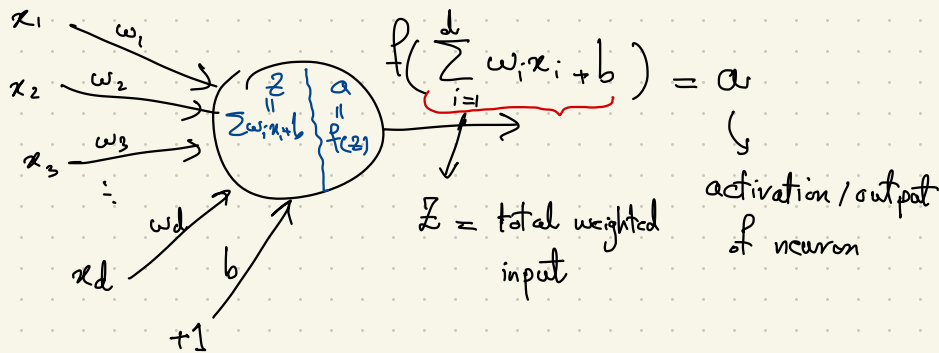


Neural Networks

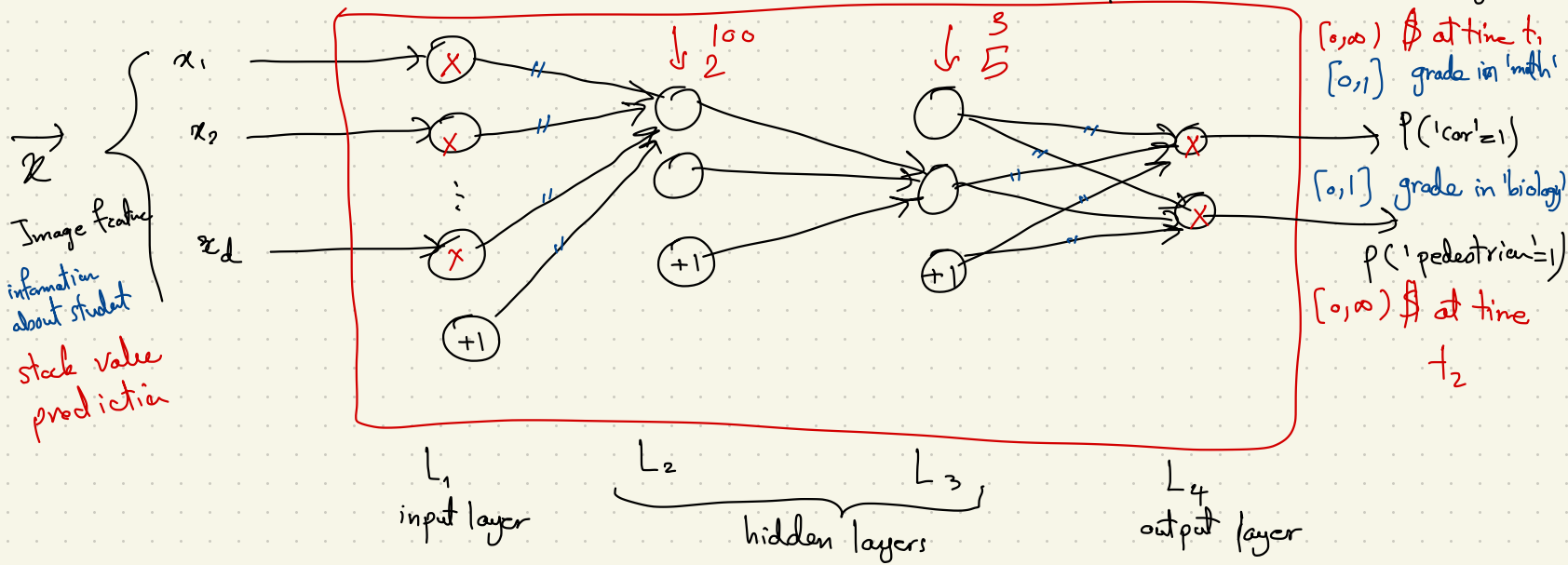
Building block : neuron



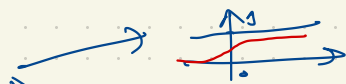
Feedforward Neural Networks :

