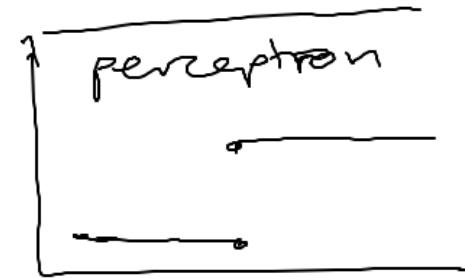
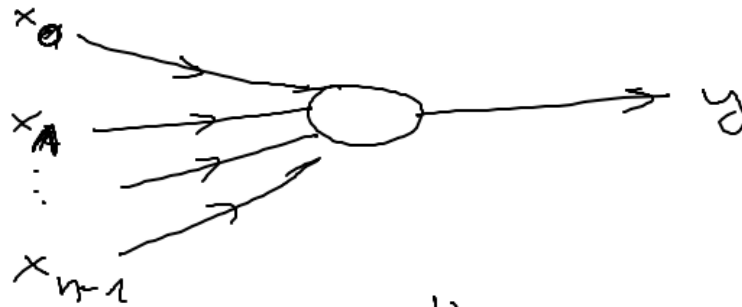


# Neural nets

Neuron

is a function  
of a form

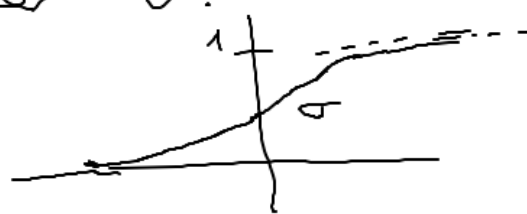


$$\mathbb{R}^n \ni (x_0, x_1, \dots, x_{n-1}) \mapsto \sigma \left( \underbrace{\alpha_0 x_0 + \dots + \alpha_{n-1} x_{n-1}}_{\text{Weights}} + \underbrace{\alpha_n}_{\text{bias}} \right) \in \mathbb{R}$$

$\begin{cases} x_n = -1 \\ \alpha_n x_n = -\alpha_n \end{cases}$

Common choices for  $\sigma$ :

(S1)  $\sigma(x) = \frac{1}{1 + e^{-x}}$   
sigmoid



ReLU

$$\sigma(x) = \max(x, 0)$$

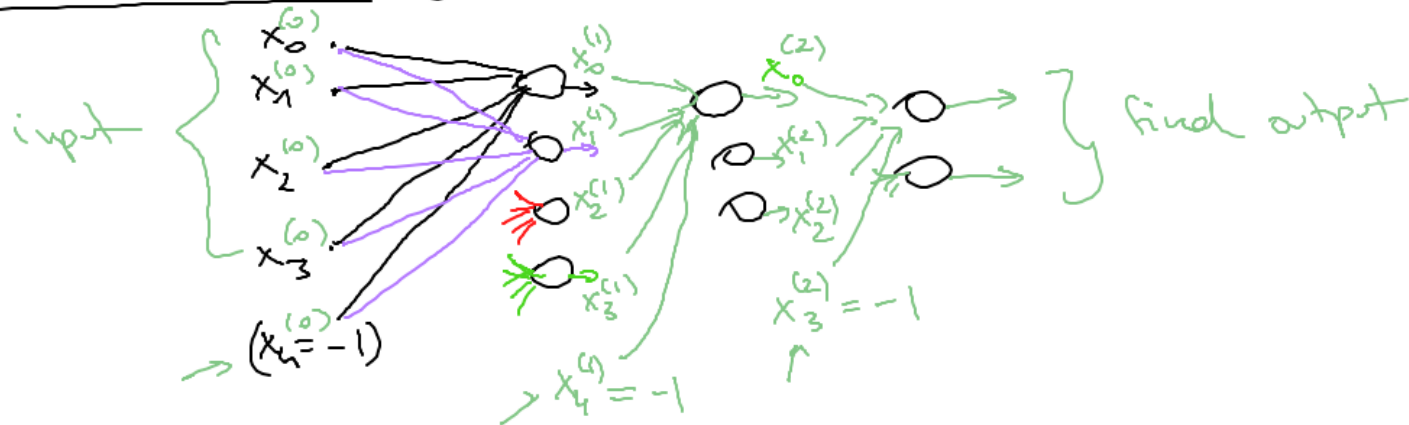


One usually takes  $x_n := -1$  and then

$$(x_0, \dots, x_{n-1}) \mapsto \sigma\left(\sum_{k=0}^n \alpha_k x_k\right) \leftarrow$$

so we get a nicer formula.

