# RECAP SOFAR
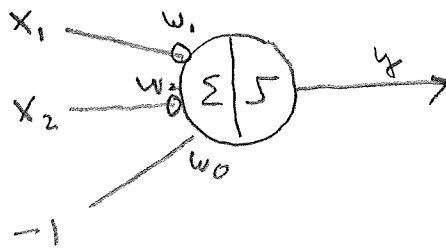
* model of a single neuron



$$y = logsig(w_1 x_1 + w_2 x_2 - w_0)$$
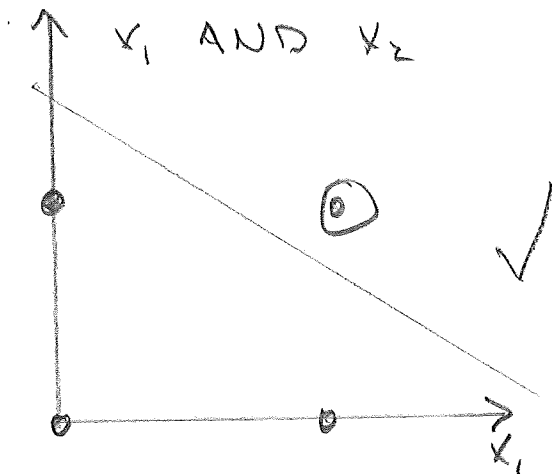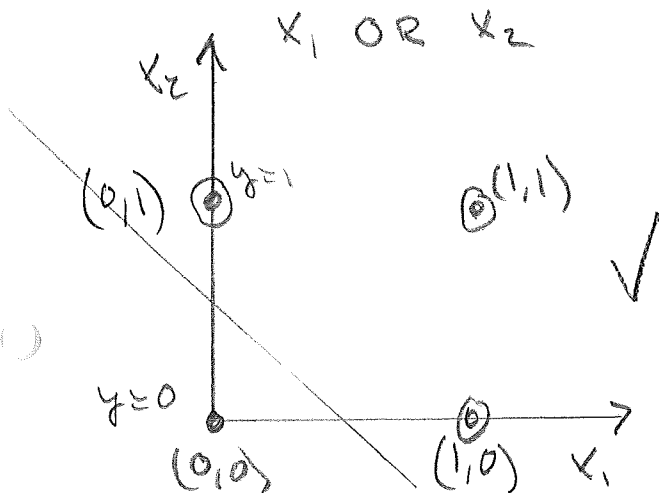
$$class = \begin{cases} 1 & \text{if } y = 0 \ (<0.5) \\ 2 & \text{if } y = 1 \ (>0.5) \end{cases}$$

* Training using gradient descent (GD)

$$w_i^{t+1} = w_i^t - \mu \frac{\partial \lambda_{MSE}}{\partial w_i} \qquad \lambda_{MSE} : loss$$
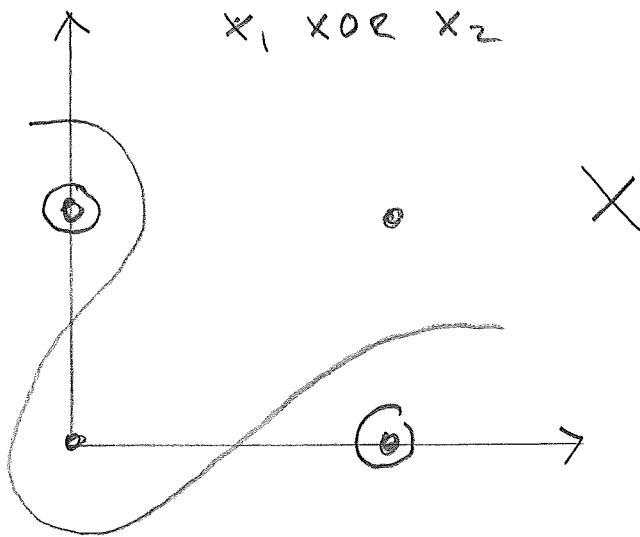
## SINGLE NEURON EXAMPLES



