

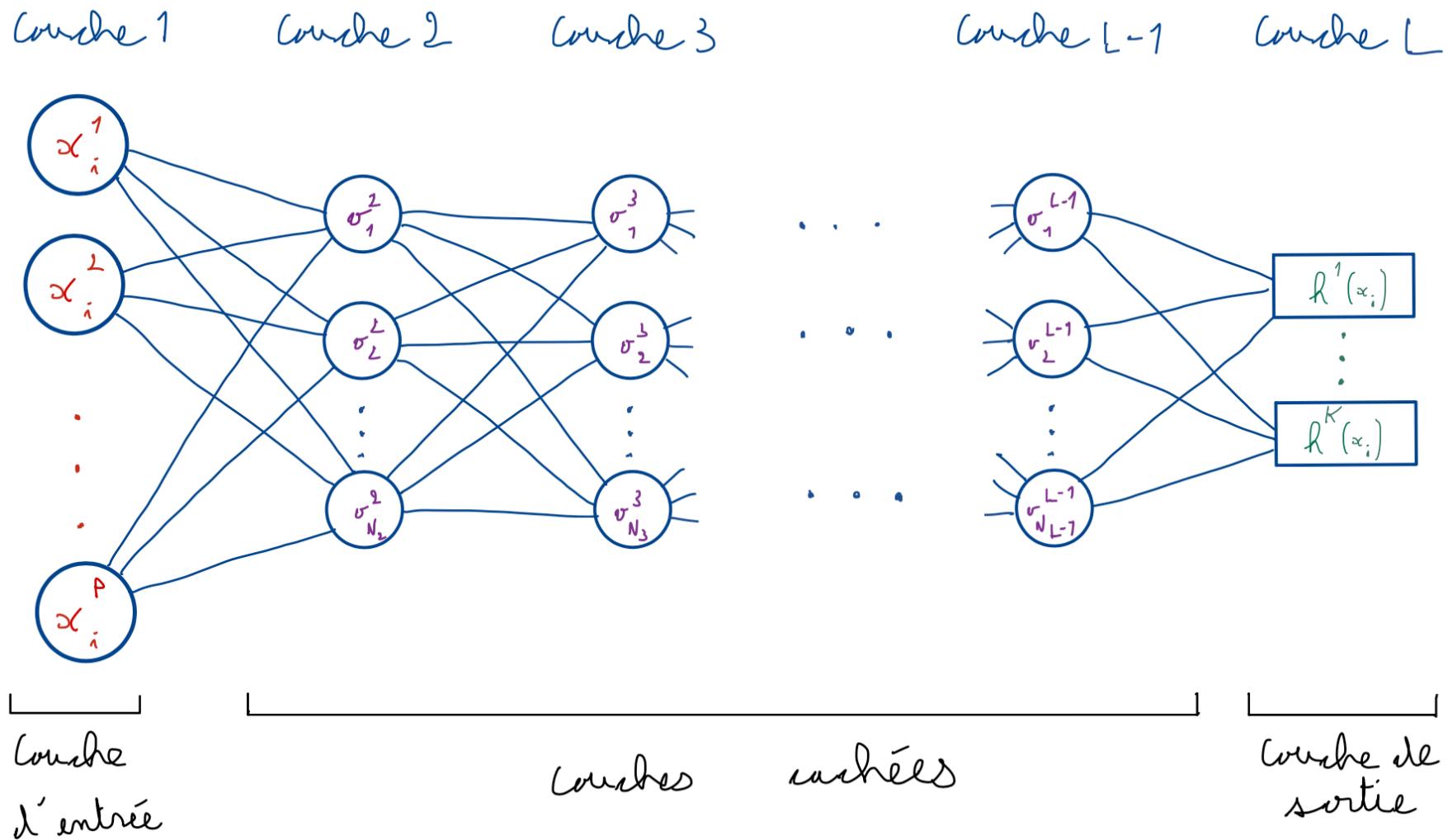
Introduction au calcul GP-GPU

Partie 3 : Application aux réseaux de neurones convolutifs

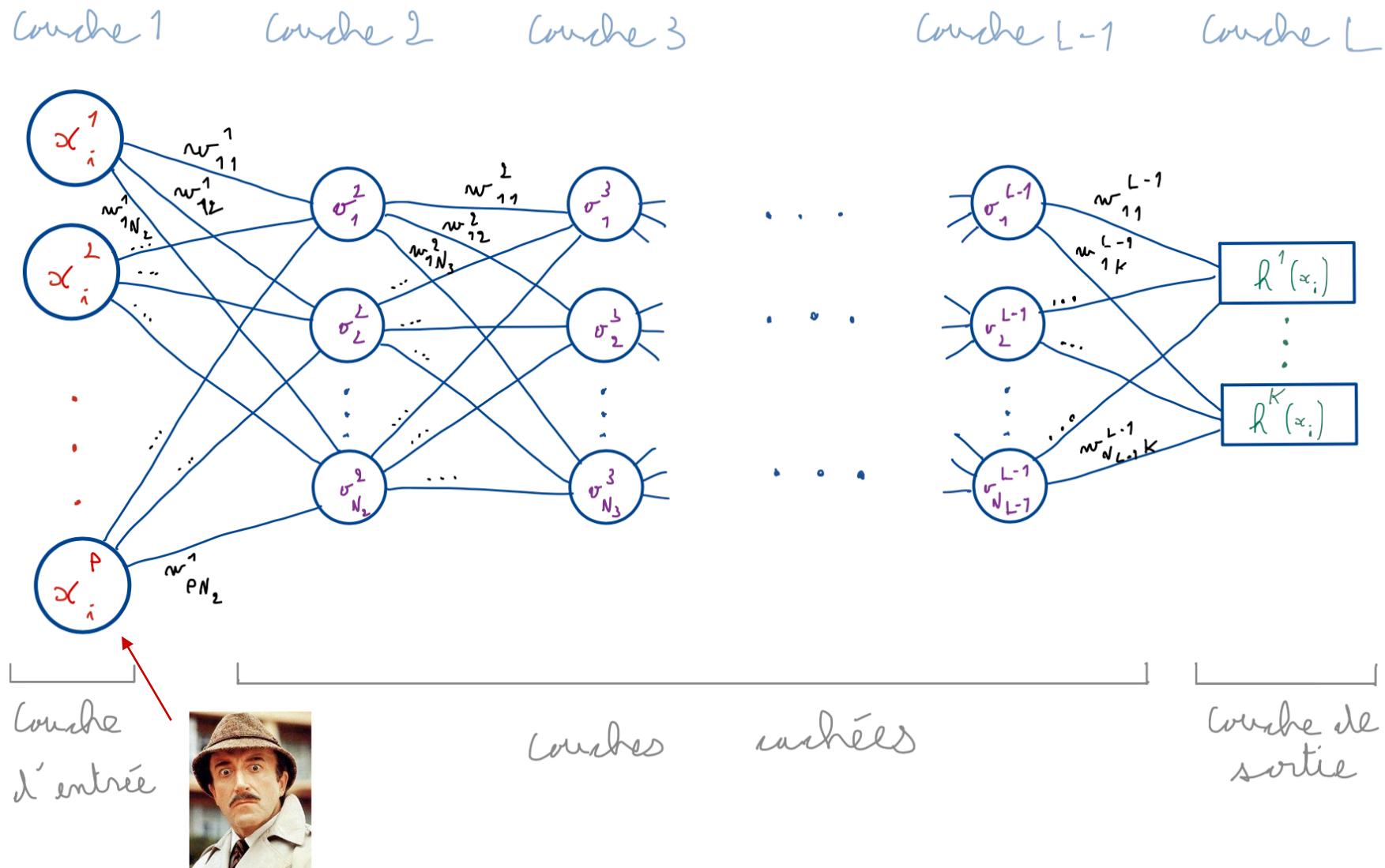
Laurent Risser

Institut de mathématiques de Toulouse
lrisser@math.univ-toulouse.fr

1) Calcul d'un gradient dans un réseau de neurones



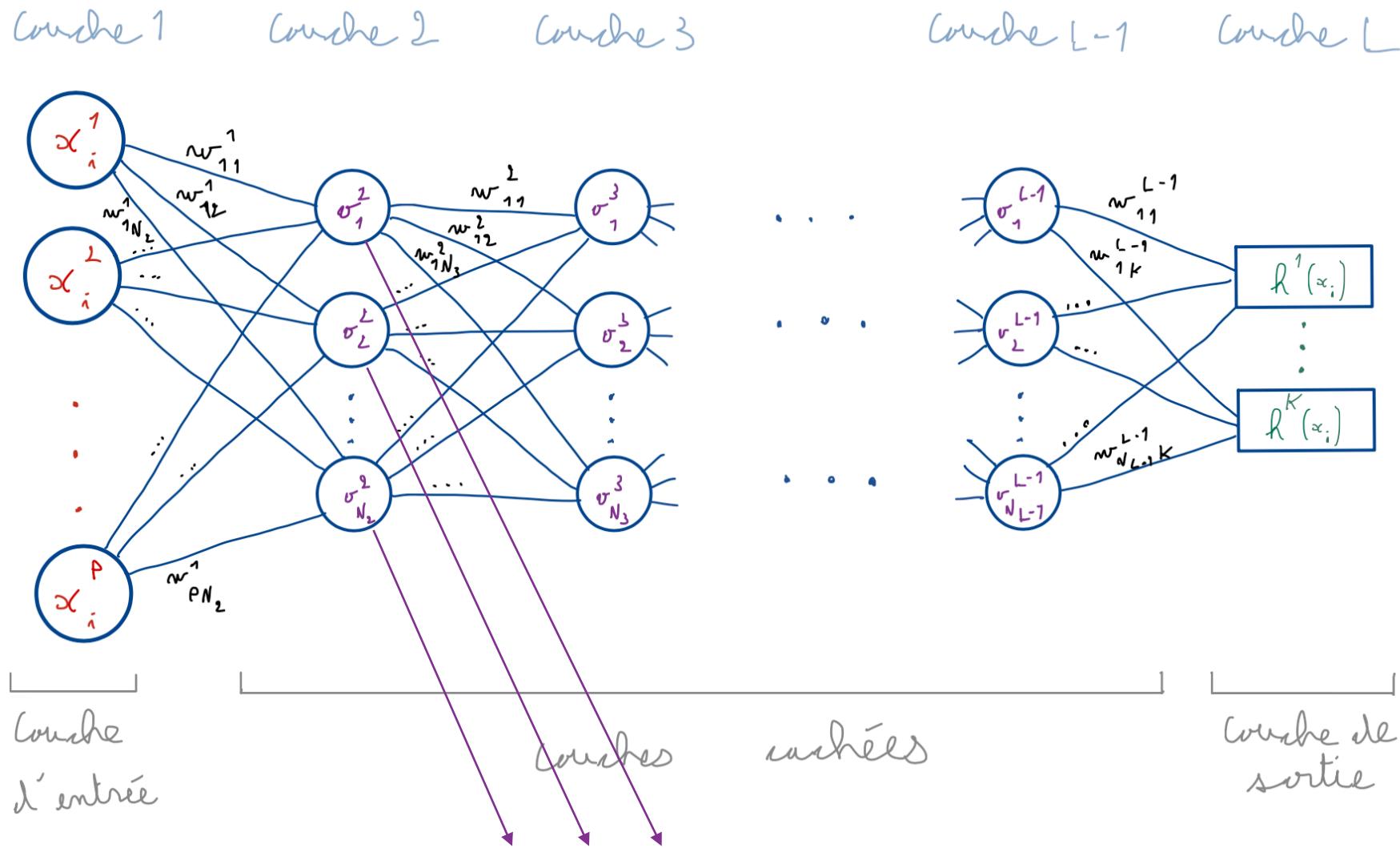
1) Calcul d'un gradient dans un réseau de neurones



