

§5. Depth Separation

Q: Why deep neural networks when 2-layer NNs are already universal approximators?

A1: \exists input distr. D , target function f_* s.t. w.r.t. D ,

no $\text{poly}(d/\epsilon)$ -width 2-layer networks can approximate f_* (ES16, SES19)

while \exists some $\text{poly}(d/\epsilon)$ -width 3-layer network can efficiently learn f_* (RZG23)

§5.1 Negative results.

1. Fourier transform.

(Assume all integrals exist and are finite)

Fourier transform. $\mathcal{F}: L^2 \rightarrow L^2$

$$f \mapsto (\mathcal{F}f)(\xi)$$

$$:= \int_{\mathbb{R}^d} f(x) \exp(-2\pi i \langle x, \xi \rangle) dx$$

$$\text{Notation: } \hat{f} := \mathcal{F}f$$

Properties: \mathcal{F} -linear, invertible, isometric,

$$\widehat{fg} = \hat{f} * \hat{g}$$

$$2. \mathbb{R}^d > 0 \text{ s.t. } \text{Vol}(\mathbb{R}^d B_d) = 1$$

\hookrightarrow unit ball in \mathbb{R}^d .

\mathcal{F} - inverse Fourier transform of $\mathbb{1}_{\mathbb{R}^d B_d}$.

\mathcal{F} - isometric $\Rightarrow \mathcal{F}^2$ is a density.

\hookrightarrow let this be the input distr.

g - target

$f(x) = \sum_{k=1}^m f_k(\langle v_k, x \rangle)$ - two-layer network,
 $(v_k)_{k=1}^m$ weights.

$$\text{Loss} = \|f - g\|_{L^2(\mathcal{F}^2)}^2 = \|f\mathcal{F} - g\mathcal{F}\|_{L^2}^2 = \|\widehat{f\mathcal{F}} - \widehat{g\mathcal{F}}\|_{L^2}^2.$$

3. Lemma. $\widehat{f\varphi}$ is supported within some tubes.

$$\text{supp}(\widehat{f\varphi}) \subset \bigcup_{k=1}^m (\text{span}(V_k) + R\alpha B_d) =: T$$

Pf. (Informal). $\widehat{f\varphi} = \sum_{k=1}^m \widehat{f_k \varphi} = \sum_{k=1}^m \widehat{f_k} * \widehat{\varphi} = \sum_{k=1}^m \widehat{f_k} * \mathbb{1}_{R\alpha B_d}$

$$\Rightarrow \text{supp}(\widehat{f\varphi}) \subset \bigcup_{k=1}^m (\text{supp}(\widehat{f_k}) + R\alpha B_d).$$

Claim. $\text{supp} \widehat{f_k} \subseteq \text{span } V_k$.

