

# STAT 946 - Topics in Probability and Statistics: Mathematical Foundations of Deep Learning

## Lecture 7

Lucas Noritomi-Hartwig  
University of Waterloo

September 24, 2025 from 16h00 to 17h20 in M3 3103

October 6th: Majid (Deep RL), Marty (MF Transformers) presenting articles.

October 20th: ? ?

October 8th: Last day to decide on a topic and direction for a project.

Recall

$$\begin{aligned} f(x^\alpha; \theta) &= \frac{1}{\sqrt{n}} W_d \varphi(h_d) \\ h_{l+1}^\alpha &= \frac{1}{\sqrt{n}} W_l \varphi(h_l) \\ h_1^\alpha &= \frac{1}{\sqrt{n_0}} W_0 x^\alpha \end{aligned}$$

Backward neuron:

$$\begin{aligned} g_l^\alpha &= \sqrt{n} \frac{\partial f^\alpha}{\partial h_l^\alpha} \\ z_l^\alpha &= \frac{1}{\sqrt{n}} W_l^\top g_{l+1}^\alpha \\ &= \frac{1}{\sqrt{n}} W_l^\top \underbrace{\text{diag}(\varphi'(h_{l+1}^\alpha))}_{=: D_{l+1}^\alpha} z_{l+1}^\alpha \\ \mathcal{F}_l^z &= \sigma(\{h_k^\alpha\}_{\alpha, k}, \{z_k\}_{k \geq l}) \end{aligned}$$

$$\begin{aligned} W|\sigma(W\varphi) &\stackrel{\text{d}}{=} WP_\varphi + \underbrace{\tilde{W}}_{\text{indep. copy}} P_\varphi^\perp \\ z_l^\alpha | \mathcal{F}_l^z &= \frac{1}{\sqrt{n}} \left( P_{\varphi_l} W_l^\top + P_{\varphi_l}^\perp \tilde{W}_l^\top \right) D_{l+1}^\alpha z_{l+1}^\alpha | \mathcal{F}_l^z & (h_{l+1} = W_l \varphi(h_l)) \\ &= P_{\varphi_l} z_l^\alpha + P_{\varphi_l}^\perp \underbrace{\frac{1}{\sqrt{n}} \tilde{W}_l^\top g_{l+1}^\alpha}_{\tilde{z}_l^\alpha \sim \mathcal{N}(0, G_{l+1} \otimes I)} | \mathcal{F}_l^z \end{aligned}$$

$$\begin{aligned}
G_l^{\alpha\beta} &= \frac{1}{n} \left\langle g_l^\alpha, g_l^\beta \right\rangle \\
&= \frac{1}{n} \left\langle D_l^\alpha z_l^\alpha, D_l^\beta z_l^\beta \right\rangle \\
&= \frac{1}{n} \left\langle D_l^\alpha (P_{\varphi_l} z_l^\alpha + P_{\varphi_l}^\perp \tilde{z}_l^\alpha), D_l^\beta (P_{\varphi_l} z_l^\beta + P_{\varphi_l}^\perp \tilde{z}_l^\beta) \right\rangle \\
&= \frac{1}{n} \left\langle D_l^\alpha P_{\varphi_l} z_l^\alpha, D_l^\beta P_{\varphi_l} z_l^\beta \right\rangle \tag{1}
\end{aligned}$$

$$+ \frac{1}{n} \left\langle D_l^\alpha P_{\varphi_l}^\perp \tilde{z}_l^\alpha, D_l^\beta P_{\varphi_l} z_l^\beta \right\rangle \tag{2}$$

$$+ \frac{1}{n} \left\langle D_l^\alpha P_{\varphi_l} z_l^\alpha, D_l^\beta P_{\varphi_l}^\perp \tilde{z}_l^\beta \right\rangle \tag{3}$$

$$+ \frac{1}{n} \left\langle D_l^\alpha P_{\varphi_l}^\perp \tilde{z}_l^\alpha, D_l^\beta P_{\varphi_l}^\perp \tilde{z}_l^\beta \right\rangle \tag{4}$$

$$(1) = \frac{1}{n} \sum_{j=1}^n \underbrace{\varphi' (h_{l,j}^\alpha)}_{\in \Theta(1)} (P_{\varphi_l} z_l^\alpha)_j$$

$\|P_{\varphi_l} z_l^\alpha\|^2$  is a projection on the column space of  $\varphi_l$ , namely  $\text{col}(\varphi_l)$ , where  $\varphi_l$  is  $n \times m$ .