$\forall k = 1, \dots, p :$

$$o_k^1 = x_i^k$$

1) Calcul d'un gradient dans un réseau de neurones

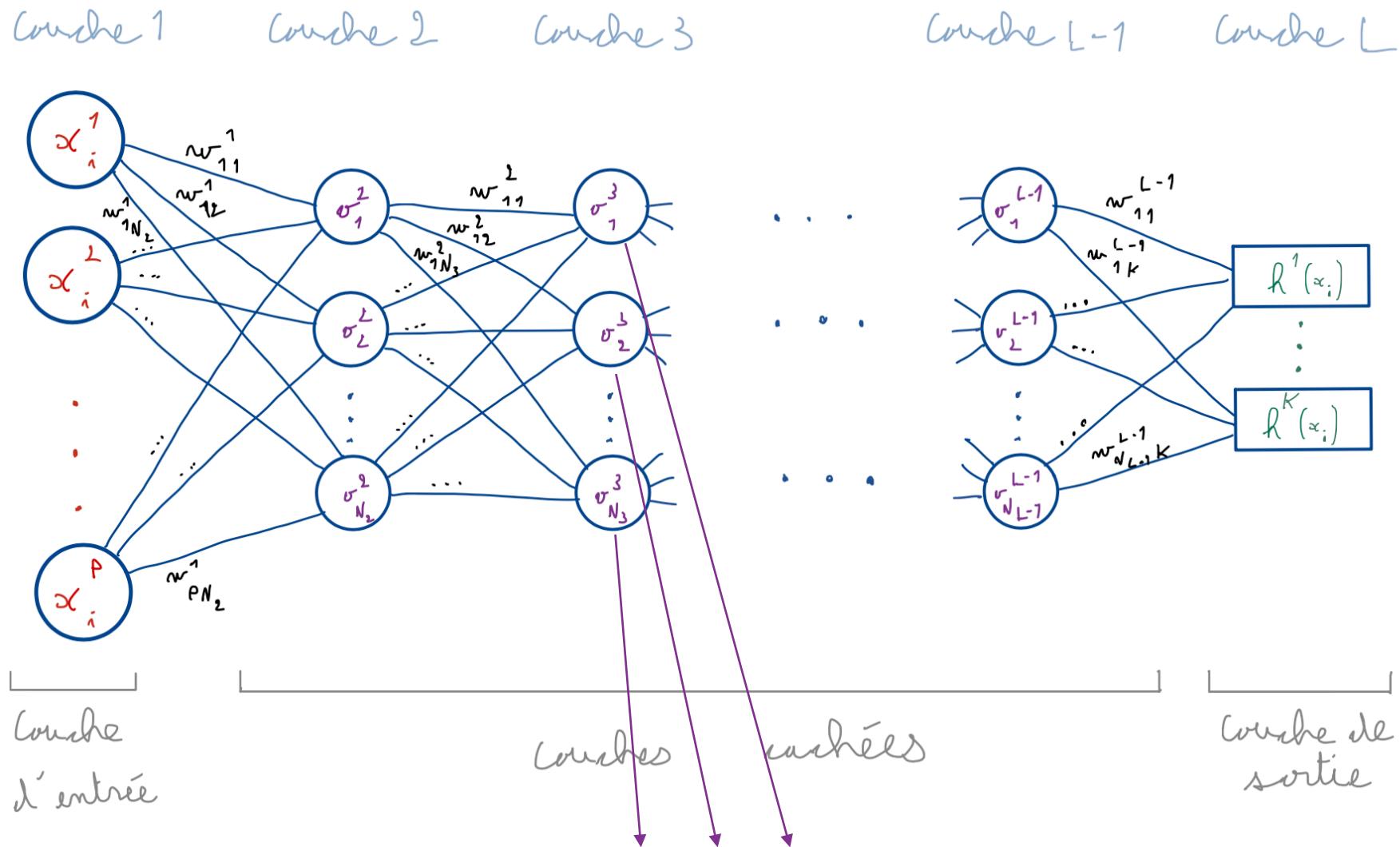


$\forall k = 1, \dots, p :$
 $o_k^1 = x_i^k$

$$\begin{cases} s_k^{l+1} = \sum_{p \in \{1, \dots, N_l\}} w_{pk}^l o_p^l \\ o_k^{l+1} = \varphi(s_k^{l+1}) \end{cases}$$

$\varphi(\cdot)$ non-linéaire \rightarrow ici $\varphi(s_k^{l+1}) = \text{ReLU}(s_k^{l+1}) = \max(0, s_k^{l+1})$

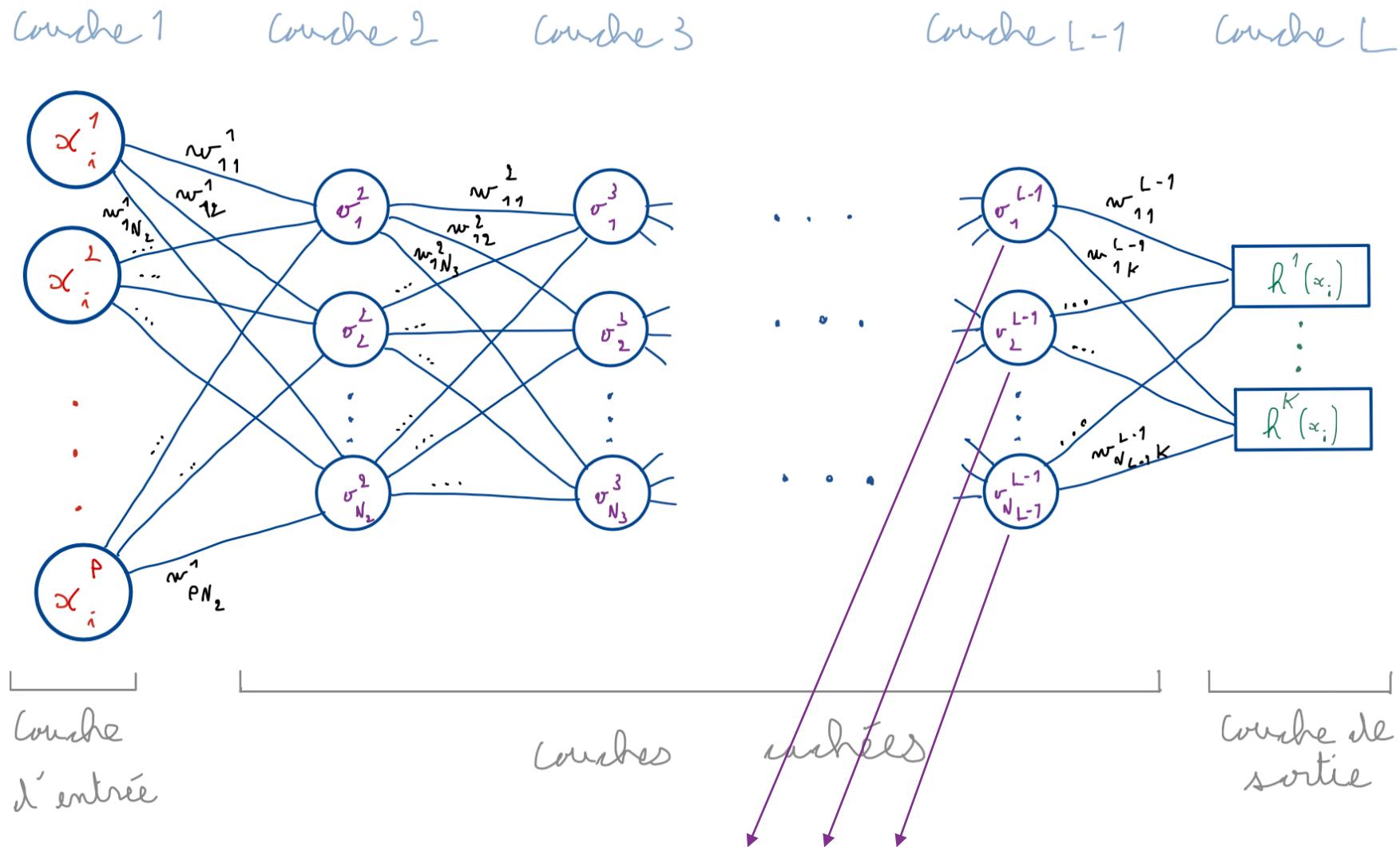
1) Calcul d'un gradient dans un réseau de neurones



$$\begin{cases} s_k^{l+1} = \sum_{p \in \{1, \dots, N_l\}} w_{pk}^l o_p^l \\ o_k^{l+1} = \varphi(s_k^{l+1}) \end{cases}$$

$\varphi(\cdot)$ non-linéaire \rightarrow ici $\varphi(s_k^{l+1}) = \text{ReLU}(s_k^{l+1}) = \max(0, s_k^{l+1})$

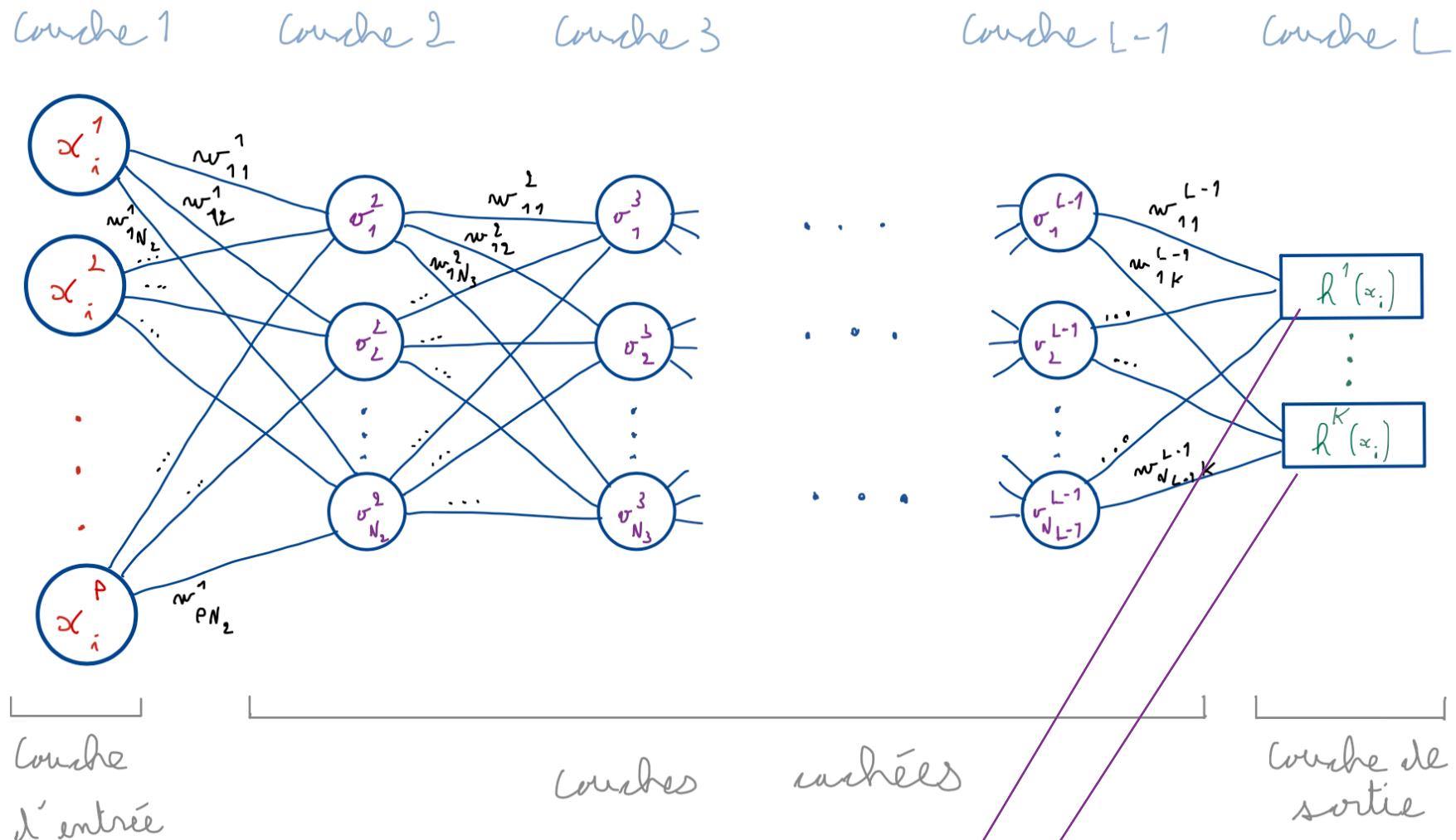
1) Calcul d'un gradient dans un réseau de neurones



$$\begin{cases} s_k^{l+1} = \sum_{p \in \{1, \dots, N_l\}} w_{pk}^l o_p^l \\ o_k^{l+1} = \varphi(s_k^{l+1}) \end{cases}$$

