

Tutorial 7

26 February 2022 11:00

* ℓ_p are convex but not differentiable everywhere.

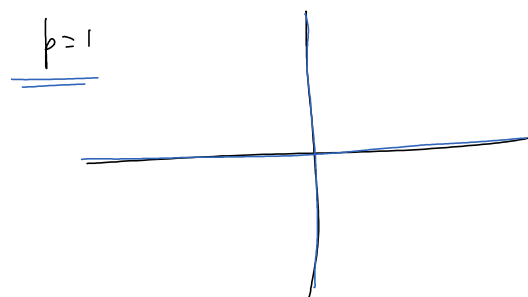
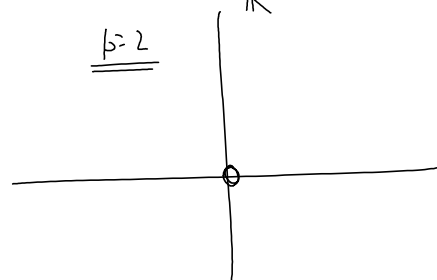
$$f(x) = \|x\|_p := \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} \quad x \in \mathbb{R}^n, p \geq 1$$

* $p > 1$: convex but not differentiable at $x = \underline{0}$

* $p = 1$

$$f(x) = \sum_{i=1}^n |x_i|$$

convex but not differentiable at any x such that $x_i = 0$ for some i .



Example $p=2$

$$f(x) = \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

$$\nabla f(x) = \begin{bmatrix} \frac{2x_1}{2\|x\|_2} \\ \vdots \\ \frac{2x_n}{2\|x\|_2} \end{bmatrix} = \underline{\underline{\frac{x}{\|x\|_2}}}$$

defined everywhere except
 $x = \underline{0}$

Convexity

$$\begin{aligned} f(y) - f(x) &\stackrel{?}{\geq} \nabla f(x)^T (y-x) \\ \|y\|_2 - \|x\|_2 &\stackrel{?}{\geq} \frac{x^T (y-x)}{\|x\|_2} = \frac{x^T y}{\|x\|_2} - \frac{\|x\|_2^2}{\|x\|_2} \\ &\stackrel{?}{\geq} \frac{x^T y}{\|x\|_2} - \|x\|_2 \end{aligned}$$

$$\|y\|_2 \geq \frac{x^T y}{\|x\|_2}$$

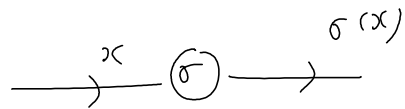
$$\|x\|_2 \|y\|_2 \geq x^T y \quad \checkmark$$

Neural Networks

Neuron

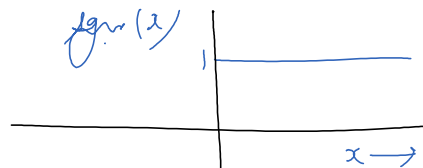
$$\sigma: \mathbb{R} \rightarrow \mathbb{R}$$

(a non linear function)