$$\varphi_l = \begin{bmatrix} \varphi \left( \underbrace{h_l^{(1)}}_{n \times 1} \right) & \varphi \left( \underbrace{h_l^{(2)}}_{n \times 1} \right) & \dots & \varphi \left( \underbrace{h_l^{(m)}}_{n \times 1} \right) \end{bmatrix}$$

In the trivial case, i.e.,  $m = 1$ ,

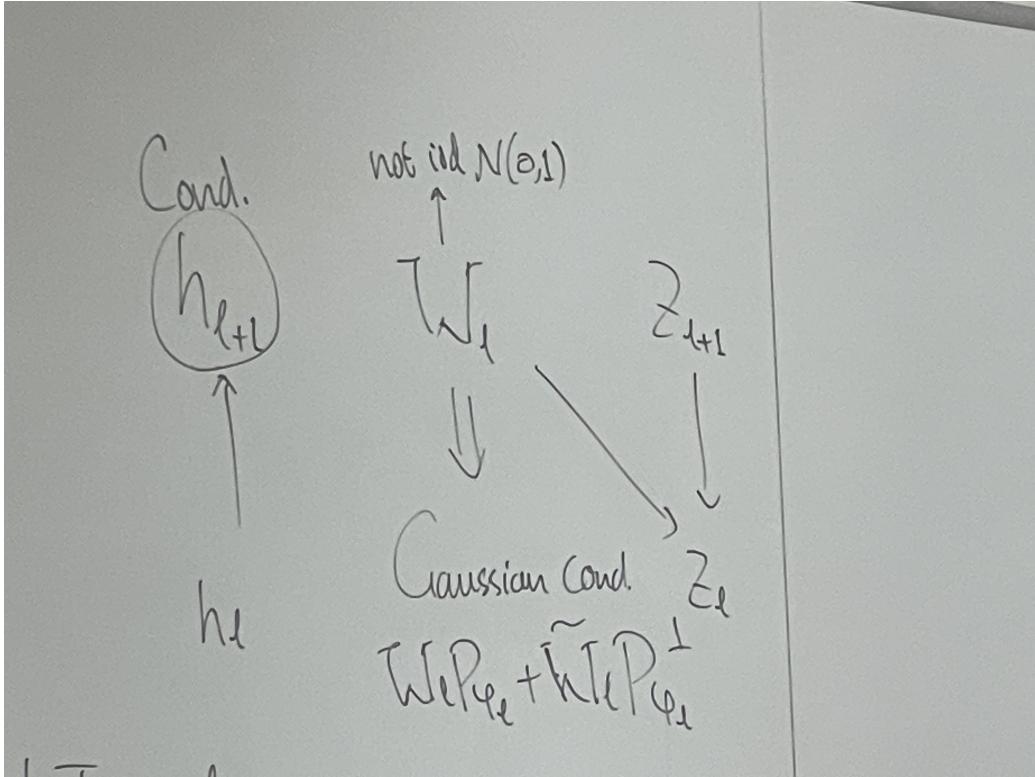
$$\begin{aligned}
\varphi_l &= e_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^n \\
P_{\varphi_l} z_l^\alpha &= \begin{bmatrix} z_{l,1}^\alpha \\ 0 \\ \vdots \\ 0 \end{bmatrix}
\end{aligned}$$

So as  $n \rightarrow \infty$ ,  $\|P_{\varphi_l} z_l^\alpha\|^2 \lesssim \frac{m}{n} \|z_l^\alpha\|^2$

$$(2) + (3) = \frac{1}{n} \sum_i \dots \left( \underbrace{P_{\varphi_l}^\perp \tilde{z}_l^\alpha}_{\text{zero mean, indep.}} \right)_i \dots \tag{\dagger}$$

$$\tilde{z}_l^\alpha = \frac{1}{\sqrt{n}} \underbrace{\tilde{W}_l^\top}_{\text{indep. copy of } W_l^\top} g_{l+1}^\alpha$$

$$\begin{aligned}
(\dagger) &= \frac{1}{\sqrt{n}} \Theta(1) \rightarrow \text{CLT scaling} \\
&\rightarrow 0
\end{aligned}$$



$$(4) = \frac{1}{n} \sum_{i=1}^n \varphi'(h_{l,i}^\alpha) \tilde{z}_{l,i}^\alpha \varphi'(h_{l,i}^\beta) \tilde{z}_{l,i}^\beta - \dots \underbrace{(P_{\varphi_l} \tilde{z}_l^\alpha)}_{\Theta(m,n) \rightarrow 0} \dots (P_{\varphi_l} \tilde{z}_{l,i}^\beta)$$

