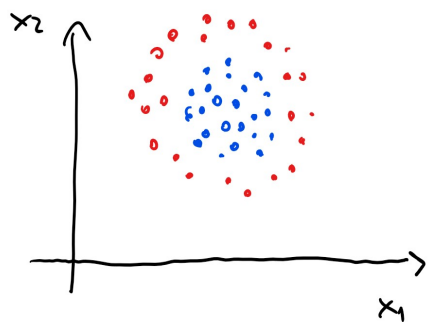


## NEURAL NETWORKS

ASSUME I'VE THE FOLLOWING CLASSIFICATION PROBLEM



IT IS AN EASY PROBLEM, THE CLASS DEPENDS ON THE DISTANCE FROM THE MEAN OF THE DATASET, HOWEVER THE DATA IS NOT LINEARLY SEPARABLE (AT LEAST NOT IN THIS INPUT SPACE). SO I'VE NO HOPE TO SOLVE THE PROBLEM WITH A LINEAR MODEL (WITHOUT USING THE KERNEL TRICK)

WE'VE SEEN THAT KERNEL TRICK ALLOWS TO TRANSFORM THE INPUT SPACE INTO A FEATURE SPACE IN WHICH THE DATA MAY BE LINEARLY SEPARABLE, HOWEVER THIS TRANSFORMATION IS HAND CRAFTED, SO I MUST CHOOSE THE KERNEL THAT TRANSFORMS THE DATA INTO THE RIGHT FEATURE SPACE.

NEURAL NETWORKS ALLOW TO DO SOMETHING AMAZING: THEY LEARN AN INTERMEDIATE REPRESENTATION OF DATA (CALLED LATENT REPRESENTATION) THAT ALLOWS TO SOLVE THE SPECIFIC TASK.

THIS LATENT REPRESENTATION CAN BE VERY COMPLEX (I.E. A VERY NONLINEAR TRANSFORMATION) BUT IT IS OBTAINED STACKING MULTIPLE SIMPLE TRANSFORMATIONS THAT CONSIST IN A LINEAR MAP AND A NONLINEAR FUNCTION APPLIED ELEMENTWISE.

LET'S DESIGN A NEURAL NETWORK WITH 3 LAYERS THAT SOLVES THE CLASSIFICATION PROBLEM DESCRIBED ABOVE.

WE HAVE OUR DATA MATRIX

$$\mathbf{X} \in \mathbb{R}^{N \times 2}$$

WE FIRST APPLY A LINEAR MAP THAT MAPS THE DATA INTO A HIGHER DIMENSIONAL SPACE, FOR EXAMPLE  $\mathbb{R}^5$ , WE ALSO APPLY A NONLINEARITY ELEMENTWISE, SUCH THAT THE MAPPING IS NONLINEAR, THEN WE APPLY A LINEAR MODEL (LOGISTIC REGRESSION) TO THE POINTS IN THE HIDDEN SPACE

$\bar{X} \in \mathbb{R}^{N \times 2}$   
↓  
DATA MATRIX

$W' \in \mathbb{R}^{2 \times 5}$   
↓  
LINEAR MAP

$q: \mathbb{R} \rightarrow \mathbb{R}$   
↓  
NONLINEAR FUNCTION

$W^2 \in \mathbb{R}^{5 \times 1}$   
↓  
LINEAR MAP

$\delta: \mathbb{R} \rightarrow \mathbb{R}$   
↓  
SIGMOID FUNCTION

- FIRST I COMPUTE THE HIDDEN REPRESENTATION

$$H = q(\bar{X} W')$$

$H \in \mathbb{R}^{N \times 5}$

- THEN I APPLY LOGISTIC REGRESSION TO THE DATA IN THE HIDDEN SPACE

$$O = \delta(H W^2)$$

$O \in \mathbb{R}^{N \times 1}$

- THEN I CAN GROUP THE LOSS

$$L = - \sum_{i=1}^N y[i] \log(O[i]) + (1 - y[i]) \log(1 - O[i])$$

- I CAN ALSO WRITE THE OUTPUT AS A FUNCTION OF THE INPUT AS

$$O = \delta(q(\bar{X} W') W^2)$$

### BACKPROPAGATION

NOW THE TRICKY PART, I'VE TO COMPUTE THE GRADIENT OF THE LOSS WITH RESPECT TO ALL THE MODEL'S PARAMETERS

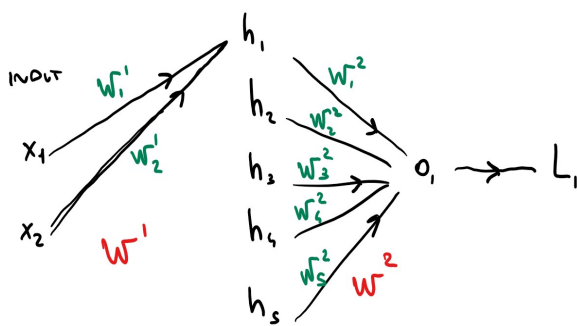
$$\nabla_{\theta} L = ?$$

I CAN DO IT APPLYING THE CHAIN RULE

ASSUME I'VE  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  AND  $g: \mathbb{R}^2 \rightarrow \mathbb{R}^2$   
 $(x, y) \mapsto z = f(x, y)$  AND  $(s, t) \mapsto (x = g_1(s, t), y = g_2(s, t))$   
 AND LET  
 $F(s, t) = (f \circ g)(s, t) = f(g_1(s, t), g_2(s, t))$  THEN

$$\frac{\partial F}{\partial s} = \frac{\partial f}{\partial x} \cdot \frac{\partial g_1(s, t)}{\partial s} + \frac{\partial f}{\partial y} \cdot \frac{\partial g_2(s, t)}{\partial s}$$

$$\frac{\partial F}{\partial t} = \frac{\partial f}{\partial x} \cdot \frac{\partial g_1(s, t)}{\partial t} + \frac{\partial f}{\partial y} \cdot \frac{\partial g_2(s, t)}{\partial t}$$