Shallow NN \leftarrow 1 hidden

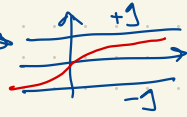
Deep NN \leftarrow ≥ 2 hidden layers



$f(\cdot)$ = activation function



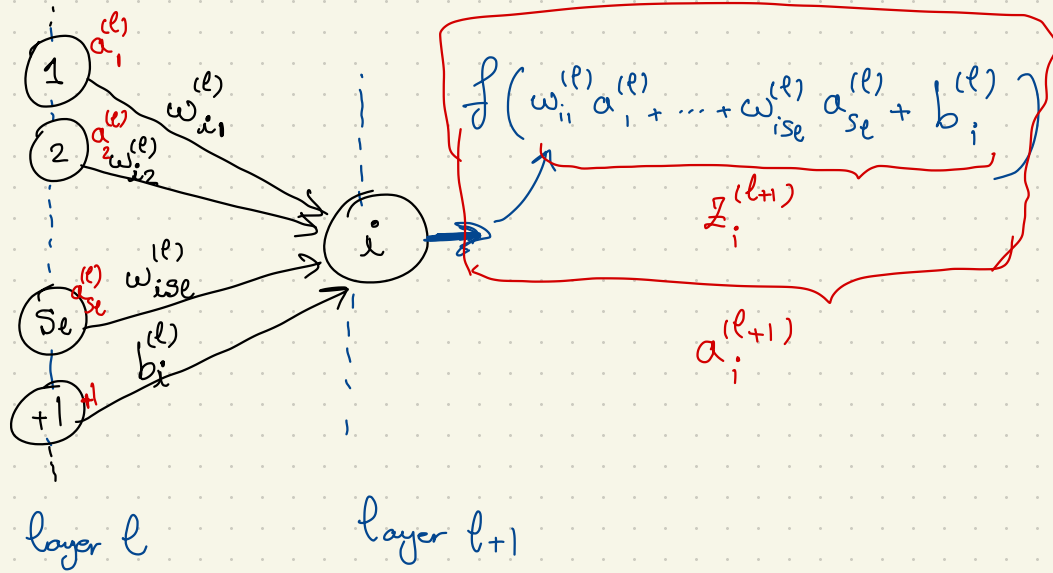
Sigmoid \Rightarrow grade $[0, 1]$

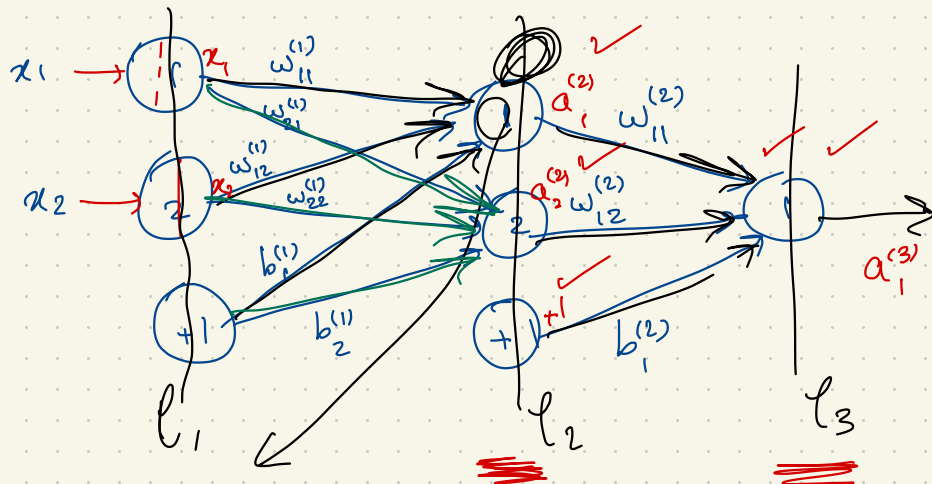


Tanh \Rightarrow class ± 1



Relu \Rightarrow stock value \$





$$z_1^{(1)} = a_1^{(1)} = x_1$$

$$z_2^{(1)} = a_2^{(1)} = x_2$$

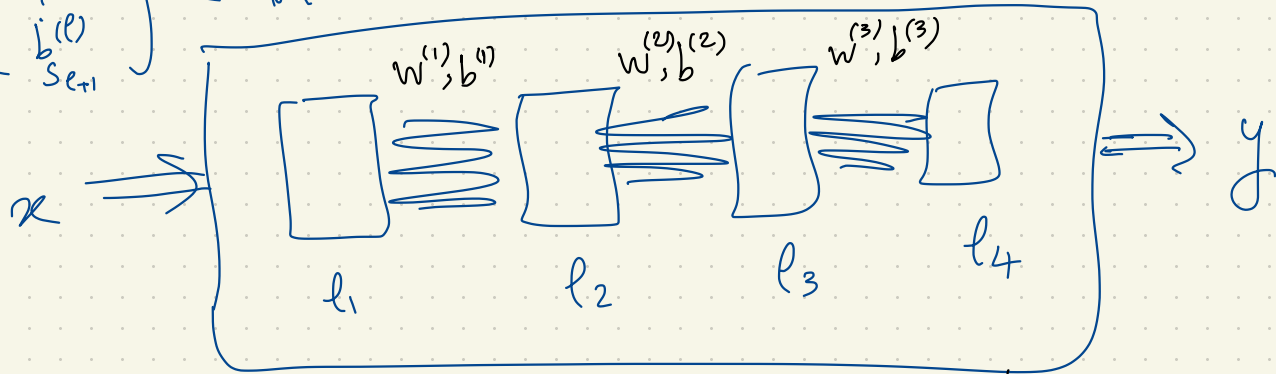
$$\begin{cases} z_1^{(2)} = w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + b_1^{(1)} \\ a_1^{(2)} = f(z_1^{(2)}) \end{cases}$$

$$\begin{cases} z_2^{(2)} = w_{21}^{(1)} x_1 + w_{22}^{(1)} x_2 + b_2^{(1)} \\ a_2^{(2)} = f(z_2^{(2)}) \end{cases}$$

$$\begin{cases} z_1^{(3)} = w_{11}^{(2)} a_1^{(2)} + w_{12}^{(2)} a_2^{(2)} + b_1^{(2)} \\ a_1^{(3)} = f(z_1^{(3)}) \end{cases}$$

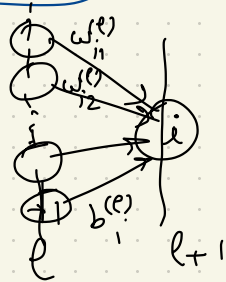
$$W^{(l)} = \begin{bmatrix} w_{11}^{(l)} & w_{12}^{(l)} & \dots & w_{1s_l}^{(l)} \\ \vdots & \vdots & & \vdots \\ w_{s_{l+1}1}^{(l)} & w_{s_{l+1}2}^{(l)} & \dots & w_{s_{l+1}s_l}^{(l)} \end{bmatrix} \in \mathbb{R}^{s_{l+1} \times s_l} \quad s_l = \# \text{ nodes in layer } l \text{ except } (+1).$$

$$b^{(l)} = \begin{bmatrix} b_1^{(l)} \\ \vdots \\ b_{s_{l+1}}^{(l)} \end{bmatrix} \in \mathbb{R}^{s_{l+1}}$$



