

Part 1

Question 1

Is it possible to define a function $f: \{-1\}^d \rightarrow \{-1\}$ that the lower bound on B is larger than 2^{d-1} ?

(In order to let a shallow network with width B realizes f)

Answer: no! it's not possible.

Lemma: We can realize each f with $\leq 2^{d-1}$ AND gates

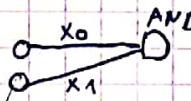
Proof that represent $\{x | f(x) = 1\}$

We will prove it by induction on d .

base: $d=2$

Let $f: \{-1\}^2 \rightarrow \{-1\}$. We will show that a network with width $2^{d-1} = 2^{2-1} = 2$ can realize f . Define $A := \{x | f(x) = 1\}$

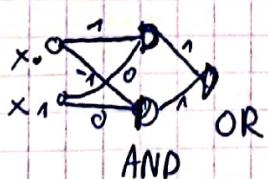
① If $|A| = j \leq 2$, then we can build network with width j , s.t. each AND gate represents some $x \in A$

[if $x = x_0 x_1$ then  will be the weights to this gate] and the weights after each AND gate

will be 1. In that case, $N(x) = 1$ if and only if one of the AND gates outputs 1 if and only if $x \in A$ by definition of the gates $\Rightarrow N$ realizes f with at most 2 width.

② if $|A| = 4$, then $f \equiv 1$, the network N can be

$N:$



So, for every $x = x_0 x_1 \in \{-1\}^2$, $x_0 = 1$ or $x_0 = -1$

So one of the AND gates outputs 1 as desired (then the OR gate outputs 1). Thus N realizes f with 2 width.

③ If $|A|=3$, then by let's sign $A = \{y_1, y_2, y_3\}$,
 $y_i = x_0^i x_i$.

By the pigeonhole principle there must exist y_i, y_j

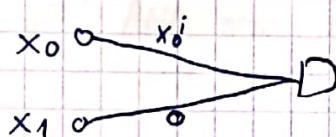
$i \neq j$ s.t. $x_0^i = x_0^j$ (or $x_i = x_j$) and

$x_i^i \neq x_j^j$ (or $x_0^i \neq x_0^j$)

W.l.o.g we'll assume $x_0^i = x_0^j$ and $x_i^i \neq x_j^j$ then

the AND gate which will be defined now, can represents

y_i, y_j and only them:



Thus, we can build network with such AND gate and another AND gate for y_k ($k \neq i, k \neq j$, the left item on A) and as before the weights after the AND gates are 1.

So, we've built a network with 2 AND gates that realizes f .

Thus, the induction base for $d=2$ is correct.

Step:

Let's assume that we can build a network with width 2^{d-2} such that can realize any function $g: \{-1, 1\}^{d-1} \rightarrow \{-1, 1\}$.

Let $f: \{-1, 1\}^d \rightarrow \{-1, 1\}$ be a function.

We will prove that we can build a network with at most 2^{d-1} AND gates to realize f .

Let's sign, as in the induction base, $A := \{y \mid f(y)=1\}$

and define $B = \{y = x_1 x_2 \dots x_d \mid y \in A \text{ and } x_1 = 1\}$

$C = \{y = x_1 x_2 \dots x_d \mid y \in A \text{ and } x_1 = -1\}$

Of course, it holds $A = B \cup C$

Lemma

- ① We can build $\leq 2^{d-2}$ AND gates which represent B 's items
- ② We can build $\leq 2^{d-2}$ AND gates which represent C 's items.

If the lemma is true, then we can build a network with $\leq 2^{d-2} + 2^{d-2} = 2^{d-1}$ AND gates represents $B \cup C = A$ as desired.

Proof of lemma

We'll prove ① and the proof for ② will be the same.

$$B = \{y = x_1 x_2 \dots x_d \mid u \in A, x_1 = 1\}$$

$$\text{define } g: \{\pm 1\}^{d-1} \rightarrow \{\pm 1\} \quad g(x_2 \dots x_d) = f(1x_2 \dots x_d)$$

Then by the induction we can build $\leq 2^{d-2}$ AND gates that represent $\{x_2 \dots x_d \mid g(x_2 \dots x_d) = 1\}$, thus, if

for each gate we'll add a weight of 1 on x_1 (in the

sense of $f: \{\pm 1\}^d \rightarrow \{\pm 1\}$) then these $\leq 2^{d-2}$ AND

gates will assure $\{x_2 \dots x_d \in \{x_2 \dots x_d \mid g(x_2 \dots x_d) = 1\} \text{ and } x_1 = 1\}$

which means, by definition of g , they will assure and

represent the items in $B = \{1x_2 \dots x_d \mid f(1x_2 \dots x_d) = g(x_2 \dots x_d) = 1\}$

as desired.

□

Foundations of Deep Learning – Homework Assignment #2

Adi Alm & Tomer Epshteyn

Part 1: (2)

Question:

Prove that with polynomial width \bar{B} the deep network cannot realize all possible functions. In particular derive an exponential in d lower bound on \bar{B} required in order for its hypothesis space to be y^X .

Solution:

Denote by $\bar{H}_{\bar{B}}$ the class of functions realizable by deep AND-OR networks with width \bar{B} (as defined in class). We will prove that $\bar{H}_{\bar{B}} \neq y^X$ by combinatorial reasoning for $\bar{B} \in \mathcal{O}(\text{poly}(d))$. And that for $\bar{H}_{\bar{B}} = y^X$ we need \bar{B} to be exponential in d .

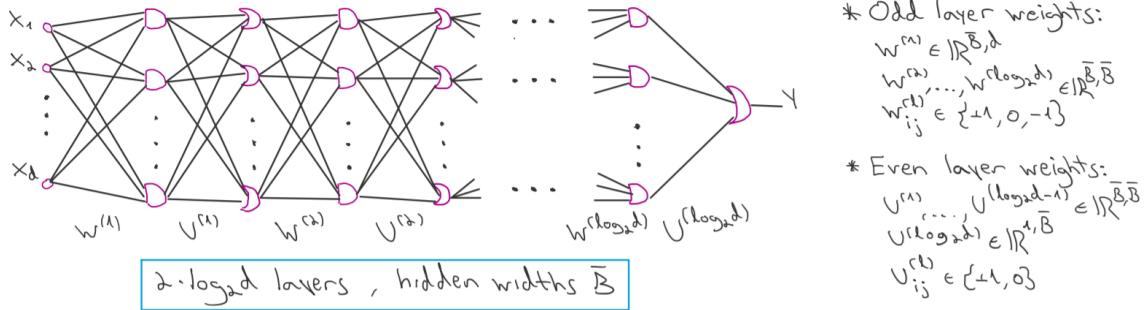
- $|y^X| = |y|^{|X|} = 2^{2^d}$
- $|\bar{H}_{\bar{B}}| = ?$

$$|\bar{H}_{\bar{B}}| \leq \#\{\text{Parameter combinations in AND - OR deep nets of width } \bar{B}\}$$

Let's count the number of parameter combinations. Reminder:

Deep network

Denote by $\bar{H}_{\bar{B}}$ ($\bar{B} \in \mathbb{N}$) the mappings realizable by the deep network:



- $W^{(1)}$ has $\bar{B} \times d$ parameters from $\{-1, 0, +1\}$, yielding: $3^{\bar{B} \times d}$ combinations
- $W^{(j)}$ has \bar{B}^2 parameters from $\{-1, 0, +1\}$, for $j \in \{2, \dots, \log_2(d)\}$, yielding: $3^{(\log_2(d)-1) \times \bar{B}^2}$ combinations
- $U^{(\log_2(d))}$ has \bar{B} parameters from $\{0, +1\}$, yielding: $2^{\bar{B}}$ combinations
- $U^{(j)}$ has \bar{B}^2 parameters from $\{0, +1\}$, yielding: $2^{(\log_2(d)-1) \times \bar{B}^2}$ combinations.

Bringing it all together:

$$|\bar{H}_{\bar{B}}| \leq 3^{\bar{B} \times d} \times 3^{(\log_2(d)-1) \times \bar{B}^2} \times 2^{\bar{B}} \times 2^{(\log_2(d)-1) \times \bar{B}^2}$$

Let $\bar{B} \in \mathcal{O}(\text{poly}(d))$.

$$\begin{aligned}
 |\bar{H}_{\bar{B}}| &\leq 3^{\bar{B} \times d} \times 3^{(\log_2(d)-1) \times \bar{B}^2} \times 2^{\bar{B}} \times 2^{(\log_2(d)-1) \times \bar{B}^2} = 3^{d\bar{B} + (\log_2(d)-1)\bar{B}^2} 2^{\bar{B} + (\log_2(d)-1)\bar{B}^2} \leq \\
 &\leq 4^{d\bar{B} + (\log_2(d)-1)\bar{B}^2} 4^{\bar{B} + (\log_2(d)-1)\bar{B}^2} = 4^{d\bar{B} + (\log_2(d)-1)\bar{B}^2 + \bar{B} + (\log_2(d)-1)\bar{B}^2} = \\
 &\stackrel{(*)}{\leq} 4^{(d+1)\bar{B} + 2(\log_2(d)-1)\bar{B}^2} = 4^{(d+1)\bar{B}^2 + 2(d+1)\bar{B}^2} = 4^{3(d+1)\bar{B}^2} = 2^{6(d+1)\bar{B}^2}
 \end{aligned}$$

(Where (*) holds because $\bar{B} \leq \bar{B}^2$, $\log_2(d) - 1 \leq d + 1$, and monotonicity of $x \mapsto 4^x$)

For any $\bar{B} \in \mathcal{O}(\text{poly}(d))$:

$2^{6(d+1)\bar{B}^2} = 2^{\mathcal{O}(\text{poly}(d))}$ and in particular $2^{6(d+1)\bar{B}^2} < 2^{2^d}$ so $|\bar{H}_{\bar{B}}| < |y^X|$.

So,

$\bar{H}_{\bar{B}} \neq y^X$ for $B \in \mathcal{O}(\text{poly}(d))$.

We will now derive an exponential lower bound for \bar{B} for us to achieve $\bar{H}_{\bar{B}} = y^X$.

We saw

$$|\bar{H}_{\bar{B}}| \leq 2^{6(d+1)\bar{B}^2}$$

For $\bar{H}_{\bar{B}} = y^X$ we must have

$$2^{6(d+1)\bar{B}^2} \geq 2^{2d}$$

\Rightarrow

$$6(d+1)\bar{B}^2 \geq 2^d \Rightarrow \bar{B}^2 \geq \frac{2^d}{6(d+1)}.$$

So for $\bar{H}_{\bar{B}} = y^X$, we must have $\bar{B}^2 \geq \frac{2^d}{6(d+1)}$. And in particular $\bar{B} \geq \frac{2^{\frac{d}{2}}}{\sqrt{6(d+1)}}$.

