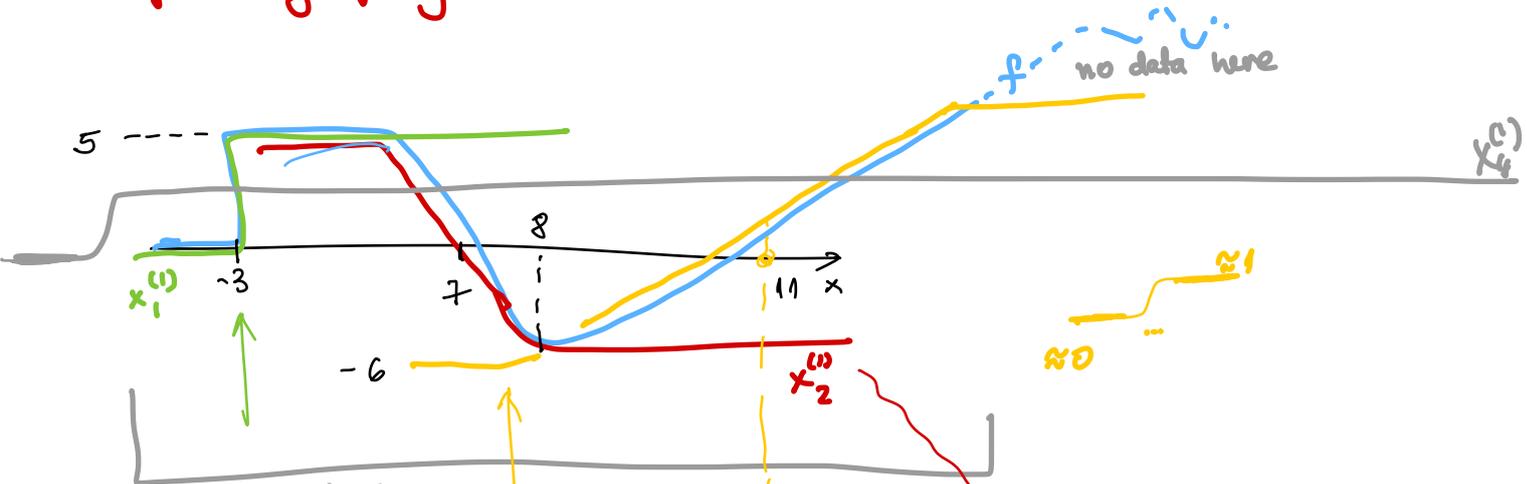


# Lecture 11

- How a nn represents an  $f$
- Analysis & practical issues

Q1 ✓  
HW5 Wed 2/25  
train nn!  
Lec notes  
TB Posted

# Representing $f$ by a nn.



input  $x$  in this range

$$x_1^{(1)} = x$$

$$(x-11) \cdot 0.5 = 0.5x - 5.5 = z_2^{(1)}$$

$$x_2^{(1)} = \varphi(z_2^{(1)})$$

$$(x-3) \cdot 1000 = \underbrace{1000x}_{W_{11}^{(1)}} - \underbrace{3000}_{W_{10}^{(1)}} = z_1^{(1)}$$