FOR EXAMPLE LET US COMPUTE  $\frac{\partial L}{\partial w_1^2}$  APPLYING THE CHAIN RULE

$$\frac{\partial L}{\partial w_1^2} = \left( \frac{\partial L}{\partial o} \right) \frac{\partial o}{\partial w_1^2}$$

NOW I CAN COMPUTE  $\frac{\partial L}{\partial w_1^1}$

$$\frac{\partial L}{\partial w_1^1} = \left( \frac{\partial L}{\partial o} \right) \left( \frac{\partial o}{\partial h_1} \right) \frac{\partial h_1}{\partial w_1^1} + \left( \frac{\partial L}{\partial o} \right) \left( \frac{\partial o}{\partial h_2} \right) \frac{\partial h_2}{\partial w_1^1} + \left( \frac{\partial L}{\partial o} \right) \left( \frac{\partial o}{\partial h_3} \right) \frac{\partial h_3}{\partial w_1^1} + \left( \frac{\partial L}{\partial o} \right) \left( \frac{\partial o}{\partial h_4} \right) \frac{\partial h_4}{\partial w_1^1} + \left( \frac{\partial L}{\partial o} \right) \left( \frac{\partial o}{\partial h_5} \right) \frac{\partial h_5}{\partial w_1^1}$$

NOTE THAT THIS TERM IS REPEATED OVER AND OVER, I CAN COMPUTE IT JUST ONCE !!

IN THE SAME WAY I CAN COMPUTE  $\frac{\partial L}{\partial w_2^1}$

$$\frac{\partial L}{\partial w_2^1} = \left( \frac{\partial L}{\partial o} \right) \left( \frac{\partial o}{\partial h_1} \right) \frac{\partial h_1}{\partial w_2^1} + \left( \frac{\partial L}{\partial o} \right) \left( \frac{\partial o}{\partial h_2} \right) \frac{\partial h_2}{\partial w_2^1} + \left( \frac{\partial L}{\partial o} \right) \left( \frac{\partial o}{\partial h_3} \right) \frac{\partial h_3}{\partial w_2^1} + \left( \frac{\partial L}{\partial o} \right) \left( \frac{\partial o}{\partial h_4} \right) \frac{\partial h_4}{\partial w_2^1} + \left( \frac{\partial L}{\partial o} \right) \left( \frac{\partial o}{\partial h_5} \right) \frac{\partial h_5}{\partial w_2^1}$$

WE'VE A LOT OF REPEATED TERMS !!

### IDEA

- COMPUTE THE OUTPUT FOR THE INPUT GOING FORWARD (FOLLOWING THE ARROWS OF THE COMPUTATIONAL GRAPH)

- COMPUTE THE GRADIENTS GOING BACKWARD

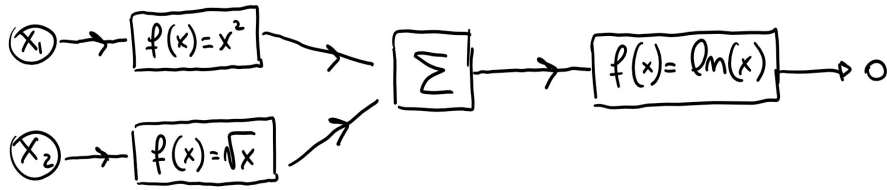
→ IN PROGRAMMING ENVIRONMENTS THAT SUPPORT AUTOMATIC DIFFERENTIATION (AUTOGRAD) EACH OPERATION HAS A FORWARD AND BACKWARD FUNCTIONS, THE FORWARD OPERATION COMPUTES THE OUTPUT W.R.T. THE INPUT, WHILE THE BACKWARD FUNCTION COMPUTES THE GRADIENT OF THE OUTPUT W.R.T. THE INPUT

# EXAMPLE

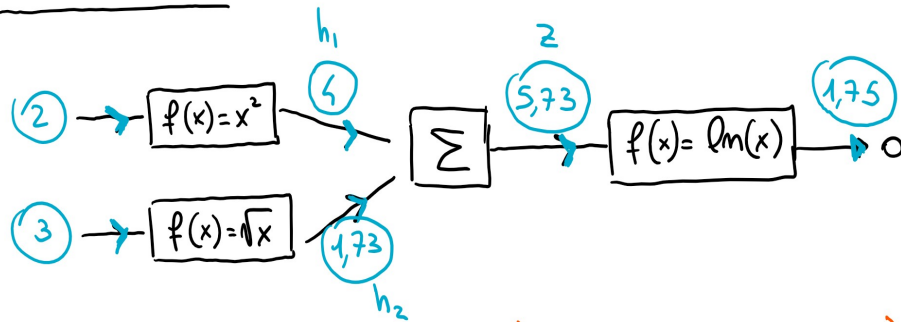
## COMPUTATIONAL GRAPH

○ DATA

□ OPERATION



## FORWARD PASS



## BACKWARD PASS

