

Depth Separation in Learning via Representation Costs

Suzanna Parkinson* Greg Ongie** Rebecca Willett* Ohad Shamir† Nathan Srebro‡

*University of Chicago **Marquette University †Weizmann Institute of Science ‡Toyota Technical Institute at Chicago

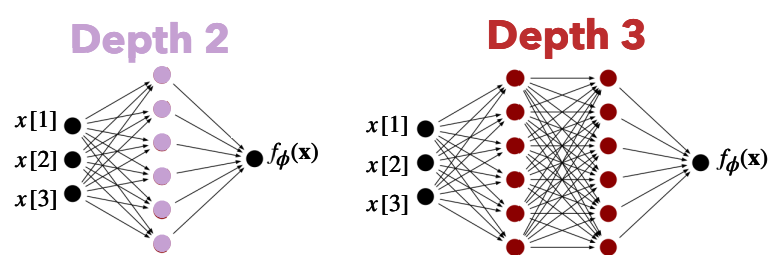
Are **depth-2** or **depth-3** neural networks better at **learning**?

Depth- L Neural Networks

$$\phi = (\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_{L-1}, \mathbf{a})$$

$$\sigma(x) = \text{ReLU}(x) = \max(0, x)$$

$$f_\phi(\mathbf{x}) = \mathbf{a}^\top \sigma \left(\mathbf{W}_{L-1} \cdot \sigma \left(\dots \sigma \left(\mathbf{W}_2 \sigma \left(\mathbf{W}_1 \mathbf{x} \right) \right) \right) \right)$$



PAC Learning

- The output of a learning rule \mathcal{A} trained with m samples is (ϵ, δ) **-Probably Approximately Correct** if with probability $1 - \delta$ over the training samples $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$, the **generalization error** is less than ϵ :

$$\mathcal{L}_{\mathcal{D}}(\mathcal{A}(S)) := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[(\mathcal{A}(S)(\mathbf{x}) - f(\mathbf{x}))^2 \right] < \epsilon.$$

- If our learning rule \mathcal{A} gives a model that is (ϵ, δ) -Probably Approximately Correct using $m(\epsilon, \delta)$ samples, then we say that we can **learn** with **sample complexity** $m(\epsilon, \delta)$.

Controlling Generalization Error

- We often end up with error bounds like this:

$$\underbrace{\mathcal{L}_{\mathcal{D}}(\mathcal{A}(S))}_{\text{Generalization Error}} \leq \underbrace{\inf_{g \in \mathcal{G}} \mathcal{L}_{\mathcal{D}}(g)}_{\text{Approximation Error}} + 2 \underbrace{\sup_{g \in \mathcal{G}} |\mathcal{L}_S(g) - \mathcal{L}_{\mathcal{D}}(g)|}_{\text{Estimation Error}}$$

- Approximation error:** Need existence of **one** good approximator $g \in \mathcal{G}$.¹² Both depth **2** and **3** networks of arbitrary width are universal approximations of continuous functions.
- Estimation error:** Controlled using **size** of \mathcal{G} , here analyzed in terms of **Rademacher complexity**.³⁴ Naively, depth **3** networks have more parameters and so form a bigger model class

What if we measure model **size** in terms of **norm** of parameters instead of **number** of parameters?⁴⁵

Weight Decay & Representation Cost

$$\hat{\phi}_S \in \arg \min_{\phi} \mathcal{L}_S(f_\phi) + \lambda C_L(\phi) \text{ where } C_L(\phi) = \frac{1}{L} \left(\sum_{\ell=1}^{L-1} \|\mathbf{W}_\ell\|_F^2 + \|\mathbf{a}\|_2^2 \right)$$



Weight Decay Cost

$$\mathcal{A}_L(S) \in \arg \min_{g \in \mathcal{N}_L} \mathcal{L}_S(g) + \lambda R_L(g) \text{ where } R_L(g) = \inf_{\phi} C_L(\phi) \text{ s.t. } f_\phi = g$$

Representation Cost

Understanding **representation costs** across different depths helps us understand gaps in **learning** capabilities

Depth Separation in Approximation

$\exists f$ that requires **exponential width** (in dimension) with depth **2** but only **polynomial width** with depth **3** to be **approximated**.⁶⁷⁸

Depth Separation in Learning

- $\mathbf{x} \sim \text{Unif}(\mathbf{S}^{d-1} \times \mathbf{S}^{d-1}), f(\mathbf{x}) \in [-1, 1]$
- Depth-**2** vs. Depth **3** learning rules:

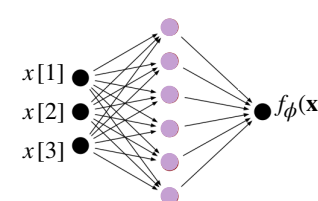
$$\mathcal{A}_2(S) \in \arg \min_{g \in \mathcal{N}_2} \mathcal{L}_S(g) + \lambda_2 R_2(g) \text{ vs.}$$

$$\mathcal{A}_3(S) \in \arg \min_{g \in \mathcal{N}_3} \mathcal{L}_S(g) + \lambda_3 R_3(g)$$

$\exists f$ that requires **exponential sample complexity** with depth **2** but only **polynomial sample complexity** with depth **3** to **learn**.

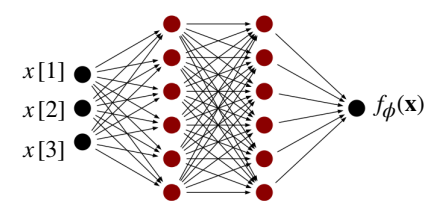
Key Idea: Choose f so that...

Large **representation cost** with **Depth 2**



Expensive

Small **representation cost** with **Depth 3**

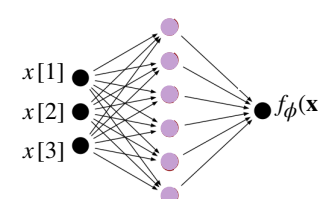


Cheap

$\forall f$ that can be **learned** with **polynomial sample complexity** with depth **2** can also be **learned** with **polynomial sample complexity** with depth **3**.

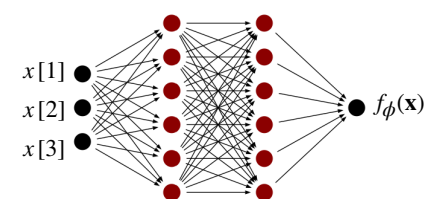
Key Idea:

Small **representation cost** with **Depth 2**



Cheap

Small **representation cost** with **Depth 3**



Cheap

Easy with depth 2

Easy with depth 3

Open Questions & Extensions

- Depth separation between **other depths**?
- Other architectures** beyond MLPs?
- We've implicitly assumed that we're **close to global minima** of our objective. How does **optimization** and the **loss-landscape** affect learning at different depths?

¹ Hornik (1991)

² Shen et al. (2022)

³ Bartlett & Mendelson (2001)

⁴ Neyshabur et al. (2015)

⁵ Bartlett (1996)

⁶ Eldan & Shamir (2016)

⁷ Daniely (2017)

⁸ Safran et al. (2021)

