

# Lecture 10: architectural benefits

Announcements:

- \* Next 5 lectures on zoom.
- \* Hw2 probably out Oct 2.

Original material, now "open questions"

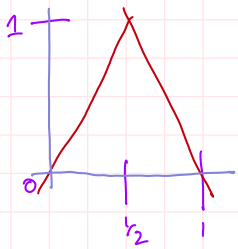
- ① Architectural benefits of convolution & attention layers.
- ② Benefits of activation choice
- ③ Modeling distributions
- ④ low norm approximation
- ⑤ effect of depth [non-algorithms]
  - (i) polynomial networks (sum-product networks)
  - (ii) many layers to fewer layers
  - (iii) depth 3 vs depth 2

Many-to-few layers

We'll construct efficient approximations w. k many layers; warning: algorithm's relevance unknown.

All based on "triangle mapping"  $\Delta: \mathbb{R} \rightarrow \mathbb{R}^n$

also: poor scaling w/ dimension

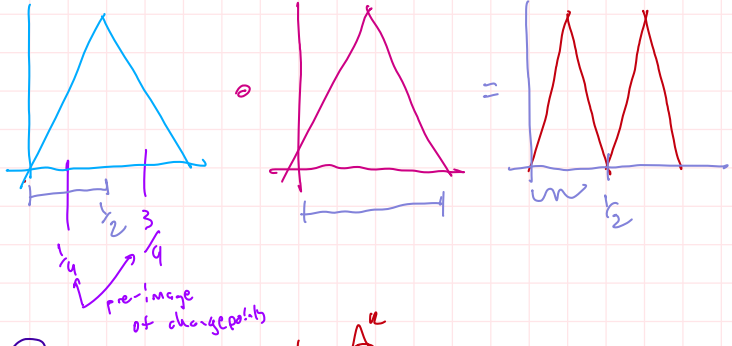


$$\Delta(x) = \begin{cases} 2x & x \in [0, \frac{1}{2}) \\ 2(1-x) & x \in [\frac{1}{2}, 1) \\ 0 & \text{otherwise} \end{cases}$$

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

[2-layer, 3-node ReLU network.]

$$(\Delta \circ \Delta)(x) = \Delta(\Delta(x))$$



Remark: Approximating this with shallow network needs exponentially many nodes, and we used 3k nodes.

One attempt to capture fractal structure:

Proposition: "fractal property"  $\Delta^L(x) = \Delta(2^L x - \lfloor 2^L x \rfloor) = \Delta(\langle 2^L x \rangle)$  (2^L side by side then squeezed.)

Proof: Induction on #layers i:  $(\forall x \in [0,1] \text{ not on outside } \Delta^i(x) = \Delta(\langle 2^i x \rangle)) \Rightarrow (\forall x \in [0,1] \Delta^{i+1}(x) = \Delta(\langle 2^{i+1} x \rangle))$

Base case  $\Delta^1(x) = \Delta(\langle 2^1 x \rangle) \checkmark$

I.H.: Let  $x \in [0,1]$  & target depth  $i+1$  be given, assume I.H.

$$\Delta^{i+1}(x) = \Delta^i(\Delta(x)) = \Delta^i(2x) = \Delta(\langle 2^{i+1} \cdot 2x \rangle) = \Delta(\langle 2^i \cdot 2x \rangle)$$

$$= \Delta(\langle 2^i x \rangle) \quad \left\{ \begin{array}{l} x \in [\frac{1}{2}, 1] \text{ if } i=1, \Delta^1(x) = x \\ \Delta^i(x) = \Delta^i(\Delta(\Delta(x))) = \Delta^i(\Delta(2(1-x))) \\ \text{(using } \Delta(y) = \Delta(1-y)) = \Delta^i(\Delta(1-2(1-x))) = \Delta^i(\Delta(2x-1)) \\ = \Delta^i(2x-1) = \Delta(\langle 2^i(2x-1) \rangle) = \Delta(\langle 2^i x \rangle) \end{array} \right.$$

Applications of  $\Delta^L$ :