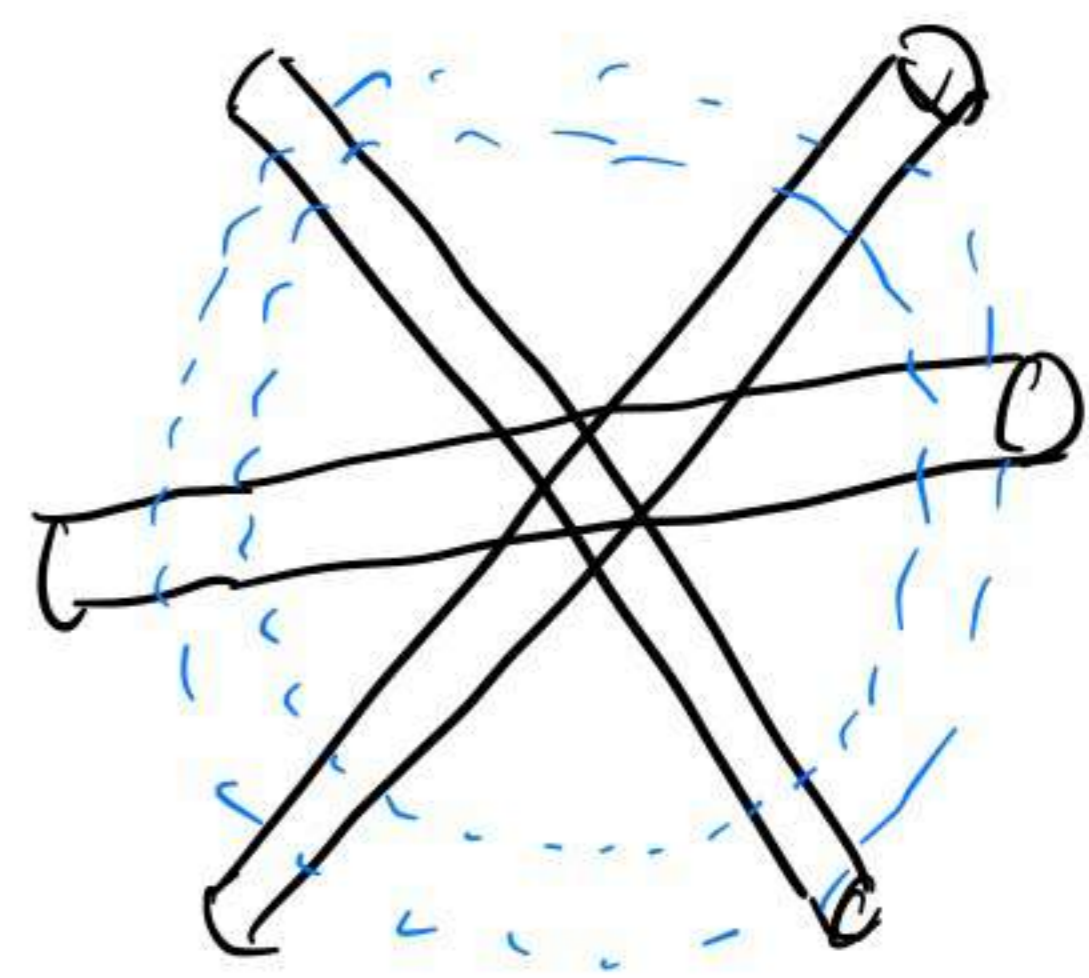
$$\xi \in \mathbb{R}^d. \quad \xi = aV_k + u, \quad a \in \mathbb{R}, u \perp V_k.$$

$$\begin{aligned} \widehat{f_k}(\xi) &= \int f_k(V_k \cdot x) \exp(-2\pi i \langle x, \xi \rangle) dx \\ &= \int f_k(V_k \cdot x) \exp(-2\pi i a \langle x, V_k \rangle) \exp(-2\pi i \langle x, u \rangle) dx =: \widehat{f_k}(a, u). \end{aligned}$$

Symmetry $\Rightarrow \widehat{f_k}(a, u) = \widehat{f_k}(a, -u)$

$$\Rightarrow \widehat{f_{ic}}(\xi) = \frac{1}{2} (\widehat{f_k}(a, u) + \widehat{f_k}(a, -u)) = \int f_k(V_k \cdot x) \exp(-2\pi i a \langle x, V_k \rangle) \cos(-2\pi \langle x, u \rangle) dx$$

$$\Rightarrow = 0 \text{ if } u \neq 0.$$



← T
sparse
when d
is large.

$\widehat{g\varphi}$ radial.

has some mass away
from the origin.