$$z_i^{(l+1)} = \sum_{j=1}^{s_l} w_{ij}^{(l)} a_j^{(l)} + b_i^{(l)}$$

$$Z_i^{(l+1)} = \begin{bmatrix} w_{i1}^{(l)} & w_{i2}^{(l)} & \dots & w_{i s_l}^{(l)} \end{bmatrix} \begin{bmatrix} a_1^{(l)} \\ \vdots \\ a_{s_l}^{(l)} \end{bmatrix} + b_i^{(l)}$$



$$\underbrace{\begin{bmatrix} \vec{z}_1^{(l+1)} \\ \vdots \\ \vec{z}_i^{(l+1)} \\ \vdots \\ \vec{z}_{s_{l+1}}^{(l+1)} \end{bmatrix}}_{\vec{z}^{(l+1)} \in \mathbb{R}^{s_{l+1}}} = \underbrace{\begin{bmatrix} w_{11}^{(l)} & w_{12}^{(l)} & \dots & w_{1s_l}^{(l)} \\ w_{i1}^{(l)} & w_{i2}^{(l)} & \dots & w_{is_l}^{(l)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{s_{l+1}1}^{(l)} & w_{s_{l+1}2}^{(l)} & \dots & w_{s_{l+1}s_l}^{(l)} \end{bmatrix}}_{W^{(l)} \in \mathbb{R}^{s_{l+1} \times s_l}} \underbrace{\begin{bmatrix} a_1^{(l)} \\ \vdots \\ a_{s_l}^{(l)} \end{bmatrix}}_{\vec{a}^{(l)} \in \mathbb{R}^{s_l}} + \underbrace{\begin{bmatrix} b_1^{(l)} \\ \vdots \\ b_i^{(l)} \\ \vdots \\ b_{s_{l+1}}^{(l)} \end{bmatrix}}_{\vec{b}^{(l)} \in \mathbb{R}^{s_{l+1}}}$$

Forward
Propagation

$$\begin{aligned}
 \vec{z}^{(l+1)} &= \underline{W}^{(l)} \vec{a}^{(l)} + \vec{b}^{(l)} \\
 \vec{a}^{(l+1)} &= f(\vec{z}^{(l+1)}) =
 \end{aligned}$$

$\vec{a}^{(n)} = \text{output}$

$$\vec{z}^{(1)} = \vec{a}^{(1)} = \vec{x}$$

$$\begin{aligned}
 l=1 &\rightarrow \begin{pmatrix} \vec{z}^{(2)} \\ \vec{a}^{(2)} \end{pmatrix} \\
 l=2 &\rightarrow \begin{pmatrix} \vec{z}^{(3)} \\ \vec{a}^{(3)} \end{pmatrix} \\
 &\vdots
 \end{aligned}$$

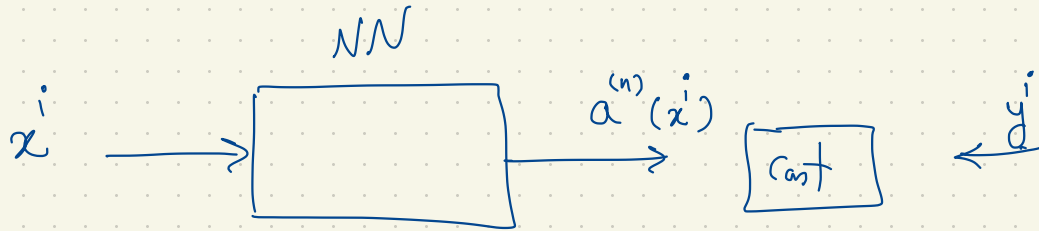
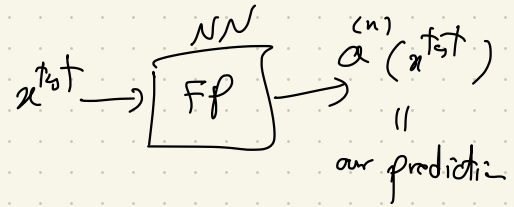
We learn weights and biases $\{w^{(l)}, b^{(l)}\}_{l=1}^{n-1}$ using training data.