Neural network (dense)



A neural network is a function which is a composition of neurons

# Universality Theorem (there are many versions)

Thm. (Cybenko)

Let  $\sigma$  - sigmoidal, continuous function.  
( $\sigma(-\infty) = 0, \sigma(\infty) = 1$ )

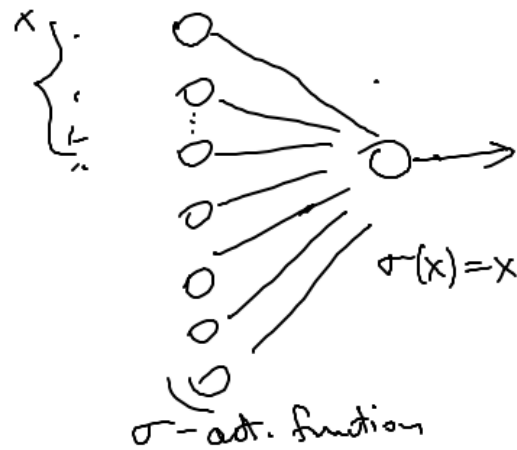
Then the finite sums of the form

$$G(x) = \sum_{j=1}^n \alpha_j \sigma(\underbrace{y_j^T x}_{\text{weights}} + \underbrace{\theta_j}_{\text{bias}})$$

$$, x \in \underline{[0,1]^n}$$

are dense in  $(C([0,1]^n), \|\cdot\|_{\text{sup}})$ .

{ Proof is short, but some background in Functional Analysis is needed.



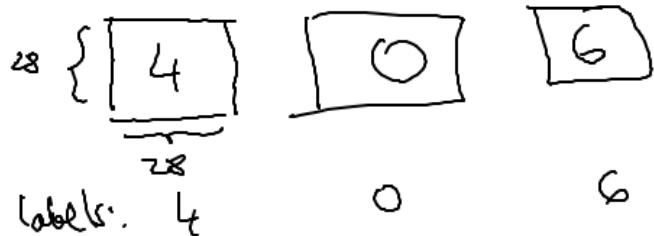
Neural nets are often used in classification problems:

the input is of one of the given  $N$  classes

the aim: given input, decide to what class it belongs to

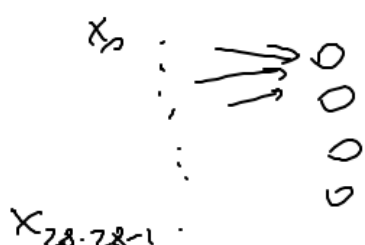
MNIST database of hand-written digits:

the picture is a point in  $\mathbb{R}^{28 \cdot 28}$



60k pictures

expected output for digit "7"



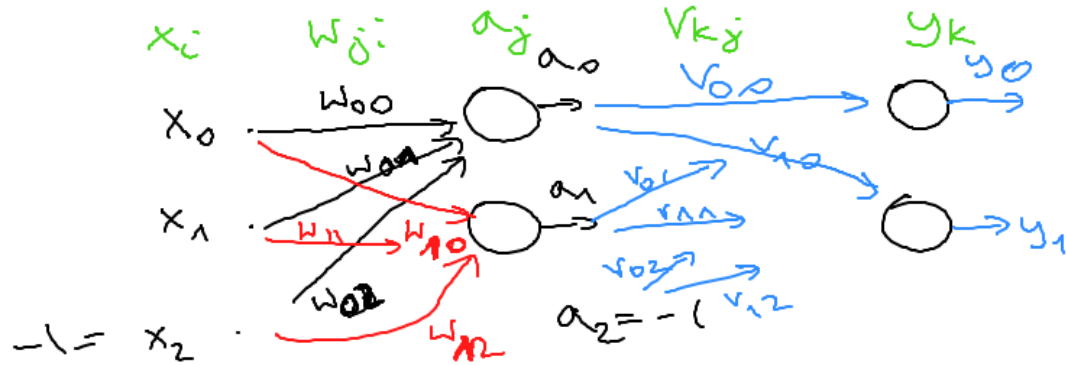
MNIST: last layer

00 → (0,1)	0.1	
01 → (0,1)	0.05	
⋮	⋮	
07	0.8	answer: "7"
09 → (0,1)	0.2	

last layer is of size  $N$   
and sigmoid is used as an activation function

arg max

Neural net with 2 layers with 2 neurons in each layer.



