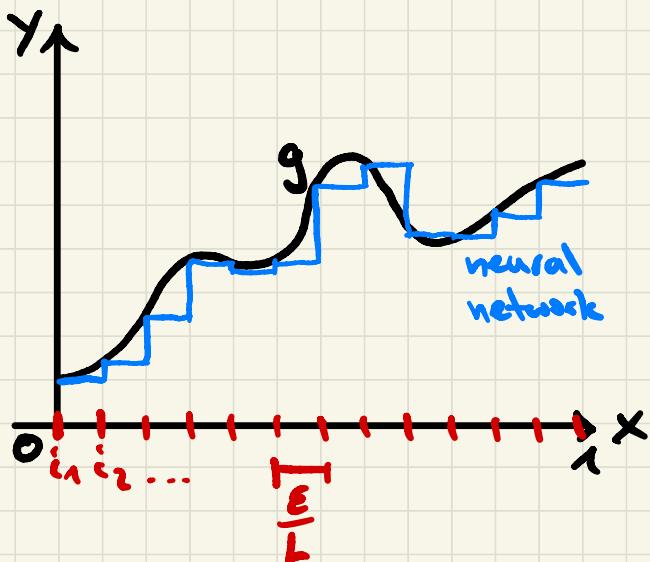


## Proof for L-Lipschitz functions

Assume  $g: [0,1] \rightarrow \mathbb{R}$  is L-Lipschitz.

Then we can construct an appropriate 2-layer neural network using two steps:



1. Partition the x-axis into intervals of size  $\frac{\epsilon}{L}$ .

2. For each interval  $k$ , construct an ANN  $f_k$  such that

$$f_k(x) = \begin{cases} g(i_k) & \text{if } x \in [i_k, i_{k+1}] \\ 0 & \text{else} \end{cases}$$

1. is a trivial step. For 2., we choose the following activation function:

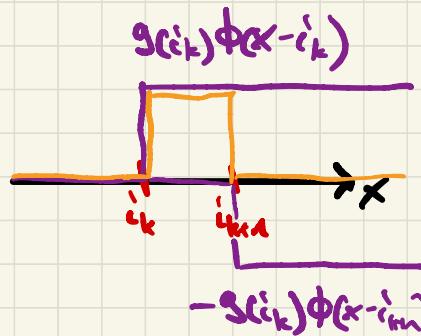
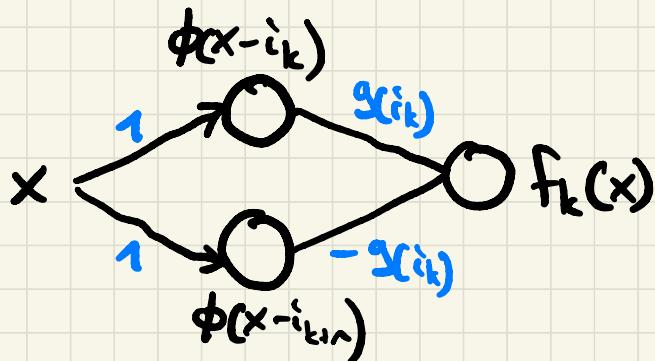
$$\phi(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{else} \end{cases}$$

Then the following construction works:

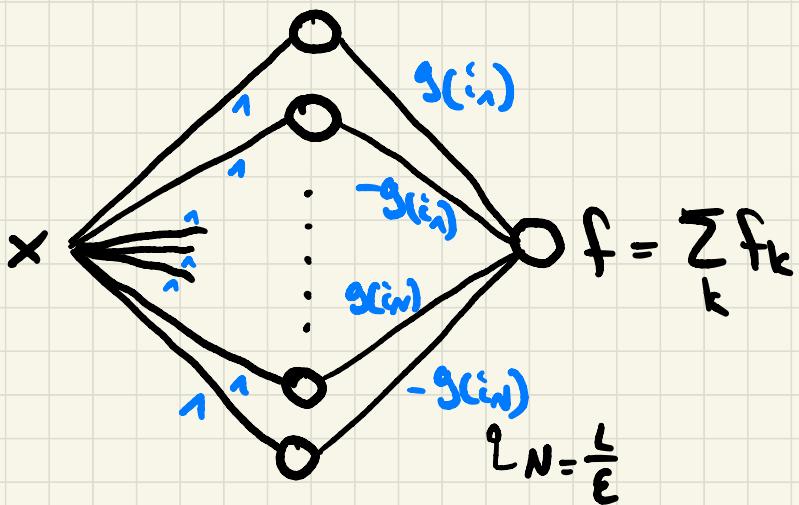
$$f_k(x) = g(i_k) \cdot [\phi(x - i_k) - \phi(x - i_{k+1})]$$

