

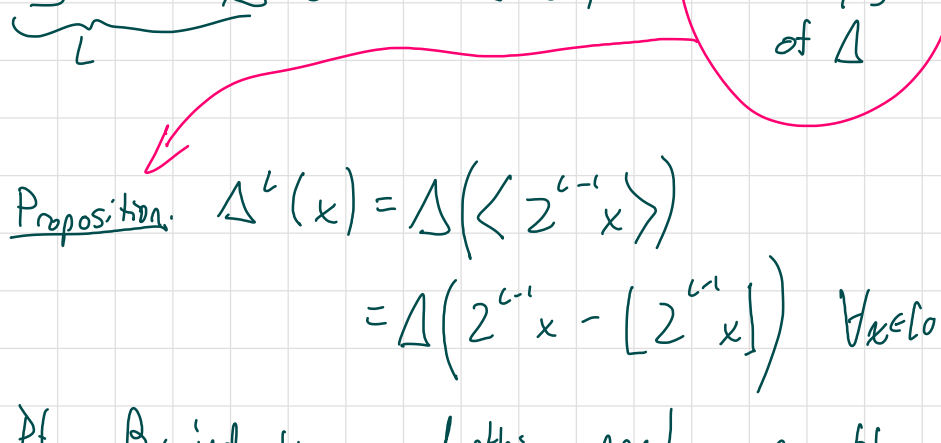
# Lec 8

\* hwk 1 due Wednesday. (No late hwk.)

## Recap Apx:

- \* Elementary/constructive approximations
  - \* "universal approximation" (Stone-Weierstrass) (one hidden layer)
  - \* Fourier/Barron (one hidden layer)
  - \* Neural tangent kernel
- } near initialization

\* Deep vs shallow (non-algorithmic)



$\Delta \circ \dots \circ \Delta (x) = \Delta^L(x)$  2<sup>L-1</sup> regular copies of Δ

Proposition.  $\Delta^L(x) = \Delta(\langle 2^{L-1}x \rangle)$   
 $= \Delta(2^{L-1}x - \lfloor 2^{L-1}x \rfloor) \quad \forall x \in [0, 1]$

Pf. By induction on depth; peel one off every iteration; details in notes.

## Theorem (Telgarsky '15).

$\forall$  depth  $L \geq 2$   
 $f(x) := \Delta^{L^2+2}(x)$  use  $3(L^2+2)$  ReLUs total in  $2(L^2+2)$  layers

$\forall g: \mathbb{R} \rightarrow \mathbb{R} \quad \leq 2^L$  ReLUs in  $\leq L$  layers.

$\int_0^1 |g(x) - f(x)| \geq \frac{1}{32}$ .

## Remark (Why L'?).

Seems to imply existence of natural-ish production problems when risk is large.

## Proof plan.

- ① Characterize  $f(x) := \Delta^{L^2+2}(x) = \Delta(\langle 2^{L^2+2}x \rangle)$  ✓
- ② Define & bound "complexity" of  $g$ . 0 ?
- ③ 1 & 2 together imply theorem ?

Remark Effective justifications of depth & architectural choices are all open.

## ② Define & bound "complexity"

Definition. Given  $f: \mathbb{R} \rightarrow \mathbb{R}$ , define "number of affine pieces"  $N_A(f)$  as

$$N_A(f) = \begin{cases} \infty & f \text{ is not piecewise affine} \\ N & \text{cardinality of smallest partition of } \mathbb{R} \text{ such that } f \text{ is affine restricted to any partition element} \end{cases}$$

} into intervals

Remark.  $\exists$  multivariate version; similar reasoning to VC dimension in later lectures.

Lemma. Given  $f: \mathbb{R} \rightarrow \mathbb{R}$ , ReLU network with  $m$  ReLUs,  $L$  layers,  $(m_1, \dots, m_L)$  with  $\sum_{i=1}^L m_i = m$ .

\*  $N_A(f) \leq \left(\frac{2^m}{L}\right)^L$ .

\* For any node in layer  $i$ , letting  $g: \mathbb{R} \rightarrow \mathbb{R}$  denote its output as a function of whole network input,

$N_A(g) \leq 2^i \prod_{j=1}^i m_j$ .

Remark. "Complexity"  $N_A$  grows exponentially in  $L$ , polynomially in  $m$ .