$$\rightarrow \mathbb{E} \left[ \varphi'(h_{l,i}^\alpha) \varphi'(h_{l,i}^\beta) \mathbb{E} \left[ \underbrace{\tilde{z}_{l,i}^\alpha \tilde{z}_{l,i}^\beta}_{\rightarrow G_{l+1}^{\alpha\beta}} | \mathcal{F}_{l+1}^z \right] \right]$$

$$P_{\varphi_l}^\perp \tilde{z}_l^\alpha = P_{\varphi_l}^\perp \frac{1}{\sqrt{n}} \tilde{W}_l^\top g_{l+1}^\alpha \sim \mathcal{N} \left( 0, G_{l+1}^{\alpha\beta} \otimes \underbrace{P_{\varphi_l}^\perp}_{=I-P_{\varphi_l}} \right)$$

Aside:

$$\sum_{i=1}^n X_i \rightarrow X \implies \sum_{i=1, i \neq 2}^n X_i \rightarrow X,$$

i.e., finitely many random variables contribute infinitesimally small amounts to the limit.

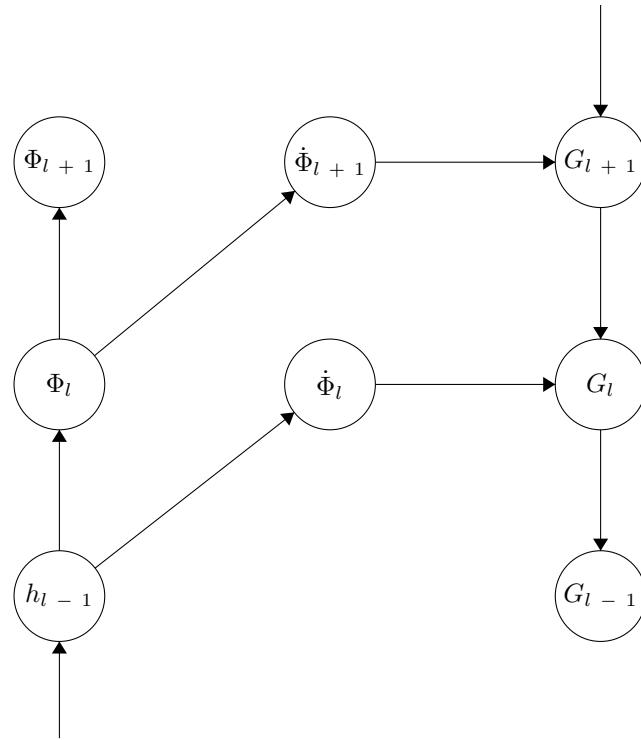
$$\dot{\Phi}_l^{\alpha\beta} := \mathbb{E} \left[ \varphi'(h_{l,i}^\alpha) \varphi'(h_{l,i}^\beta) \right]$$

$$G_l^{\alpha\beta} = \dot{\Phi}_l^{\alpha\beta} \odot G_{l+1}^{\alpha\beta}$$

where  $A \odot B = [a_{ij} b_{ij}]_{ij}$  is the Hadamard product of  $A$  and  $B$ .

NTK:

$$\begin{aligned}
 K^{\alpha\beta} &= \sum_{l=0}^d \langle \nabla_{W_l} f^\alpha, \nabla_{W_l} f^\beta \rangle \\
 &\sum_{l=0}^d \frac{1}{n} \langle g_{l+1}^\alpha, g_{l+1}^\beta \rangle \frac{1}{n} \langle \varphi_l^\alpha, \varphi_l^\beta \rangle \\
 &= \sum_{l=0}^d G_{l+1}^{\alpha\beta} \Phi_l^{\alpha\beta}
 \end{aligned}$$



### Theorem

$$\begin{cases} \Phi_{l+1} = f_1(\Phi_l) \\ G_l = \dot{\Phi}_l \odot G_{l+1} \\ \Phi_0 = \left[ \frac{1}{n_0} \langle x^\alpha, x^\beta \rangle \right]_{\alpha, \beta} \\ G_{d+1} = 11^\top \end{cases}$$

$$K = \sum_{l=0}^d \Phi_l \odot G_{l+1}$$

$$\dot{\Phi}_{l+1} = f_2(\Phi_l)$$

$$\dot{\Phi}_{l+1} = \mathbb{E} [\varphi'(h^\alpha) \varphi'(h^\beta)]$$

$$\begin{bmatrix} h^\alpha \\ h^\beta \end{bmatrix} \sim \mathcal{N} \left( 0, \begin{bmatrix} \Phi_l^{\alpha\alpha} & \Phi_l^{\alpha\beta} \\ \Phi_l^{\beta\alpha} & \Phi_l^{\beta\beta} \end{bmatrix} \right)$$

$$z_l^\alpha = P_{\varphi_l} z_l^\alpha + P_{\varphi_l}^\perp \tilde{z}_l^\alpha$$

## Generalization

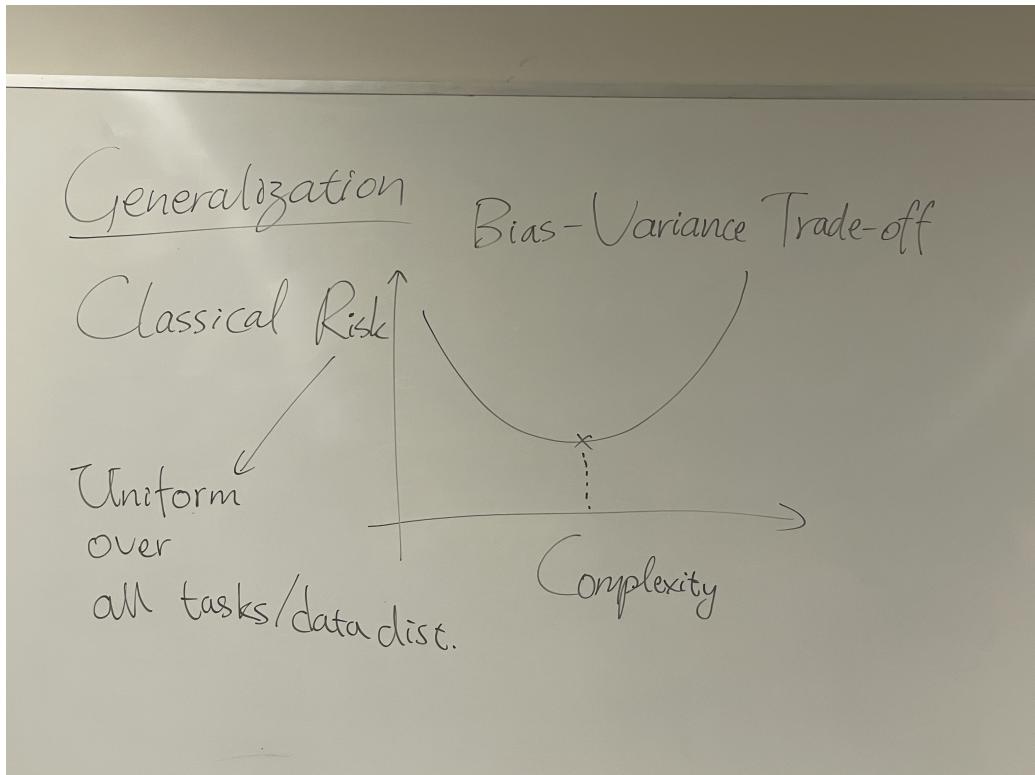
Something of which we have a completely new understanding since 2018-2019.

Classical view: bias-variance tradeoff, i.e., in order to perform better (have low risk), we must regularize our model's complexity.

What is risk? In classical theories, “risk” is defined in a very restrictive way.

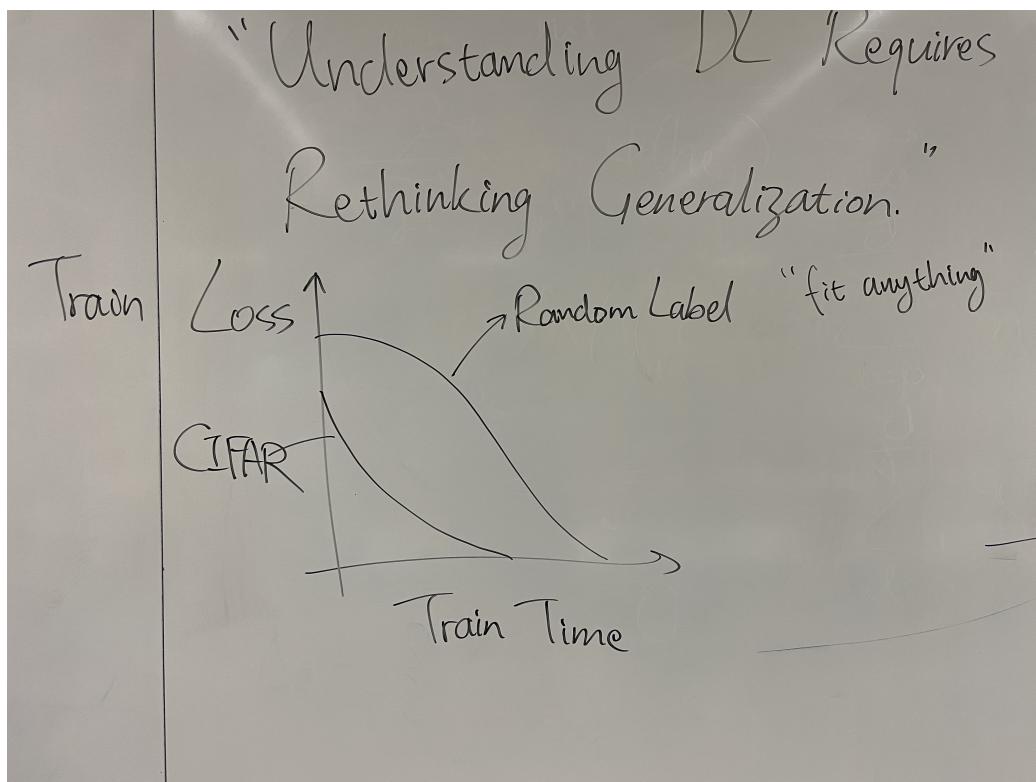
### No Free Lunch Theorem

Informally: If algorithm  $A$  can fit any data, then  $\exists$  a task (data distribution) on which  $A$  fails, i.e., there is some fixed amount of risk under which  $A$  cannot output a model that performs better than the risk.

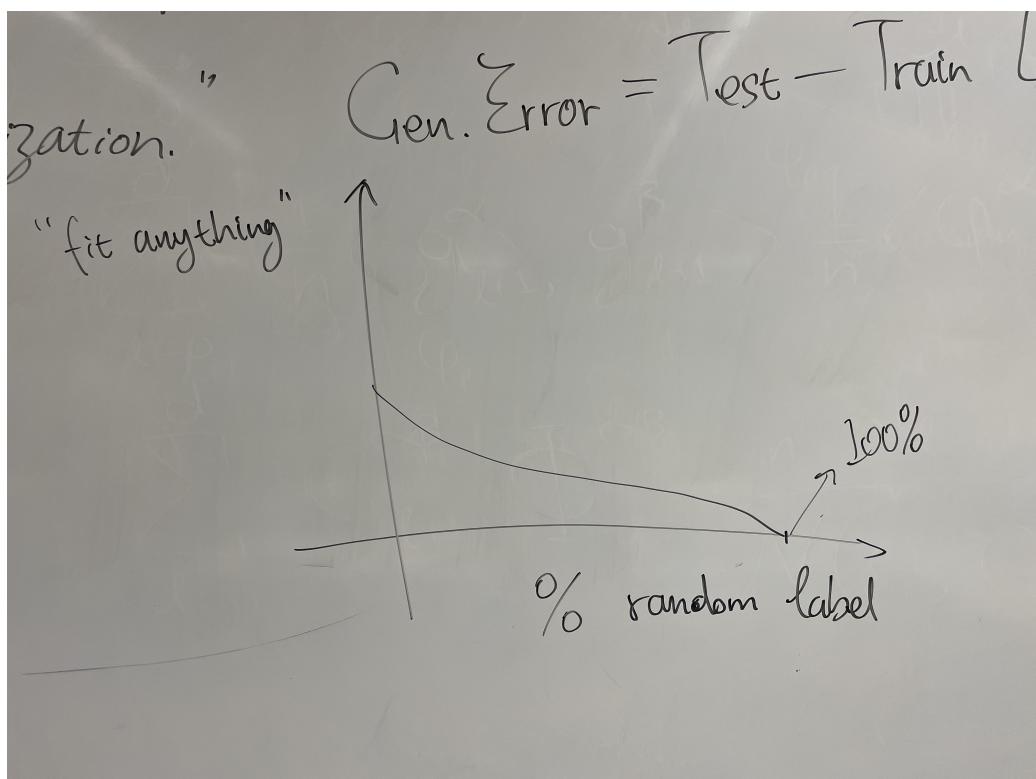


“Understanding Deep Learning Requires Rethinking Generalization” - Zhang et al (2016)

First plot: Training loss



Second plot: Generalization error (Test - Training loss)



Belkin et al. (2018) "Double Descent"