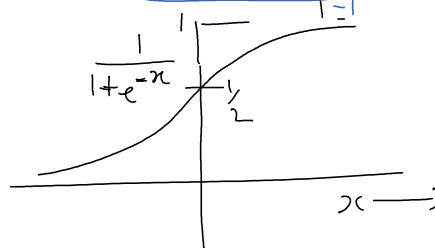


Examples

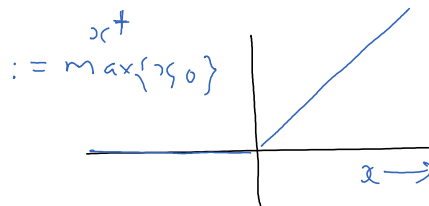
① Sign function



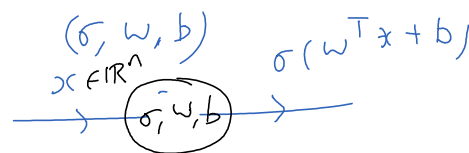
② Sigmoid



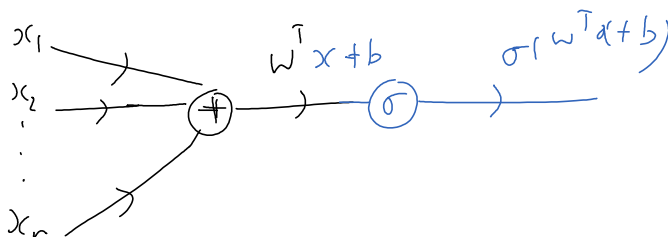
③ ReLU



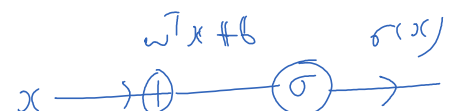
Neurons with vector inputs



$$\sigma: \mathbb{R} \rightarrow \mathbb{R}$$



what does a neuron do



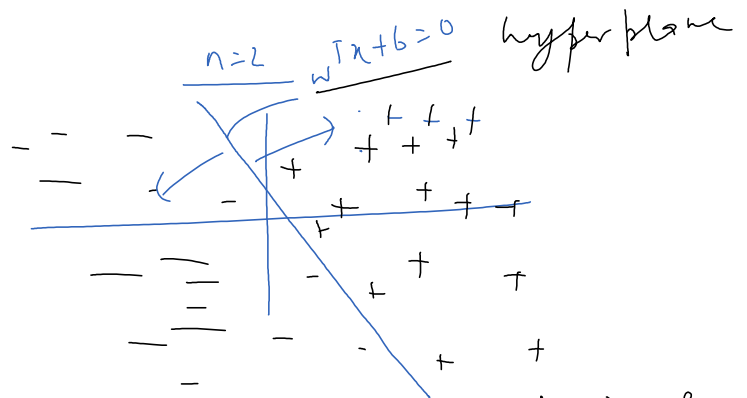
what does a neuron do

(σ, w, b)

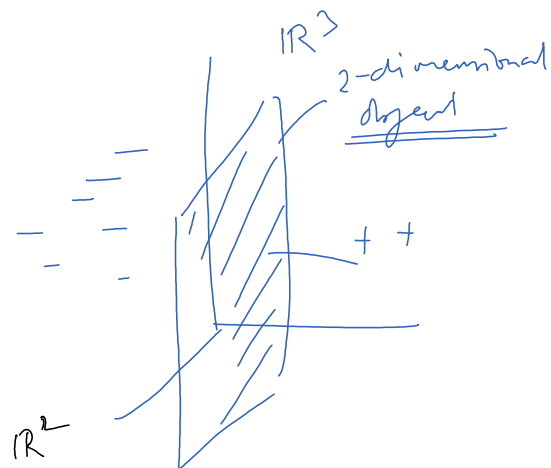
$\text{sgn}(\cdot)$

$$\sigma_{w,b}(x) = \begin{cases} 1 & \text{if } w^T x \geq -b \\ -1 & \text{if } w^T x < -b \end{cases}$$

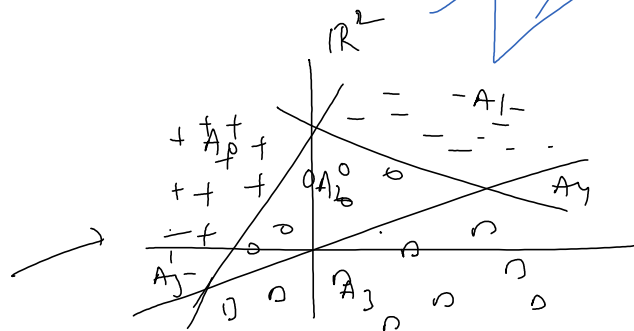
$x \in \mathbb{R}^n$



A single neuron has divided \mathbb{R}^n into two $(n-1)$ dimensional spaces by drawing a hyperplane through it.



what we want



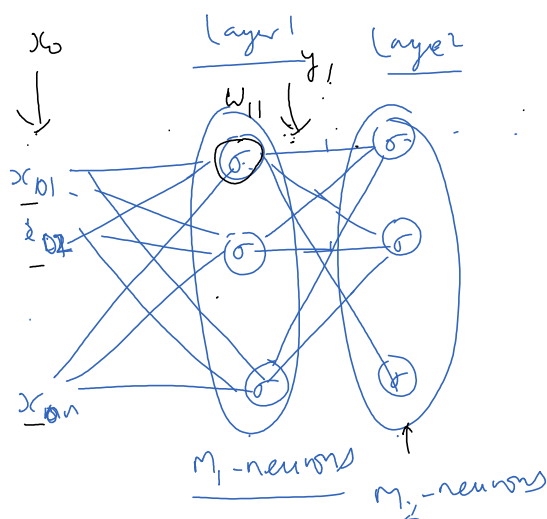
equivalent to implementing

$$f(x) = \begin{cases} 1 & \text{if } x \in A^+ \\ -1 & \text{otherwise} \end{cases} \quad \forall i = 0, 1, \dots, l$$

multilayer NNS allow us to divide n -dimensional space in more complicated ways.