$$= g(i_k) \phi(1 \cdot x - i_k) - g(i_k) \phi(1 \cdot x - i_{k+1})$$

This is a 2-layer neural network!



Our final function is  $f(x) = \sum_k f_k(x)$ , which is simply a wider 2-layer neural network.



With this construction, we can complete the proof:

$$\|f(x) - g(x)\|_\infty = \left\| \sum_j f_j(x) - g(x) \right\|_\infty$$

pick interval with highest deviation  $\Rightarrow \max_{i_k} \|g(i_k) - g(x)\|_\infty$

$g$  is Lipschitz  $\leq L \cdot \|i_k - x\|_\infty$

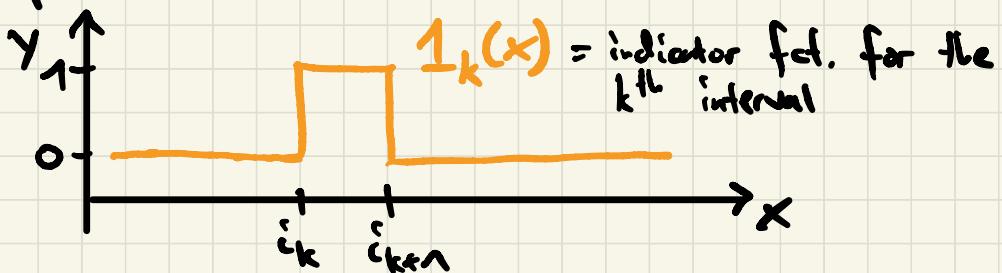
$$\leq L \cdot \frac{\epsilon}{L}$$

$$= \epsilon$$

□

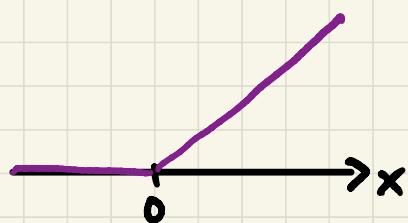
## Case of ReLU act. fcts.:

Our proof relied on the fact that we can construct indicator functions for every interval



We achieved this using threshold functions as act. fcts. But: in modern ANNs, these are rarely used. Far more famous are ReLUs:

$$\text{ReLU}(x) = \max(0, x)$$

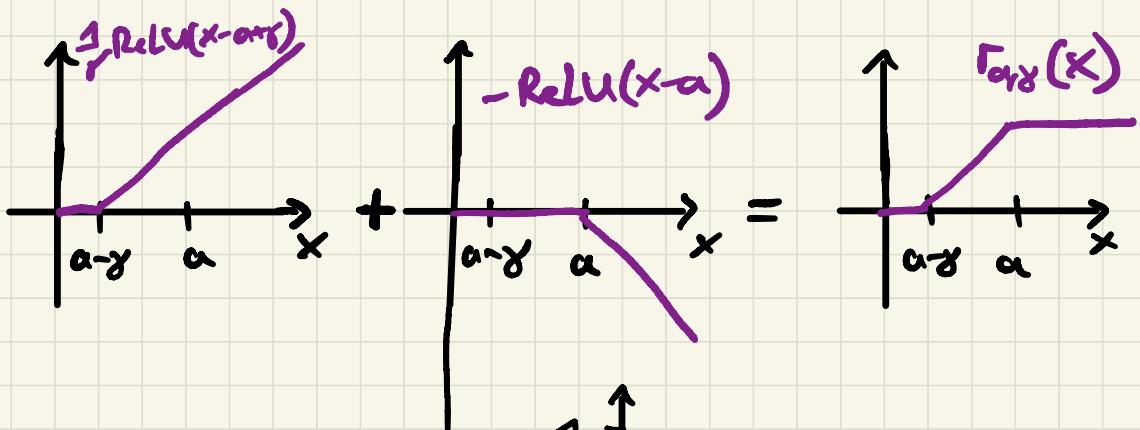


To show that the theorem also holds (up to constants) for ReLUs, we show that  $1_k$  can be constructed using ReLUs!