- ① Approximate  $x^2$  and  $x \cdot y$  up to accuracy  $\epsilon$  with  $1/\epsilon^2$  nodes & layers.  $\Rightarrow$  polynomials  $\Rightarrow$  smooth functions.
- ② Parity function on  $d$  bits:  $x \in \{\pm 1\}^d, \prod_{i=1}^d x_i = \Delta^L(\frac{d + \sum x_i}{2d})$ . [d=2^L]
- ③ If  $f(z) = f(1-z) \forall z \in [0,1]$ , then  $f(\Delta^L(x)) = f(\Delta(\langle 2^L x \rangle))$ ,  $2^{L+1}$  copies of  $f$ ! "viral fractal property"
- ④ In particular, can use ③ to extract bits

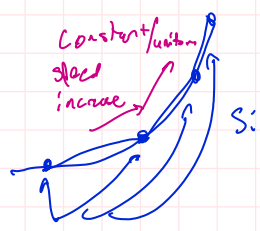
Approximating  $x^2$  via  $\Delta$

$\forall \epsilon > 0 \quad x^2 = \int_0^{\infty} 2\sigma(x-b) db$  : places poles uniformly;

need of enough uniformity for  $\Delta^L$ .

Consider progressive affine approximations of  $x^2$ :

$\forall i \quad S_i = \left\{ \frac{0}{2^i}, \frac{1}{2^i}, \frac{2}{2^i}, \dots, \frac{2^i}{2^i} \right\}, \quad h_i :=$  affine interpolation of  $x^2$  along  $S_i$



= given  $x \in [0, 1]$ ,  $\exists \tau \in [0, 1]$   
 &  $k \in \{0, \dots, 2^i - 1\}$  s.t.  $x = \frac{k+\tau}{2^i}$   

$$h_i(x) = (1-\tau) \left( \frac{k}{2^i} \right)^2 + \tau \left( \frac{k+1}{2^i} \right)^2$$

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

$$= x + \sum_{j=1}^i (h_{j+1}(x) - h_j(x)).$$

Consider  $h_{j+1}(x) - h_j(x)$

\* case 1  $x \in S_j$  : Since  $S_{j+1} \supseteq S_j$ ,  $h_{j+1}(x) - h_j(x) = x^2 - x^2 = 0$ .

\* case 2  $x \in S_{j+1} \setminus S_j$  : pick  $k$  integer s.t.  $x = \frac{2k+1}{2^{j+1}}$ ,

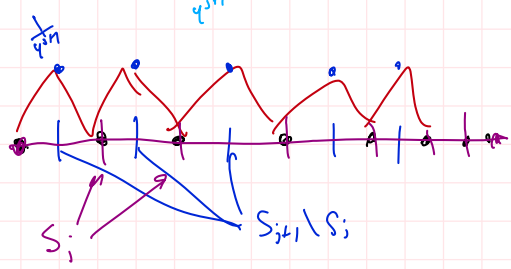
$$h_{j+1}(x) - h_j(x) = \left( \frac{2k+1}{2^{j+1}} \right)^2 - \frac{1}{2} \left( \left( \frac{2k}{2^{j+1}} \right)^2 + \left( \frac{2k+2}{2^{j+1}} \right)^2 \right)$$

$$= \frac{1}{4^{j+1}} \left( \underbrace{4k^2 + 4k + 1}_{\square} - \frac{1}{2} \left( \underbrace{4k^2}_{\square} + \underbrace{4k^2 + 8k + 4}_{\square} \right) \right)$$

$$= -\frac{1}{4^{j+1}}$$

\* case 3  $x \in [0, 1] \setminus S_{j+1}$ :

$$h_{j+1}(x) - h_j(x) = \frac{-\Delta^j(x)}{4^{j+1}}$$

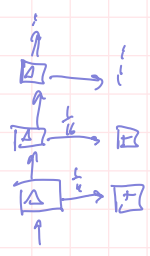


$$\Rightarrow h_i(x) = x - \sum_{j=0}^{i-1} \frac{\Delta^j(x)}{4^{j+1}}$$

Theorem: Let  $h_i$  be the piecewise affine interpolation of  $x^2$  along  $S_i$ .

①  $h_i$  can be written as Roll network with  $S_i$  nodes in  $2^i$  layers.

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



Proof ① what we had before, but reuse  $\Delta$ :

② [Basically same proof as  $h_i(x) - h_{i+1}(x)$ , see notes.]

Remark: ① need  $\lceil \ln(1/\epsilon) \rceil$  layers & nodes to get  $\epsilon$  accuracy

② Multiplication via polarization identity:  $xy = \frac{1}{2} \left( (x+y)^2 - x^2 - y^2 \right)$ .

$\Rightarrow$  poly results  $\Rightarrow \dots$