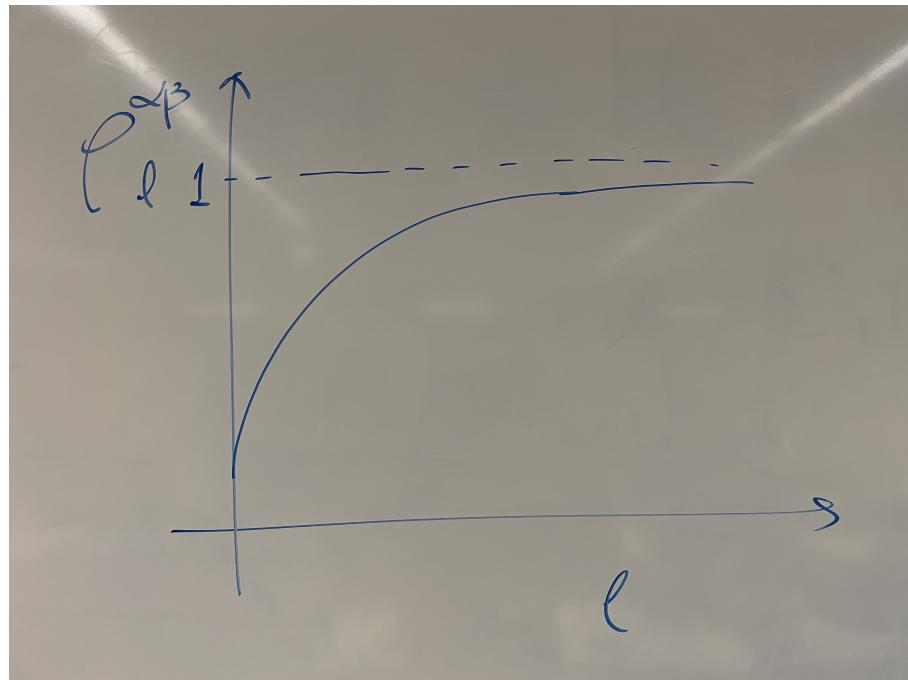
$$\begin{aligned}
\mathbb{E} [\Phi_{l+1} | \mathcal{F}_l] &= \mathbb{E} \left[\frac{c}{n} \sum_{j=1}^n \varphi(h_{l+1,j})^2 | \mathcal{F}_l \right] \\
&= \mathbb{E} [c \varphi(h_{l+1,j})^2 | \mathcal{F}_l] \\
&= 2 \mathbb{E} [\varphi(w)^2], \quad w \sim \mathcal{N}(0, \Phi_l) \\
&= 2 \frac{1}{2} \Phi_l \\
\implies \Phi_{l+1} &= \Phi_l + \frac{1}{\sqrt{n}} \underbrace{\frac{1}{\sqrt{n}} \sum_{j=1}^n (c \varphi(h_{l+1,j})^2 - \Phi_l)}_{\xi_l}
\end{aligned}$$

$$\begin{aligned}
\mathbb{E} [\xi_l^2 | \mathcal{F}_l] &= \mathbb{E} \left[(c \varphi(h_{l+1,j})^2 - \Phi_l)^2 | \mathcal{F}_l \right] \\
&= c^2 \mathbb{E} \left[\underbrace{\varphi(h_{l+1,j})^4}_{h_{l+1,j}^4 \mathbb{1}_{\{h_{l+1,j}>0\}}} | \mathcal{F}_l \right] - \Phi_l^2 \\
&= \underbrace{\frac{c^2}{2}}_{=2} \mathbb{E} \left[\underbrace{h_{l+1,j}^4}_{3\Phi_l^2} | \mathcal{F}_l \right] - \Phi_l^2 \\
&= 5\Phi_l^2
\end{aligned}$$

$$\begin{aligned}
\Phi_{l+1} &= \Phi_l + \sqrt{5}\Phi_l \varepsilon_l, \quad \varepsilon_l \sim \mathcal{D}(\mathbb{E}[\varepsilon_l] = 0, \mathbb{E}[\varepsilon_l^2] = 1) \\
&\rightarrow d\Phi_l = \sqrt{5}\Phi_l dB_\tau, \quad \frac{d}{n} \rightarrow \bar{\tau} \\
\implies \frac{\Phi_l}{\Phi_0} &\sim \exp\left(\mathcal{N}\left(-\frac{5}{2}\bar{\tau}, 5\bar{\tau}\right)\right)
\end{aligned}$$

In the case of $m \geq 2$,

$$\begin{aligned}
h_{l+1}^\alpha &= \sqrt{\frac{c}{n}} W_l \varphi(h_l^\alpha) \\
\rho_l^{\alpha\beta} &= \frac{\Phi_l^{\alpha\beta}}{\sqrt{\Phi_l^{\alpha\alpha} \Phi_l^{\beta\beta}}} \in [-1; 1] \\
\rho_{l+1}^{\alpha\beta} &= f(\rho_l^{\alpha\beta}) + \dots \\
\rho_l^{\alpha\beta} &\rightarrow 1 \text{ as } l \rightarrow \infty
\end{aligned}$$



- $h_l^\alpha \approx h_l^\beta$ constant
- Transformers with $\rho \approx 1$
- \implies vanishing gradients (Noci et al., 2022)