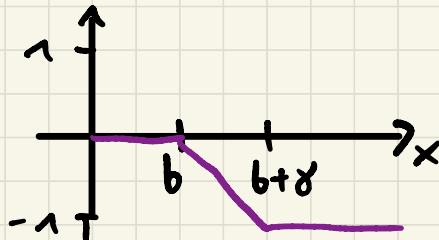
Indicator fct.: Assume  $a \in \mathbb{R}$ ,  $\gamma > 0$ . Then

$$\Gamma_{a,\gamma}(x) = \frac{1}{\gamma} [\text{ReLU}(x - a + \gamma) - \text{ReLU}(x - a)]$$

looks like:



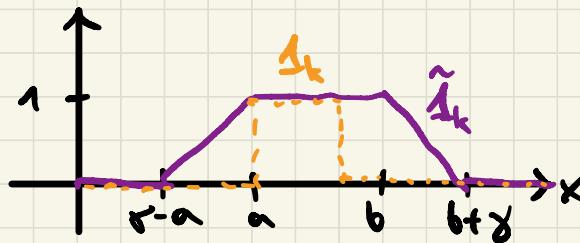
Also note that  $\Gamma_{b,-\gamma}$  is:



Thus, we can approximate  $\mathbf{1}_k$ , with  $a = i_k$ ,  $b = i_{k+1}$ ,

using

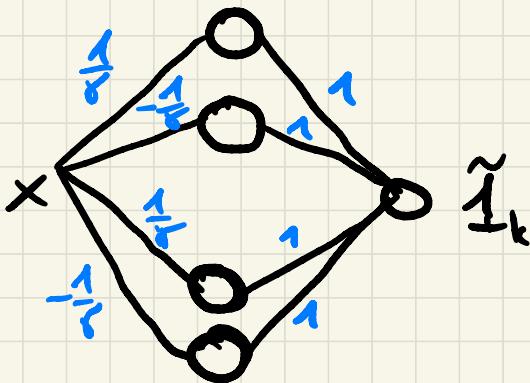
$$\tilde{\chi}_k(x) = \Gamma_{a,\gamma}(x) + \Gamma_{b,-\gamma}(x)$$



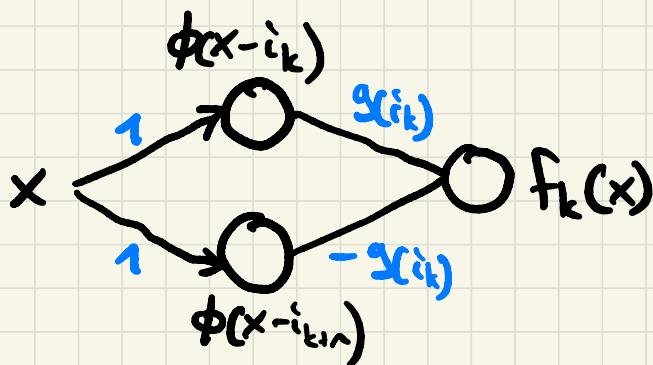
For smaller values of  $\gamma$ , we get closer to the actual indicator fct.! In fact, we can find a  $\gamma$  such that

$$\|\mathbf{1}_k(x) - \tilde{\mathbf{1}}_k(x)\|_\infty < \epsilon$$

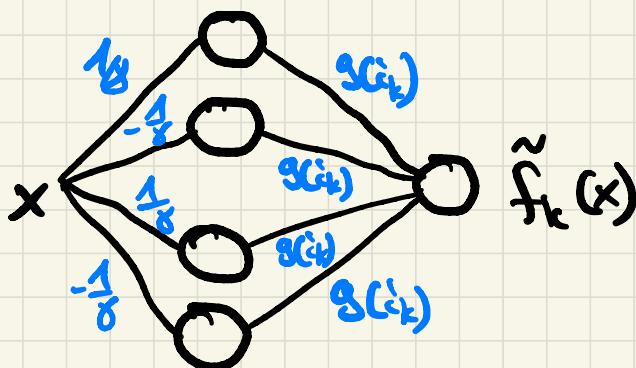
Summary: we can approximate the indicator fct. to given error  $\epsilon$  using a 2-layer ReLU neural network!