$$a_0 = \sigma(\underbrace{w_{00}x_0 + w_{01}x_1 + w_{02}x_2}_{\text{net}_0})$$

$$a_1 = \sigma(\underbrace{w_{10}x_0 + w_{11}x_1 + w_{12}x_2}_{\text{net}_1})$$

$$W = \begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \end{bmatrix}$$

$$W \cdot \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} \text{net}_0 \\ \text{net}_1 \end{bmatrix}$$

$$y_0 = \sigma(v_{00}a_0 + v_{01}a_1 + v_{02}a_2)$$

Let's suppose that for some input  $\begin{bmatrix} x_0 \\ x_1 \end{bmatrix}$  we expect the net to output  $\begin{bmatrix} t_0 \\ t_1 \end{bmatrix}$  but actually it outputs  $\begin{bmatrix} y_0 \\ y_1 \end{bmatrix}$

we introduce a loss function

$$L(t, y) = \frac{1}{2} \|t - y\|_2^2 = \frac{1}{2} \sum_{k=0}^{N-1} (t_k - y_k)^2$$

↑            ↑  
target    actual output

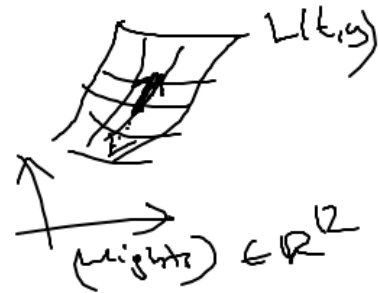
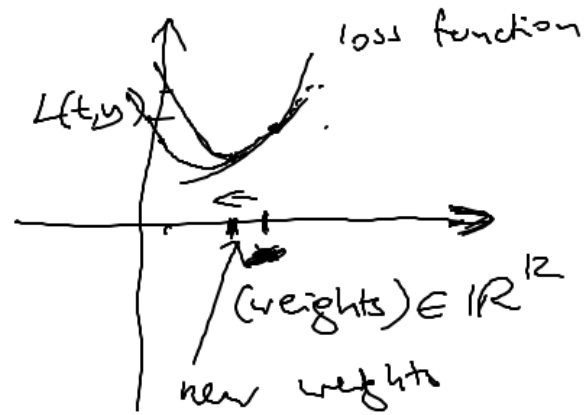
We want to adjust the weights so that the loss function gets smaller.

general idea: calculate the gradient of  $L$   
with respect to the weights  $(w_{00}, \dots, v_{12})$

and then:

$$\rightarrow (\text{new weights}) = (\text{old weights}) - c \left( \frac{\partial L}{\partial \text{weights}} \right)$$

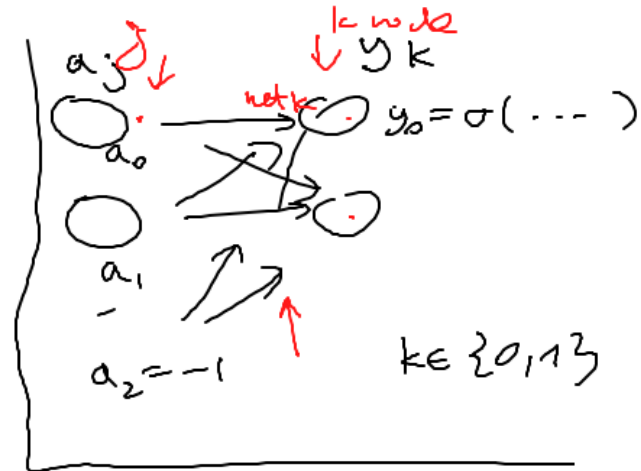
↑  
learning rate,  
e.g.  $c = 0.01$



$$(\text{new } v_{kj}) = (\text{old } v_{kj}) - c \frac{\partial L}{\partial v_{kj}}$$