Layer 2

$$x^{(2)} = \varphi_{\text{out}}(z^{(2)})$$

$$z^{(2)} = 5 \cdot x_1^{(1)} + 11 \cdot x_2^{(1)} + 12 \cdot x_3^{(1)} - 11$$

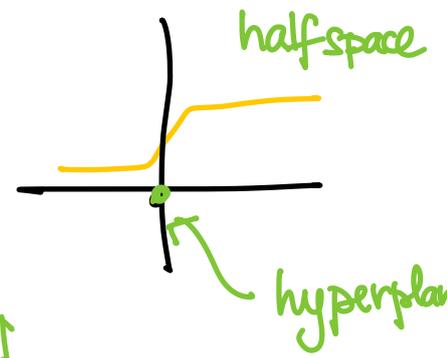
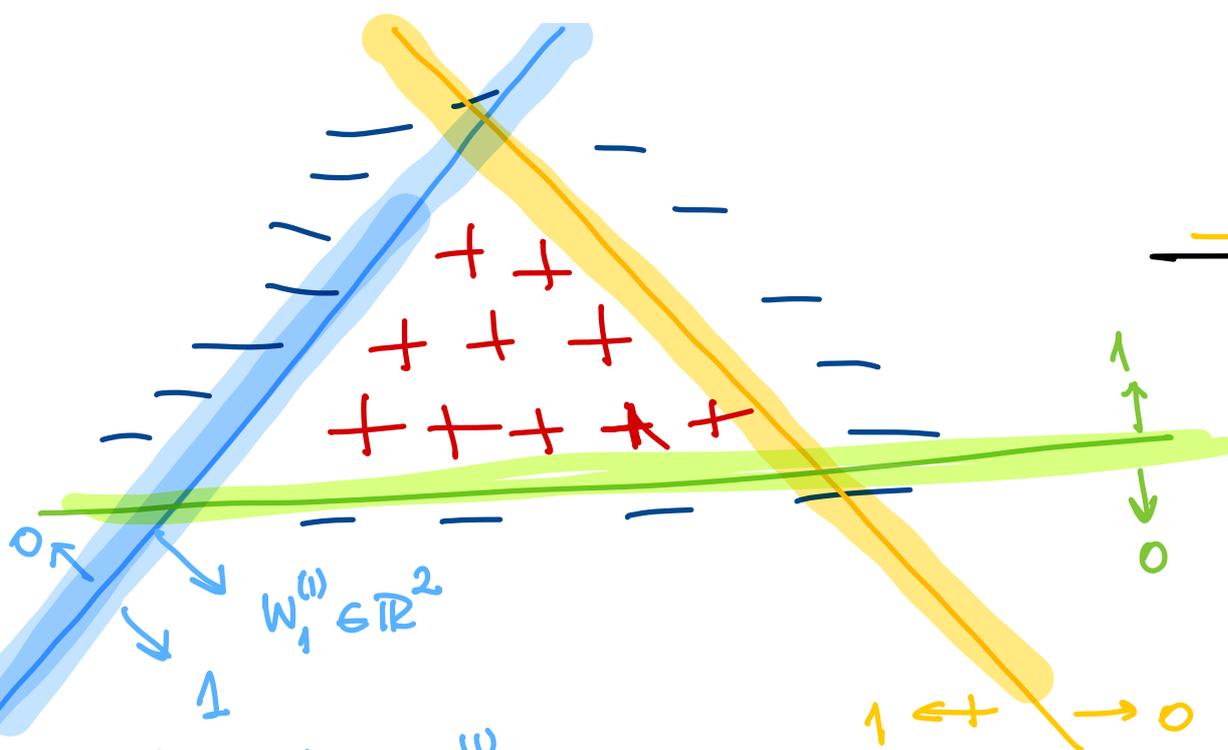
$\varphi_{\text{out}}(z) = z$   
 $W_2^{(2)}$ ,  $W_3^{(2)}$ ,  $W_0^{(2)}$

$$z_2^{(1)} = (x-7)(-1) = 7-x$$

$$x_2^{(1)} = \varphi(7-x)$$

adding a constant

$$\text{OR } \left. \begin{aligned} z_4^{(1)} &= x + 100000 \\ x_4^{(1)} &= \varphi(z_4^{(1)}) \end{aligned} \right\}$$



$$z_1^{(1)} = W_1^{(1)} x + W_0^{(1)}$$

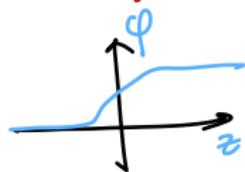
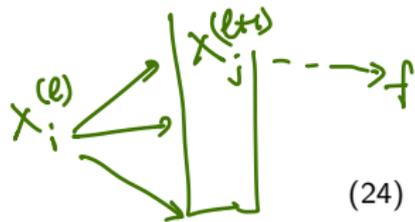
$$= \sum_1^{(1)} \begin{bmatrix} x \\ 1 \end{bmatrix}$$

$$x_1^{(1)} = \varphi(z_1^{(1)})$$

$$x^{(2)} = \text{Agg}_{\varphi_{\text{out}}}(z^{(2)})$$

$$z^{(2)} = 1 \cdot x_1^{(1)} + 1 \cdot x_2^{(1)} + 1 \cdot x_3^{(1)} - 2.5 \cdot x_4^{(1)}$$

$$x_1^{(1)} = \text{ct}$$

From 2 layers to  $L$  layers
 $|z| \text{ large} \Rightarrow \phi' \approx 0$  "saturation"


(24)

$$\frac{\partial x_i^{(l)}}{\partial z_i^{(l)}} = \phi'(z_i^{(l)})$$

$$\frac{\partial x_i^{(l)}}{\partial W_i^{(l)}} = \phi'(z_i^{(l)}) x^{(l-1)}$$

$$z_i^{(l)} = W_i^{(l)} x^{(l-1)} \quad (25)$$

(26)

$$\left\{ \begin{array}{l} \frac{\partial x_i^{(l)}}{\partial x_i^{(l-1)}} = \phi'(z_i^{(l)}) W_i^{(l)} \end{array} \right.$$

(27)

$$\frac{\partial f}{\partial x_i^{(l)}} = \sum_{j=1}^{\text{next}} \frac{\partial f}{\partial x_j^{(l+1)}} \frac{\partial x_j^{(l+1)}}{\partial x_i^{(l)}}$$