An exponential in d lower bound on \bar{B} , as desired.

■

Part 2

Question 1

proposition: A shallow network of width $B \geq 2$ can realize any piecewise linear mapping with $\leq B$ pieces.

Conversely, any mapping realizable by such network is piecewise linear with $\leq B+1$ pieces.

Proof

First, I will explain why the condition $B \geq 2$ is necessary.

Explanation:

By definition of a piecewise linear function, we only have to explain why the proposition doesn't hold for the case $B=1$.

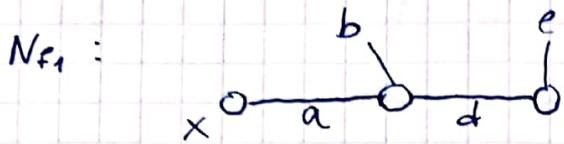
In that case, it means that the function is linear in $(-\infty, \infty)$ which means, each function f which is piecewise linear with 1 piece is $f(x) = \alpha_f x + \beta_f \quad \forall x \in (-\infty, \infty)$

Let f_1 be a linear function $f_1(x) = \alpha_1 x + \beta_1$
s.t. $\alpha_1 \neq 0$

Assume, in contradiction, that the proposition holds for the case $B=1$.

Then, by the proposition, there exists a shallow network with width $B=1$ that can realize f_1 .

Let's sign this network as N_{f_1} and the weights as:



which means there exist $a, b, d, e \in \mathbb{R}$ s.t

$$N_{f_1}(x) = d \cdot \text{RELU}(ax + b) + e$$

Of course $a \neq 0$, otherwise N_{f_1} would output the same output for each input, in contradiction that it realizes f_1 which outputs other values for each two inputs ($f_1(x) = \alpha_1 \cdot x + \beta_1$, and $\alpha_1 \neq 0$).

It holds that $ax + b = 0$ for $x = -\frac{b}{a}$.

If $a < 0$ then $ax + b \leq 0$ for $-\frac{b}{a}, -\frac{b}{a} + 1$.

If $a > 0$ then $ax + b \leq 0$ for $-\frac{b}{a}, -\frac{b}{a} - 1$.

In both cases we found 2 different points which $ax + b \leq 0$, and thus $\text{RELU}(ax + b) = \max(0, ax + b) = 0$, thus $N_{f_1}(x)$ outputs $(d \cdot 0 + e = e)$ for both points in contradiction that f_1 never outputs the same value for 2 different points, so N_{f_1} doesn't realize f_1 , in contradiction.

\Rightarrow For our proposition, it must hold $B \geq 2$.

□

Part 2 - question 1

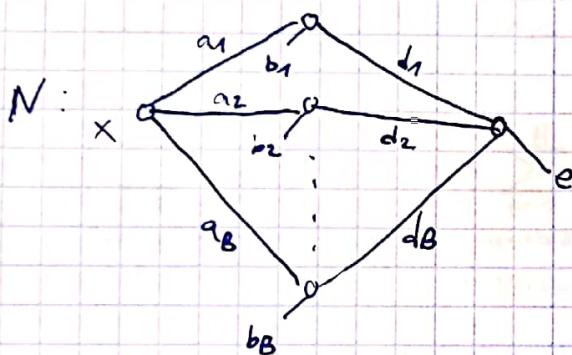
Now, we will prove the following statement/proposition:

Any mapping realizable by such network is piecewise linear with $\leq B+1$ pieces.

proof

Let sign such network as N and its weights as

$$\{a_i, b_i, d_i\}_{i=1}^B, e$$



$$N(x) = \sum_{i=1}^B \text{RELU}(a_i x + b_i) \cdot d_i + e$$

For each $1 \leq i \leq B$ s.t $a_i \neq 0$ it holds that

at $c_i := \frac{-b_i}{a_i}$: $a_i c_i + b_i = 0$. We will sign these

points as $\{c_1, \dots, c_j\}$ for $j = |\{a_i | a_i \neq 0\}| \leq B$

and assume w.l.o.g that $c_1 < c_2 < \dots < c_j$ (otherwise we would rename them) and $c_0 = -\infty$, $c_{j+1} = \infty$.

Then, by definition of these points, we conclude that

for each interval (c_i, c_{i+1}) $[0 \leq i \leq j]$ it holds

for each $1 \leq i \leq B$ $a_i x + b_i \geq 0$ or $a_i x + b_i \leq 0$, for all $x \in (c_i, c_{i+1})$.

We have got that each interval (c_i, c_{i+1}) holds that in it $a_i x + b_i \geq 0$ or $a_i x + b_i \leq 0$ for each $i \in \mathbb{N}$.

Thus, For a specific interval (c_i, c_{i+1}) :

$$N(x) = \sum_{i=1}^B \text{RELU}(a_i x + b_i) d_i + e$$

$$\begin{aligned} \text{when } a_i x + b_i \leq 0 & \quad \checkmark = \sum_{i=1}^B (a_i x + b_i) d_i + e = \sum_{\substack{i=1 \\ a_i x + b_i \geq 0 \\ \text{in } (c_i, c_{i+1})}}^B a_i d_i x + b_i d_i + e \\ \text{then } \text{RELU}(1) = 0 & \quad a_i x + b_i \geq 0 \\ & \quad \text{in } (c_i, c_{i+1}) \end{aligned}$$

$$= \left(\sum_{\substack{i=1 \\ a_i x + b_i \geq 0 \\ \text{in } (c_i, c_{i+1})}}^B a_i d_i \right) x + \left(\sum_{\substack{i=1 \\ a_i x + b_i \geq 0 \\ \text{in } (c_i, c_{i+1})}}^B b_i d_i + e \right)$$

which means that N is linear in (c_i, c_{i+1}) , but that's true for each $0 \leq i \leq j$ and thus we conclude that N realizes a piecewise linear function with $j+1 \leq B+1$ pieces as we wished.

□

Now, we will prove the last part of the question:

A shallow network of width $B \geq 2$ can realize any piecewise linear mapping with $\leq B$ pieces.

Proof

We will prove the proposition for linear mappings with B pieces and if there would be less than B pieces, we would use the same solution ~~but another way~~, but put in the last neurons of the network zeros on their weights.

Let f be a linear mapping with B pieces.

By definition, there exist:

$$-\infty = c_0 < c_1 < c_2 < \dots < c_{B-1} < c_B = \infty$$

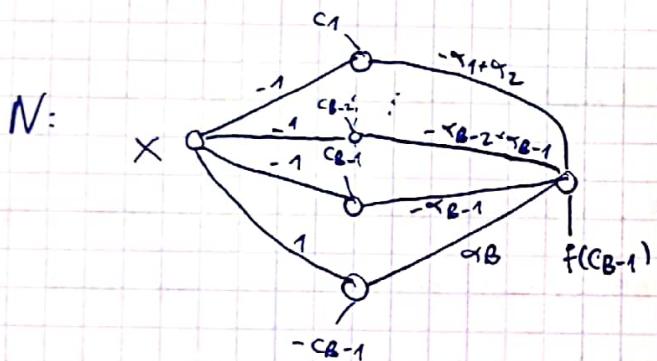
and $\{\alpha_i, \beta_i\}_{i=1}^B$

such that:

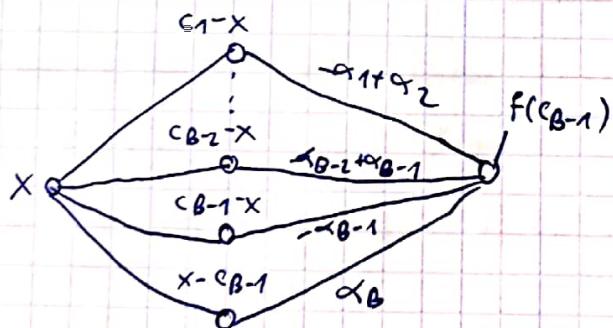
$$\forall 0 \leq i \leq B-1: \forall x \in (c_i, c_{i+1}) \quad f(x) = \alpha_{i+1}x + \beta_{i+1}.$$

We will build a network with width B that realizes f :

We will define the network N as the following:



more clearly:



We'll prove that $N(x) = f(x)$ for every x and we will finish.

① We will prove that $\forall x \geq c_{B-1} \quad N(x) = f(x)$ and then

② we will prove ~~that~~ $\forall 0 \leq i \leq B-2$ that $N(x) = f(x)$ on each (c_i, c_{i+1})

$0 \leq i \leq B-2$

① Let $x \geq c_{B-1}$, then $\forall 1 \leq i \leq B-1 \quad c_i - x \leq 0$,

$$\text{thus } \text{RELU}(c_i - x) = 0$$

and we conclude that

$$N(x) = \sum_{i=1}^{B-1} \text{RELU}(c_i - x) \cdot d_i + \text{RELU}(x - c_{B-1}) \cdot \alpha_B + f(c_{B-1})$$

$$= \text{RELU}(x - c_{B-1}) \cdot \alpha_B + f(c_{B-1})$$

$$x - c_{B-1} \geq 0$$

$$= (x - c_{B-1}) \alpha_B + f(c_{B-1})$$

$$[\text{Note that at } x = c_{B-1} \quad N(x) = f(c_{B-1})]$$

$$= \alpha_B \cdot x - \alpha_B c_{B-1} + f(c_{B-1})$$

Thus, $N(x)$ is linear on $x \geq c_{B-1}$, and has the same slope

as f has on $x \geq c_{B-1}$ (the slope is α_B) and also $N(c_{B-1}) = f(c_{B-1})$

Then, we conclude that they represent the same line on

$$x \geq c_{B-1} \Rightarrow N(x) = f(x) \quad \forall x \geq c_{B-1}. \quad \square$$

② Let $0 \leq i \leq B-2$. We'll show that $N(x) = f(x)$

on (c_i, c_{i+1}) . It holds for every $x \in (c_i, c_{i+1})$:

$$N(x) = \sum_{j=1}^{B-1} \text{RELU}(c_j - x) d_j + \text{RELU}(x - c_{B-1}) \alpha_B + f(c_{B-1})$$

$$= \sum_{j=i+1}^{B-1} (c_j - x) d_j + f(c_{B-1})$$

$$= -\left(\sum_{j=i+1}^{B-1} d_j\right)x + \sum_{j=i+1}^{B-1} c_j d_j + f(c_{B-1})$$

Thus N is linear on (c_i, c_{i+1}) and the slope there is

$$\begin{aligned} -\left(\sum_{j=i+1}^{B-1} d_j\right) &= -\left(d_{i+1} + d_{i+2} + \dots + d_{B-1}\right) \\ &= -\left(\overbrace{-\alpha_{i+1} + \alpha_{i+2}}^{\alpha_{i+1}}, \overbrace{\alpha_{i+2} + \alpha_{i+3}}^{\alpha_{i+2}}, \dots, \overbrace{+\alpha_{B-1} - \alpha_{B-1}}^{\alpha_{B-1}}\right) \\ &= -(-\alpha_{i+1}) = \alpha_{i+1} \end{aligned}$$

Until now, we have got that for each $0 \leq i \leq B-2$
 N acts on (c_i, c_{i+1}) as a linear function with
 the same slope which f acts on that interval,
 that means:

$$f(x) = \alpha_{i+1}x + \beta_{i+1} \quad (\text{by definition})$$

$$N(x) = \gamma_{i+1}x + \delta_{i+1}$$

$$\text{where } \delta_{i+1} = \sum_{j=1+1}^{B-1} c_j d_j + f(c_{B-1})$$

So it only remains for us to explain $\gamma_{i+1} = \beta_{i+1}$

so we conclude that N realizes the function f .

