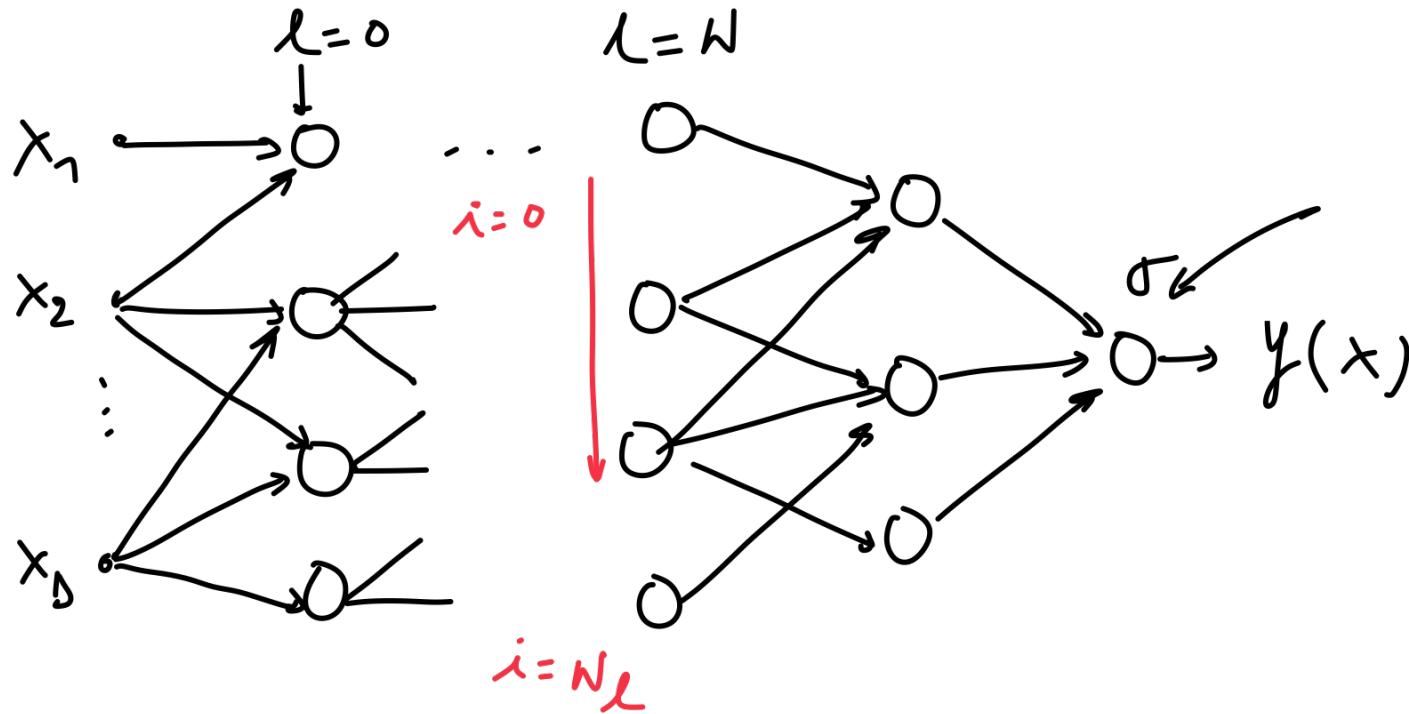


Réseau de Neurones



→ Apprentissage du réseau via Maximum de vraisemblance

$$p(\{t(x^{(i)}) = t^{(i)}\}_{i=1}^N | \{x^{(i)}\}) = \prod_{i=1}^N p(t(x^{(i)}) = t^{(i)} | x^{(i)})$$

Cette fois-ci c'est le réseau lui-même qui est utilisé pour représenter $p(t(x^{(i)}) = t^{(i)} | x^{(i)})$

$$p(t(x^{(i)}) = 1 | x^{(i)}) = y(x^{(i)})$$

$$p(t(x^{(i)}) = 0 | x^{(i)}) = 1 - y(x^{(i)})$$

$$p(t(x^{(i)}) = t^{(i)} | x^{(i)}) = y(x^{(i)})^{t^{(i)}} (1 - y(x^{(i)}))^{1 - t^{(i)}}$$