$X_1$ OR $X_2$

$X_1$ AND $X_2$

$$logsig(-10) \sim 0, \ logsig(+10) \sim 0$$

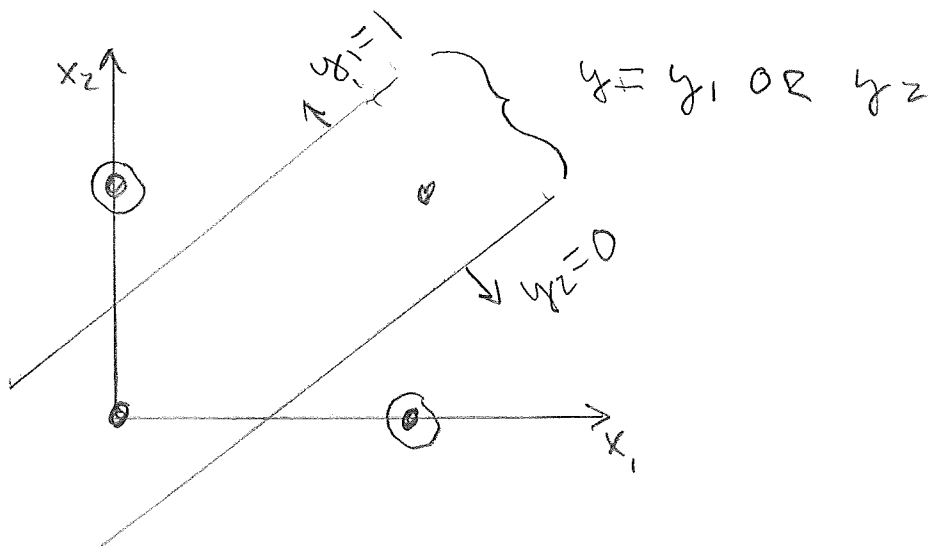$$\underbrace{w_1 x_1}_{30} + \underbrace{w_2 x}_{30} - \underbrace{w_0}_{20}$$

<NOTEBOOK>

$$\underbrace{w_1 x_1}_{15} + \underbrace{w_2 x_2}_{15} - \underbrace{w_0}_{20}$$

<NOTEBOOK>

$X_1$ XOR $X_2$



$\Rightarrow$ The First AI Winter

## Neural Network (multi-layer Perceptron)



$y = y_1$ OR $y_2$

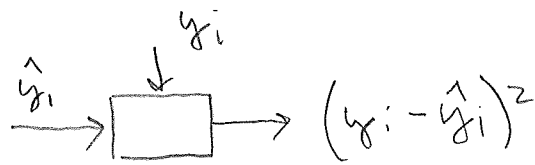$y_1:$ $\underbrace{W_{11}X_1}_{-20} + \underbrace{W_{12}X_2}_{30} - \underbrace{W_{10}}_{20}$

$y_2:$ $\underbrace{W_{21}X_1}_{30} + \underbrace{W_{22}X_2}_{-20} \div \underbrace{W_{20}}_{20}$

$y:$ $\underbrace{W_1 y_1}_{30} + \underbrace{W_2 y_2}_{30} - \underbrace{W_0}_{20}$

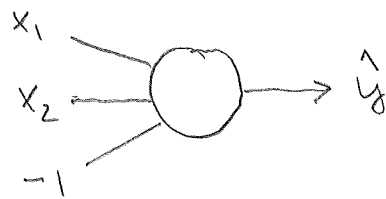〈NOTEBOOK〉

# Gradient Descent for MLP

## Loss function is minimized



Forward: $\mathcal{L}_{MSE} = (y_i - \hat{y}_i)^2$

Backward (to input): $\dfrac{\partial \mathcal{L}_{MSE}}{\partial \hat{y}} = 2(y - \hat{y}) \cdot -1$