We have already proved that $N(c_{B-1}) = f(c_{B-1})$

and since f and N are continuous (N is composition of continuous functions, and f by definition is continuous)

then it means that N acts on $[c_{B-2}, c_{B-1}]$ exactly

like f does \rightarrow which means $\gamma_{B-1} = \beta_{B-1}$.

Same Slope

and same value at one point. The last also means that $N(c_{B-2}) = f(c_{B-2})$, and we can yield that $\gamma_{B-2} = \beta_{B-2}$.

Iteratively, we can yield that $\gamma_{i+1} = \beta_{i+1} \quad \forall 0 \leq i \leq B-2$

Thus, N realizes the function f as desired.

□

Question 2

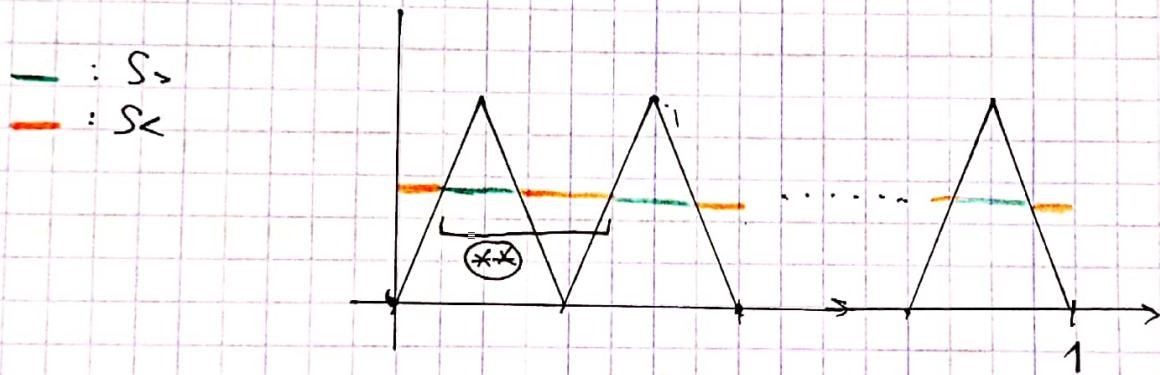
In the context of the expressive efficiency with inapproximability analysis delivered in class - derive a lower bound on the number of intervals in $S_>$ and $S_<$ that a piecewise linear mapping with $\leq B+1$ pieces must miss.

Solution

We'll remind the definitions for $S_>$, $S_<$:

$$S_> := \{x \in [0,1] : g^{\circ L-1}(x) > \frac{1}{2}y\}$$

$$S_< := \{x \in [0,1] : g^{\circ L-1}(x) < \frac{1}{2}y\}$$



A function misses interval in $S_>$ if $f \leq \frac{1}{2}$ on that interval.

" " " " " $S_<$ " " $f \geq \frac{1}{2}$ " " " .

We proved that the number of intervals is $2^{L-1} + 1$ - there are $2^{L-2} + 1$ intervals in $S_<$ and 2^{L-2} intervals in $S_>$.

We'll look at the first interval of $S_<$, we'll sign it $I_1^<$.

Now we'll define 2^{L-2} tuples such the i 'th tuple $1 \leq i \leq 2^{L-2}$ consists of the i 'th interval of $S_>$ and the following interval in $S_<$. For example, $\textcircled{**}$ will be the first tuple.

Lemma

Let f be a piecewise linear with $\leq B+1$ pieces.

We'll prove that f can avoid missing at most

$$\lceil \frac{1}{2}(2^{L-1}+1) + \frac{1}{2}(B+1) \rceil$$
 intervals.

Proof of lemma

Since f is piecewise linear with $\leq B+1$ pieces, it has at most B "breaking points" where the slope is change there.

We'll look at the j breaking points which inside one of our tuple $[\{x \mid \begin{array}{l} x \text{ is breaking point of } f \\ \text{and } x \in [\frac{1}{4} \cdot 2^{L-2}, 1] \end{array}\}, y]$.

of course $j \leq B$ (there are at most B breaking points on all \mathbb{R}).

We'll count how many ~~of the~~ pairs of intervals which are the tuples we've defined - f can avoid missing both intervals in the tuple (the one in S_S , and the one in S_L).

*

Note that between each two breaking points ~~f~~ is linear and thus if there exists a tuple of intervals (S_S^i, S_L^i) that f didn't miss, then it means that after this tuple (or inside him, depends where the ^{next} breaking point is) $f < \frac{1}{2}$ and thus there cannot be another tuple which avoids missing both intervals there.

Thus, between each two breaking points (or the first breaking point and the point 1), f can avoid missing at most 1 tuple (which means f can avoid missing only once a tuple, two intervals)

Therefore, from all the tuples we have, f can avoid missing at most j tuples.

Now, there are two cases:

- ① For any tuple until the first breaking point, f misses one of the intervals in that tuple:

In that case we've got that at most j tuples f avoided missing the intervals in them (from the 2^{L-2} tuples), and of course that if missed a tuple (which means that f wouldn't be able to avoid missing both intervals in the tuple) than f avoided missing one of them and missed the second.

In total we look at $2^{L-2} \cdot 2 = 2^{L-1}$ intervals, and we know that f avoided missing at most j tuples, thus

f avoided missing at most $\lceil \frac{1}{2} \cdot 2^{L-1} + \frac{1}{2} \cdot j \rceil$ intervals from the 2^{L-1} intervals. In addition we have the first interval $I_1^<$ which maybe f avoided missing it too, thus in total

f avoided missing at most

$$\lceil \frac{1}{2} \cdot 2^{L-1} + \frac{1}{2} \cdot j \rceil + 1 = \lceil \frac{1}{2} \cdot (2^{L-1} + 1) + \frac{1}{2} (j+1) \rceil$$

- ② There exists a tuple which f avoided missing both its intervals [the tuple is until the first breaking point]:

In that case we get that f avoided $j+1$ tuples at most but in that case f must miss the interval $I_1^<$ because we assume f is linear until the first breaking point and because it avoided a tuple it means that $f > \frac{1}{2}$ in the beginning of the tuple, and f has negative slope, thus oh $I_1^<$ $f > \frac{1}{2}$ necessarily.

Thus, from this explanation + the explanation of ① :

$$f \text{ avoided missing at most } \lceil \frac{1}{2} \cdot 2^{L-1} + \frac{1}{2} (j+1) \rceil$$

We've got that f avoided missing at most

$$\left\lceil \frac{1}{2} \cdot (2^{L-1} + 1) + \frac{1}{2} (j+1) \right\rceil \leq \left\lceil \frac{1}{2} \cdot (2^{L-1} + 1) + \frac{1}{2} (B+1) \right\rceil$$

$j \leq B$ intervals.

Thus, because the number of local intervals is $2^{L-1} + 1$
it means that f must miss at least:

$$\left\lfloor \frac{1}{2} \cdot (2^{L-1} + 1) - \frac{1}{2} (B+1) \right\rfloor$$

□

Part 2 - Question 3

We have to modify the universality and expressive efficiency analysis given in class so that they apply to leaky RELU activation $\sigma(z) = \max\{\alpha z, z\}$ for $\alpha \in (0, 1)$ instead of RELU.

Proof

Notice that in the regular proof (with RELU) we only used the proposition \otimes from question 1:

A shallow network of width $B \geq 2$ can realize any piecewise linear mapping with $\leq B$ pieces.

Thus, it's enough for us to prove it now with shallow network that uses leaky RELU and we'll conclude universality also for networks that use leaky RELU.

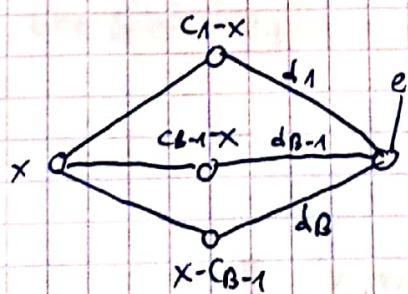
Proof of proposition \otimes

Let f be a piecewise linear with $\leq B$ pieces,

defined by $-\infty = c_0 < c_1 < c_2 < \dots < c_{B-1} < c_B = \infty$

$$f(x) = \alpha_{i+1}x + \beta_{i+1} \quad 0 \leq i \leq B-1$$

We'll build a network N which realizes f :



where we haven't defined d_0, \dots, d_{B-1}, e yet.

We've already proved that on each interval (c_i, c_{i+1}) , N is linear (in question 1), and we would like to see d_1, \dots, d_B, e s.t. $N(x) = f(x)$ on each interval. [The proof that N is linear on each interval was with RELU but it is exactly the same for leaky RELU - instead looking at the line $y=0$ for $x \leq 0$, the line will be $y = \alpha x$]

$$N(x) = \sum_{i=1}^{B-1} \text{LReLU}(c_i - x) \cdot d_i + \text{LReLU}(x - c_{B-1}) d_B + e$$

1. $x \geq c_{B-1}$: on that part, we will get $\text{LReLU}(x - c_{B-1}) = x - c_{B-1}$ and $\text{LReLU}(c_i - x) = \alpha(c_i - x) \quad \forall 1 \leq i \leq B-1$

thus

$$\begin{aligned} N(x) &= \sum_{i=1}^{B-1} \alpha(c_i - x) d_i + d_B(x - c_{B-1}) + e \\ &= (\sum_{i=1}^{B-1} -\alpha d_i + d_B) x + \sum_{i=1}^{B-1} \alpha c_i d_i - c_{B-1} + e \\ \Rightarrow \text{we would want } \sum_{i=1}^{B-1} -\alpha d_i + d_B &= \alpha_B \quad (\text{slope}) \end{aligned}$$

2. $c_{B-2} \leq x \leq c_{B-1}$:

With the same process we'll get:

$$\sum_{i=1}^{B-2} -\alpha d_i - d_{B-1} + \alpha d_B = \alpha_{B-1}$$

⋮

We'll get the B 'th equation

$$\sum_{i=1}^{B-1} -d_i + \alpha d_B = \alpha_1$$

Thus, we've got B linear equation.