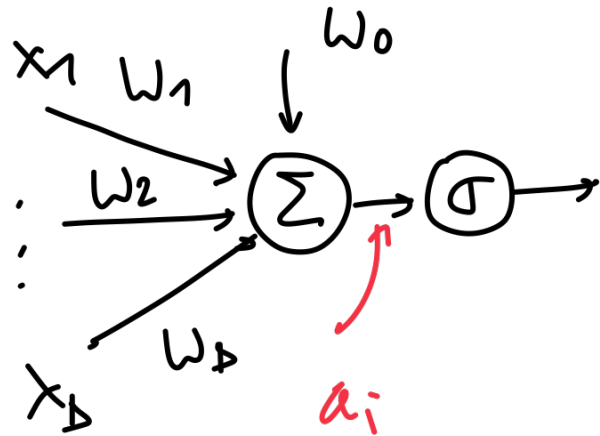
→ On peut en fait entraîner le réseau en maximisant la log vraisemblance,

$$\max \log \prod_{i=1}^N p(t(x^{(i)}) = t^{(i)} | x^{(i)})$$

$$\max \log \prod_{i=1}^N y(x^{(i)}) t^{(i)} (1 - y(x^{(i)}))^{1 - t^{(i)}}$$

$$\min J = - \sum_{i=1}^N t^{(i)} \log \underbrace{y(x^{(i)})}_{\sigma(a_{out})} + (1 - t^{(i)}) \log \underbrace{(1 - y(x^{(i)}))}_{\sigma(a_{out})}$$

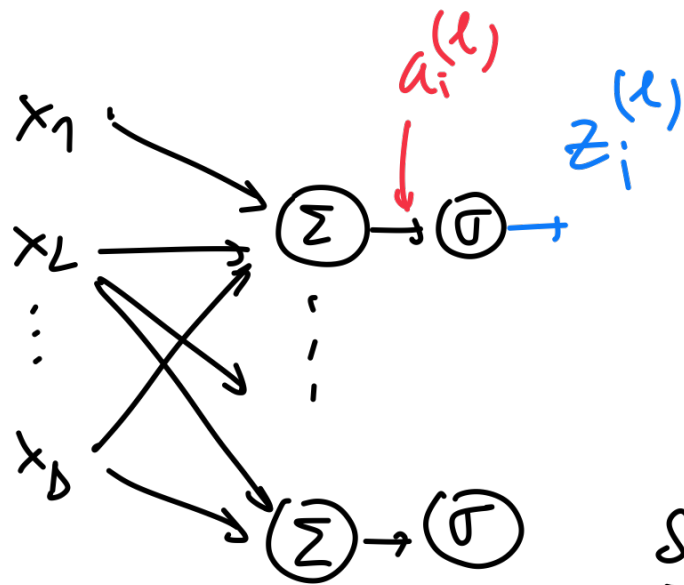
ENTROPIE BINAIRE CROISÉE = J



$a_i^{(l)}$ = practivation associé au neurone i de la couche l

$$\underbrace{a_i^{(l)}} = \sum_{j=1}^{N_l} \underbrace{w_{ij}^{(l)}} \underbrace{z_j^{(l-1)}} + w_{i0}^{(l)}$$

$$\underbrace{z_j^{(l-1)}} = \underbrace{\sigma(a_j^{(l-1)})}$$



→ Objectif : obtenir les gradients

$$\frac{\partial J}{\partial w_{ij}^{(l)}}$$

Step 1
$$\frac{\partial J}{\partial w_{ij}^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} \cdot \frac{\partial a_i^{(l)}}{\partial w_{ij}^{(l)}}$$