$\varphi(\cdot)$ non-linéaire \rightarrow ici $\varphi(s_k^{l+1}) = \text{ReLU}(s_k^{l+1}) = \max(0, s_k^{l+1})$

1) Calcul d'un gradient dans un réseau de neurones

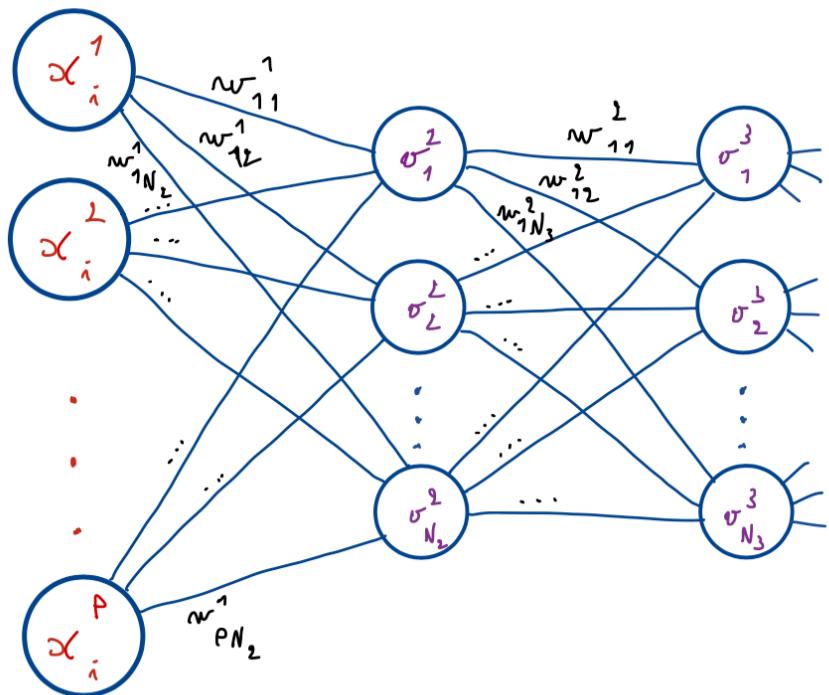


$$\begin{cases} s_k^{l+1} = \sum_{p \in \{1, \dots, N_l\}} w_{pk}^l o_p^l \\ o_k^{l+1} = \varphi(s_k^{l+1}) \end{cases}$$

$\varphi(\cdot)$ non-linéaire \rightarrow ici $\varphi(s_k^{l+1}) = \text{logistic}(s_k^{l+1}) = 1/(1 + e^{-s_k^{l+1}})$

$$\forall k = 1, \dots, K : \quad h^k(x_i) = o_k^L$$

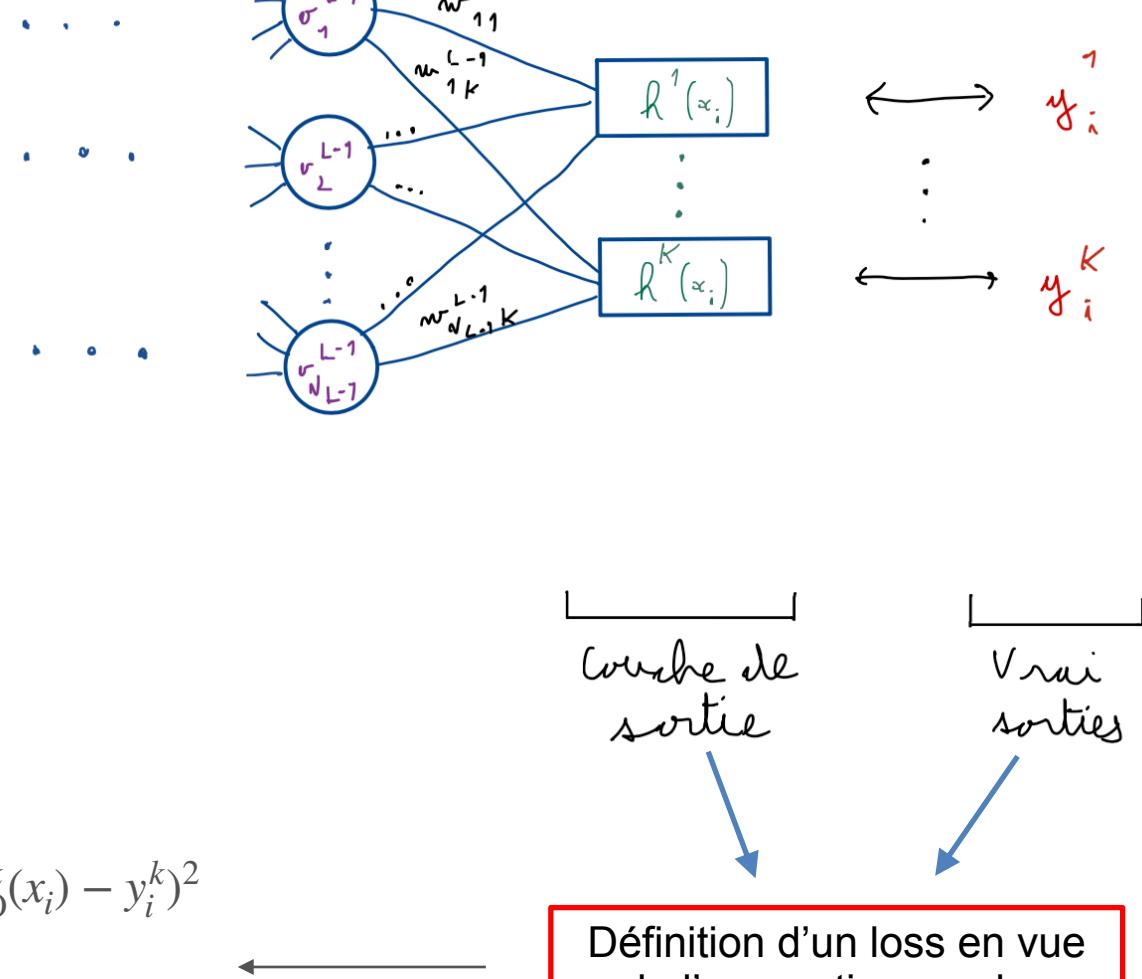
1) Calcul d'un gradient dans un réseau de neurones



$$\text{Loss possible : } loss \left(h_{\Theta}(x_i), y_i \right) = \sum_{k=1}^K (h_{\Theta}^k(x_i) - y_i^k)^2$$

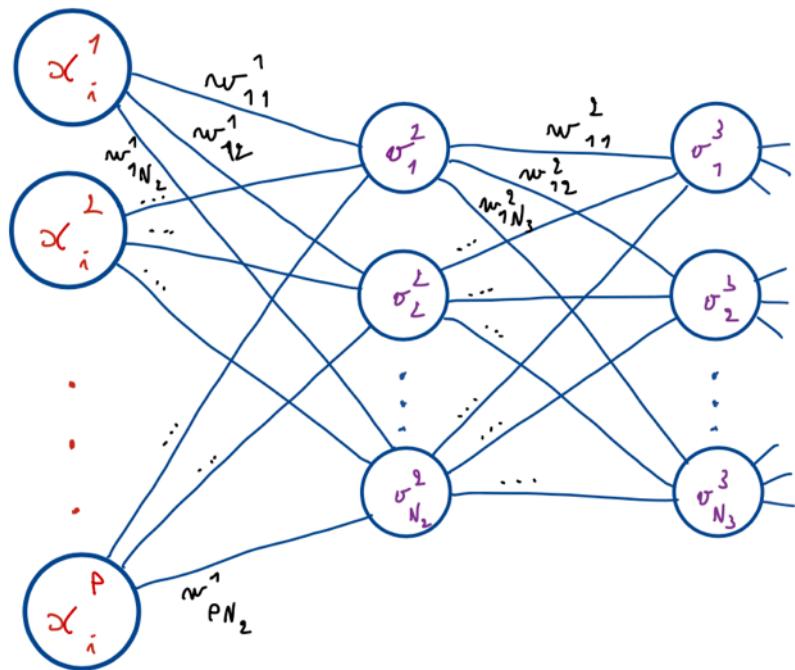
Remarque :

- Il existe d'autres losses comme la *Cross-entropy*.
 - Le loss doit être dérivable par rapport aux $h_{\Theta}(x_i)$

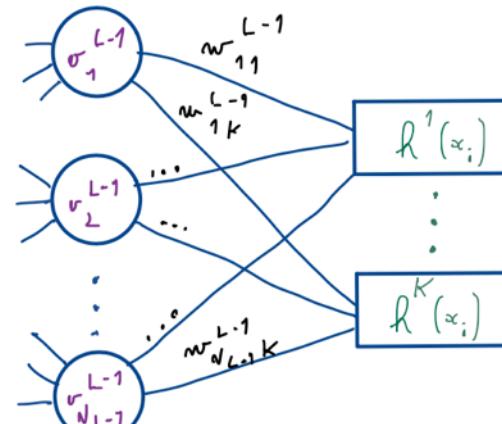


Définition d'un loss en vue de l'apprentissage des paramètres $\Theta = \{w_{pk}^l\}_{p,k,l}$

1) Calcul d'un gradient dans un réseau de neurones



...
...
...



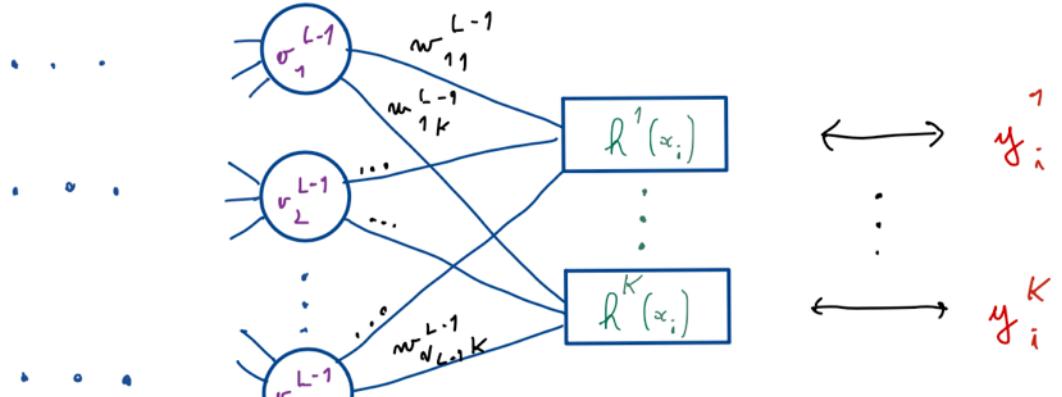
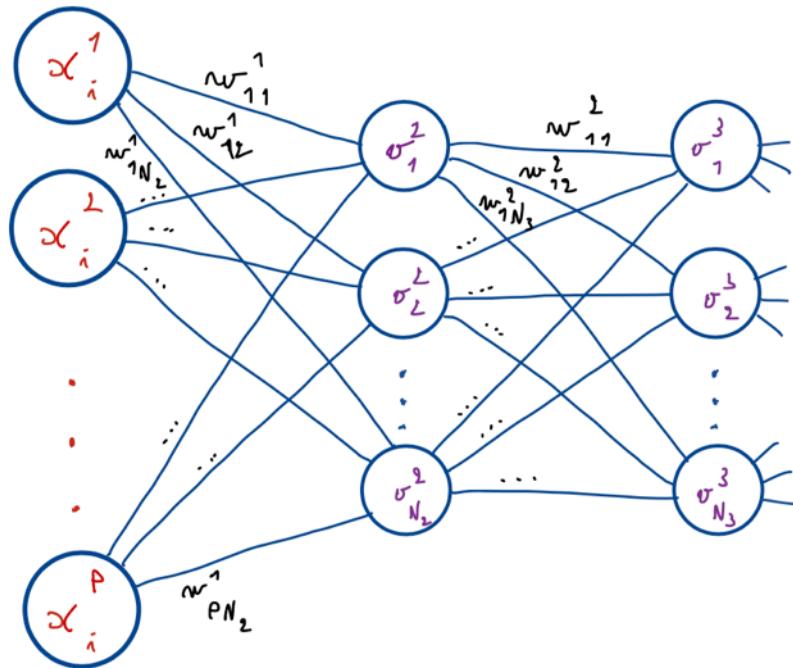
y_i^1
⋮
 y_i^K