We'll write our B equations as matrix:

$$\begin{bmatrix} -a & -a & \dots & -a & 1 \\ ; & -a & \dots & -1 & a \\ ; & ; & \dots & -1 & ; \\ ; & ; & \dots & -1 & ; \\ -a & \dots & \dots & -1 & a \end{bmatrix} \begin{pmatrix} d_1 \\ d_B \end{pmatrix} = \begin{pmatrix} \alpha_B \\ \alpha_{B-1} \\ ; \\ ; \\ \alpha_1 \end{pmatrix}$$

Let's sign the matrix as A.

Then if A is invertible, then there exists a solution

for $\begin{pmatrix} d_1 \\ d_B \end{pmatrix}$ and all the equations we want to exist

will exist. In that case, we will assign the left variable e to be such $N(c_{B-1}) = f(c_{B-1})$
and then by the explanation in question 1 we will
get that $N(x) = f(x)$ [the explanation will be the same
because LReLU is continuous function and thus N is
still continuous because it is composition of Continuous functions]

Now, it's easy to see that A is invertible because by
B-1 actions on the first B-1 rows ($r_i \leftarrow r_i - \alpha r_B$)

we'll get the matrix:

$$\begin{bmatrix} -a & -a & \dots & -a & 1 \\ ; & -a & \dots & -1 & a \\ ; & ; & \dots & -1 & ; \\ ; & ; & \dots & -1 & ; \\ -a & \dots & \dots & -1 & a \end{bmatrix} \rightarrow \begin{bmatrix} \textcircled{O} & & & 1-a^2 \\ & -1+a & & * \\ & & -1+a & \\ & & & -1 \end{bmatrix} = B$$

where B is invertible because $\det(B) = (-1)^{B+1} (1-a^2) \cdot (-1+a)^{B-2} \cdot (-1) \neq 0$
(the only permutation which doesn't have the value 0) $a \in (0, 1)$
 $\Rightarrow A$ is invertible $[\det(A) \neq 0]$

Now, let's focus on the proof of expressive efficiency.

In that part, we have used the proposition from question 1 (not only the first part of it, we have proved in the Universality with leaky RELU)

Note the second part of the proposition: $\star\star$

Any mapping realizable by shallow network with width $B \geq 2$ is piecewise linear with $\leq B+1$ pieces.

Then the proof was given to this part in question 1 with regular RELU holds exactly for networks with leaky RELU [instead $y=0$ when $a_i x + b_i \leq 0$ we'll get $a_i(a_i x + b_i)$]

Moreover

~~Especially~~, when we proved that $\exists \bar{h} \in H_B$ for $B = O(1)$

s.t. $\bar{h} \notin H_B$ unless $B = \exp(L)$

We built $= \bar{h} = g^{\circ k}$ when we prove that

$$g(x) = \text{RELU}[2x] - \text{RELU}[4x-2] + \text{RELU}[2x-2]$$

Now, from the recitation, we proved that:

$$\text{RELU}[wx+b] = \alpha \text{RELU}[wx+b] + \beta \text{RELU}[-wx-b]$$

Thus, we can write: for fixed $\alpha, \beta \in \mathbb{R}$

$$\alpha = \frac{1}{1-q^2} \quad \beta = \frac{q}{1-q^2}$$

$$g(x) = (\alpha \text{RELU}[2x] + \beta \text{RELU}[-2x])$$

$$- (\alpha \text{RELU}[4x-2] + \beta \text{RELU}[-4x+2])$$

$$+ (\alpha \text{RELU}[2x-2] + \beta \text{RELU}[-2x+2])$$

And therefore we can write g as sum of $\bar{B}=6 \in O(1)$

leaky RELU, thus by composing on it k times $g^{\circ k}$ can still be implemented in the network with $\bar{B}=6$ and (leaky RELU).

The other parts of the proof remain the same.

Question 4:

Prove that the shallow and deep networks are \mathcal{F} -universal in the sense of $d(\cdot, \cdot)$, where $\mathcal{F} \subseteq \mathbb{R}^{\mathbb{R}}$ is the set of Riemann integrable functions and $d: \mathcal{F} \times \mathcal{F} \rightarrow \mathbb{R}_{\geq 0}$ is defined by:

$$d(f_1, f_2) = \int_0^1 |f_1(x) - f_2(x)| dx$$

Proof

We would like to prove the following proposition: \otimes

$$\forall \epsilon > 0 \quad \forall f \in \mathcal{F} \quad \exists B \exists N \in H_B : d(f, N) \leq \epsilon$$

That will prove that the shallow network is \mathcal{F} -universal and since we've already proved that $H_B \subseteq \bar{H}_B$ where $\bar{B} = B$ that will prove that the deep network is \mathcal{F} -universal too.

Lemma: Each $f \in \mathcal{F}$ is bounded: $\exists M \in \mathbb{R} : |f(x)| \leq M \quad \forall x \in [0, 1]$

proof:

Recall that a function is supposed to be Riemann integrable if for all $\epsilon > 0$ there exists a partition s.t. $U - L < \epsilon$

$$[\text{when } U = \sum_{i=0}^n \sup_{t \in [x_i, x_{i+1}]} f(t) \cdot (x_{i+1} - x_i)]$$

$$[L = \sum_{i=0}^n \inf_{t \in [x_i, x_{i+1}]} f(t) \cdot (x_{i+1} - x_i)]$$

Suppose f becomes unbounded (say, unbounded above) near the point x_0 . Then the definition of U isn't defined well because in the interval where x_0 belongs to, the supremum is not defined well, contradiction. [the same for L]

Thus, there exists M s.t. $|f(x)| \leq M \quad \forall x \in [0, 1]$

Proof of proposition \otimes :

Let $\epsilon > 0$, $f \in \mathcal{F}$. By the lemma, there exists $M \in \mathbb{R}$ s.t. $|f(x)| \leq M$.

In addition, $f \in \mathcal{F}$ and therefore there exists a set of intervals $[x_0, x_1], [x_1, x_2], \dots, [x_n, x_{n+1}]$ s.t.: $x_0 = 0, x_{n+1} = 1$,
 $U = \sum_{i=0}^n \sup_{t \in [x_i, x_{i+1}]} f(t) \cdot (x_{i+1} - x_i)$, $L = \sum_{i=0}^n \inf_{t \in [x_i, x_{i+1}]} f(t) \cdot (x_{i+1} - x_i)$

and $U - L < \frac{\epsilon}{2}$.

Let's define $U: [0, 1] \rightarrow \mathbb{R}$ $U(x) = \sup_{\substack{y \in [x_i, x_{i+1}] \\ i \text{ s.t. } x \in [x_i, x_{i+1}]}} f(y)$
 $L: [0, 1] \rightarrow \mathbb{R}$ $L(x) = \inf_{\substack{y \in [x_i, x_{i+1}] \\ i \text{ s.t. } x \in [x_i, x_{i+1}]}} f(y)$

Then - by definition $L(x) \leq f(x) \leq U(x) \quad \forall x \in [0, 1]$

and thus:

$$\begin{aligned} \int_0^1 |U(x) - f(x)| dx &\leq \int_0^1 |U(x) - L(x)| dx \leq \int_0^1 |U(x) - L(x)| dx \\ &= \int_0^1 |U(x)| dx - \int_0^1 |L(x)| dx = U - L \leq \frac{\epsilon}{2} \end{aligned}$$

Thus, we've got $\int_0^1 |U(x) - f(x)| dx \leq \frac{\epsilon}{2}$

Now we'll define $N(x) \in \mathcal{H}_B$ for $B = 2n-1$ (which means a function that we can represent with shallow network with $B = 2n-1$ width), we will prove that

$$\int_0^1 |N(x) - U(x)| dx \leq \frac{\epsilon}{2}.$$

Say it is true, then we'll get:

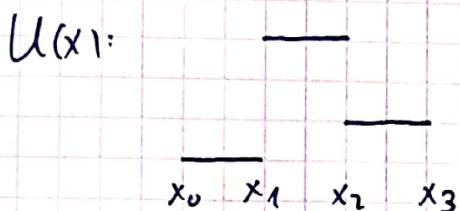
$$\begin{aligned} d(N, f) &= \int_0^1 |N(x) - f(x)| dx = \int_0^1 |N(x) - U(x) + U(x) - f(x)| dx \leq \int_0^1 |N(x) - U(x)| dx + \int_0^1 |U(x) - f(x)| dx \\ &\leq \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon \end{aligned}$$

as we wished.

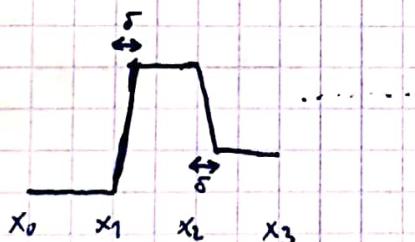
Defining the function $N(x)$:

We'll define a function which is piecewise linear with $B=2n-1$ pieces, and from question 1 we already know that a shallow network of width B can realize it.

We'll draw the function $U(x)$ which was defined earlier:



$U(x)$ is composed from set of intervals s.t. on each interval U is const. We'll define $N(x)$ by the following [N is almost like U]:



$$\text{where } \delta = \frac{\varepsilon}{2n \cdot M}$$

[From each interval $[x_i, x_{i+1}]$ $i \geq 1$, we take δ from the beginning and make a line to the end of the previous interval]

Of course $N(x)$ is refined well and has $n + (n-1) = 2n-1$ pieces which compose the function as we wished
So a network with width $B = 2n-1$ can implement it.

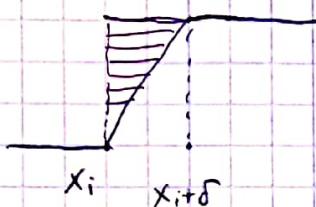
Now, we'll prove that $\int_0^1 |N(x) - U(x)| \leq \frac{\varepsilon}{2}$

Notice that N and U are only differ on each beginning of size δ for each interval and thus:

$$\int_0^1 |N(x) - U(x)| dx = \sum_{i=1}^n \int_{x_i}^{x_i+\delta} |N(x) - U(x)| dx$$

Now for each $i \geq 1$, $\int_{x_i}^{x_i+\delta} |N(x) - U(x)| dx$ is actually

the area of the following triangle:



As we know, $|f(x)| \leq M$ and therefore for every two different points x, y $|f(x) - f(y)| \leq 2M$.