$$= \delta_i^{(l)} z_j^{(l-1)}$$

gradient
par rapport
aux paramètres
du neurone i de
la couche l

Step 2
$$\delta_{out} = \frac{\partial J}{\partial a_{out}}$$

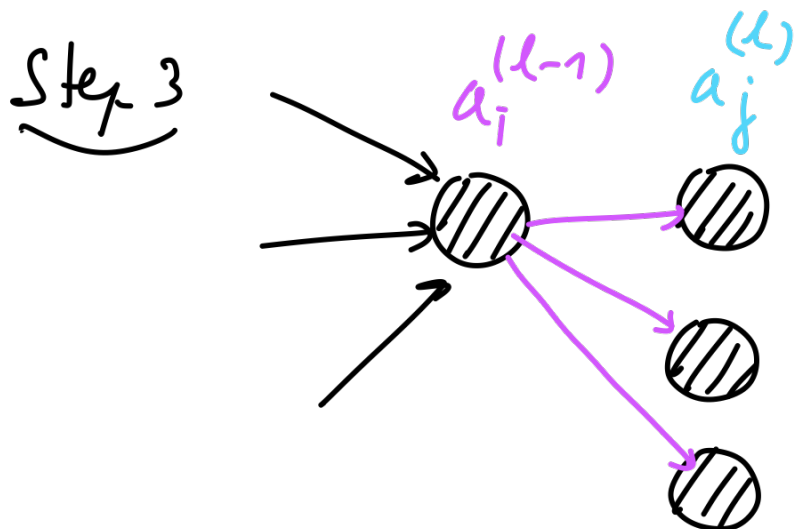
$$\delta_{out} = \frac{\partial}{\partial a_{out}} \left\{ -t^{(l)} \log(\sigma(a_{out})) - (1-t^{(l)}) \log(1-\sigma(a_{out})) \right\}$$

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

$$\delta_{out} = -t^{(L)} \frac{\sigma'(a_{out})}{\sigma(a_{out})} - (1 - t^{(L)}) \frac{(-\sigma'(a_{out}))}{1 - \sigma(a_{out})}$$

$$\sigma'(x) = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1}{1 + e^{-x}} \right) \left(1 - \frac{1}{1 + e^{-x}} \right) = \sigma(x)(1 - \sigma(x))$$

$$\begin{aligned} \delta_{out} &= -t^{(L)} (1 - \sigma(a_{out})) + (1 - t^{(L)}) \sigma(a_{out}) \\ &= \sigma(a_{out}) - t^{(L)} \end{aligned}$$



Chaque $a_i^{(l-1)}$ est utilisé par tous les neurones de la couche suivante (couche l)

$$\frac{\partial J}{\partial a_i^{(l-1)}} = \delta_i^{(l-1)} = \sum_{j=1}^{N_l} \underbrace{\frac{\partial J}{\partial a_j^{(l)}}}_{\delta_j^{(l)}} \cdot \underbrace{\frac{\partial a_j^{(l)}}{\partial a_i^{(l-1)}}}_{w_{ji}^{(l)} \sigma'(a_i^{(l-1)})}$$

pour le second facteur, en utilisant la définition des préactivations

$$a_j^{(l)} = \sum_{i=1}^{N_l} w_{ji}^{(l)} \sigma(a_i^{(l-1)}) + w_{j0}^{(l)}$$

Step 3: à partir des $\delta_i^{(l)}$ obtenir les $\delta_i^{(l-1)}$ via

$$\delta_i^{(l-1)} = \sum_{j=1}^{N_l} \delta_j^{(l)} w_{ji}^{(l)} \sigma'(a_i^{(l-1)}) \quad (*)$$

Back propagation: (démarer avec $w_{ij}^{(l)}$ aléatoires)

- Pour un exemple $\{x^{(i)}, t^{(i)}\}$, propager l'exemple à travers tout le réseau et générer toutes les valeurs de pré et post activations $a_i^{(l)}, z_j^{(l)}$ pour toute couche l et tout neurone i

→ Calculer $\delta_{out} = y(x^{(i)}) - t^{(i)}$

→ Dériver tous les $\delta_i^{(l)}$ à partir de δ_{out} via (*),

→ Dériver les gradients via

$$\frac{\partial J}{\partial w_{ij}^{(l)}} = \delta_i^{(l)} z_j^{(l-1)}$$

→ Réaliser un pas dans la direction opposée au gradient

$$w_{ij}^{(L)} \leftarrow w_{ij}^{(L)} - \eta \frac{\partial J}{\partial w_{ij}^{(L)}}$$