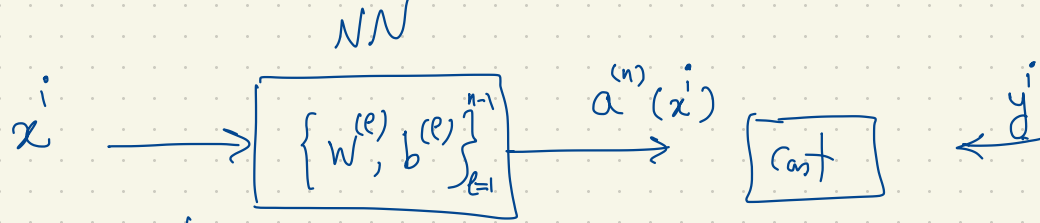
$D = \{(x^i, y^i)\}_{i=1}^N$ training data ($x^i \in \mathbb{R}^d, y^i \in \mathbb{R}^m$)

$\begin{pmatrix} \text{car} \\ \text{pedestrian} \end{pmatrix} \in \mathbb{R}^2 \quad \begin{pmatrix} 101 \\ 131 \\ 191 \end{pmatrix} \in \mathbb{R}^{10}$

Once learned, classification of a test sample x^{test}

is just applying FP to x^{test}





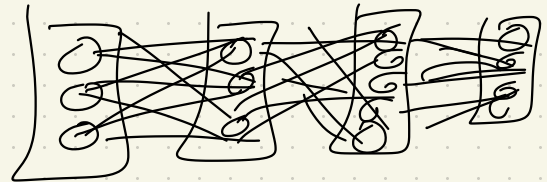
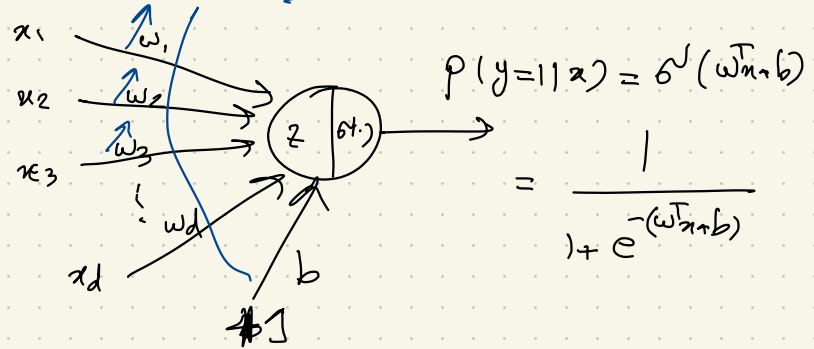
$$\min_{\{w^{(e)}, b^{(e)}\}_{e=1}^{n-1}} \frac{1}{N} \sum_{i=1}^N \text{cost}(a^{(n)}(x^i), y^i) + \lambda \underline{\text{regularizer}}$$

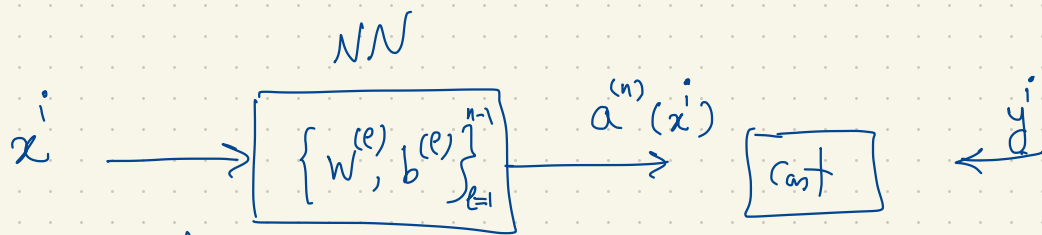
Logistic regression $(+ \lambda \|w\|_2^2)$

Simplest NN with only

1 neuron whose

activation is sigmoid func.