Now we've got:

$$\begin{aligned} \int_{x_i}^{x_i+\delta} |N(x) - U(x)| dx &= \text{Area of triangle } \triangle = \frac{\delta \cdot |f\left(\frac{x_i+x_{i+1}}{2}\right) - f\left(\frac{x_i+x_{i-1}}{2}\right)|}{2} \\ &\leq \frac{\delta \cdot 2M}{2} = \frac{\frac{\epsilon}{2mn} \cdot 2M}{2} = \frac{\epsilon}{2m} \\ &\text{We chose } \delta = \frac{\epsilon}{2mn} \end{aligned}$$

Therefore,

$$\int_0^1 |N(x) - U(x)| dx = \sum_{i=1}^n \int_{x_i}^{x_i+\delta} |N(x) - U(x)| dx \leq \sum_{i=1}^n \frac{\epsilon}{2n} = n \cdot \frac{\epsilon}{2n} = \frac{\epsilon}{2}$$

as we wished.



Foundations of Deep Learning – Homework Assignment #2

Adi Alm & Tomer Epshtain

Part 3: (1)

Prove (Lemma): (“matricization of outer product = Kronecker product of matricizations”)

Let T and \bar{T} be tensors of order n and \bar{n} respectively. Let $I \subset [n + \bar{n}]$ and denote by $I - n$ the set obtained by subtracting n from each element in I . Then:

$$[|T \otimes \bar{T}|]_I = [|T|]_{I \cap [n]} \odot [|T|]_{I - n \cap [\bar{n}]}$$

Proof:

Let $T \in \mathbb{R}^{m_1 \times \dots \times m_n}$ be a tensor of order n and $\bar{T} \in \mathbb{R}^{\bar{m}_1 \times \dots \times \bar{m}_{\bar{n}}}$ be a tensor of order \bar{n} .

Denote $I = \{i_1, \dots, i_{|I|}\} \subset [n + \bar{n}]$ where $1 \leq i_1 < \dots < i_{|I|} \leq n + \bar{n}$.

Are the two sides of the target equation equal?

1. Sanity check – they have equal dimensions... Both sides of the equation are matrices of dimensions:

$$\left(\prod_{i \in I \cap [n]} m_i \right) \cdot \left(\prod_{i \in I - n \cap [\bar{n}]} \bar{m}_i \right) \times \left(\prod_{i \in I \cap [\bar{n}]} m_i \right) \cdot \left(\prod_{i \in I - n \cap [\bar{n}]} \bar{m}_i \right)$$

2. We'll use linearity to make our proof less complicated.

- a. Linearity of tensor outer product:

$$\begin{aligned} (\alpha A + \beta B) \otimes C &= \alpha(A \otimes C) + \beta(B \otimes C) \\ C \otimes (\alpha A + \beta B) &= \alpha(C \otimes A) + \beta(C \otimes B) \end{aligned}$$

For any scalars $\alpha, \beta \in \mathbb{R}$ and any tensors A, B, C such that A and B have equal dimensions.

- b. Linearity of Matricization:

$$[|\alpha A + \beta B|]_I = \alpha [|A|]_I + \beta [|B|]_I$$

For any scalars $\alpha, \beta \in \mathbb{R}$ and any tensors A, B such that A and B have equal dimensions.

- c. Linearity of Kronecker Product:

$$\begin{aligned} (\alpha A + \beta B) \odot C &= \alpha(A \odot C) + \beta(B \odot C) \\ C \odot (\alpha A + \beta B) &= \alpha(C \odot A) + \beta(C \odot B) \end{aligned}$$

For any scalars $\alpha, \beta \in \mathbb{R}$ and any tensors A, B, C such that A and B have equal dimensions.

Any tensor $T \in \mathbb{R}^{m_1 \times \dots \times m_n}$ can be represented as a linear combinations of tensors $E_i \in \mathbb{R}^{m_1 \times \dots \times m_n}$ with scalars $\alpha_i \in \mathbb{R}$ $i \in [m_1 \cdot \dots \cdot m_n]$, such that E_i is a tensor with ‘1’ in a single entry and ‘0’ in all other entries.

$$T = \sum_{i=1}^{m_1 \cdot \dots \cdot m_n} \alpha_i E_i$$

3. Claim:

Let $E_i \in \mathbb{R}^{m_1 \times \dots \times m_n}$ be some tensor with ‘1’ in a single entry and ‘0’ in all other entries,

Let $\bar{E}_i \in \mathbb{R}^{\bar{m}_1 \times \dots \times \bar{m}_{\bar{n}}}$ be some tensor with ‘1’ in a single entry and ‘0’ in all other entries.

Let $I \subset [n + \bar{n}]$ and denote by $I - n$ the set obtained by subtracting n from each element in I .

Then:

$$[|E_i \otimes \bar{E}_i|]_I = [|E_i|]_{I \cap [n]} \odot [|E_i|]_{I - n \cap [\bar{n}]}$$

4. Assume we proved Claim (3).

We'll first show this completes our proof:

Let $T \in \mathbb{R}^{m_1 \times \dots \times m_n}$ be a tensor of order n and $\bar{T} \in \mathbb{R}^{\bar{m}_1 \times \dots \times \bar{m}_{\bar{n}}}$ be a tensor of order \bar{n} .

By (2) we have $T = \sum_{i=1}^{m_1 \dots m_n} \alpha_i E_i$ and $\bar{T} = \sum_{i=1}^{\bar{m}_1 \dots \bar{m}_{\bar{n}}} \bar{\alpha}_i \bar{E}_i$

$$\begin{aligned}
[|T \otimes \bar{T}|]_I &= \left[\left| \sum_{i=1}^{m_1 \dots m_n} \alpha_i E_i \otimes \sum_{i=1}^{\bar{m}_1 \dots \bar{m}_{\bar{n}}} \bar{\alpha}_i \bar{E}_i \right| \right]_I = \\
&\quad \{ \text{Linearity of tensor outer product} \} \\
&= \left[\left| \sum_{i=1}^{m_1 \dots m_n} \sum_{\bar{i}=1}^{\bar{m}_1 \dots \bar{m}_{\bar{n}}} \alpha_i \bar{\alpha}_{\bar{i}} \cdot (E_i \otimes \bar{E}_{\bar{i}}) \right| \right]_I = \\
&\quad \{ \text{Linearity of matricization} \} \\
&= \sum_{i=1}^{m_1 \dots m_n} \sum_{\bar{i}=1}^{\bar{m}_1 \dots \bar{m}_{\bar{n}}} \alpha_i \bar{\alpha}_{\bar{i}} \cdot [|E_i \otimes \bar{E}_{\bar{i}}|]_I = \\
&\quad \{ (3) \} \\
&= \sum_{i=1}^{m_1 \dots m_n} \sum_{\bar{i}=1}^{\bar{m}_1 \dots \bar{m}_{\bar{n}}} \alpha_i \bar{\alpha}_{\bar{i}} \cdot [|E_i|]_{I \cap [n]} \odot [|E_{\bar{i}}|]_{I - n \cap [\bar{n}]} = \\
&\quad \{ \text{Linearity of kronecker product} \} \\
&= \sum_{i=1}^{m_1 \dots m_n} \alpha_i [|E_i|]_{I \cap [n]} \odot \sum_{\bar{i}=1}^{\bar{m}_1 \dots \bar{m}_{\bar{n}}} \bar{\alpha}_{\bar{i}} [|E_{\bar{i}}|]_{I - n \cap [\bar{n}]} = \\
&\quad \{ \text{Linearity of matricization} \} \\
&= \left[\left| \sum_{i=1}^{m_1 \dots m_n} \alpha_i E_i \right| \right]_{I \cap [n]} \odot \left[\left| \sum_{i=1}^{\bar{m}_1 \dots \bar{m}_{\bar{n}}} \bar{\alpha}_i \bar{E}_i \right| \right]_{I - n \cap [\bar{n}]} = [|T|]_{I \cap [n]} \odot [|\bar{T}|]_{I - n \cap [\bar{n}]}
\end{aligned}$$

Which gives us:

$$[|T \otimes \bar{T}|]_I = [|T|]_{I \cap [n]} \odot [|\bar{T}|]_{I - n \cap [\bar{n}]}$$

As desired.

It remains to prove claim 3:

Claim 3:

Let $E_i \in \mathbb{R}^{m_1 \times \dots \times m_n}$ be some tensor with '1' in a single entry and '0' in all other entries,

Let $\bar{E}_{\bar{i}} \in \mathbb{R}^{\bar{m}_1 \times \dots \times \bar{m}_{\bar{n}}}$ be some tensor with '1' in a single entry and '0' in all other entries.

Let $I \subset [n + \bar{n}]$ and denote by $I - n$ the set obtained by subtracting n from each element in I .

Then:

$$[|E_i \otimes \bar{E}_{\bar{i}}|]_I = [|E_i|]_{I \cap [n]} \odot [|\bar{E}_{\bar{i}}|]_{I - n \cap [\bar{n}]}$$

Proof of claim 3:

Let $E_i \in \mathbb{R}^{m_1 \times \dots \times m_n}$ be some tensor with '1' in a single entry and '0' in all other entries,
i.e. $(E_i)_{d_1, \dots, d_n} = 1$ for some $d_k \in [m_k]$, $k \in [n]$. And E_i is 0 for all other entries.

Let $\bar{E}_{\bar{i}} \in \mathbb{R}^{\bar{m}_1 \times \dots \times \bar{m}_{\bar{n}}}$ be some tensor with '1' in a single entry and '0' in all other entries.

i.e. $(\bar{E}_{\bar{i}})_{\bar{d}_1, \dots, \bar{d}_{\bar{n}}} = 1$ for some $\bar{d}_k \in [\bar{m}_k]$, $k \in [\bar{n}]$. And $\bar{E}_{\bar{i}}$ is 0 for all other entries.

Start by viewing LHS: $[|E_i \otimes \bar{E}_{\bar{i}}|]_I$

- **Tensor outer product:**

$E_i \otimes \bar{E}_{\bar{i}} \in \mathbb{R}^{m_1 \times \dots \times m_n \times \bar{m}_1 \times \dots \times \bar{m}_{\bar{n}}}$ is a tensor defined by:

$$(E_i \otimes \bar{E}_{\bar{i}})_{l_1, \dots, l_n, \bar{l}_1, \dots, \bar{l}_{\bar{n}}} = E_{i, l_1, \dots, l_n} \cdot \bar{E}_{\bar{i}, \bar{l}_1, \dots, \bar{l}_{\bar{n}}}$$

So this produces a tensor with '1' in a single entry. The entry with index:

$$(d_1, \dots, d_n, \bar{d}_1, \dots, \bar{d}_{\bar{n}})$$

And '0' in every other entry.