Fix  $k_0, j_0$  and let  $\underline{v} = v_{k_0 j_0}$

$$L(t, y) = \frac{1}{2} \sum_k (t_k - y_k)^2 = \frac{1}{2} \sum_k (t_k - \varphi(\sum_j v_{kj} a_j))^2$$



$$\frac{\partial L}{\partial v} = \sum_k (t_k - \varphi(\sum_j v_{kj} a_j)) \cdot \frac{\partial}{\partial v} (-\varphi(\sum_j v_{kj} a_j)) =$$

$$= -\sum_k (t_k - \varphi(\sum_j v_{kj} a_j)) \cdot \varphi'(\sum_j v_{kj} a_j) \cdot \frac{\partial}{\partial v} (\sum_j v_{kj} a_j) =$$

$$= (\varphi(\sum_j v_{kj} a_j) - t_{k_0}) \varphi'(\sum_j v_{kj} a_j) \cdot a_{j_0}$$

"error" at node  $k_0$ 
for sigmoid

$\varphi'(\text{net}_{k_0}) = y_k(1-y_k)$

Fix  $i, j_0$ , let  $\underline{w} = w_{j_0 i_0}$

$$a_j = \sigma(\sum_i w_{ji} x_i) \leftarrow$$

$$\frac{\partial L}{\partial w} = \sum_k (t_k - \varphi(\sum_j v_{kj} a_j)) \cdot \frac{\partial}{\partial a_{j_0}} (-\varphi(\sum_j v_{kj} a_j)) \cdot \frac{\partial a_{j_0}}{\partial w}$$

$$= \sum_k (t_k - \varphi(\sum_j v_{kj} a_j)) \cdot \varphi'(\sum_j v_{kj} a_j) \cdot v_{k j_0} \cdot \sigma'(\sum_i w_{j_0 i} x_i) \cdot x_{i_0}$$