$$\min_{\{w^{(e)}, b^{(e)}\}_{e=1}^{n-1}} \frac{1}{N} \sum_{i=1}^N \text{cost}(a^{(n)}(x^i), y^i) + \lambda \text{ regularizer}$$

$$\min_{\{w^{(e)}, b^{(e)}\}_{e=1}^{n-1}} \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \| a^{(n)}(x^i) - y^i \|_2^2 + \frac{\lambda}{2} \sum_{e=1}^{n-1} \| w^{(e)} \|_F^2$$

Frobenius norm

$$\| A \|_F^2 \triangleq \sum_{i=1}^m \sum_{j=1}^n a_{ij}^2$$

$m \times n$

↓
train a NN

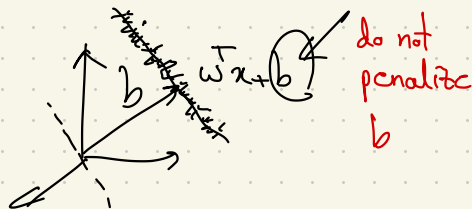
$$\min_{\{w^{(e)}, b^{(e)}\}_{e=1}^{n-1}} \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \|a^{(n)}(x^i) - y^i\|_2^2 + \frac{\lambda}{2} \sum_{e=1}^{n-1} \|w^{(e)}\|_F^2$$

cost

Solve via G.D.

$$\frac{\partial \text{cost}}{\partial w^{(e)}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \frac{1}{2} \|a^{(n)}(x^i) - y^i\|_2^2}{\partial w^{(e)}} + \lambda w^{(e)}$$

$$\frac{\partial \text{cost}}{\partial b^{(e)}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \frac{1}{2} \|a^{(n)}(x^i) - y^i\|_2^2}{\partial b^{(e)}}$$



* Initialize $\{w^{(l)}, b^{(l)}\}_{l=1}^{n-1}$ (e.g. $\mathcal{N}(0, \epsilon^2)$, $\epsilon \approx 0.01$)

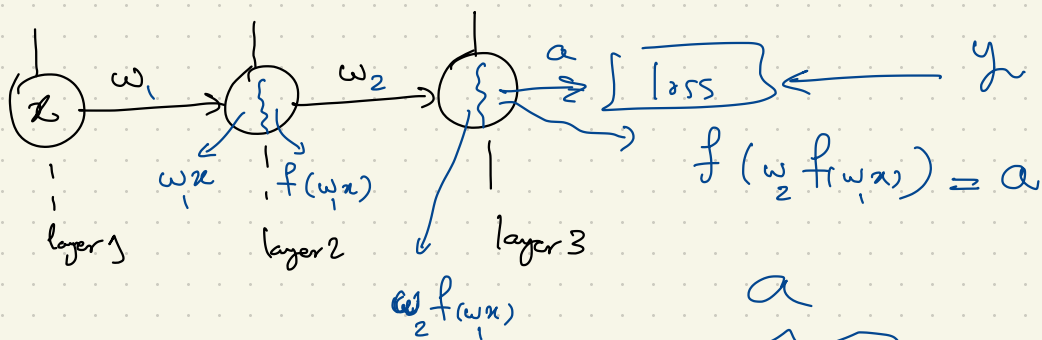
* For $t=1, 2, \dots, T$ (until convergence)

$$\begin{cases} w^{(l)} \leftarrow w^{(l)} - \eta \frac{\partial \text{cost}}{\partial w^{(l)}} \\ b^{(l)} \leftarrow b^{(l)} - \eta \frac{\partial \text{cost}}{\partial b^{(l)}} \end{cases} \quad \forall l=1, \dots, n-1$$

$\{w^{(l)}, b^{(l)}\}$

NWs have a mechanism to efficiently compute derivative wrt weights and biases : Back propagation

