

Lec 9

Ann.

- * Homework due tomorrow
- * Maybe hybrid Thursday

Recap

* Today is final approx lecture

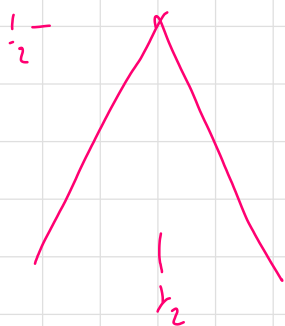
$$\inf_{f \in \mathcal{F}} \mathcal{R}(f) \quad \text{vs} \quad \inf_{g \in \text{anyhow}} \mathcal{R}(g) \text{ else}$$

* Various approx of continuous function (2-hidden & 1-hidden)

* Former approx (connections to ntk, non-worst-case estimates.)

* NTK (near initialization, largewidth network \approx first order Taylor)

* Depth helps



$$\Delta(x) = 2\sigma(x) - 4\sigma(x - \frac{1}{2}) + 2\sigma(x-1) = \begin{cases} 2x & x \in [0, \frac{1}{2}] \\ 2(1-x) & x \in (\frac{1}{2}, 1] \\ 0 & \text{o.v.} \end{cases}$$

$$N_A \left(\begin{array}{l} \text{ReLU network} \\ \text{with } \leq L \text{ layers} \\ \leq m \text{ nodes} \end{array} \right) \leq \left(\frac{2m}{L} \right)^L$$

$$\left[2^L \prod_{j=L} m_j \right]$$

Approximating x^2 .

* Why?

"polarization identity"

$$xy = \frac{1}{2}((x+y)^2 - x^2 - y^2)$$

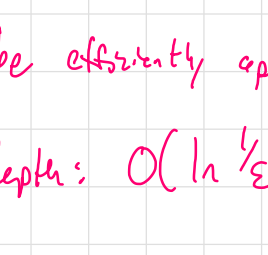
efficient $\Delta^L \Rightarrow$ eff. x^2

\Rightarrow products \Rightarrow multivariable polynomials \Rightarrow Taylor expansions
 \Rightarrow "Sobolev space"

"smalish" deep networks

$$x^2 = \int_0^1 2b \mathbb{1}\{x \geq b\} db$$

$$= \int_0^1 2 \sigma_b(x-b) db$$



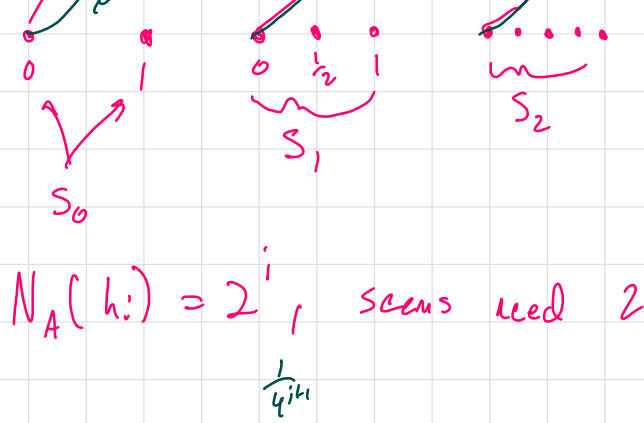
constant; hit can be efficiently approx.

Homework: $\left\lfloor \frac{2}{\sqrt{\epsilon}} \right\rfloor$ with depth: $O(\ln \frac{1}{\epsilon})$.

(Figured out by Dmitry Yarotsky)