□

4. Lemma. $g, w, \|g\|_{L^2} = \|w\|_{L^2} = 1.$

$\text{supp } g \subseteq T.$

w -radial, $\int_{(2Rd)B_d} w^2 \geq \delta, \delta \in [0, 1]$

\hookrightarrow has some mass away from the origin.

$$\Rightarrow \langle g, w \rangle_{L^2} \leq 1 - \frac{\delta}{2} + \underbrace{m \exp(-cd)}_{\rightarrow 0 \text{ as } d \rightarrow \infty}$$

Pf. Define $A = (2Rd)B_d^c$ and

$$h(r) = \frac{\text{Vol}(rS^{d-1} \cap T)}{\text{Vol}(rS^{d-1})}$$

claim (without pf). $h \downarrow, h(r) \leq m \exp(-cd).$

$$\begin{aligned} \text{Write } \langle g, w \rangle_{L^2} &= \int_{A^c} g w + \int_A g w \\ &\leq \|g\|_{L^2} \|w \mathbb{1}_A\|_{L^2} + \int_A g w \\ &\leq \sqrt{1-\delta} + \int_A g w. \end{aligned}$$

claim. $\int_A g w \leq m \exp(-cd/2).$

$$\text{Pf. } \int_A g w = \int_A \tilde{g} w \leq \|w\|_{L^2} \|\tilde{g} \mathbb{1}_A\|_{L^2}$$

$$= \sqrt{\int_{2Rd}^{\infty} \tilde{g}^2(r) \text{Vol}(rS^{d-1}) dr}$$

$$\begin{aligned} \text{Note. } \tilde{g}(r) &= \frac{\int_{T \cap rS^{d-1}} g}{\text{Vol}(rS^{d-1})} = h(r) \frac{\int_{T \cap rS^{d-1}} g}{\text{Vol}(T \cap rS^{d-1})} \\ &\leq h(r) \sqrt{\frac{\int_{T \cap rS^{d-1}} g^2}{\text{Vol}(T \cap rS^{d-1})}} \end{aligned}$$

$$\begin{aligned}
\int_A g^w &\leq \sqrt{\int_{2R_A B_d} h^2(r) \frac{\int_{r \cap \mathbb{S}^{d-1}} g^2}{\text{Vol}(r \cap \mathbb{S}^{d-1})} \text{Vol}(r \cap \mathbb{S}^{d-1})} \\
&= \sqrt{\int_{2R_A B_d} h(r) \int_{\mathbb{S}^{d-1}} g^2} \\
&\leq \sqrt{h(2R_A B_d) \underbrace{\int_{2R_A B_d} \int_{r \cap \mathbb{S}^{d-1}} g^2}_{\leq \|g\|_2^2}} \\
&\leq \sqrt{h(2R_A B_d)} \leq \text{mexp}(-c d/2). \quad \square
\end{aligned}$$

[SS17]: 2-layer NN cannot efficiently approximate ball indicators

[SES19]: $\sim \text{ReLU}(1 - \|x\|)$

5. Lemma. $\exists \tilde{g}(x) = \sum_{i=1}^N \varepsilon_i g_i(\|x\|)$, where $\varepsilon_i \in \{\pm 1\}$, g_i is the indicator function of some interval Δb_i . s.t. $w = \hat{g}g / \|\hat{g}g\|_2$ satisfies the condition of the previous lemma. for some universal constant δ .

Corollary. $\|f - \tilde{g}\|_{L^2(g^2)} \geq \Omega(\delta)$.

General message: high-frequency \Rightarrow harder to approximate/learn.

Hierarchical structure?

[Abbe, Boix-Adsera, Misiakiewicz, 22]. Merged-staircase property

[ABM 23] follow up.

[Allen-Zhu, Li 21]. Backward feature correction.

§5.2 Positive results.

Target $f_*(x) = \sigma(1 - \|x\|)$. σ - ReLU.

learner $\left\{ \begin{array}{l} F(x; \mu_1) = \mathbb{E}_{w_1 \sim \mu_1} \|w\| \sigma(v \cdot x) \quad \text{first layer} \\ f(x; \mu_2, \mu_1) = \mathbb{E}_{\substack{(w_2, b_2) \\ \sim \mu_2}} \sigma(w_2 F(x; \mu_1) + b_2) \quad \text{second layer} \end{array} \right.$

Loss. $\mathcal{L}(\mu_1, \mu_2) = \text{MSE} = \frac{1}{2} \mathbb{E}_X \left[(f_*(x) - f(x; \mu_2, \mu_1))^2 \right]$.