- **Matricization:**

$[(E_i \otimes \bar{E}_i)]_I \in \mathbb{R}^{(\prod_{i \in I \cap [n]} m_i) \cdot (\prod_{i \in I - n \cap [\bar{n}]} \bar{m}_i) \times (\prod_{i \in I \cap [n]} m_i) \cdot (\prod_{i \in I \cap [\bar{n}]} \bar{m}_i)}$ is a matrix

Defined such that:

$$([(E_i \otimes \bar{E}_i)]_I)_{u,v} = (E_i \otimes \bar{E}_i)_{d_1, \dots, d_n, \bar{d}_1, \dots, \bar{d}_{\bar{n}}}$$

for every $d_k \in [m_k], k \in [n]$ and $\bar{d}_k \in [\bar{m}_{\bar{k}}], \bar{k} \in [\bar{n}]$

Where:

$$u = 1 + \sum_{t=1}^{|I \cap [n]|} (d_{i_t} - 1) \prod_{t'=t+1}^{|I \cap [n]|} m_{i_{t'}} + \sum_{t=1}^{|I - n \cap [\bar{n}]|} (\bar{d}_{i_t} - 1) \prod_{t'=t+1}^{|I - n \cap [\bar{n}]|} \bar{m}_{i_{t'}} \prod_{l=1}^{|I \cap n|} m_{i_l}$$

$$v = 1 + \sum_{t=1}^{|I \cap [n]|} (d_{i_t} - 1) \prod_{t'=t+1}^{|I \cap [n]|} m_{i_{t'}} + \sum_{t=1}^{|I - n \cap [\bar{n}]|} (\bar{d}_{i_t} - 1) \prod_{t'=t+1}^{|I - n \cap [\bar{n}]|} \bar{m}_{i_{t'}} \prod_{l=1}^{|I \cap n|} m_{i_l}$$

i.e. We have a matrix with a single '1' entry in index (u, v) corresponding to index $(d_1, \dots, d_n, \bar{d}_1, \dots, \bar{d}_{\bar{n}})$ in $E_i \otimes \bar{E}_i$.

Simply putting it, if we denote the matrix in LHS as

$$A \in \mathbb{R}^{(\prod_{i \in I \cap [n]} m_i) \cdot (\prod_{i \in I - n \cap [\bar{n}]} \bar{m}_i) \times (\prod_{i \in I \cap [n]} m_i) \cdot (\prod_{i \in I \cap [\bar{n}]} \bar{m}_i)}$$

We have a matrix with a single '1' entry at (u_A, v_A) and all other entries are '0' where

$$u_A = 1 + \sum_{t=1}^{|I \cap [n]|} (d_{i_t} - 1) \prod_{t'=t+1}^{|I \cap [n]|} m_{i_{t'}} + \sum_{t=1}^{|I - n \cap [\bar{n}]|} (\bar{d}_{i_t} - 1) \prod_{t'=t+1}^{|I - n \cap [\bar{n}]|} \bar{m}_{i_{t'}} \prod_{l=1}^{|I \cap n|} m_{i_l}$$

$$v_A = 1 + \sum_{t=1}^{|I \cap [n]|} (d_{i_t} - 1) \prod_{t'=t+1}^{|I \cap [n]|} m_{i_{t'}} + \sum_{t=1}^{|I - n \cap [\bar{n}]|} (\bar{d}_{i_t} - 1) \prod_{t'=t+1}^{|I - n \cap [\bar{n}]|} \bar{m}_{i_{t'}} \prod_{l=1}^{|I \cap n|} m_{i_l}$$

Let's move on to viewing RHS: $[(E_i)]_{I \cap [n]} \odot [(\bar{E}_i)]_{I - n \cap [\bar{n}]}$

- **Matricizations:**

- We have an order n tensor $E_i \in \mathbb{R}^{m_1 \times \dots \times m_n}$, so its matricization w.r.t $I \cap [n]$ is:

$$[(E_i)]_{I \cap [n]} \in \mathbb{R}^{(\prod_{i \in I \cap [n]} m_i) \times (\prod_{i \in I \cap [n]} m_i)}$$

Is defined such that:

$$([(E_i)]_{I \cap [n]})_{u,v} = E_{i, d_1, \dots, d_n}$$

for every $d_k \in [m_k], k \in [n]$.

Where:

$$u = 1 + \sum_{t=1}^{|I \cap [n]|} (d_{i_t} - 1) \prod_{t'=t+1}^{|I \cap [n]|} m_{i_{t'}} \\ v = 1 + \sum_{t=1}^{|I \cap [n]|} (d_{i_t} - 1) \prod_{t'=t+1}^{|I \cap [n]|} m_{i_{t'}}$$

Simply putting it, this yields a matrix $B \in \mathbb{R}^{(\prod_{i \in I \cap [n]} m_i) \times (\prod_{i \in I \cap [n]} m_i)}$ with a single '1' entry at (u_B, v_B) and all other entries are '0'. Where:

$$u_B = 1 + \sum_{t=1}^{|I \cap [n]|} (d_{i_t} - 1) \prod_{\substack{t' = t+1 \\ |I^C \cap [n]|}}^{I \cap [n]} m_{i_{t'}} \\ v_B = 1 + \sum_{t=1}^{|I^C \cap [n]|} (d_{i_t} - 1) \prod_{\substack{t' = t+1 \\ |I \cap [n]|}}^{I^C \cap [n]} m_{i_{t'}}$$

- And an order \bar{n} tensor $\bar{E}_{\bar{i}} \in \mathbb{R}^{\bar{m}_1 \times \dots \times \bar{m}_{\bar{n}}}$, so its matricization w.r.t $I - n \cap [\bar{n}]$ is:

$$[\bar{E}_{\bar{i}}]_{I-n \cap [\bar{n}]} \in \mathbb{R}^{(\prod_{i \in I-n \cap [\bar{n}]} \bar{m}_i) \times (\prod_{i \in I^C-n \cap [\bar{n}]} \bar{m}_i)}$$

Is defined such that:

$$([\bar{E}_{\bar{i}}]_{I-n \cap [\bar{n}]})_{u,v} = \bar{E}_{\bar{i} \bar{d}_1, \dots, \bar{d}_{\bar{n}}}$$

for every $d_k \in [\bar{m}_k]$, $k \in [\bar{n}]$.

Where:

$$u = 1 + \sum_{\substack{t=1 \\ |I^C-n \cap [\bar{n}]|}}^{|I-n \cap [\bar{n}]|} (\bar{d}_{i_t} - 1) \prod_{\substack{t' = t+1 \\ |I^C-n \cap [\bar{n}]|}}^{I-n \cap [\bar{n}]} \bar{m}_{i_{t'}} \\ v = 1 + \sum_{\substack{t=1 \\ |I^C-n \cap [\bar{n}]|}}^{|I-n \cap [\bar{n}]|} (\bar{d}_{i_t} - 1) \prod_{\substack{t' = t+1 \\ |I \cap [n]|}}^{I \cap [n]} \bar{m}_{i_{t'}}$$

Simply putting it, this yields a matrix $C \in \mathbb{R}^{(\prod_{i \in I-n \cap [\bar{n}]} \bar{m}_i) \times (\prod_{i \in I^C-n \cap [\bar{n}]} \bar{m}_i)}$ with a single '1' entry at (u_C, v_C) and all other entries are '0'. Where:

$$u_C = 1 + \sum_{\substack{t=1 \\ |I^C-n \cap [\bar{n}]|}}^{|I-n \cap [\bar{n}]|} (\bar{d}_{i_t} - 1) \prod_{\substack{t' = t+1 \\ |I^C-n \cap [\bar{n}]|}}^{I-n \cap [\bar{n}]} \bar{m}_{i_{t'}} \\ v_C = 1 + \sum_{\substack{t=1 \\ |I^C-n \cap [\bar{n}]|}}^{|I-n \cap [\bar{n}]|} (\bar{d}_{i_t} - 1) \prod_{\substack{t' = t+1 \\ |I \cap [n]|}}^{I \cap [n]} \bar{m}_{i_{t'}}$$

Kronecker Product:

We now have two matrices:

$$B \in \mathbb{R}^{(\prod_{i \in I \cap [n]} m_i) \times (\prod_{i \in I^C \cap [n]} m_i)}, C \in \mathbb{R}^{(\prod_{i \in I-n \cap [\bar{n}]} \bar{m}_i) \times (\prod_{i \in I^C-n \cap [\bar{n}]} \bar{m}_i)}$$

Each matrix with a single '1' entry and all other entries being '0'.

B 's '1' entry: (u_B, v_B) , C 's '1' entry: (u_C, v_C) .

B & C 's Kronecker product yields a matrix $D \in \mathbb{R}^{(\prod_{i \in I \cap [n]} m_i) \cdot (\prod_{i \in I-n \cap [\bar{n}]} \bar{m}_i) \times (\prod_{i \in I^C \cap [n]} m_i) \cdot (\prod_{i \in I^C-n \cap [\bar{n}]} \bar{m}_i)}$ with a single '1' entry at index (u_D, v_D) where:

$$u_D = \left(\prod_{i \in I-n \cap [\bar{n}]} \bar{m}_i \right) \cdot (u_B - 1) + u_C \\ v_D = \left(\prod_{i \in I^C-n \cap [\bar{n}]} \bar{m}_i \right) \cdot (v_B - 1) + v_C$$

Bringing it ALL together...

- On the LHS we got a matrix A :

$$A \in \mathbb{R}^{(\prod_{i \in I \cap [n]} m_i) \cdot (\prod_{i \in I-n \cap [\bar{n}]} \bar{m}_i) \times (\prod_{i \in I^C \cap [n]} m_i) \cdot (\prod_{i \in I^C-n \cap [\bar{n}]} \bar{m}_i)}$$

with a single '1' entry at (u_A, v_A) and all other entries are '0' where

$$u_A = 1 + \sum_{t=1}^{|I \cap [n]|} (d_{i_t} - 1) \prod_{t'=t+1}^{|I \cap [n]|} m_{i_{t'}} + \sum_{t=1}^{|I-n \cap [\bar{n}]|} (\bar{d}_{i_t} - 1) \prod_{t'=t+1}^{|I-n \cap [\bar{n}]|} \bar{m}_{i_{t'}} \prod_{l=1}^{|I \cap n|} m_{i_l}$$

$$v_A = 1 + \sum_{t=1}^{|I^C \cap [n]|} (d_{i_t} - 1) \prod_{t'=t+1}^{|I^C \cap [n]|} m_{i_{t'}} + \sum_{t=1}^{|I^C - n \cap [\bar{n}]|} (\overline{d}_{i_t} - 1) \prod_{t'=t+1}^{|I^C - n \cap [\bar{n}]|} \overline{m}_{i_{t'}} \prod_{l=1}^{|I^C \cap n|} m_{i_l}$$

- On the RHS we got a matrix D :

$$D \in \mathbb{R}^{(\prod_{i \in I \cap [n]} m_i) \cdot (\prod_{i \in I - n \cap [\bar{n}]} \overline{m}_i) \times (\prod_{i \in I^C \cap [n]} m_i) \cdot (\prod_{i \in I^C - n \cap [\bar{n}]} \overline{m}_i)}$$

with a single '1' entry at (u_D, v_D) and all other entries are '0' where

$$u_D = \left(\prod_{i \in I - n \cap [\bar{n}]} \overline{m}_i \right) \cdot (u_B - 1) + u_C$$

$$v_D = \left(\prod_{i \in I^C - n \cap [\bar{n}]} \overline{m}_i \right) \cdot (v_B - 1) + v_C$$