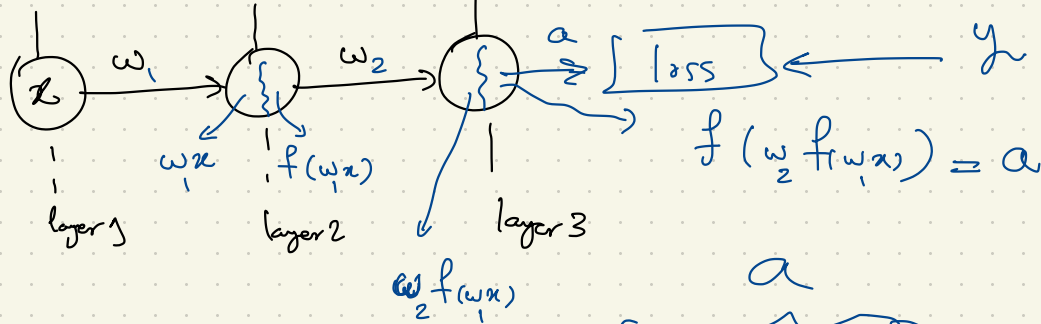
$$\left\{ \begin{array}{c} \frac{\partial \text{cost}}{\partial w^{(n-1)}} \\ \frac{\partial \text{cost}}{\partial b^{(n-1)}} \end{array} \right\} \longrightarrow \begin{array}{c} \frac{\partial \text{cost}}{\partial w^{(n-2)}} \\ \frac{\partial \text{cost}}{\partial b^{(n-2)}} \end{array} \longrightarrow \dots \longrightarrow \begin{array}{c} \frac{\partial \text{cost}}{\partial w^{(1)}} \\ \frac{\partial \text{cost}}{\partial b^{(1)}} \end{array} \right\}$$



$$\text{Cost} = \frac{1}{2} (a - y)^2 = \frac{1}{2} \left(\underbrace{f(w_2 f(w_1 x))}_a - y \right)^2$$

To compute $\frac{\partial \text{Cost}}{\partial w_1 \text{ or } w_2}$ need to compute $\frac{\partial a}{\partial w_1 \text{ or } w_2}$

$$\begin{aligned} \frac{\partial \text{Cost}}{\partial w_1} &= \frac{\partial \text{Cost}}{\partial a} \times \frac{\partial a}{\partial w_1} \\ &= \frac{\partial \frac{1}{2} (a - y)^2}{\partial a} \times \frac{\partial a}{\partial w_1} = (a - y) \times \frac{\partial a}{\partial w_1} \end{aligned}$$



$$\text{Cost} = \frac{1}{2} (a - y)^2 = \frac{1}{2} \left(\underbrace{f(w_2 f(w_1 x))}_a - y \right)^2$$

$$\frac{\partial a}{\partial w_2} = \frac{\partial f(w_2 f(w_1 x))}{\partial w_2} = a f'(w_2 a)$$

$$\begin{aligned} \frac{\partial a}{\partial w_1} &= \frac{\partial a}{\partial (w_2 f(w_1 x))} \times \frac{\partial (w_2 f(w_1 x))}{\partial w_1} \\ &= \frac{\partial a}{\partial (w_2 f(w_1 x))} \times f(w_1 x) \times f'(w_2 f(w_1 x)) \\ &= f'(w_2 f(w_1 x)) \times w_2 \times \frac{\partial f(w_1 x)}{\partial w_1} \times f'(w_1 x) \end{aligned}$$

Pytorch :

• Backward ()