Recherche des paramètres optimaux $\hat{\Theta}$ tels que :

$$\begin{aligned}\hat{\Theta} &= \arg \min_{\Theta=\{w_0, \dots, w_p\}} R_{\Theta} ((x_i, y_i)_{i=1, \dots, n}) \\ &= \arg \min_{\Theta=\{w_0, \dots, w_p\}} \frac{1}{n} \sum_{i=1}^n \text{loss}(h_{\Theta}(x_i), y_i)\end{aligned}$$

→ descente de gradient en grande dimension

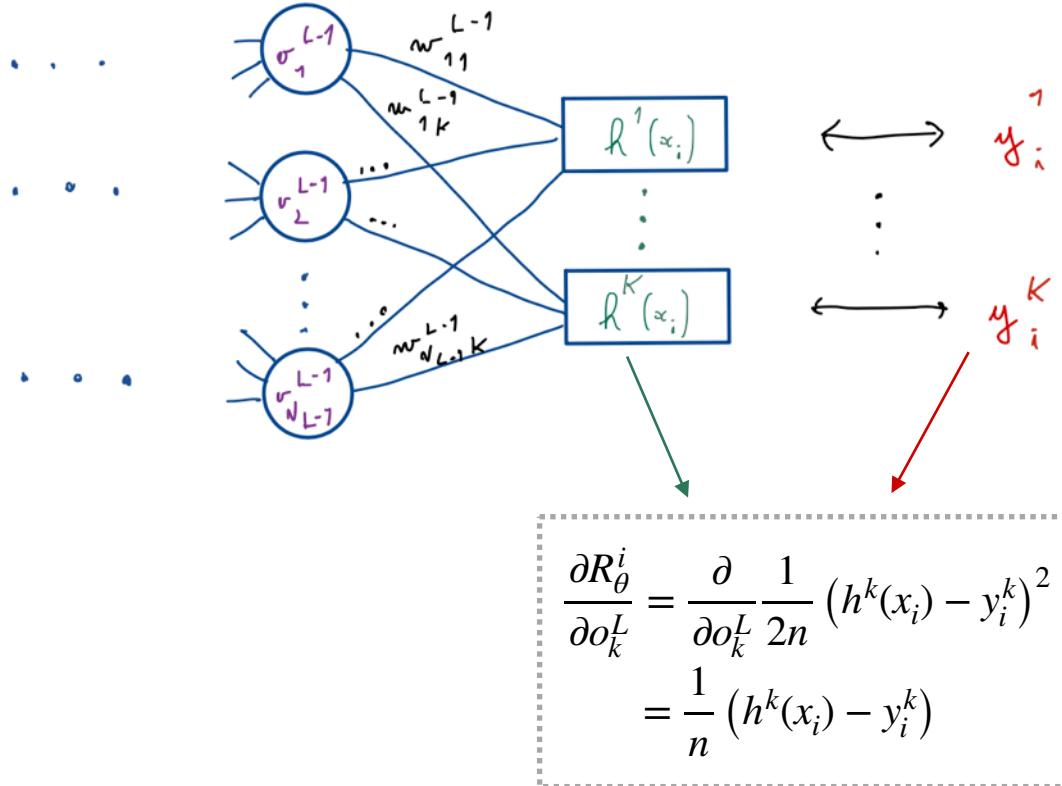
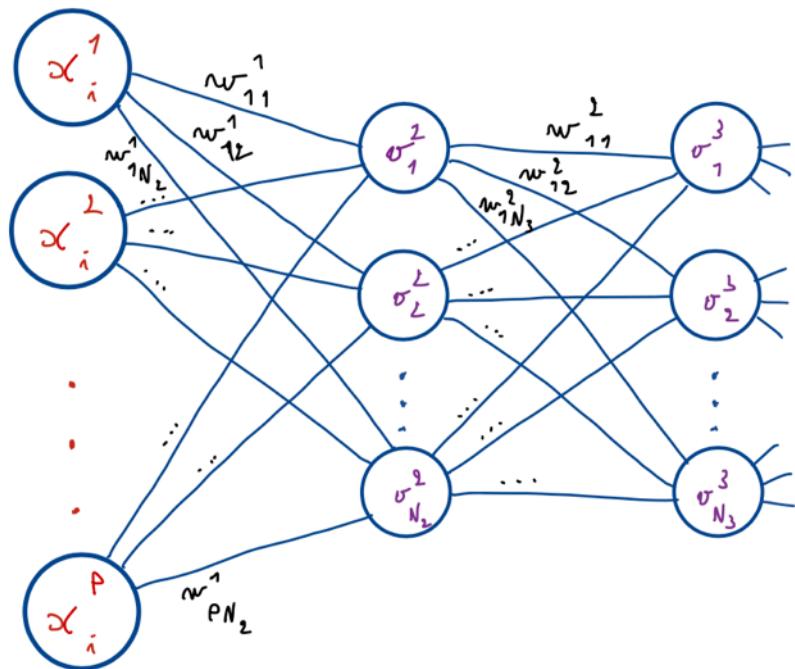
1) Calcul d'un gradient dans un réseau de neurones



$$\begin{aligned}
\nabla R_{\Theta}((x_i, y_i)_{i=1, \dots, n}) &= \left(\frac{\partial R_{\Theta}(\dots)}{\partial w_{1,1}^1}, \frac{\partial R_{\Theta}(\dots)}{\partial w_{1,2}^1}, \dots, \frac{\partial R_{\Theta}(\dots)}{\partial w_{N_{L-1},K}^{N_{L-1}}} \right)^T && \text{(Gradient du risque empirique)} \\
&= \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \left(\frac{\partial \text{loss}(h^k(x_i), y_i^k)}{\partial w_{1,1}^1}, \frac{\partial \text{loss}(h^k(x_i), y_i^k)}{\partial w_{1,2}^1}, \dots, \frac{\partial \text{loss}(h^k(x_i), y_i^k)}{\partial w_{N_{L-1},K}^{N_{L-1}}} \right)^T \\
&= \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K \left(\frac{\partial (h^k(x_i) - y_i^k)^2}{\partial w_{1,1}^1}, \frac{\partial (h^k(x_i) - y_i^k)^2}{\partial w_{1,2}^1}, \dots, \frac{\partial (h^k(x_i) - y_i^k)^2}{\partial w_{N_{L-1},K}^{N_{L-1}}} \right)^T
\end{aligned}$$

Calcul analytique possible mais les $h^k(x_i)$ sont issus de fonctions composées couches après couches
→ dérivation avec le théorème des fonctions composées (*chain rule*) !

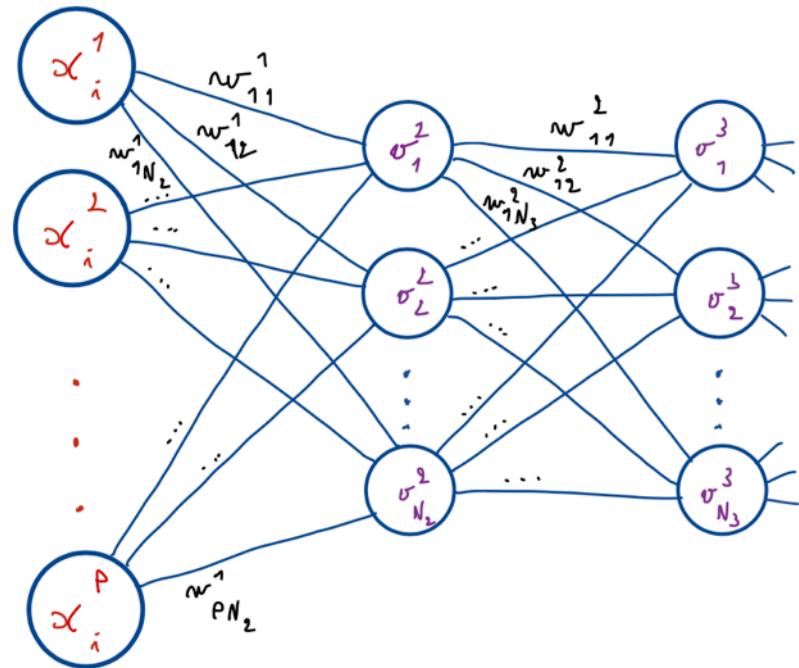
1) Calcul d'un gradient dans un réseau de neurones



Remarques :

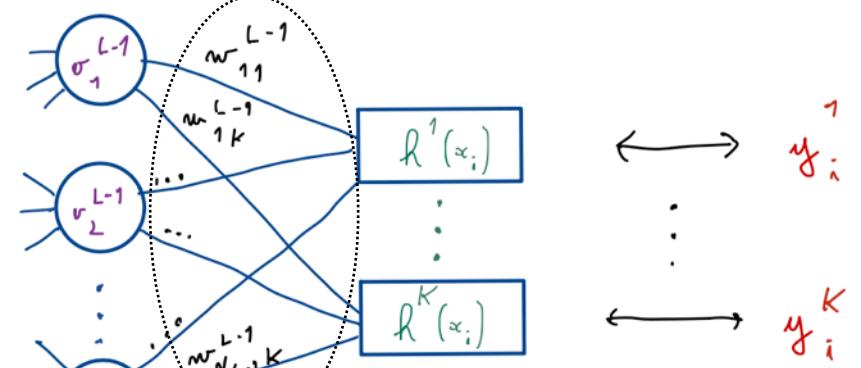
- Notation simplifiée $R_\theta^i = \text{loss}(h(x_i), y_i)$
 - $h^k(x_i)$ est la sortie o_k^L pour l'observation i

1) Calcul d'un gradient dans un réseau de neurones



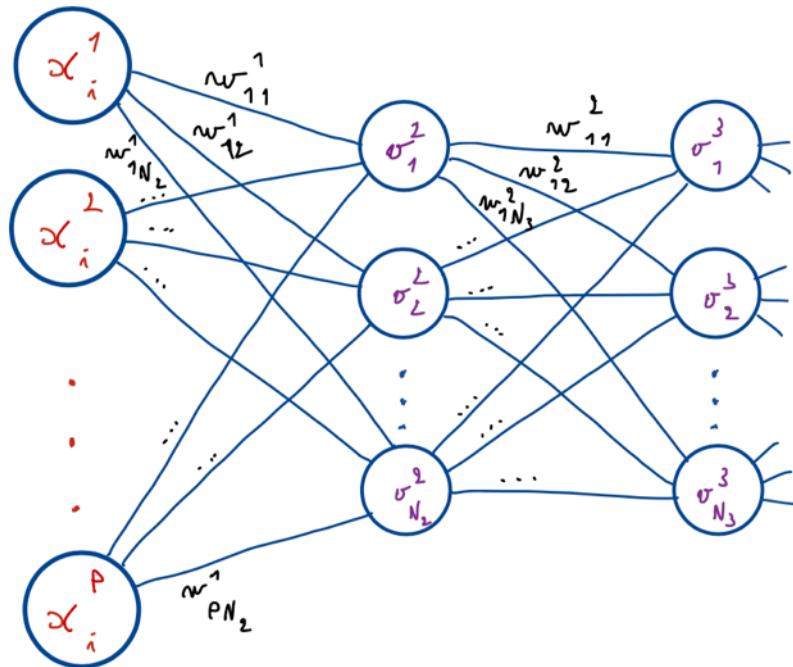
$$\frac{\partial R_\Theta^i}{\partial w_{pk}^l} = \frac{\partial R_\Theta^i}{\partial o_k^{l+1}} \frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l}$$

$$l = L - 1$$

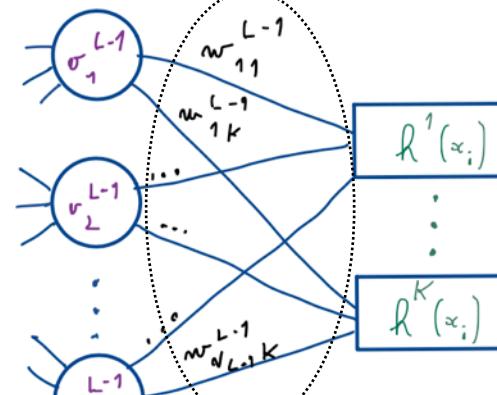


$$\begin{aligned} \frac{\partial R_\Theta^i}{\partial o_k^L} &= \frac{\partial}{\partial o_k^L} \frac{1}{2n} (h^k(x_i) - y_i^k)^2 \\ &= \frac{1}{n} (h^k(x_i) - y_i^k) \end{aligned}$$