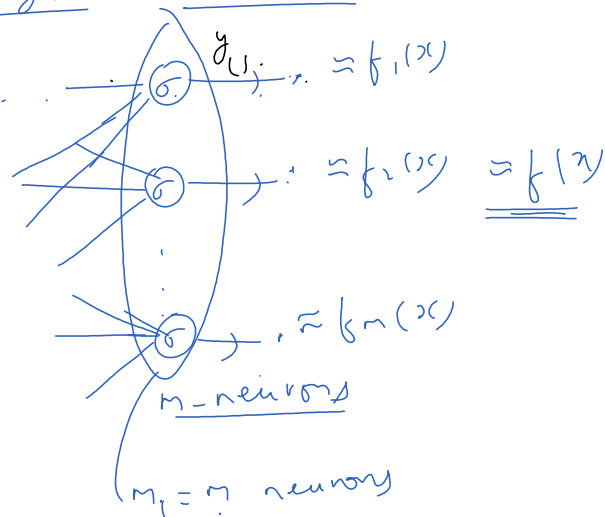
in more complicated ways.

Layers



σ, w, b
 $m_1, \dots, m_L = m$

$x \in \mathbb{R}^n$
 $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$
 - goal is: realize f .



- what is in our control

w_{11} - weight vector of the first neuron of the first layer
 b_{11} - bias of the first neuron of the first layer

- m_1, \dots, m_L
 - $w_{li}, 1 \leq l \leq L$
 $1 \leq i \leq m_l \forall l$
 $b_{li}, 1 \leq l \leq L$
 $1 \leq i \leq m_l \forall l$

we can configure these

- we assume that m_1, \dots, m_{L-1} are given

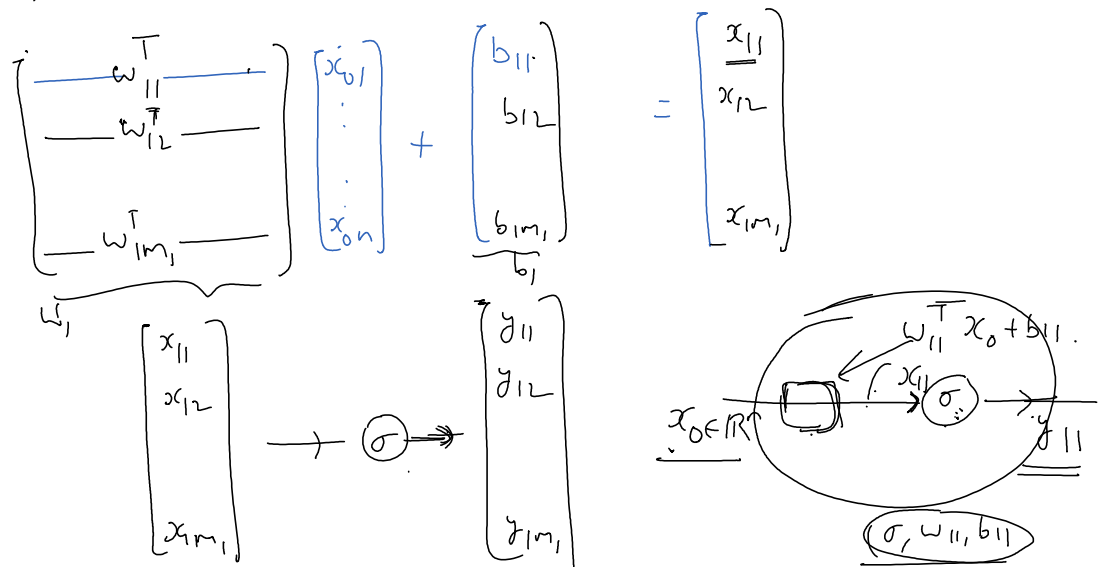
- w_{li}, b_{li} are configured using gradient descent