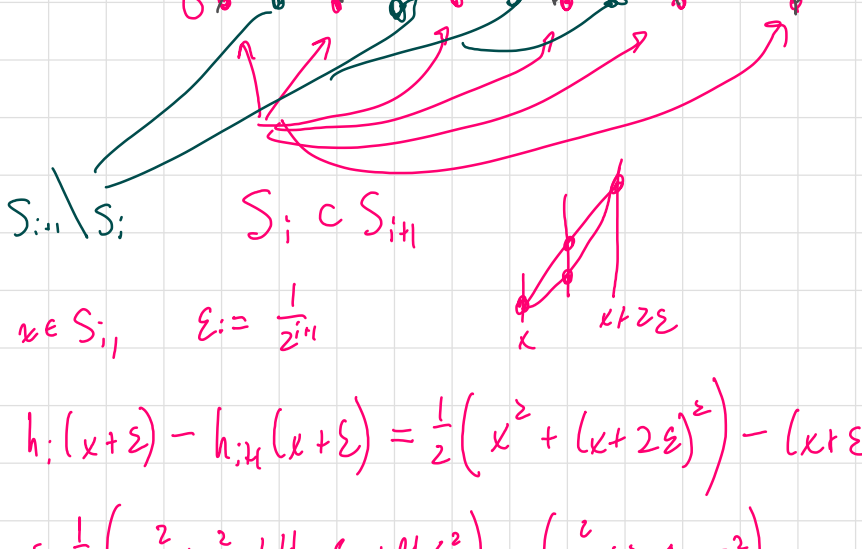
$S_i := \left\{ \frac{0}{2^i}, \frac{1}{2^i}, \dots, \frac{2^i}{2^i} \right\}$ grid $[0,1]$ at scale $\frac{1}{2^i}$

$h_i :=$ secant approximation of x^2 at points in S_i



" h_i is affine interpolation of x^2 along S_i "

$N_A(h_i) = 2^i$, seems need 2^i ReLUs



$$h_i(x+\epsilon) - h_{i+1}(x+\epsilon) = \frac{1}{2}(x^2 + (x+2\epsilon)^2) - (x+\epsilon)^2$$

$$= \frac{1}{2}(x^2 + x^2 + 4x\epsilon + 4\epsilon^2) - (x^2 + 2x\epsilon + \epsilon^2)$$

$$= \epsilon^2$$

By induction
$$h_i(x) = h_0(x) + \sum_{j=1}^i (h_{j+1}(x) - h_j(x))$$

$$= x - \sum_{j=1}^i \frac{\Delta^j(x)}{4^j}$$

Theorem (Yarotsky).

Let h_i be as above. $\delta := \frac{1}{4^{i+1}}$

① h_i can be written with $\leq 4^i$ ReLUs in $\leq 2^i$ layers

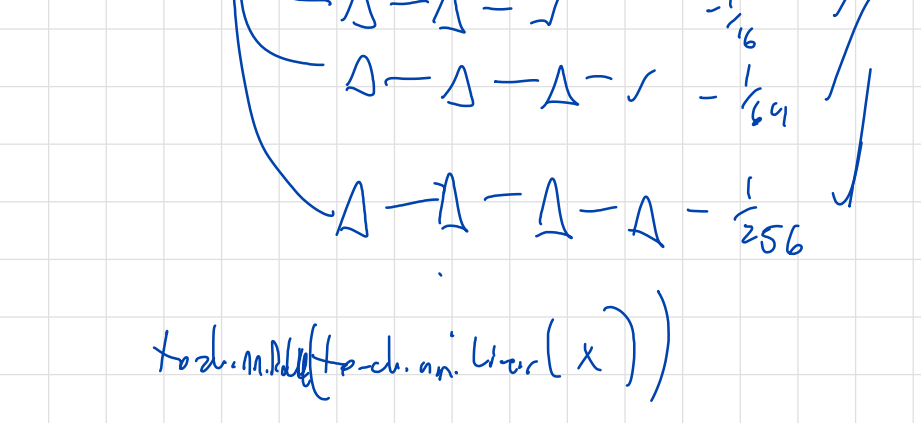
② $\sup_{x \in [0,1]} |h_i(x) - x^2| \leq \frac{1}{4^{i+1}}$

③ Any ReLU network $g: \mathbb{R} \rightarrow \mathbb{R}$ with $\leq L$ layers $\leq m$ nodes

satisfies
$$\int_0^1 (g(x) - x^2)^2 dx \geq \frac{1}{5760 \left(\frac{2m}{L}\right)^L}$$

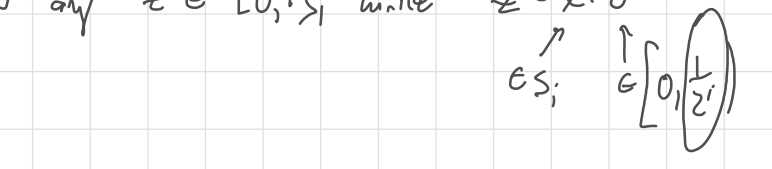
$\geq N_A(g)$

Proof.