GF. $v_i \in \mu_1$ - $\dot{v}_i = \mathbb{E}_X \left\{ S_2(x) (v_i \sigma(v_i \cdot x) + \|v_i\| \sigma'(v_i \cdot x) x) \right\}$.

where $S_2(x) = (f_*(x) - f(x)) \mathbb{E}_{w, b_2} \left\{ \sigma'(w_2 F(x) + b_2) w_2 \right\}$

(the dynamics of 2nd layer neurons are not very important)

1. The infinite-width dynamics are much simpler than the finite-width ones.

Lemma. μ spherically symmetric. $\Rightarrow \int_{\mathbb{R}^d} \|w\| \delta(w \cdot x) = C_T \frac{\int_{\mathbb{R}^d} \|w\|^2}{\sqrt{d}} \|x\|$. $C_T = \Theta(d)$.

Pf. $\int_{\mathbb{R}^d} \|w\| \delta(w \cdot x) = \int_{\mathbb{R}^d} \|w\|^2 \delta(w \cdot x) = \left(\int_{\mathbb{R}^d} \|w\|^2 \right) \left(\int_{\mathbb{R}^d} \delta(w \cdot x) \right)$ □

\downarrow \downarrow
 \mathbb{I} -homogeneity spherical symmetry $\Rightarrow \frac{C_T}{\sqrt{d}} \|x\|$

• Define $\alpha = \frac{C_T}{\sqrt{d}} \int_{\mathbb{R}^d} \|w\|^2$. (\leftarrow well-defined for discrete μ too).

• In the infinite-width limit $F(x; \mu_t) = \alpha \|x\|$. (1)

• Lemma. Spherically symmetric g . $\Rightarrow \int_{\mathbb{R}^d} g(x) \delta(v \cdot x) = \frac{G}{\sqrt{d}} \int_{\mathbb{R}^d} \{g(w) \|x\|\} \|v\|$, $\int_{\mathbb{R}^d} g(x) \delta'(v \cdot x) x = \frac{G}{\sqrt{d}} \int_{\mathbb{R}^d} \{g(w) \|x\|\} \bar{v}$.

$\Rightarrow \frac{d}{dt} \alpha = \int_{\mathbb{R}^d} \frac{\partial \alpha}{\partial w_i} \frac{\partial w_i}{\partial t} = \dots = \frac{4G^2}{d} \int_{\mathbb{R}^d} \{S(x) \|x\|\} \int_{\mathbb{R}^d} \|w\|^2 = \frac{4G}{\sqrt{d}} \int_{\mathbb{R}^d} \{S(x) \|x\|\} \alpha$. (2)

(1), (2) $\Rightarrow \alpha_t \in \mathbb{R}$ characterize the first layer (in the infinite-width limit)

Next goal: poly(d)-width discretization.

2. Symmetrization.

Given $g: \mathbb{R}^d \rightarrow \mathbb{R}$, define its symmetrization $\tilde{g}(x) := \mathbb{E}_{x' \in \|x\| S^{d-1}} g(x')$

Observation - $\tilde{F}(x) = \mathbb{E}_{w \sim \mathcal{D}_1} \left[\mathbb{E}_{x' \sim \|x\| S^{d-1}} \left[\sum_i \sigma(w_i \cdot x') \right]^2 \right]$ $= \frac{G}{\sqrt{d}} \mathbb{E} \|w\|^2 \|x\| = \alpha \|x\|$

not necessarily spherically symmetric $\rightarrow \frac{G}{\sqrt{d}} \|x\|$

i.e. the infinite-width network \rightarrow the symmetrization of the finite-width network.