Neural Network implements a for $\tilde{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$

goal: $\tilde{f}(\cdot) \approx f(\cdot)$ via tuning w_{li} and b_{li}

understanding $\tilde{f}(\cdot)$

$$\tilde{f} = w_{11}^T \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} + b_{11} \begin{bmatrix} x_1 \end{bmatrix}$$



x_0 : input to the first layer

y_1 : output of the first layer

$x_1 = w_1 x_0 + b_1$
 $y_1 = \sigma(x_1)$ } mathematical representation of the first layer.

following similar steps
 $\forall 1 \leq k \leq L$

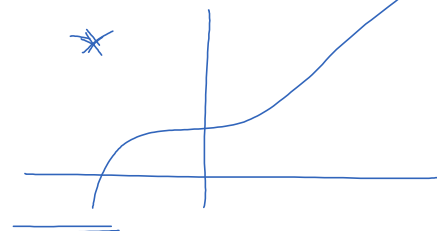
$$x_k = w_k y_{k-1} + b_k$$

$$y_k = \sigma(x_k)$$

$$\tilde{f}(x_0) = y_L$$

How are we conveyed f !

$$f(x) = \quad x \in \mathbb{R}^n$$



We are given n -input-output pairs

$$\left. \begin{array}{l} x_0^{(1)} \\ x_0^{(2)} \\ \vdots \\ x_0^{(n)} \end{array} \right\} \begin{array}{l} f(x_0^{(1)}) \\ f(x_0^{(2)}) \\ \vdots \\ f(x_0^{(n)}) \end{array} \quad \text{training data}$$

$x_0^{(n)}$ $f(x_0^{(n)})$
 * we want outputs $y_1^{(n)}, y_2^{(n)}, \dots, y_L^{(n)}$ to be close to $f(x_0^{(n)})$...
 $\dots b(x_0^{(n)})$
 $J(\underline{w_1, \dots, w_L, b_1, \dots, b_L}) = \sum_{n=1}^N \underbrace{\ell(y_L^{(n)}, f(x_0^{(n)}))}_{\text{total discrepancy}}$

$$\sum_{i=1}^N f_i(x)$$

Objective would be to minimize this

Gradient descent start with $w_1^{(0)}, \dots, w_L^{(0)}, b_1^{(0)}, \dots, b_L^{(0)}$

and

$$\begin{cases} w_{\lambda}(t+1) = w_{\lambda}(t) - \epsilon_t \nabla_{w_{\lambda}} J(w_1(t), \dots, w_L(t), b_1(t), \dots, b_L(t)) \\ b_{\lambda}(t+1) = b_{\lambda}(t) - \epsilon_t \nabla_{b_{\lambda}} J(w_1(t), \dots, w_L(t), b_1(t), \dots, b_L(t)) \end{cases}$$

- we will focus on gradient of one term
- will drop the superscript (n) for convenience.

Systematic gradient computation is called backpropagation

$$\frac{d}{dx} f(g(x)) = g'(x) f'(g(x))$$

Back. propagation

① compute $\frac{\partial \tilde{L}}{\partial y_i}$ $i=1, \dots, m_L$

② For $\ell=L, L-1, \dots, 1$ compute

for $\ell=L$ we

$$\frac{\partial \tilde{L}}{\partial w_{\ell+1, i}} = \frac{\partial \tilde{L}}{\partial y_i} \sigma'(x_{\ell+1, i}) \frac{\partial y_{\ell+1, j}}{\partial w_{\ell+1, i}}$$

$i=1, \dots, m_{\ell}$
 $j=1, \dots, m_{\ell-1}$

r^m
 $l = L$
 w_e
 can compare

$$\frac{\partial L}{\partial w_{ji}} = \frac{\partial L}{\partial y_{ji}} \sigma'(x_{ji}) w_{l-1,j}$$

$$j = 1, \dots, m_{l-1}$$

$$\frac{\partial \tilde{L}}{\partial w_{ei}}$$

$$\frac{\partial \tilde{L}}{\partial y_{ji}} \sigma'(x_{ji})$$

$$i = 1, \dots, m_l$$

$$\left\{ \frac{\partial \tilde{L}}{\partial y_{l-1,i}} \right\}$$

$$\sum_{k=1}^{m_l} \frac{\partial \tilde{L}}{\partial y_{ek}} \sigma'(x_{ek}) w_{ek,i} \quad i = 1, \dots, m_{l-1}$$

$$\frac{\partial \tilde{L}}{\partial y_{(l-1)i}}$$