With threshold act. fct., our block for each interval looked like:



With Relu, it now looks like:



With  $\tilde{f}(x) = \sum_k \tilde{f}_k$ , we get:

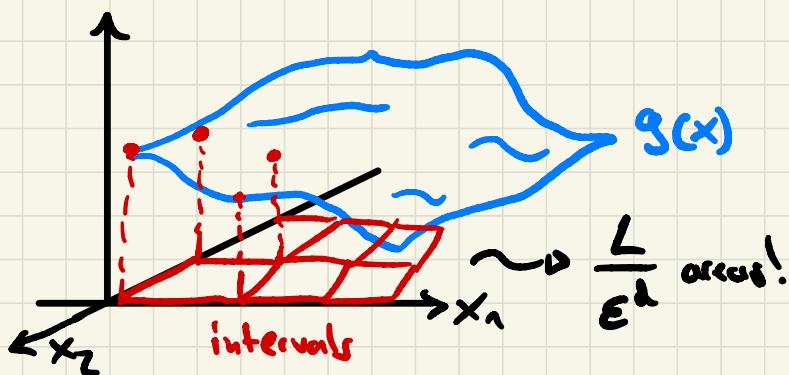
$$|\tilde{f} - g|_\infty \leq \underbrace{|\tilde{f} - f|_\infty}_{\text{set } g \text{ s.t.}} + \underbrace{|f - g|_\infty}_{< \epsilon} < 2\epsilon \quad \square$$

$< \epsilon$   
 "approx. indicators using Relu's"      "approx. g using indicator fcts."

# Multivariate functions

This is pretty much the same idea as in 1D:  
we construct d-dimensional indicator regions!

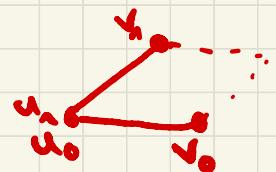
Example in 2D:



So how can we construct a fct. that is 1 on only one tile and 0 otherwise?

Assume our tile has corners  $\{(u_j^i, v_j^i) \mid j \in [1, d]\}$

We first construct indicators for every dimension using threshold act. fcts.  $\phi$ :



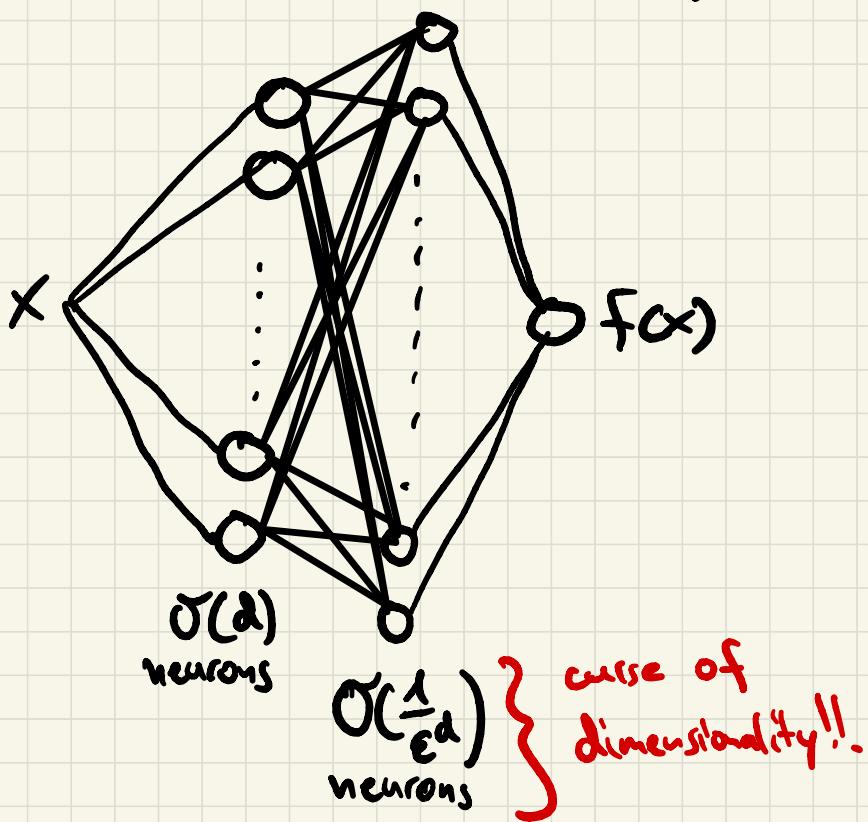
$$h_j^i(x) = \phi(x - u_j^i) - \phi(x - v_j^i)$$