Remark Eldan & Shamir '15 poly(d) nodes

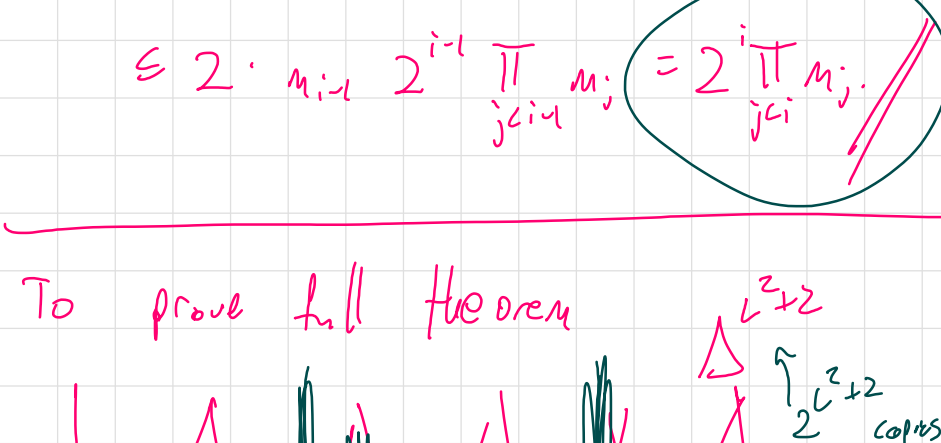
$\exists f: \mathbb{R}^d \rightarrow \mathbb{R}$  with 2 hidden layers, s.t.  $\forall g: \mathbb{R}^d \rightarrow \mathbb{R}$  with 1 hidden layer and  $o(2^d)$  nodes s.t.

$\|f - g\|_u \geq \text{constant}$ .

Lemma. Let  $\{g_1, \dots, g_k, f, a_1, \dots, a_k, b\}$  be given.

- ①  $N_A(f + g_i) \leq N_A(f) + N_A(g_i)$  [width additively increases complexity]
  - ②  $N_A(\sum_{i=1}^k a_i g_i + b) \leq \sum_{i=1}^k N_A(g_i)$
  - ③  $N_A(f \circ g) \leq N_A(f) N_A(g)$  [depth multiplicatively increases complexity]
  - ④  $N_A(x \mapsto f(\sum_{i=1}^k a_i g_i + b)) \leq N_A(f) \sum_{i=1}^k N_A(g_i)$ .
- e.g., one node

## Proof.



Sort the changepoints of the two partitions; within any adjacent pair, both  $f$  &  $g$  are affine; thus  $f+g$  also affine in that interval, thus

$N_A(f+g) \leq N_A(f) + N_A(g)$ .

② (Sketch.) Remark 6 (translation), use ① inductively.

③  $N_A(f \circ g)$ .



Let  $\mathcal{P}_A(g)$  denote the corresponding partition. For any  $u \in \mathcal{P}_A(g)$ ,  $g(u)$  is an interval, thus  $N_A(f|_{g(u)}) \leq N_A(f)$ .

Given any  $T \in \mathcal{P}_A(f|_{g(u)})$ , then  $f|_{g(u)}$  is affine in  $T$ , and  $g$  is affine in  $u \cap g^{-1}(T)$ .

$$N_A(f \circ g) \leq \sum_{u \in \mathcal{P}_A(g)} N_A((f \circ g)|_u) \leq \sum_{u \in \mathcal{P}_A(g)} \sum_{T \in \mathcal{P}_A(f|_{g(u)})} N_A((f \circ g)|_{u \cap g^{-1}(T)}) \leq N_A(g) \cdot N_A(f)$$



Remark. For images of  $\mathcal{P}_A(g)$  to hit all pieces of  $\mathcal{P}_A(f)$  is delicate;  $\Delta \circ \Delta$  approximately meets this bound.

## To prove

$N_A(g) \leq 2^i \cdot \prod_{j=1}^i m_j$

$$N_A(g) = N_A(x \mapsto \left( \sum_{i=1}^{\min} a_i h_i(x) + b \right)) \leq 2 \cdot \sum_{i=1}^{\min} N_A(h_i) \leq 2 \cdot \min_i 2^{i-1} \prod_{j=1}^{i-1} m_j = 2^i \prod_{j=1}^i m_j$$

## To prove full theorem