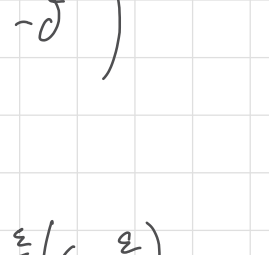
$h_i(x) = x - \sum_{j=1}^i \frac{\Delta^j(x)}{4^j}$

$O(i)$ and $O(i^2)$



total number of nodes $\sim L \cdot m$

② Consider any $z \in [0,1]$, write $z = x + \delta$



$$h_i(x+\delta) - (x+\delta)^2 = \frac{\epsilon-\delta}{\epsilon} x^2 + \frac{\delta}{\epsilon} (x+\delta)^2 - (x+\delta)^2$$

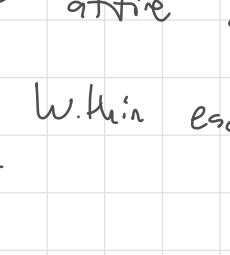
$$= x^2 - \frac{\delta}{\epsilon} x^2 + \frac{\delta}{\epsilon} (x^2 + 2x\delta + \delta^2) - (x^2 + 2x\delta + \delta^2)$$

$$= x^2 \left(1 - \frac{\delta}{\epsilon} + \frac{\delta}{\epsilon} - 1\right) + \delta(2\delta - 2\delta) + (\delta\delta - \delta^2)$$

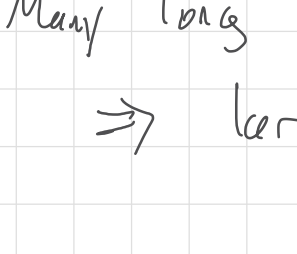
$$= \delta(\delta - \delta)$$

$$\leq \sup_{\delta \in [0, \epsilon]} \delta(\delta - \delta) = \frac{\epsilon}{2} \left(\frac{\epsilon}{2} - \frac{\epsilon}{2}\right) = \frac{\epsilon^2}{4} = \frac{1}{4^{i+1}}$$

③ $N_A(g) \leq \left(\frac{2m}{L}\right)^L$



Constant fraction of affine pieces have length $\frac{1}{2N_A(g)}$. Within each piece:



Many long intervals with large error \Rightarrow large error over $[0,1]$.

Approximating functions with "nice Taylor expansions" (Sobolev space)

Showing \Rightarrow products

$$xy = \frac{1}{2}((x+y)^2 - x^2 - y^2)$$

\nearrow \uparrow \uparrow \uparrow \uparrow

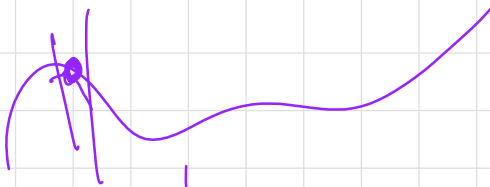
$$\frac{1}{2} \left(h_i(x+y) - h_i(x) - h_i(y) \right)$$

\downarrow

$$4h_i\left(\frac{x+y}{2}\right)$$

$xy \rightarrow \prod_{i \in S} x_i$ Monomials

\rightarrow Polynomial



P_1	P_2	P_3	P_4
P_5	P_6	P_7	P_8
P_9			

"approximate partition of unity"

Lemma 5.5

Office hours

defn "rotational invariance"

M rot matrix

$$M^T M = I$$

$$Mx \stackrel{d.}{=} x$$

↑
"equiv in dist"

$$\mathbb{E} f(x) = \mathbb{E} f(Mx)$$

$$\int \ell(y, f(x)) d_p(x, y)$$

$$\int_{\mathbb{C}_0, \mathbb{I}^d} \ell(f(x)) h(x) dx$$

↑

density

$$\int_{\mathbb{C}_0, \mathbb{I}^d} h(x) dx = 1$$

$$h \geq 0$$

$$\int_{\mathbb{C}_0, \mathbb{I}^d} h = 1$$

$$\int \ell(y, f(x)) d_p(x, y) = \mathbb{E}_{x, y} \ell(y, f(x))$$