We then check if all interval conditions are satisfied as follows:

$$h(x) = \phi\left(\sum_{i=1}^d h_i(x_i) - (d-1)\right)$$

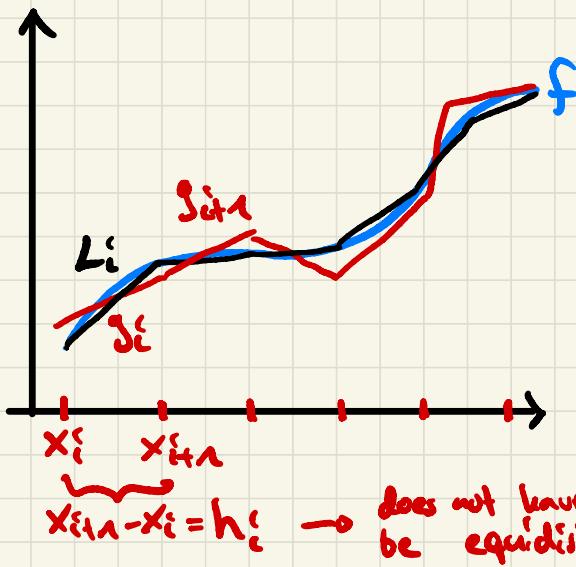
only = 1 if all  
 $h_i(x_i) = 1,$   
0 otherwise.

Our final network then has 3 layers:



Proof:  $\frac{1}{p_2}$ -bound

Let's start with a picture.



$f$ : the function we want to fit

$g_i$ : the pieces of our network

$x_i^c$ : segmentation of input

$L_i$ : linear interpolation of  $f$  using  $x_i^c$

We assume that the  $x_i^c$  were chosen such that

$$(A) \quad \|f(x) - g_i(x)\| \leq \varepsilon \text{ for } x \in [x_i^c, x_{i+1}^c] \forall i.$$

The proof consists of two parts:

1) Show that  $\varepsilon$  is related to the interpol.

error of  $L$ :  $\|e_i\| = \|f(x) - L_i(x)\| \leq 2\varepsilon$

2) Show that the max $|e_i| = \frac{h_i^c}{8} \|f''(m_i)\| + O(h_i^c)$   
with  $m_i = \frac{x_i^c + x_{i+1}^c}{2}$ .

From 1) and 2), we get our result:

$$\|e_i\| \leq 2\varepsilon \text{ for all } x \in [x_i, x_{i+1}]$$

this includes  
 $\Leftrightarrow$   
 the max.  
 error

$$\max \|e_i\| = \frac{h_i^2}{8} \|f''(m_i)\| \leq 2\varepsilon \quad \begin{matrix} \text{in the} \\ (\text{left}) \\ h_i \rightarrow 0 \\ \hat{\triangleq} \varepsilon \rightarrow 0! \end{matrix}$$

take  $\sqrt{\cdot}$   
 $\Leftrightarrow$

$$\frac{h_i}{4} \sqrt{\|f''(m_i)\|} \leq \sqrt{\varepsilon}$$

sum over  
 $\Leftrightarrow$   
 all segments

$$\sum_i \frac{h_i}{4} \sqrt{\|f''(m_i)\|} \leq p \sqrt{\varepsilon} \quad \begin{matrix} \text{we sum over} \\ p \text{ pieces} \end{matrix}$$

$$\rightarrow \varepsilon \geq \frac{C}{p^2}$$

in the limit,  
 this becomes  
 $\int_a^b \sqrt{|f''(x)|} dx$

Now lets first prove 1):

Our interpolation is given by:

$$L_i(x) = f(x_i) + \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i} (x - x_i)$$

check:  $L_i(x_i) = f(x_i)$ ,  $L_i(x_{i+1}) = f(x_{i+1}) \checkmark$

Then:  $\|e_i\| = \|f(x) - L_i(x)\| \text{ for } x \in [x_i, x_{i+n}]$

$$\stackrel{\text{add}}{=} \|f(x) - g_i(x) + g_i(x) - L_i(x)\|$$

$$|a-b| \leq |a-c| + |c-b|$$

$$\Delta\text{-ineq.} \leq \underbrace{|f - g_i|}_{\substack{\text{using (A)} \\ \text{this is}}} + \underbrace{|g_i - L_i|}$$

these are two lines on finite intervals!  
 $\Rightarrow$  They are furthest apart at the ends



$$\Rightarrow |g_i - L_i| = f(x)$$

$$\leq \max \left\{ |g_i(x_i) - L_i(x_i)|, |g_i(x_{i+n}) - L_i(x_{i+n})| \right\}$$

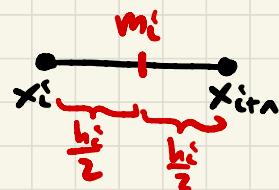
$$\leq |g_i(x_i) - f(x_i)| = f(x_{i+n})$$

using (A)  $\leq \epsilon$

$$\Rightarrow \|e_i\| \leq 2\epsilon \quad \square$$

Now the second part. For this we Taylor expand all terms in  $L_i(x)$  as well as  $f(x)$  around  $m_i = \frac{x_i + x_{i+1}}{2}$ .

Note the following:  $x_i = m_i - \frac{h_i}{2}$   $x_{i+1} = m_i + \frac{h_i}{2}$



and  $x = m_i + \alpha(x) \frac{h_i}{2}$  for  $x \in [x_i, x_{i+1}]$   
and  $\alpha \in [-1, 1]$

In particular, we get  $x - x_i = (\alpha + 1) \frac{h_i}{2}$

$$= (\alpha + 1) \frac{h_i}{2}$$

Let's look at  $L_i(x) = f(x_i) + \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i} (x - x_i)$

$$\frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i} = h_i$$

Let's Taylor:

$$f(x_i) = f(m_i - \frac{h_i}{2}) = \underline{f(m_i)} - f'(m_i) \frac{h_i}{2} + \frac{f''(m_i)}{2} \frac{h_i^2}{4} + O(h_i^3)$$

$$f(x_{i+1}) = \underline{f(m_i)} + f'(m_i) \frac{h_i}{2} + \frac{f''(m_i)}{2} \frac{h_i^2}{4} + O(h_i^3)$$

$$\Rightarrow L(x_i) = f(m_i) + \alpha \frac{h_i}{2} f'(m_i) + \frac{f''(m_i)}{8} h_i^2 + O(h_i^3)$$

Lets also Taylor expand  $f(x)$ :

$$f(x) = f(m_i + \alpha \frac{h_i}{2})$$

$$= f(m_i) + f'(m_i) \alpha \frac{h_i}{2} + f''(m_i) \frac{\alpha^2 h_i^2}{8} + O(h_i^3)$$

Combining these, we get:

$$\|f(x) - L_i(x)\| = \left| \frac{f''(m_i) h_i^2}{8} \right| \cdot |\alpha^2 - 1|$$

which has its maximum at  $\alpha = 0$  (where  $|\alpha^2 - 1|$  has its max, since  $\alpha \in [-1, 1]$  and  $|\alpha^2 - 1| \in [0, 1]$ )

□