$\mathcal{L}(\bar{y}, y)$ : loss
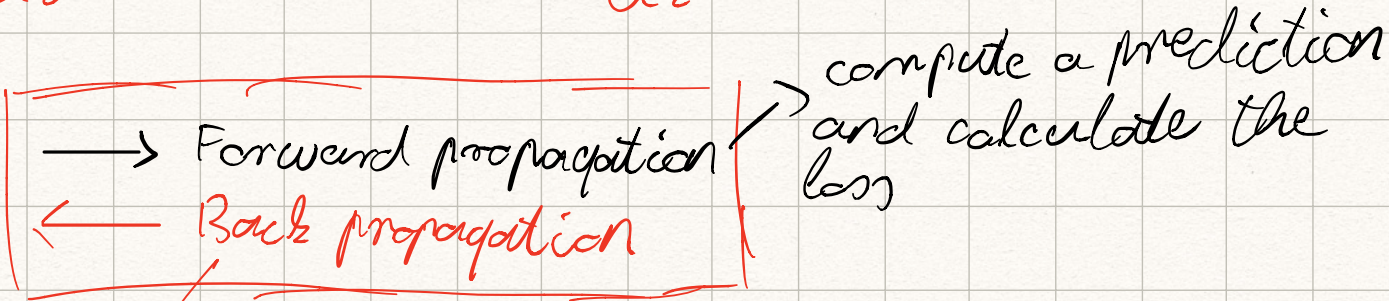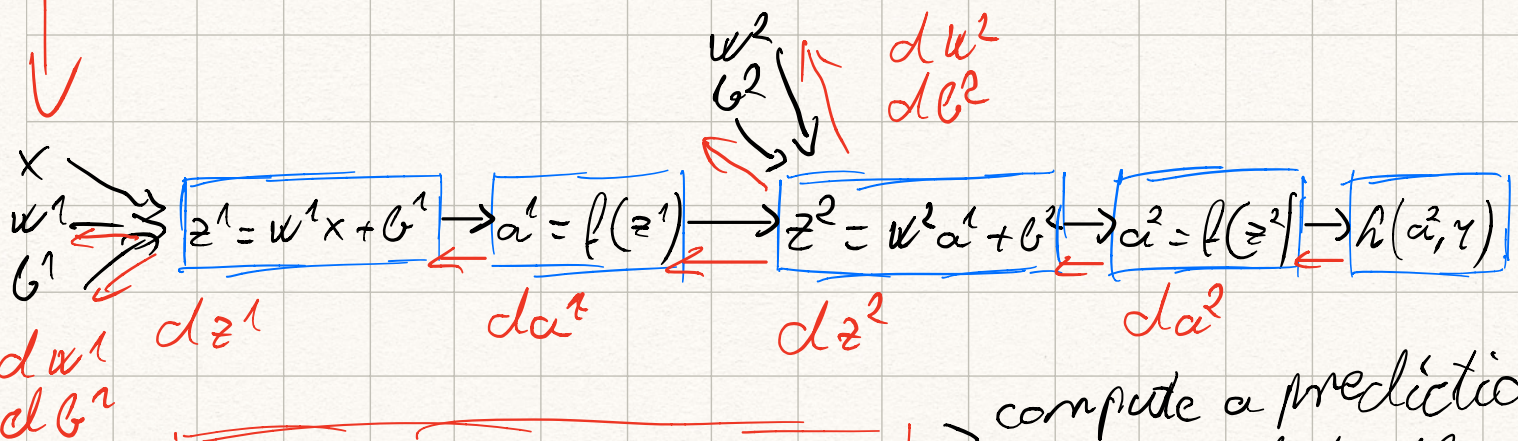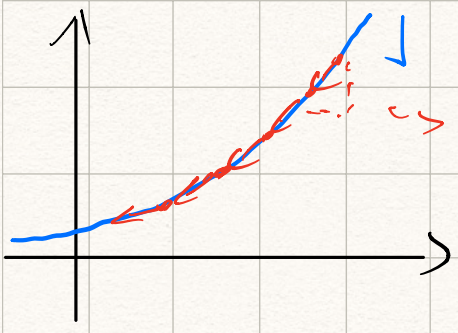
○ Training a neural network is an optimization problem. The goal is to minimize the loss by update the values of the weight matrices and biases.

$$x, w^1, 6^1 \Rightarrow \boxed{z^1 = w^1 x + 6^1} \rightarrow \boxed{a^1 = f(z^1)} \rightarrow \boxed{z^2 = w^2 a^1 + 6^2} \rightarrow \boxed{a^2 = f(z^2)} \rightarrow \boxed{\mathcal{L}(a^2, y)}$$

$w^2, 6^2 \quad dw^2 \quad d6^2$

$dw^1 \quad dz^1 \quad da^1 \quad dz^2 \quad da^2$
$d6^2$

$\longrightarrow$ Forward propagation
$\longleftarrow$ Back propagation

$\Big\}$ compute a prediction and calculate the loss

Back propagate the error through the gradients to update the weights and biases.

Cost function:

$$J(w^1, 6^1, w^2, 6^2) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\bar{y}, y)$$



$\dashrightarrow$ derivatives

$\Rightarrow \min(J(w, 6))$

Repeat until convergence {

compute forward predictions $( \overline{y}^{(i)} , i = 1 \dots n )$

$\underline{d w^1} = \frac{\partial J}{\partial w^1}$ , $dG^1 = \frac{\partial J}{\partial G^1}$ / $w^2, G^2$ is the same

$w^1 := w^1 - h \cdot d w^1$
$G^1 := G^1 - h \cdot dG^1$ / $w^2, G^2$ is the same

}

$\rightarrow$ learning rate

compute these derivative terms with

Backpropagation

$\underline{d z^2} = a^2 - Y$ , $Y = [y^1, y^2 \dots y^n]$ labeled data

$\underline{d w^2} = \frac{1}{m} d z^2 \cdot a^{1T}$

$\underline{dG^2} = \frac{1}{m} \text{sum}(d z^2, axis = 1)$ $(n, 1)$ array

$\underline{d z^1} = w^{2T} \cdot d z^2 * \ell^{1'}(z^1)$  derivative of the activation function

$\rightarrow$ element wise

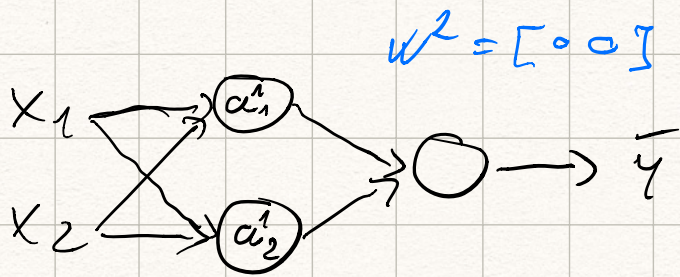$\underline{d w^1} = \frac{1}{m} d z^1 \cdot x^T$
$\rightarrow a^T$

$\underline{dG^1} = \frac{1}{m} \text{sum}(d z^1, axis = 1)$

It is based on the chain rule:

$F(x) = \ell(g(x)) \Rightarrow F'(x) = \ell'(g(x)) \cdot g'(x)$

# Random initialization

$$w^2 = [\;0 \;\; 0\;]$$



$$w^1 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \Rightarrow a_1^1 = a_2^1 \Rightarrow dz_1^1 = dz_2^1$$

Both hidden units compute the same function.
⇒ Random weight matrices initialization