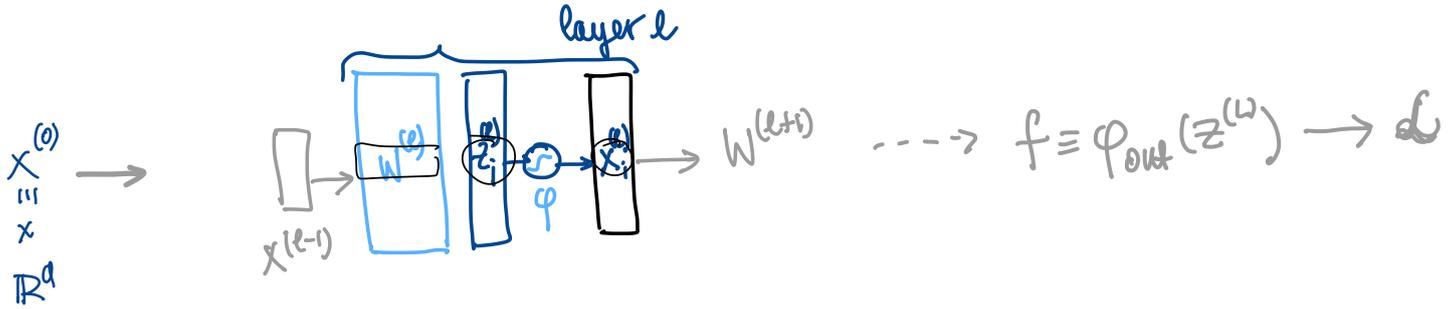
(28)

$$\frac{\partial f}{\partial W_i^{(l)}} = \frac{\partial f}{\partial x_i^{(l)}} \frac{\partial x_i^{(l)}}{\partial W_i^{(l)}}$$

$$\frac{\partial \mathcal{L}}{\partial W_i^{(l)}} = (y_* - \phi_{\text{out}}) \frac{\partial f}{\partial W_i^{(l)}}$$

(29)





for  $t = 1, 2, \dots$   
 for  $k = 1:n$  (examples)

for  $l = 1:L$  store!  
 compute  $x^{(l)}, z^{(l)}$  from  $x^{(l-1)}, W^{(l)}$

FORWARD  
 (Prediction)

$f = \varphi_{\text{out}}(z^{(L)})$ ,  $\mathcal{L}^{\text{train}} += \frac{1}{n} d(y^i, f(x^i))$  }  $\rightarrow$  used for stopping  
not  $\frac{\partial \mathcal{L}}{\partial W}$

BACK PROP.  
 for  $l = L:-1:1$   
 compute  $\frac{\partial \mathcal{L}}{\partial W_i^{(l)}}$  from  $x^{(l)}, \frac{\partial \mathcal{L}}{\partial x^{(l)}}, W^{(l)}, \frac{\partial x^{(l+1)}}{\partial x^{(l)}} \frac{\partial \mathcal{L}}{\partial f}$   
 $i = 1:m_e$   
 $\frac{\partial x^{(l)}}{\partial x^{(l-1)}} \rightarrow$  for  $\frac{\partial \mathcal{L}}{\partial W^{(l-1)}}$

# Issues with Backprop

## Practical

I. ops/iteration

$p = \# \text{parameters}$

ops/  
1 iteration =  $O(pn)$

T. Convergence  $T = ?$  # steps

S. Saturation  $\varphi' \approx 0$

L. Local minima

O. Overfitting

$$W = \{ W^{(l)} \in \mathbb{R}^{m_l \times m_{l-1}} \}_{l=1:L}$$

Ex 1)  $L=2$

$$d = 10^6 = m_0$$

$m_1 = 1000$  hidden

$m_2 = 10$  classification

$$p = 10^9 + \dots$$

Ex 2)

$$L = 100$$

$$m_0 = m_1 = \dots = m_{L-1} = 100$$

$$p \approx 100^3 \approx 10^6$$

# Convergence (T = # iteration)

## Theory

GD is "order 1"

$$\text{err}^t = \|w^t - w^p\| \quad \text{wanted } \leq \underline{\underline{\text{tol}}}$$

$$\frac{\text{err}^{t+1}}{\text{err}^t} \leq \beta < 1$$

$$\text{err}^t \sim \beta^t \xrightarrow{\text{Ex}} T = \log \text{tol}$$

## Practice

L not convex

$$\Rightarrow \frac{\partial L}{\partial W} \approx 0$$

↑  
 $w_{ij}^{(t)}$

but  
W not  
converged !!

$$H = \left[ \frac{\partial^2 L}{\partial w_i \partial w_j} \right]_{p \times p}$$