Observation - $\mathbb{E}_x \left\{ (f_{*}(x) - f(x))^2 \right\} = \mathbb{E}_x \left\{ (f_{*}(x) - \tilde{f}(x))^2 \right\} + \mathbb{E}_x \left\{ (\tilde{f}(x) - f(x))^2 \right\} \quad (3)$

$- 2 \mathbb{E}_x \left\{ (f_{*}(x) - \tilde{f}(x)) (\tilde{f}(x) - f(x)) \right\}$

$\approx \mathbb{E}_x \left\{ (f_{*}(x) - \tilde{f}(x))^2 \right\} + \frac{\overline{W}^2}{2} \mathbb{E}_x \left\{ (F(x) - \tilde{F}(x))^2 \right\} = 0.$

3. poly(d)-width discretization under symmetry.

Q: How to characterize the discretization error? (In general, making $W_2(\mu, \hat{\mu}) \leq \epsilon$ requires exponentially many samples.)

A1: Sample m neurons $\hat{\mu}_0$ from μ_0 . move them according to the infinite-/finite-width dynamics. characterize the deviation.

A2: Characterize the deviation of the relevant function. (In our case, $\bar{F} = F/\alpha$ vs. $\|\cdot\|_2$)

Lemma. $\frac{d}{dt} F(x) \approx -\frac{\bar{W}_0^2}{2} \mathbb{E}_{w_i} \left\{ \langle \nabla_{w_i} F(x), \nabla_{w_i} \mathbb{E}(\tilde{F}(x) - F(x))^2 \rangle \right\}$.

pf sketch. $\frac{d}{dt} \bar{F}(x) = \frac{\frac{d}{dt} F(x)}{\alpha} - F(x) \frac{\dot{\alpha}}{\alpha} = -\frac{1}{\alpha} \mathbb{E}_{w_i} \langle \nabla_{w_i} F(x), \nabla_{w_i} \mathbb{1} \rangle + F(x) \frac{1}{\alpha} \mathbb{E}_{w_i} \langle \nabla_{w_i} \alpha, \nabla_{w_i} \mathbb{1} \rangle$

Recall (3). $\mathbb{1} = \mathbb{1}_1 + \mathbb{1}_2$. Claim. $\nabla \mathbb{1}_1$ contributes little.

By symmetry, $\nabla_{w_i} \mathbb{1}_1 = \beta v_i$ for some $\beta \in \mathbb{R}$.

$$\langle \nabla_{w_i} F(x), v_i \rangle = \langle \nabla_{w_i} (\|w\|^2 \theta(\bar{v} \cdot x)), v_i \rangle = \theta(\bar{v} \cdot x) \langle \nabla_{w_i} \|w\|^2, v_i \rangle = 2\|w\|^2 \theta(\bar{v} \cdot x)$$

$$\langle \nabla_{w_i} \alpha, v_i \rangle = \frac{G}{\sqrt{d}} \langle \nabla_{w_i} \|w\|^2, v_i \rangle = \frac{2G}{\sqrt{d}} \|w\|^2$$

$$\Rightarrow \frac{d}{dt} \bar{F}(x) \Big|_{\mathbb{1}_2} = -\frac{\beta}{\alpha} 2 \mathbb{E}_{w_i} \|w\|^2 \theta(\bar{v} \cdot x) + \bar{F}(x) \frac{\beta}{\alpha} \frac{2G}{\sqrt{d}} \mathbb{E} \|w\|^2 = -2\beta F(x) + 2\beta \bar{F}(x) = 0$$

Corollary. $\frac{d}{dt} \|F - \Pi \cdot \Pi_2\|_{L^2}^2 \lesssim 0$. In words, the discretization error barely grows.

Takeaways.

1. The infinite-width dynamics can be much simpler than the finite-width ones.
2. poly(d)-width discretization is sometimes possible.
 - * Essentially, we are Taylor expanding the dynamics around the infinite-width trajectory and showing that the first-order error terms are good. (no compounding errors).