And:

$$u_B = 1 + \sum_{t=1}^{|I \cap [n]|} (d_{i_t} - 1) \prod_{t'=t+1}^{|I \cap [n]|} m_{i_{t'}}$$

$$v_B = 1 + \sum_{t=1}^{|I^C \cap [n]|} (d_{i_t} - 1) \prod_{t'=t+1}^{|I^C \cap [n]|} m_{i_{t'}}$$

$$u_C = 1 + \sum_{t=1}^{|I - n \cap [\bar{n}]|} (\overline{d}_{i_t} - 1) \prod_{t'=t+1}^{|I - n \cap [\bar{n}]|} \overline{m}_{i_{t'}}$$

$$v_C = 1 + \sum_{t=1}^{|I^C - n \cap [\bar{n}]|} (\overline{d}_{i_t} - 1) \prod_{t'=t+1}^{|I^C - n \cap [\bar{n}]|} \overline{m}_{i_{t'}}$$

Plotting u_B, v_B, u_C , and v_C into u_D and v_D yields equality:

$$u_A = u_D \text{ and } v_A = v_D$$

i.e matrix A (=LHS) is equal to matrix D (=RHS)

■

Foundations of Deep Learning – Homework Assignment #2

Adi Almog & Tomer Epshteyn

Part 3: (2)

Problem:

Modify the expressive efficiency analysis given in class to the case in which the deep network has $\log_4(N)$ hidden layers (Assume N is a power of 4), with pooling windows of size 4 (In class we treated $\log_2(N)$ hidden layers with pooling windows of size 2).

Solution:

Assume N is a power of 4.

Denote by $\bar{H}_{\bar{B}}$, $\bar{B} \in \mathbb{N}$ the hypothesis space corresponding to a deep network with pooling windows of size 4 and with depth $\log_4(N)$.

Similarly to notation in class, we'll denote its learnable weights by:

$$\left\{ \left(a^{l,1,\gamma}, a^{l,2,\gamma}, \dots, a^{l,\frac{N}{4^l},\gamma} \right) \right\}_{\gamma=1}^{r_l} \subset (\mathbb{R}^{r_l})^{N/4^l}, l = 0, 1, \dots, L$$

Where $L = \log_4(N)$.

Let's now state the expressive efficiency statement for this hypothesis class w.r.t H_B , the hypothesis space corresponding to shallow nets:

1. $\forall B, \exists \bar{B} \in \mathcal{O}(B)$ such that $H_B \subseteq \bar{H}_{\bar{B}}$
2. $\exists \bar{h} \in \bar{H}_{\bar{B}}$ for $\bar{B} \in \mathcal{O}(M)$ such that $\bar{h} \notin H_B$ unless $B \in \mathcal{O}(\exp(N))$.

Proof 1:

Nearly identical to the proof in class for deep nets with windows of size 2 and depth $\log_2(N)$.

Let $B \in \mathbb{N}$. Choose $\bar{B} = B$. Let $h \in H_B$. I.e. h is a function realizable by a shallow network of width B . We will prove it is realizable by a deep network with window size 4, depth $\log_4(N)$ and width $\bar{B} = B$.

Denote by $\{(a_S^{z,1}, a_S^{z,2}, \dots, a_S^{z,N})\}_{z=1}^B \subset (\mathbb{R}^M)^N$ and $a_S \in \mathbb{R}^B$ the learnable weights of the shallow net realizing h .
(S – Denoting shallow net weights)

We will now define the weights for the deep net so that it realizes h :

With $r_0 = r_1 = \dots = r_{L-1} = B$:

For $l = 1, \dots, L$: We define the weights for the conv layers $\left\{ \left(a^{l,1,\gamma}, a^{l,2,\gamma}, \dots, a^{l,\frac{N}{4^l},\gamma} \right) \right\}_{\gamma=1}^{r_l}$ by:

$$a^{l,1,\gamma} = e^\gamma$$

Where $e^\gamma \in \mathbb{R}^{r_l}$ is a vector with a single entry "1" at position γ and "0" in all other entries.

Under this definition, the 1x1 conv layers act as "passthru" layers, and the $4 \rightarrow 1$ pooling layers together constitute a global pooling (multiplying 4 entries at a time).

Altering the weights in hidden layer 0: $\{(a^{0,1,\gamma}, a^{0,2,\gamma}, \dots, a^{0,N,\gamma})\}_{\gamma=1}^{r_0} \subset (\mathbb{R}^{r_0})^N$ and output weights $a^L \in \mathbb{R}^B$ we can realize h using a deep network.

Define:

$$a^{0,k,\gamma} := a_S^{\gamma,k}$$

$$a^L := a_S$$

Using above weights of layer 0 and output layer L , together with “passthru” conv layers $1, \dots, L-1$ out deep-net realizes h .

Proof 2:

For this part of the proof we will need to use a decomposition of the network’s behavior, similar to the HT decomposition we saw in class for deep nets with windows of size 2 and depth $\log_2(N)$.

In class we saw the equivalence of shallow nets to tensors defined by the *CP* decomposition, and the equivalence of deep nets (with windows of size 2 and depth $\log_2(N)$) to tensors defined by the *HT* decomposition:

$$\text{shallow-nets} \Leftrightarrow \text{CP decomp}, \quad \text{deep-nets} \Leftrightarrow \text{HT decomp}$$

In our case, deep nets with windows of size 4 and depth $\log_4(N)$ we obtain the following decomposition which we will call the STAR decomposition denoted by (*):

$$\Phi^{1,j,\gamma} = \sum_{\alpha=1}^{r_0} a_\alpha^{1,j,\gamma} \cdot (a^{0,4j-3,\alpha} \otimes a^{0,4j-2,\alpha} \otimes a^{0,4j-1,\alpha} \otimes a^{0,4j,\alpha})$$

...

$$\Phi^{l,j,\gamma} = \sum_{\alpha=1}^{r_{l-1}} a_\alpha^{l,j,\gamma} \cdot (\Phi^{l-1,4j-3,\alpha} \otimes \Phi^{l-1,4j-2,\alpha} \otimes \Phi^{l-1,4j-1,\alpha} \otimes \Phi^{l-1,4j,\alpha})$$

...

$$\Phi^{L-1,j,\gamma} = \sum_{\alpha=1}^{r_{L-2}} a_\alpha^{L-1,j,\gamma} \cdot (\Phi^{L-2,4j-3,\alpha} \otimes \Phi^{L-2,4j-2,\alpha} \otimes \Phi^{L-2,4j-1,\alpha} \otimes \Phi^{L-2,4j,\alpha})$$

$$A^{(*)} = \sum_{\alpha=1}^{r_{L-1}} a_\alpha^L \cdot (\Phi^{L-1,1,\alpha} \otimes \Phi^{L-1,2,\alpha} \otimes \Phi^{L-1,3,\alpha} \otimes \Phi^{L-1,4,\alpha})$$

As we saw in class, by viewing the matricization of the CP decomposition (\leftrightarrow shallow-nets) we saw that for any choice of learnable parameters, any shallow net is equivalent to a tensor A such that $\text{rank}[|A^{\text{CP}}|] \leq Z$ as a sum of Z rank 1 matrices.

We’ll consider a similar analysis for STAR decomposition (Denoted by (*)) corresponding to deep nets with windows of size 4 and depth $\log_4(N)$.

We’ll define the weights of the deep net as follows:

- Layer 0:

Let $r_0 \geq M$ be the width of hidden layer 0. Denote “Layer 0” weights by:

$$\{(a^{0,1,\gamma}, a^{0,2,\gamma}, \dots, a^{0,N,\gamma})\}_{\gamma=1}^{r_0} \subset (\mathbb{R}^{r_0})^N$$

And define:

$$a^{0,j,\gamma} := e^\gamma \in \mathbb{R}^{r_0}$$

for all $\gamma \in [M]$ where e^γ is “1” at entry γ and “0” elsewhere.

And:

$$a^{0,j,\gamma} := \bar{0} \in \mathbb{R}^{r_0}$$

for all $\gamma > M$, where $\bar{0}$ is the all zeros vector.

- Layer 1:

Let r_1 be the width of hidden layer 1. Denote “Layer 1” weights by:

$$\{(a^{1,1,\gamma}, a^{1,2,\gamma}, \dots, a^{1,N/4,\gamma})\}_{\gamma=1}^{r_1} \subset (\mathbb{R}^{r_1})^{N/4}$$

$$a^{1,j,1} = \bar{1} \in \mathbb{R}^{r_1}, a^{1,j,\gamma} = \bar{0} \in \mathbb{R}^{r_1}$$

for all $\gamma > 1$ where $\bar{1}$ is the all ones vector and $\bar{0}$ is the all zeros vector.

Here we achieve:

$[\Phi^{1,j,\gamma}]$ is the identity matrix for $\gamma = 1$ and the zero matrix for $\gamma > 1$.

- Layers 2, ..., (L-1):

Let r_2, \dots, r_{L-1} be the widths of layers 2, ..., $(L - 1)$. Denote their weights by:

$$\{(a^{l,1,\gamma}, a^{l,2,\gamma}, \dots, a^{0,N/4^l,\gamma})\}_{\gamma=1}^{r_l} \subset (\mathbb{R}^{r_l})^{\frac{N}{4^l}}, l = 2, \dots, (L - 1)$$

And define:

$$a^{l,j,\gamma} = \begin{cases} e^1, & \gamma = 1 \\ \bar{0}, & \gamma > 1 \end{cases}$$

For $l = 2, \dots, (L - 1), j = 1, 2, \dots, N/4^l, \gamma = 1, 2, \dots, r_l$.

- Output weights:

Define $a^L = e^1$.

Under our above assignments we achieve with the canonical decomposition:

$$\begin{aligned} [|A^{(*)}|] &= \left[\left| \sum_{\alpha=1}^{r_{L-1}} a_{\alpha}^L \cdot (\Phi^{L-1,1,\alpha} \otimes \Phi^{L-1,2,\alpha} \otimes \Phi^{L-1,3,\alpha} \otimes \Phi^{L-1,4,\alpha}) \right| \right] = \\ &\quad \text{\{Linearity of matricization\}} \\ &= \sum_{\alpha=1}^{r_{L-1}} a_{\alpha}^L \cdot [|(\Phi^{L-1,1,\alpha} \otimes \Phi^{L-1,2,\alpha} \otimes \Phi^{L-1,3,\alpha} \otimes \Phi^{L-1,4,\alpha})|] = \\ &\quad \text{\{"matricization of outer product equals Kronecker product of matricizations" property\}} \\ &= \sum_{\alpha=1}^{r_{L-1}} a_{\alpha}^L \cdot [|(\Phi^{L-1,1,\alpha})| \odot [|(\Phi^{L-1,2,\alpha})| \odot [|(\Phi^{L-1,3,\alpha})| \odot [|(\Phi^{L-1,4,\alpha})|]]] = \\ &\quad \text{\dots Keep going down the network with (*) decomposition \dots Calculations identical to those held in class \dots} \\ &\dots = Id_{M^2} \odot Id_{M^2} \odot \dots \odot Id_{M^2} = Id_{(M^2)^{\frac{N}{4}}} \end{aligned}$$

Id_{M^2} appears $N/4$ times, where Id_n is a $n \times n$ identity matrix.