1) Calcul d'un gradient dans un réseau de neurones



...
...
...



y_i^1
⋮
 y_i^K

$$\frac{\partial R_\Theta^i}{\partial w_{pk}^l} = \frac{\partial R_\Theta^i}{\partial o_k^{l+1}} \frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l}$$

$l = L - 1$

$$\begin{aligned} \frac{\partial R_\Theta^i}{\partial o_k^L} &= \frac{\partial}{\partial o_k^L} \frac{1}{2n} (h^k(x_i) - y_i^k)^2 \\ &= \frac{1}{n} (h^k(x_i) - y_i^k) \end{aligned}$$

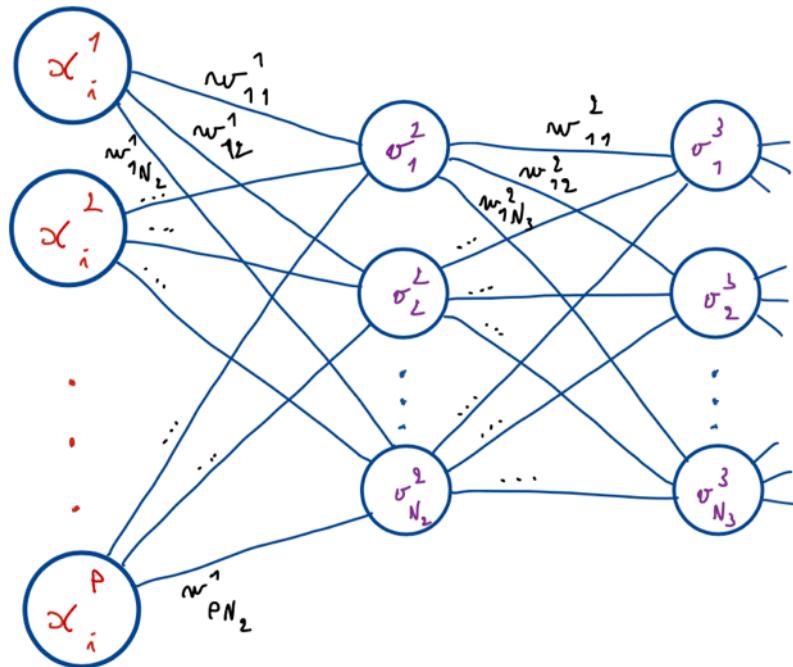
$$\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} = \frac{\partial}{\partial s_k^{l+1}} \varphi(s_k^{l+1}) \text{ [connu]}$$

$l = L - 1$

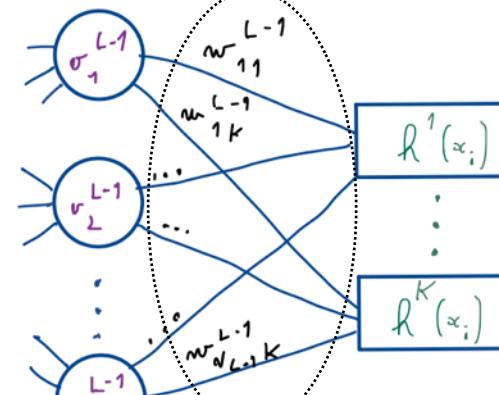
$$\begin{aligned} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l} &= \frac{\partial}{\partial w_{pk}^l} \sum_{p' \in \{1, \dots, N_l\}} w_{p'k}^l o_{p'}^l \\ &= \frac{\partial}{\partial w_{pk}^l} w_{pk}^l o_p^l \\ &= o_p^l \end{aligned}$$

$l = L - 1$

1) Calcul d'un gradient dans un réseau de neurones



...
...
...



y_i^1
⋮
 y_i^K

$$\frac{\partial R_\Theta^i}{\partial w_{pk}^l} = \frac{\partial R_\Theta^i}{\partial o_k^{l+1}} \frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l}$$

$l = L - 1$

$$\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} = \frac{\partial}{\partial s_k^{l+1}} \varphi(s_k^{l+1}) \text{ [connu]}$$

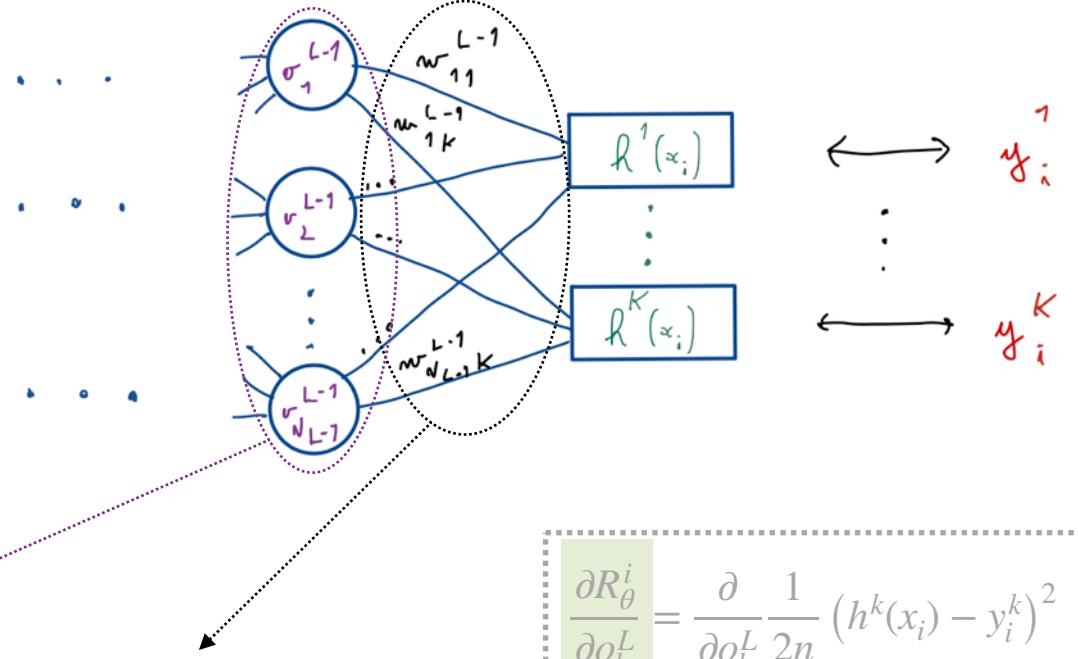
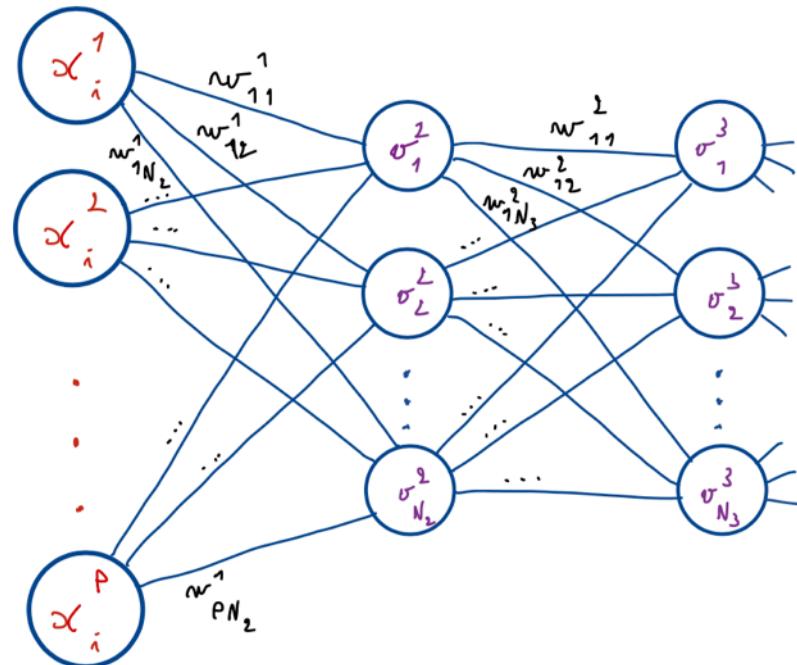
$l = L - 1$

$$\begin{aligned} \frac{\partial R_\Theta^i}{\partial o_k^L} &= \frac{\partial}{\partial o_k^L} \frac{1}{2n} (h^k(x_i) - y_i^k)^2 \\ &= \frac{1}{n} (h^k(x_i) - y_i^k) \end{aligned}$$

$$\begin{aligned} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l} &= \frac{\partial}{\partial w_{pk}^l} \sum_{p' \in \{1, \dots, N_l\}} w_{p'k}^l o_{p'}^l \\ &= \frac{\partial}{\partial w_{pk}^l} w_{pk}^l o_p^l \\ &= o_p^l \end{aligned}$$

$l = L - 1$

1) Calcul d'un gradient dans un réseau de neurones



$$\frac{\partial R_\Theta^i}{\partial o_k^l} = \sum_{p=1}^{N_{l+1}} \underbrace{\frac{\partial R_\Theta^i}{\partial o_k^{l+1}}}_{\text{Connu à la couche } l+1} \underbrace{\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}}}_{\text{Derivee de } \varphi} \underbrace{\frac{\partial s_k^{l+1}}{\partial o_k^l}}_{=w_{pk}^l}$$

$l = L - 1$

$$\frac{\partial R_\Theta^i}{\partial w_{pk}^l} = \frac{\partial R_\Theta^i}{\partial o_k^{l+1}} \frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l}$$

$l = L - 1$

$$\begin{aligned} \frac{\partial R_\Theta^i}{\partial o_k^L} &= \frac{\partial}{\partial o_k^L} \frac{1}{2n} (h^k(x_i) - y_i^k)^2 \\ &= \frac{1}{n} (h^k(x_i) - y_i^k) \end{aligned}$$

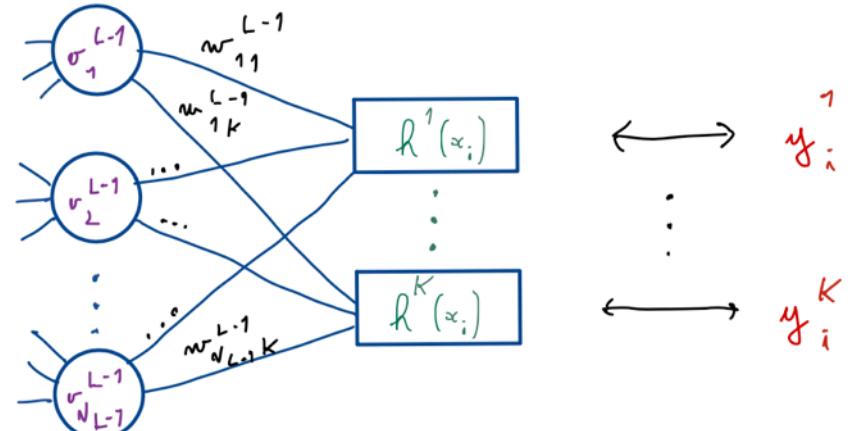
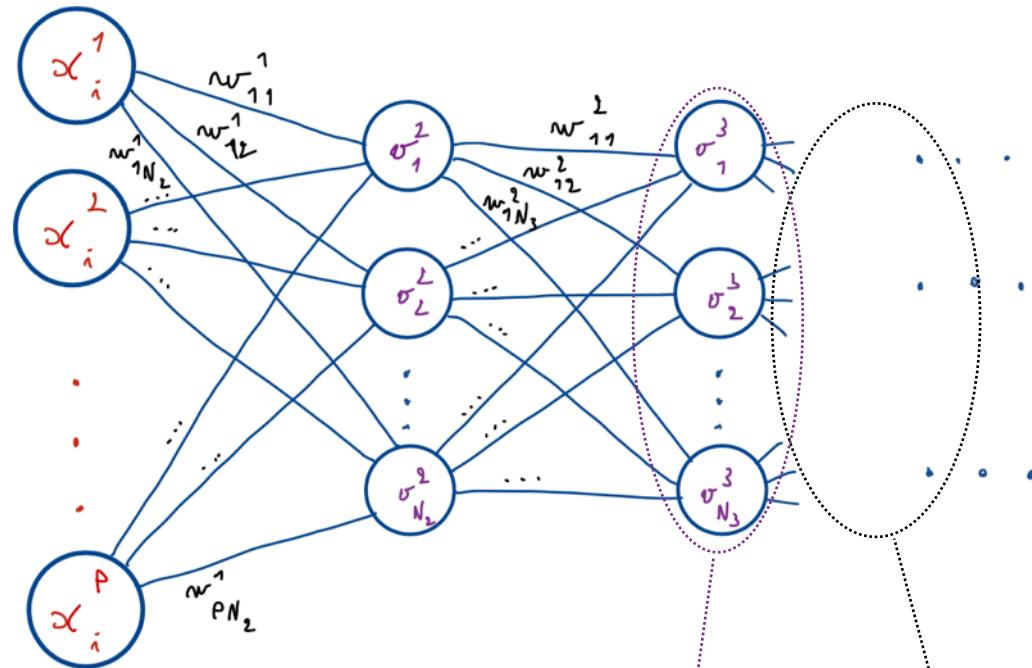
$$\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} = \frac{\partial}{\partial s_k^{l+1}} \varphi(s_k^{l+1}) \quad [\text{connu}]$$