"delta signal" (k)
net  $j_0$

(counterpart of) the "error" at node  $j_0$

ex.  
 $k_0=0, j_0=1$   
 $\sum_{k=0}^1 \sum_{j_0=0}^2 v_{kj}$

$v_{01}$

$\frac{\partial}{\partial v_{01}} v_{02} = 0$

$\frac{\partial}{\partial x} y = 0$

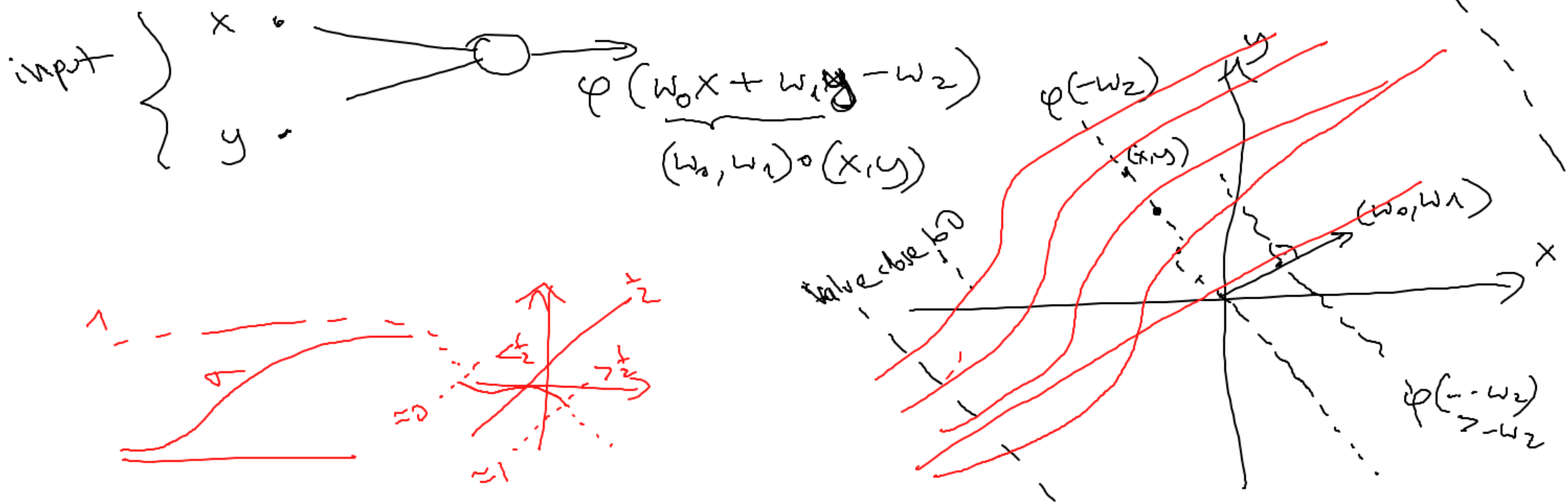
$\frac{\partial}{\partial x} x = 1$

backpropagation

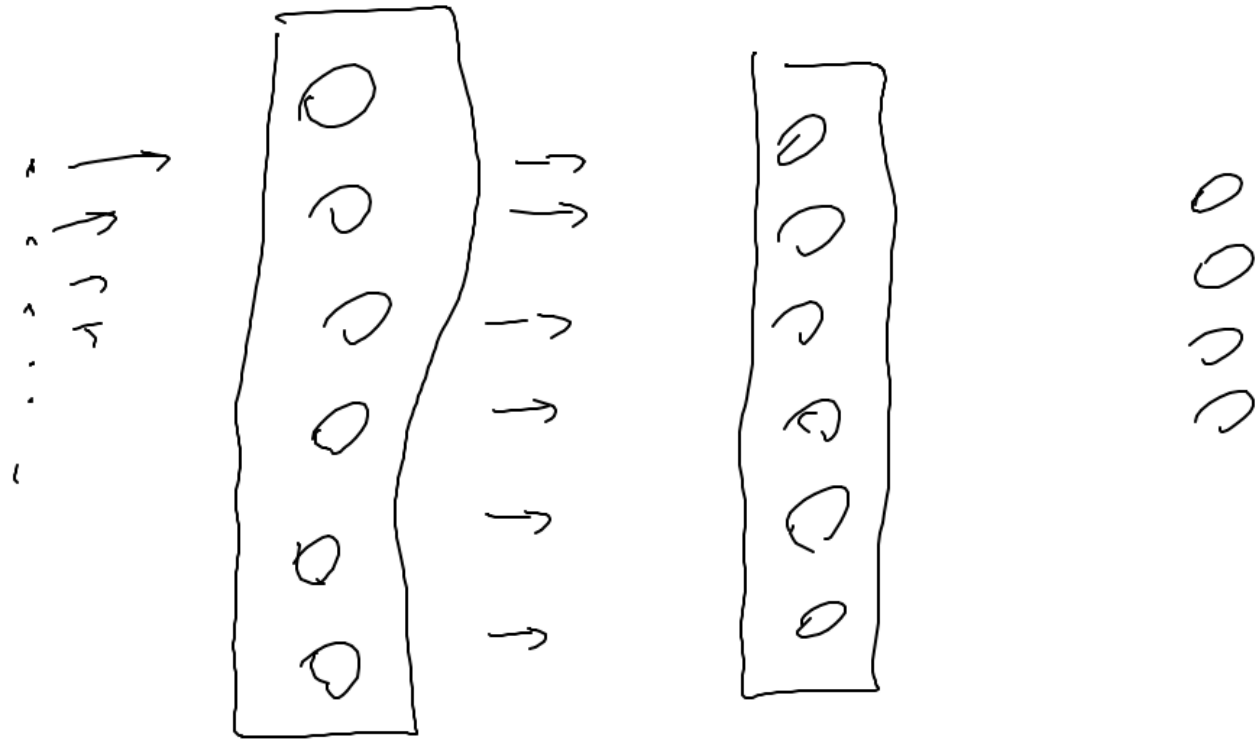
For  $\varphi(x) = \frac{1}{1+e^{-x}}$  :  $\varphi'(x) = \varphi(x)(1-\varphi(x))$

On the ~~webpage~~: an example of forward and back propagation

How does the function defined by a neuron look like?







layer is a (convenient) basic building block