In particular we showed a realization of a tensor by $(*)$ decomposition which corresponds to a deep net with windows of size 4 and depth $\log_4(N)$. The matricization via the canonical matricization

produces a matrix with $rank([|A^{(*)}|]) = (M^2)^{\frac{N}{4}} = M^{\frac{N}{2}}$.

We'll summarize: We saw in class that any realization of a shallow net produces a tensor denoted by A such that it holds the property: $\text{rank}(|A|) \leq Z$ where Z is shallow net's width.

We now saw an implementation of a deep net with windows of size 4 and depth $\log_4(N)$ such that if we denote it's corresponding tensor by $A^{(*)}$ we have $\text{rank}(|A^{(*)}|) = M^{\frac{N}{2}}$.

For this tensor to be realized by a shallow net we must have $Z \geq M^{\frac{N}{2}}$ and in particular we must have a network with exponential width (exponential in N).

■

Foundations of Deep Learning – Homework Assignment #2

Adi Alm & Tomer Epshtain

Part 3: (3)

Problem:

Consider the following “quadrant” partition of the N input elements:

$$I_{quad} = \left\{ 1, 2, \dots, \frac{1}{4}N, \frac{1}{2}N + 1, \frac{1}{2}N + 2, \dots, \frac{3}{4}N \right\}$$

$$I_{quad}^c = \left\{ \frac{1}{4}N + 1, \frac{1}{4}N + 2, \dots, \frac{1}{2}N, \frac{3}{4}N + 1, \frac{3}{4}N + 2, \dots, N \right\}$$

Prove that under this partition, the separation rank of a function realized by a deep network (with $L = \log_2(N)$ hidden layers) is no greater than $r_{L-1} \cdot r_{L-2}^2$ where r_l stands for the width of layer l .

Solution:

First, a visualization of the partition:

1		$\frac{1}{4}N + 1$		$\frac{1}{2}N + 1$		$\frac{3}{4}N + 1$	
2		$\frac{1}{4}N + 2$		$\frac{1}{2}N + 2$		$\frac{3}{4}N + 2$	
3		:		:		:	
:							
	$\frac{1}{4}N$				$\frac{3}{4}N$		N

Let $\bar{h} \in \overline{H_B}$ be a deep net. Denote by A^{HT} its corresponding tensor given by the HT decomposition:

$$\Phi^{1,j,\gamma} = \sum_{\alpha=1}^{r_0} a_{\alpha}^{1,j,\gamma} \cdot (a^{0,2j-1,\alpha} \otimes a^{0,2j,\alpha})$$

...

$$\Phi^{l,j,\gamma} = \sum_{\alpha=1}^{r_{l-1}} a_{\alpha}^{l,j,\gamma} \cdot (\Phi^{l-1,2j-1,\alpha} \otimes \Phi^{l-1,2j,\alpha})$$

...

$$\Phi^{L-1,j,\gamma} = \sum_{\alpha=1}^{r_{L-2}} a_{\alpha}^{L-1,j,\gamma} \cdot (\Phi^{L-2,2j-1,\alpha} \otimes \Phi^{L-2,2j,\alpha})$$

$$A^{HT} = \sum_{\alpha=1}^{r_{L-1}} a_{\alpha}^L \cdot (\Phi^{L-1,1,\alpha} \otimes \Phi^{L-1,2,\alpha})$$

Let's view $A^{HT'}$'s matricization w.r.t $I := I_{quad}$:

$$[|A^{HT}|]_I = \left[\left| \sum_{\alpha=1}^{r_{L-1}} a_\alpha^L \cdot (\Phi^{L-1,1,\alpha} \otimes \Phi^{L-1,2,\alpha}) \right| \right]_I =$$

{Linearity of matricization and "matricization of outer product equals Kronecker product of matricizations" property}

$$= \sum_{\alpha=1}^{r_{L-1}} a_\alpha^L \cdot \left([|\Phi^{L-1,1,\alpha}|]_{I \cap [N/2]} \odot [|\Phi^{L-1,2,\alpha}|]_{I - \frac{N}{2} \cap [N/2]} \right)$$

Let's take another step back through the HT decomposition:

$\forall I' \subset [N/2], j = 1, 2, \gamma \in [r_{L-1}]$:

$$[|\Phi^{L-1,j,\gamma}|]_{I'} = \left[\left| \sum_{\beta=1}^{r_{L-2}} a_\beta^{L-1,j,\gamma} \cdot (\Phi^{L-2,2j-1,\beta} \otimes \Phi^{L-2,2j,\beta}) \right| \right]_{I'} =$$

{Linearity of matricization and "matricization of outer product equals Kronecker product of matricizations" property}

$$\sum_{\beta=1}^{r_{L-2}} a_\beta^{L-1,j,\gamma} \cdot \left([|\Phi^{L-2,2j-1,\beta}|]_{I' \cap [N/4]} \odot [|\Phi^{L-2,2j,\beta}|]_{I' - \frac{N}{4} \cap [N/4]} \right)$$

Putting it all together:

$$\begin{aligned} [|A^{HT}|]_{I_{quad}} &= \sum_{\alpha=1}^{r_{L-1}} a_\alpha^L \cdot \left([|\Phi^{L-1,1,\alpha}|]_{I_{quad} \cap [N/2]} \odot [|\Phi^{L-1,2,\alpha}|]_{I_{quad} - \frac{N}{2} \cap [N/2]} \right) = \\ &= \sum_{\alpha=1}^{r_{L-1}} a_\alpha^L \\ &\cdot \left(\begin{array}{l} \sum_{\beta=1}^{r_{L-2}} a_\beta^{L-1,1,\alpha} ([|\Phi^{L-2,1,\beta}|]_{I_{quad} \cap [N/2] \cap [N/4]} \odot [|\Phi^{L-2,2,\beta}|]_{((I_{quad} \cap [N/2]) - \frac{N}{4}) \cap [N/4]}) \odot \\ \cdot \left(\sum_{\beta=1}^{r_{L-2}} a_\beta^{L-1,2,\alpha} ([|\Phi^{L-2,3,\beta}|]_{(I_{quad} - \frac{N}{2} \cap [N/2]) \cap [N/4]} \odot [|\Phi^{L-2,4,\beta}|]_{((I_{quad} - \frac{N}{2} \cap [N/2]) - \frac{N}{4}) \cap [N/4]}) \right) \end{array} \right) \end{aligned}$$

Let's understand that matricizations' indices:

1. $I_{quad} \cap [N/2] \cap [N/4] = [N/4]$
2. $((I_{quad} \cap [N/2]) - \frac{N}{4}) \cap [N/4] = \emptyset$
3. $(I_{quad} - \frac{N}{2} \cap [N/2]) \cap [N/4] = [N/4]$
4. $((I_{quad} - \frac{N}{2} \cap [N/2]) - \frac{N}{4}) \cap [N/4] = \emptyset$

So

$$\begin{aligned} &= \sum_{\alpha=1}^{r_{L-1}} a_\alpha^L \cdot \left(\begin{array}{l} \left(\sum_{\beta=1}^{r_{L-2}} a_\beta^{L-1,1,\alpha} ([|\Phi^{L-2,1,\beta}|]_{[\frac{N}{4}]} \odot [|\Phi^{L-2,2,\beta}|]_{\emptyset}) \right) \odot \\ \cdot \left(\sum_{\beta=1}^{r_{L-2}} a_\beta^{L-1,2,\alpha} ([|\Phi^{L-2,3,\beta}|]_{[\frac{N}{4}]} \odot [|\Phi^{L-2,4,\beta}|]_{\emptyset}) \right) \end{array} \right) = \end{aligned}$$

$$\sum_{\alpha=1}^{r_{L-1}} a_\alpha^L \cdot \left(\begin{array}{l} \sum_{\beta=1}^{r_{L-2}} a_\beta^{L-1,1,\alpha} (\text{vec}(\Phi^{L-2,1,\beta}) \cdot \text{vec}(\Phi^{L-2,2,\beta})^T) \odot \\ (\sum_{\beta=1}^{r_{L-2}} a_\beta^{L-1,2,\alpha} (\text{vec}(\Phi^{L-2,3,\beta}) \cdot \text{vec}(\Phi^{L-2,4,\beta})^T) \end{array} \right)$$

- $\text{vec}(u) \cdot \text{vec}(v)^T$ produces a rank 1 matrix
- A sum of r_{L-2} rank 1 matrix produces a $\leq r_{L-2}$ ranked matrix

Denote:

$\forall \alpha \in [r_{L-1}]$:

$$U_\alpha := \sum_{\beta=1}^{r_{L-2}} a_\beta^{L-1,1,\alpha} (\text{vec}(\Phi^{L-2,1,\beta}) \cdot \text{vec}(\Phi^{L-2,2,\beta})^T)$$

$$V_\alpha := (\sum_{\beta=1}^{r_{L-2}} a_\beta^{L-1,2,\alpha} (\text{vec}(\Phi^{L-2,3,\beta}) \cdot \text{vec}(\Phi^{L-2,4,\beta})^T))$$

So we have $\text{rank}(U_\alpha) \leq r_{L-2}$, $\text{rank}(V_\alpha) \leq r_{L-2}$ for all $\alpha \in [r_{L-1}]$.

We have:

$$[|A^{HT}|]_{I_{quad}} = \sum_{\alpha=1}^{r_{L-1}} a_\alpha^L \cdot (U_\alpha \odot V_\alpha)$$

We saw in class that $\text{rank}(A \odot B) = \text{rank}(A) \cdot \text{rank}(B)$ for any two matrices A & B .

So $\forall \alpha \in [r_{L-1}]$: $\text{rank}(U_\alpha \odot V_\alpha) = \text{rank}(U_\alpha) \cdot \text{rank}(V_\alpha) \leq r_{L-2}^2$

Therefore:

$$\text{rank}([|A^{HT}|]_{I_{quad}}) \leq r_{L-1} \cdot r_{L-2}^2$$

As the sum of r_{L-1} matrices with rank $\leq r_{L-2}^2$.

■