$l = L - 1$

$$\begin{aligned} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l} &= \frac{\partial}{\partial w_{pk}^l} \sum_{p' \in \{1, \dots, N_l\}} w_{p'k}^l o_{p'}^l \\ &= \frac{\partial}{\partial w_{pk}^l} w_{pk}^l o_p^l \\ &= o_p^l \end{aligned}$$

$l = L - 1$

1) Calcul d'un gradient dans un réseau de neurones



$$\frac{\partial R_\Theta^i}{\partial o_k^l} = \sum_{p=1}^{N_{l+1}} \underbrace{\frac{\partial R_\Theta^i}{\partial o_k^{l+1}}}_{\text{Connu à la couche } l+1} \underbrace{\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}}}_{\text{Derivee de } \varphi} \underbrace{\frac{\partial s_k^{l+1}}{\partial o_k^l}}_{=w_{pk}^l}$$

$$\frac{\partial R_\Theta^i}{\partial w_{pk}^l} = \frac{\partial R_\Theta^i}{\partial o_k^{l+1}} \frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l}$$

$l = 3$

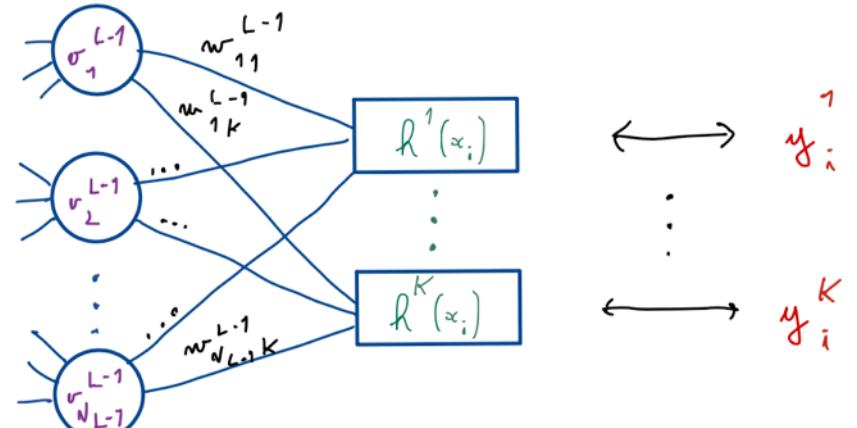
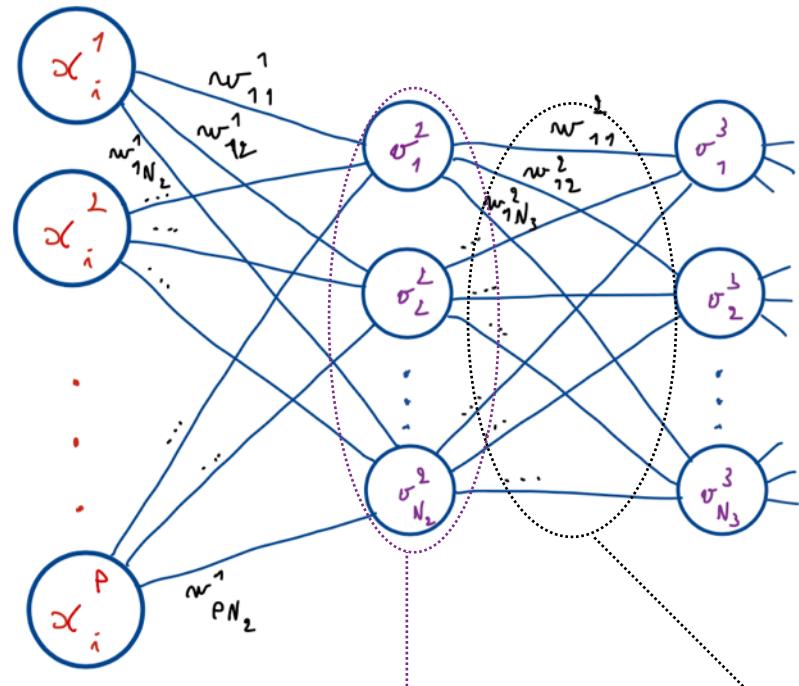
$$\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} = \frac{\partial}{\partial s_k^{l+1}} \varphi(s_k^{l+1}) \quad [\text{connu}]$$

$l = 3$

$$\begin{aligned} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l} &= \frac{\partial}{\partial w_{pk}^l} \sum_{p' \in \{1, \dots, N_l\}} w_{p'k}^l o_{p'}^l \\ &= \frac{\partial}{\partial w_{pk}^l} w_{pk}^l o_p^l \\ &= o_p^l \end{aligned}$$

$l = 3$

1) Calcul d'un gradient dans un réseau de neurones



$$\frac{\partial R_\Theta^i}{\partial o_k^l} = \sum_{p=1}^{N_{l+1}} \underbrace{\frac{\partial R_\Theta^i}{\partial o_k^{l+1}}}_{\text{Connu à la couche } l+1} \underbrace{\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}}}_{\text{Derivee de } \varphi} \underbrace{\frac{\partial s_k^{l+1}}{\partial o_k^l}}_{=w_{pk}^l}$$

$$\frac{\partial R_\Theta^i}{\partial w_{pk}^l} = \underbrace{\frac{\partial R_\Theta^i}{\partial o_k^{l+1}}}_{l=2} \underbrace{\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}}}_{l=2} \underbrace{\frac{\partial s_k^{l+1}}{\partial w_{pk}^l}}$$

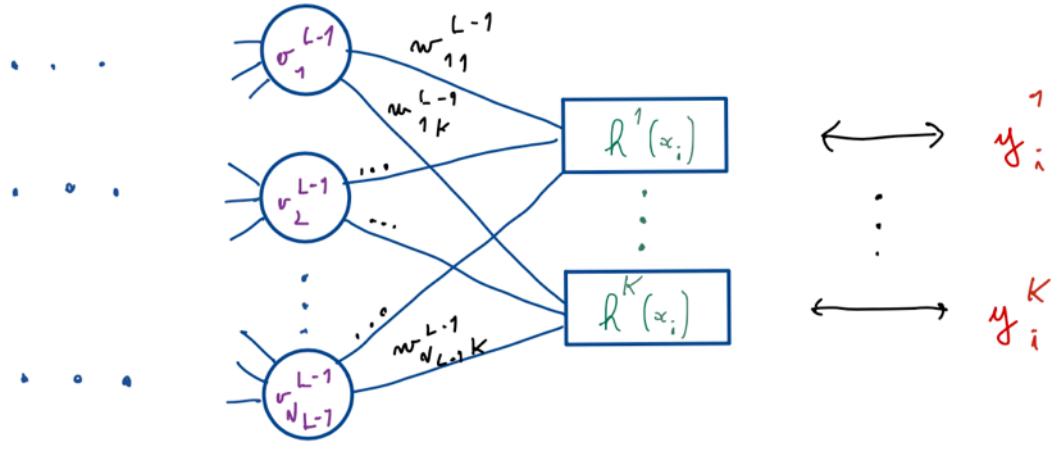
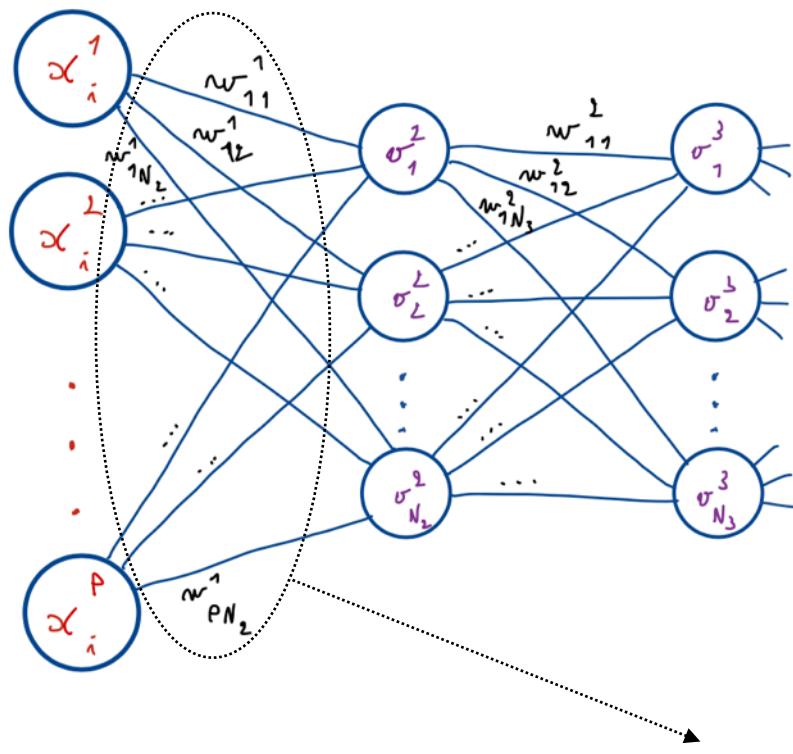
$$\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} = \frac{\partial}{\partial s_k^{l+1}} \varphi(s_k^{l+1}) \text{ [connu]}$$

$l = 2$

$$\begin{aligned} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l} &= \frac{\partial}{\partial w_{pk}^l} \sum_{p' \in \{1, \dots, N_l\}} w_{p'k}^l o_{p'}^l \\ &= \frac{\partial}{\partial w_{pk}^l} w_{pk}^l o_p^l \\ &= o_p^l \end{aligned}$$

$l = 2$

1) Calcul d'un gradient dans un réseau de neurones



$$\frac{\partial R_{\Theta}^i}{\partial w_{pk}^l} = \frac{\partial R_{\Theta}^i}{\partial o_k^{l+1}} \frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l}$$

$$l = 1$$

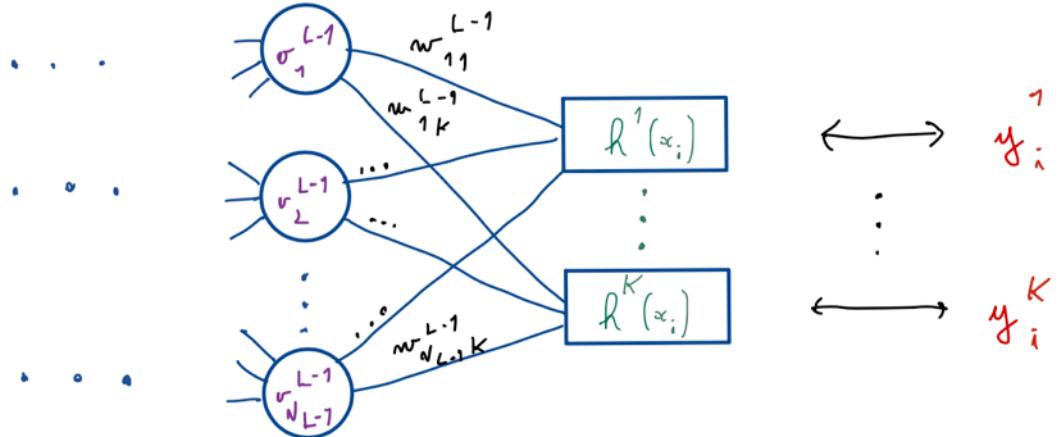
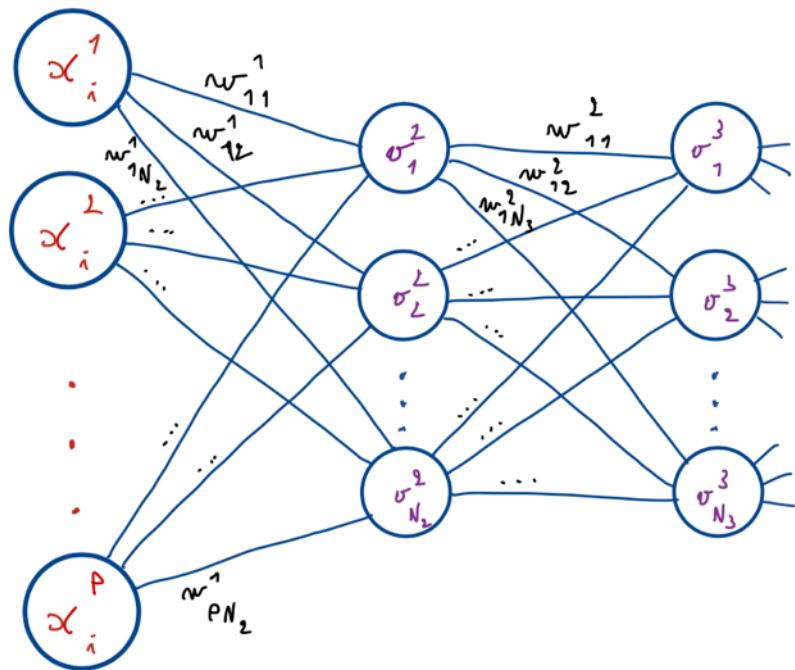
$$\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} = \frac{\partial}{\partial s_k^{l+1}} \varphi(s_k^{l+1}) \text{ [connu]}$$