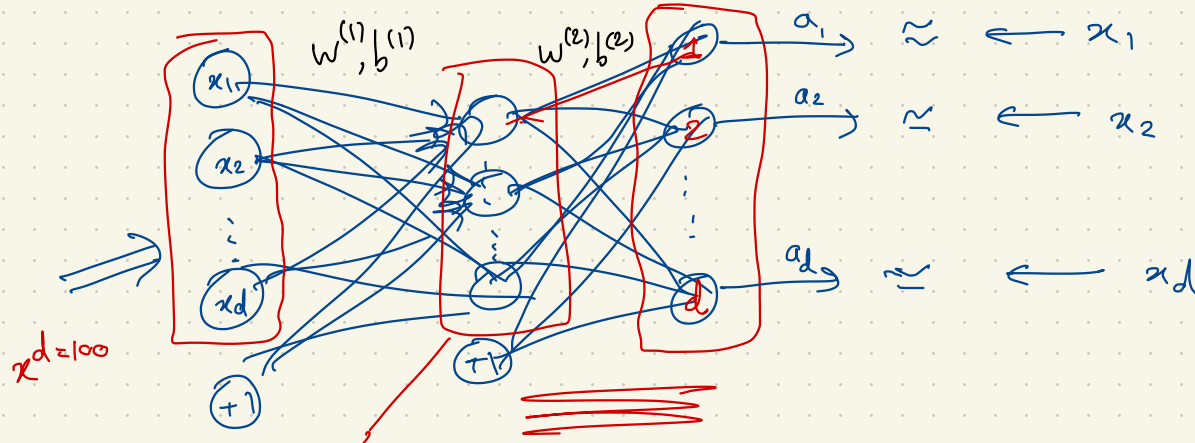


classification
or
regression

$f(z) =$  ReLU  Identity \Rightarrow regression

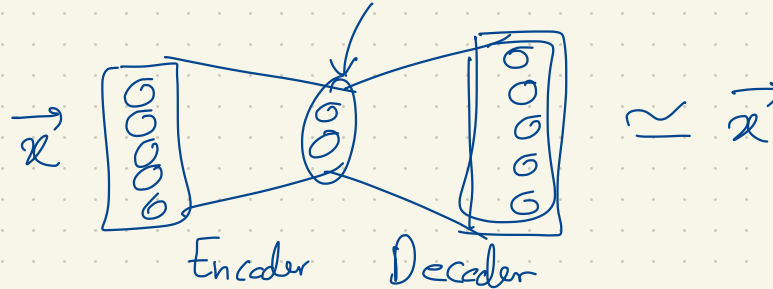
$f(z) =$ Sigmoid / tanh $\Rightarrow p(\text{'car' = 1} | x) \in [0, 1]$
[0, 1] [-1, 1] $\in [-1, 1]$
Classification

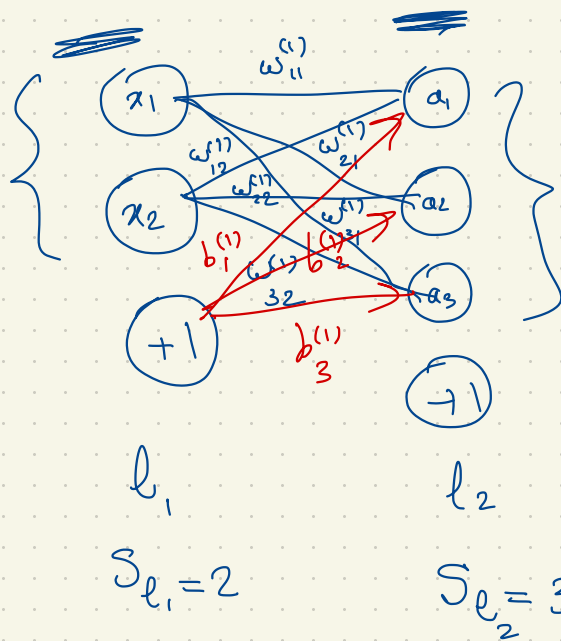
NNs can be used to reduce data dimension (dimenality reduction)



Auto encoder NN

nodes = 5





$$S_{l_1} = 2$$

$$S_{l_2} = 3$$

$$W^{(1)} = \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} \\ w_{31}^{(1)} & w_{32}^{(1)} \end{bmatrix}$$

$$3 \times 2$$

$$S_{l_2} \times S_{l_1}$$

$$b^{(1)} = \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \end{bmatrix}$$

3×1

S_{l_2}

$$\begin{cases} w^{(l)} \in \mathbb{R}^{S_{l+1} \times S_l} \\ b^{(l)} \in \mathbb{R}^{S_{l+1}} \end{cases}$$

How to select architecture:

