

Jiayuan Ye (jiayuan@comp.nus.edu.sg), Zhenyu Zhu (zhenyu.zhu@epfl.ch), Fanghui Liu (fanghui.liu@epfl.ch), Reza Shokri (reza@comp.nus.edu.sg) and Volkan Cevher (volkan.cevher@epfl.ch)



Initialization Matters: Privacy-Utility Analysis of **Overparameterized Neural Networks**

$$(\mathbf{W}_{[0:T]} \| \mathbf{W}'_{[0:T]}) = \frac{1}{2\sigma^2} \int_0^T \mathbb{E} \left[\| \nabla \mathcal{L}(\mathbf{W}_t; D) - \nabla \mathcal{L}(\mathbf{W}_t; D') \|_2^2 \right] dt \quad \text{where}$$

$$\int_0^T \mathbb{E} \left[\| \nabla \mathcal{L}(\mathbf{W}_t; D) - \nabla \mathcal{L}(\mathbf{W}_t; D') \|_2^2 \right] dt \quad \text{Initialization matters for}$$

$$\leq 2T \cdot \mathbb{E} \left[\| \nabla \mathcal{L}(\mathbf{W}_0; D) - \nabla \mathcal{L}(\mathbf{W}_0; D') \|_2^2 \right] + \frac{2c^2T}{n^2}$$

$$\text{gradient difference at initialization} \quad \text{non-smoothness cost}$$

$$+ \frac{2\beta^2}{n^2(2+\beta^2)} \left(\frac{e^{(2+\beta^2)T} - 1}{2+\beta^2} - T \right) \cdot \left(\mathbb{E} \left[\| \nabla \mathcal{L}(\mathbf{W}_0; D) \|_2^2 \right] + \sigma^2 \operatorname{rank}(M_T) + c^2 \right)$$

$$= \operatorname{randimt} \operatorname{difference} \operatorname{fluctuation} \operatorname{fluctuation} \operatorname{randimt} \operatorname{randim} \operatorname{randim}$$

$$\mathbb{E}_{\boldsymbol{W}}\left[\|\frac{\partial f(\boldsymbol{x})}{\partial \operatorname{Vec}(\boldsymbol{W})}\|_{F}^{2}\right] = \|\boldsymbol{x}\|_{2}^{2}o\left(\prod_{i=1}^{L-1}\frac{\beta_{i}m_{i}}{2}\right)\sum_{l=1}^{L}\frac{\beta_{l}}{\beta_{l}}$$



Special Case: Privacy-Utility Trade-offs for Training Linearized Network

onsider a linearized network by first-order Taylor expansion $igg| \, m{f}_{m{W}}^{lin,0}(m{x}) \equiv m{f}_{m{W}_0^{lin}}(m{x}) + rac{\partial m{f}_{m{W}}(m{x})}{\partial m{W}} \Big|_{m{W}=m{W}_0^{lin}} \left(m{W} - m{W}_0^{lin}
ight) \, ,$

nder GD, DNN can work in the **lazy training regime**, under which is linearized network well approximates DNN training

eorem: For single output linearized network with hidden layer dth m, bounded data with dimension d, under certain regularity onditions, if $d, m = \Omega(n)$ where n is size of training dataset, then

Initia	alization	Variance β_l for layer l	KL privacy bound under fixed T and σ^2	Excess Empirical risk under ε -KL privacy
	eCun	$1/m_{l-1}$	$\frac{om(L-1+\frac{d}{m})}{2^{L-1}}\cdot\frac{2T}{n^2\sigma^2}$	$ ilde{\mathcal{O}}\left(rac{1}{n^2}+\sqrt{rac{1}{2^L}arepsilon} ight)$
	He	$2/m_{l-1}$	$om(L-1+\frac{d}{m})\cdot\frac{2T}{n^2\sigma^2}$	$ ilde{\mathcal{O}}\left(rac{1}{n^2}+\sqrt{rac{1}{arepsilon}} ight)$
1	NTK	$ \begin{array}{c} 2/m_l, l < L\\ 1/o, l = L \end{array} $	$dm\left(\frac{L-1}{2}+\frac{o}{m}\right)\cdot\frac{2T}{n^2\sigma^2}$	$\tilde{\mathcal{O}}\left(\frac{1}{n^2} + \sqrt{\frac{d}{\varepsilon}}\right)$
X	lavier	$\frac{2}{m_{l-1}+m_l}$	$\frac{od(L-1+\frac{d+o}{2m})}{2^{L-3}(1+\frac{d}{m})(1+\frac{o}{m})} \cdot \frac{2T}{n^2\sigma^2}$	$\tilde{\mathcal{O}}\left(\frac{1}{n^2} + \sqrt{\frac{1}{2^L}\varepsilon}\right)$

decreases under increasing depth for $L \ge 2$

Main Takeaways

We theoretically prove and numerically show that for training DNNs with a small time, and for training linearized networks with any time

Increasing width always hurts KL privacy

•

/Increasing depth **helps** KL privacy under **certain initializations**

Under certain data regularity and large enough widths, we further prove <u>privacy-utility trade-offs</u> for training linearized networks and prove that it similarly relies on the **choice of initialization distributions**



Acknowledgements: This work was supported by Hasler Foundation Program: Hasler Responsible AI (project number 21043), the Swiss National Science Foundation (SNSF) under grant number 200021_205011, Google PDPO faculty research award, Intel within the www.private-ai.org center, Meta faculty research award, the NUS Early Career Research Award (NUS ECRA award number NUS ECRA FY19 P16), and the National Research Foundation, Singapore under its Strategic Capability Research Centres Funding Initiative.