$$l = 1$$

$$\begin{aligned}
 \frac{\partial s_k^{l+1}}{\partial w_{pk}^l} &= \frac{\partial}{\partial w_{pk}^l} \sum_{p' \in \{1, \dots, N_l\}} w_{p'k}^l o_{p'}^l \\
 &= \frac{\partial}{\partial w_{pk}^l} w_{pk}^l o_p^l \\
 &= o_p^l
 \end{aligned}$$

$$l = 1$$

1) Calcul d'un gradient dans un réseau de neurones

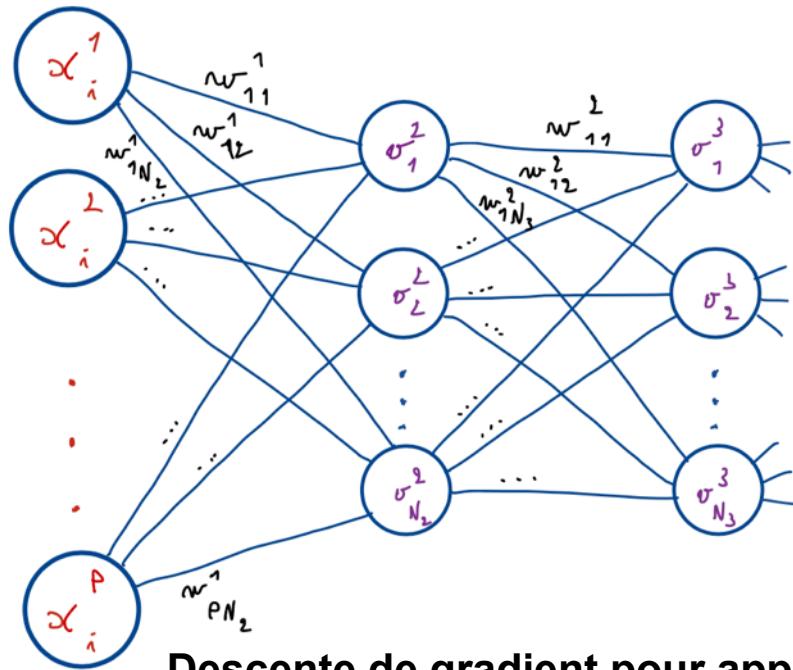


BINGO !!!

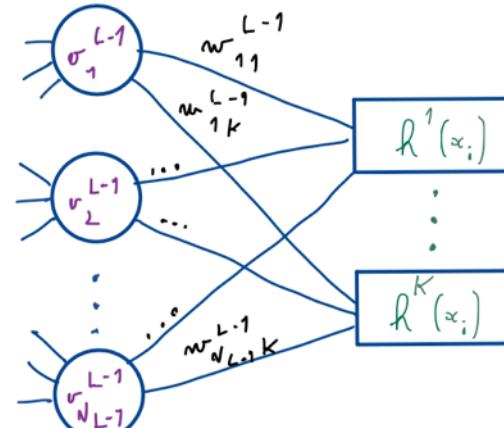
On a calculé le gradient pour une observation i :

$$\nabla R_{\Theta}^i \left((x_i, y_i)_{i=1, \dots, n} \right) = \left(\frac{\partial R_{\Theta}^i(\dots)}{\partial w_{1,1}^1}, \frac{\partial R_{\Theta}^i(\dots)}{\partial w_{1,2}^1}, \dots, \frac{\partial R_{\Theta}^i(\dots)}{\partial w_{N_{L-1},K}^{N_{L-1}}} \right)^{\top}$$

1) Calcul d'un gradient dans un réseau de neurones



...
...
...



$\leftarrow \rightarrow$
:
 y_i^1
 $\leftarrow \rightarrow$
:
 y_i^K

Descente de gradient pour apprendre les $\{w_{pk}^l\}_{p,k,l}$:

- (Forward)
 - Pour chaque x_i , calcule $h(x_i) = (h^1(x_i), h^2(x_i), \dots, h^K(x_i))^T$
 - Pour chaque x_i , calcule $\frac{\partial R_\theta^i}{\partial o_k^L}$ sur la couche de sortie en utilisant $h^k(x_i)$ et y_i^k

- (Backward)

- Pour chaque x_i et pour l de $L - 1$ à 1 : Calcule les $\frac{\partial R_\theta^i}{\partial w_{pk}^l}$ et $\frac{\partial R_\theta^i}{\partial o_p^l}$ en utilisant les $\frac{\partial R_\theta^i}{\partial o_p^{l+1}}$

- (Etape de descente de gradient)

- $\forall \{p, k, l\} : \left(w_{pk}^l \right) = \left(w_{pk}^l \right) - \lambda \frac{\partial R_\theta}{\partial w_{pk}^l} = \left(w_{pk}^l \right) - \lambda \frac{1}{n} \sum_{i=1}^n \frac{\partial R_\theta^i}{\partial w_{pk}^l}$

Remarque : Utilise des informations de la phase forward

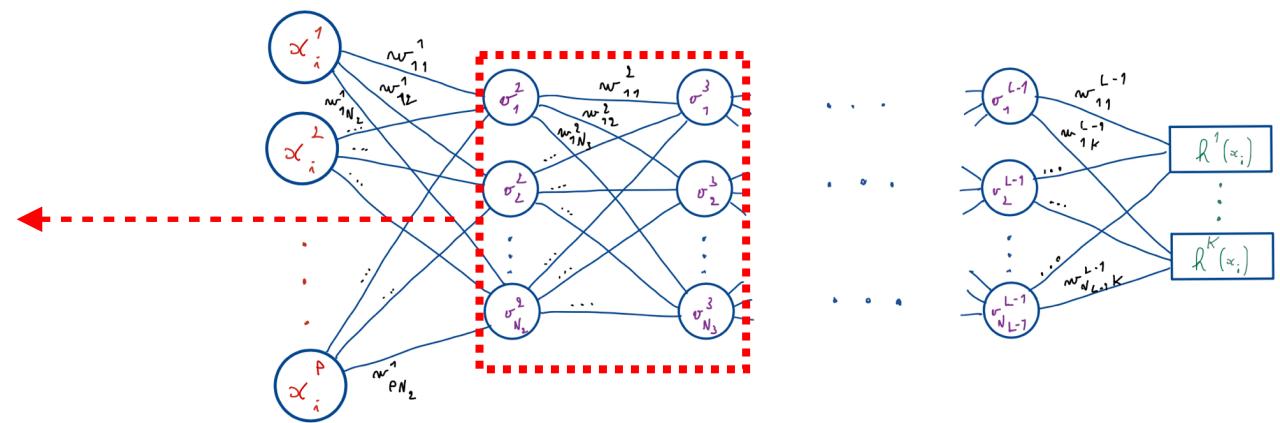
1) Calcul d'un gradient dans un réseau de neurones

Passage *forward* d'une couche dense d'un réseau de neurones :

$$\left(s_j^{l+1} \right)^b = \sum_{i=1}^{N_l} \left(o_i^l \right)^b w_{i,j}^l$$

$\forall j \in \{1, \dots, N_{l+1}\}$

$\forall b \in \{1, \dots, B\}$



$$\begin{pmatrix} \vdots \\ \left(s_j^{l+1} \right)^b \\ \vdots \\ \vdots \end{pmatrix} = \begin{pmatrix} w_{1,1}^l & w_{2,1}^l & \cdots & w_{N_1,1}^l \\ w_{1,2}^l & w_{2,2}^l & \cdots & w_{N_1,2}^l \\ \vdots & \vdots & \ddots & \vdots \\ w_{1,N_1}^l & w_{2,N_1}^l & \cdots & w_{N_1,N_1}^l \end{pmatrix} \begin{pmatrix} \left(\sigma_1^l \right)^1 & \left(\sigma_1^l \right)^2 & \cdots & \left(\sigma_1^l \right)^B \\ \left(\sigma_2^l \right)^1 & \left(\sigma_2^l \right)^2 & \cdots & \left(\sigma_2^l \right)^B \\ \vdots & \vdots & \ddots & \vdots \\ \left(\sigma_{N_1}^l \right)^1 & \left(\sigma_{N_1}^l \right)^2 & \cdots & \left(\sigma_{N_1}^l \right)^B \end{pmatrix}$$

1) Calcul d'un gradient dans un réseau de neurones

Passage *backward* d'une couche dense d'un réseau de neurones :

$$\frac{\partial R_{\Theta}^i}{\partial (o_k^l)^b} = \sum_{p=1}^{N_{l+1}} \frac{\partial R_{\Theta}^i}{\partial (o_k^{l+1})^b} \frac{\partial (o_k^{l+1})^b}{\partial (s_k^{l+1})^b} \frac{\partial (s_k^{l+1})^b}{\partial (o_k^l)^b} = \sum_{p=1}^{N_{l+1}} \frac{\partial R_{\Theta}^i}{\partial (o_k^{l+1})^b} \left(1_{(s_k^{l+1})^b \geq 0} (w_{pk})^b \right)$$

Tout juste calculé par l'algorithme de backpropagation

$$\frac{\partial R_{\Theta}^i}{\partial (w_{pk})^b} = \sum_{p=1}^{N_{l+1}} \frac{\partial R_{\Theta}^i}{\partial (o_k^{l+1})^b} \frac{\partial (o_k^{l+1})^b}{\partial (s_k^{l+1})^b} \frac{\partial (s_k^{l+1})^b}{\partial (w_{pk})^b} = \sum_{p=1}^{N_{l+1}} \frac{\partial R_{\Theta}^i}{\partial (o_k^{l+1})^b} \left(1_{(s_k^{l+1})^b \geq 0} (o_p^l)^b \right)$$