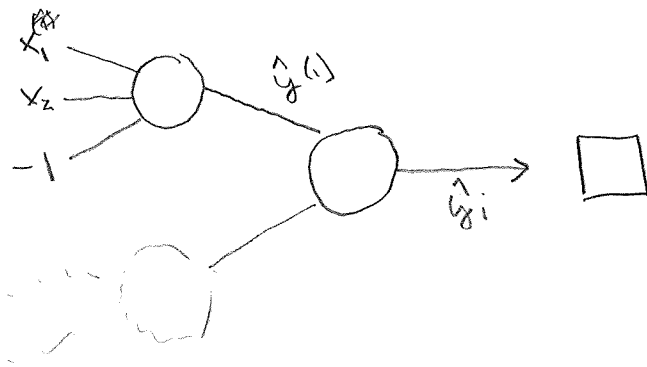$\dfrac{\partial \, y - \hat{y}}{\partial \hat{y}} = -1$

## Output neuron



Forward: $\hat{y} = logsig(w_1 x_1 + w_2 x_2 - w_0)$

Backward: $\dfrac{\partial \mathcal{L}_{MSE}}{\partial x_i} = \underbrace{\dfrac{\partial \mathcal{L}_{MSE}}{\partial \hat{y}}}_{previous} \cdot \dfrac{\partial \hat{y}}{\partial x_i} = \dfrac{\partial \mathcal{L}_{MSE}}{\partial \hat{y}} \cdot logsig(...)\,(1 - logsig(...))$

$\cdot \underbrace{w_i}_{-1 \; for \; i=0}$

update: $w_i^{(t+1)} = w_i^{(t)} - \mu \dfrac{\partial \mathcal{L}_{MSE}}{\partial w_i}$

$= w_i^{(t)} - \mu \cdot \dfrac{\partial \mathcal{L}_{MSE}}{\partial \hat{y}} \cdot logsig(...)\,(1 - logsig(...))$

$\cdot x_i \; (-1 \; for \; i=0)$

# hidden layer neuron



Forward: $\hat{y}^{(1)} = \text{logsig}\left(w_1^{(1)} x_1 + w_2^{(1)} x_2 - w_0\right)$

Backward: $\dfrac{\partial \mathcal{L}_{MSE}}{\partial x_i} = \underbrace{\dfrac{\partial \mathcal{L}_{MSE}}{\partial \hat{y}} \cdot \dfrac{\partial \hat{y}}{\partial \hat{y}^{(1)}}}_{\text{previous}} \cdot \dfrac{\partial \hat{y}^{(1)}}{\partial x_i}$

Update: $\dfrac{\partial \mathcal{L}_{MSE}}{\partial w_i^{(1)}} = \dfrac{\partial \mathcal{L}_{MSE}}{\partial \hat{y}} \cdot \dfrac{\partial \hat{y}}{\partial \hat{y}^{(1)}} \cdot \dfrac{\partial \hat{y}^{(1)}}{\partial w_i^{(1)}}$

$$w_i^{(t+1),(1)} = w_i^{(t+1),(1)} - \mu \dfrac{\partial \mathcal{L}_{MSE}}{\partial w_i^{(1)}}$$

$\Rightarrow$ No matter how many layers and neurons all computation is local:

1. Forward pass: compute output for given inputs

2. Backward pass: compute gradient, concatenate w/ received gradient and pass forward

3. Update: update weights given the current gradient

$$\frac{\partial f(g(h(x)))}{\partial x} = f'(g(h(x))) \cdot \frac{\partial g(h(x))}{\partial x}$$

$$= f'(g(h(x))) \cdot g'(h(x)) \cdot h'(x)$$

$f'$: loss bw

$g(h(x))$: output fw

$g'$: output bw

$h(x)$: hidden fw

$h'$: hidden bw

$x$: input fw



$h_1(x_1, x_2)$

$x_1$

$-1$

$x_2$

$h_2(x_1, x_2)$

$-1$

$g(h_1, h_2)$

$y$

$f(g(\cdot), y)$