Tout juste calculé par l'algorithme de backpropagation

Doit être stocké lors de la phase forward

Doit être stocké lors de la phase forward

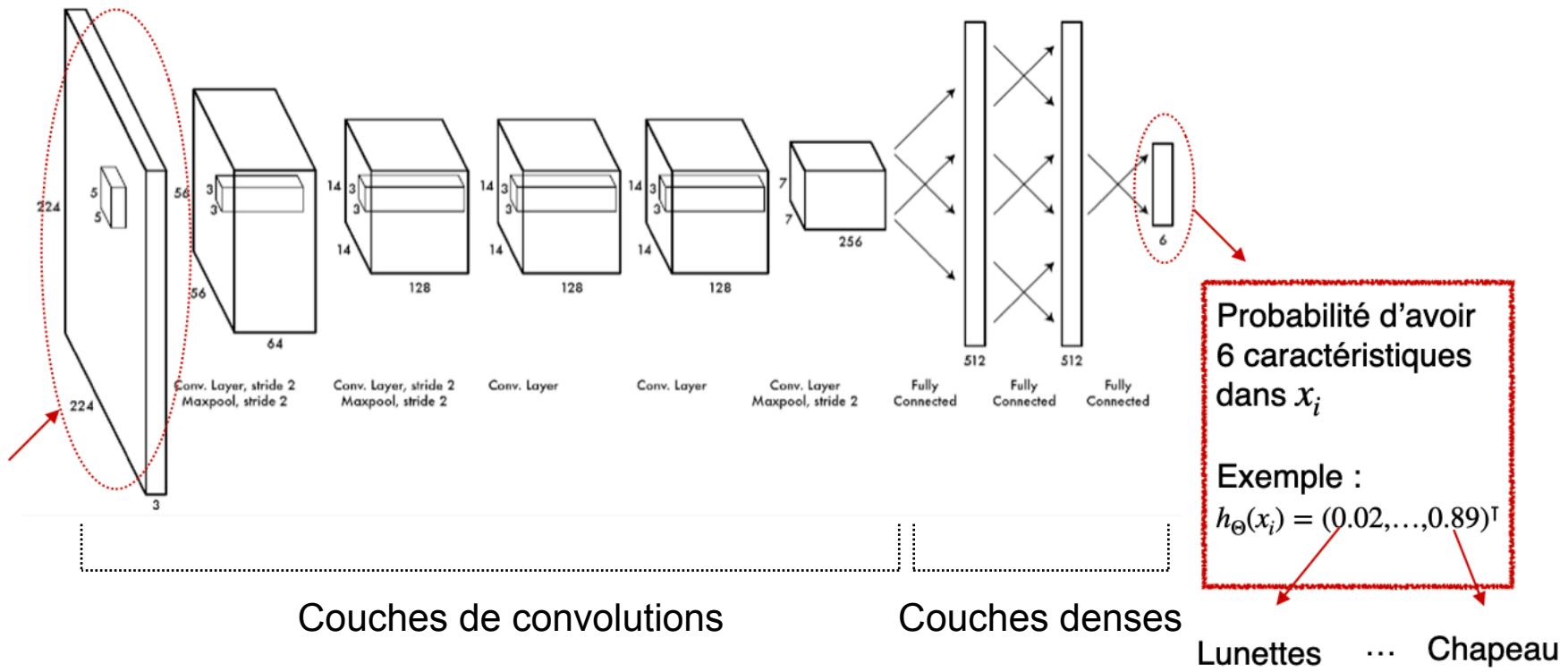
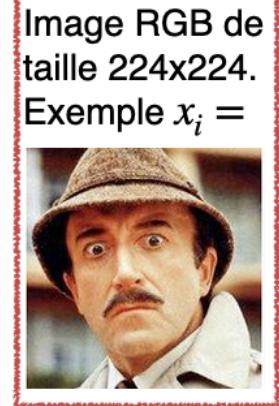
Paramètre du réseau

Doit être stocké lors de la phase forward

Remarques :

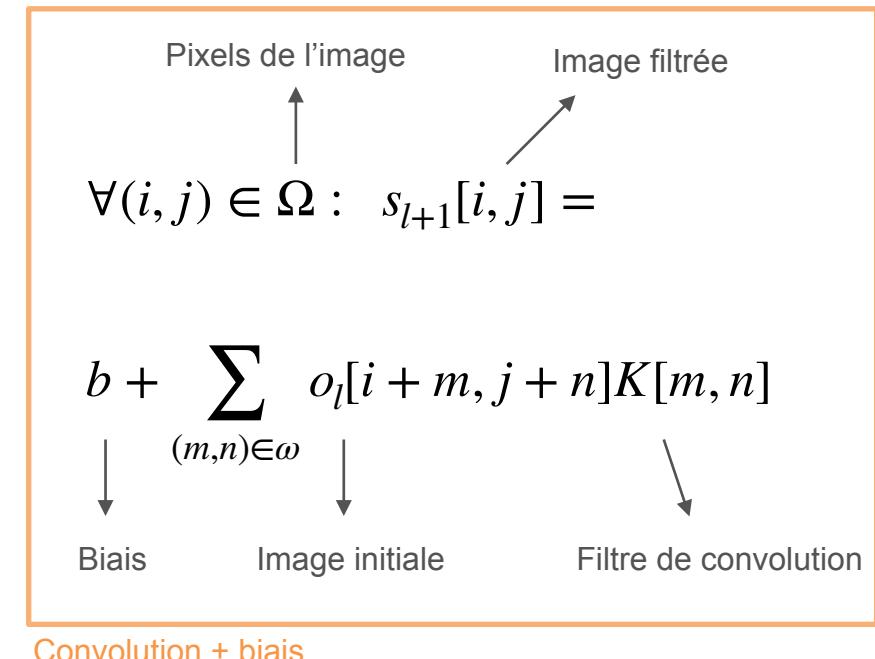
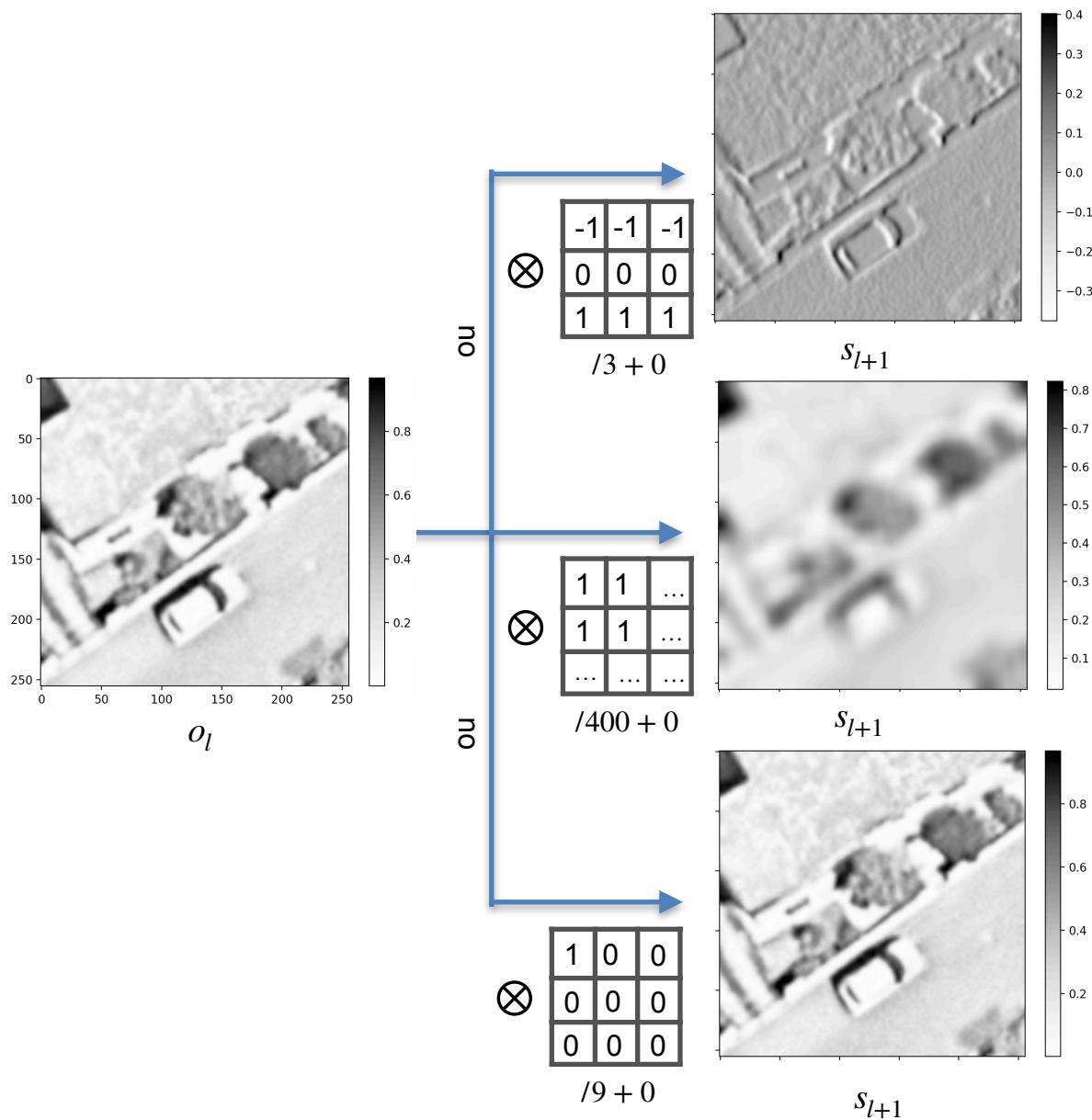
- Multiplications matrices / matrices
- Nécessite des informations stockées lors de la phase *forward*

2) Réseaux convolutifs



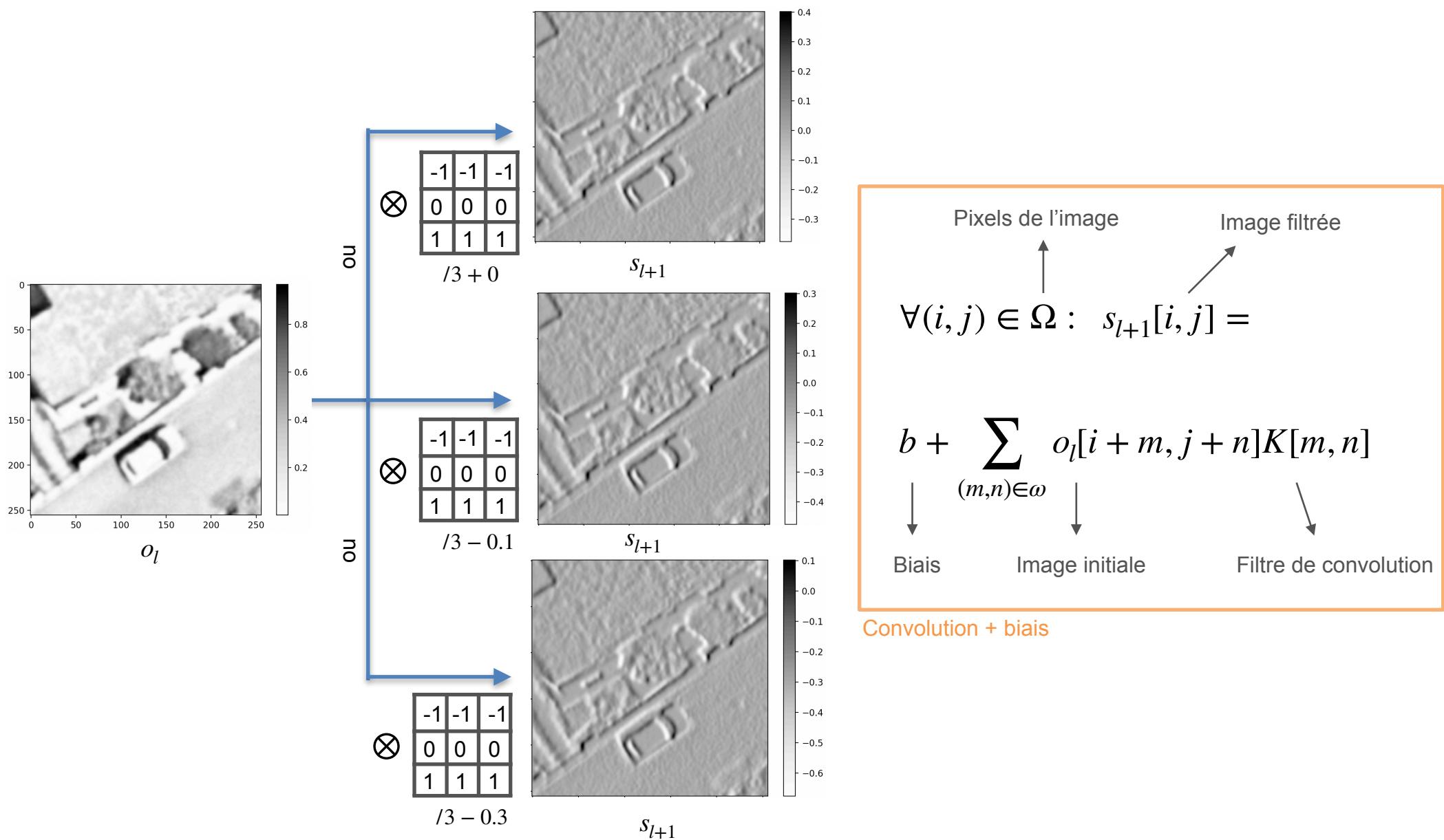
2) Réseaux convolutifs

Briques de base : Convolution



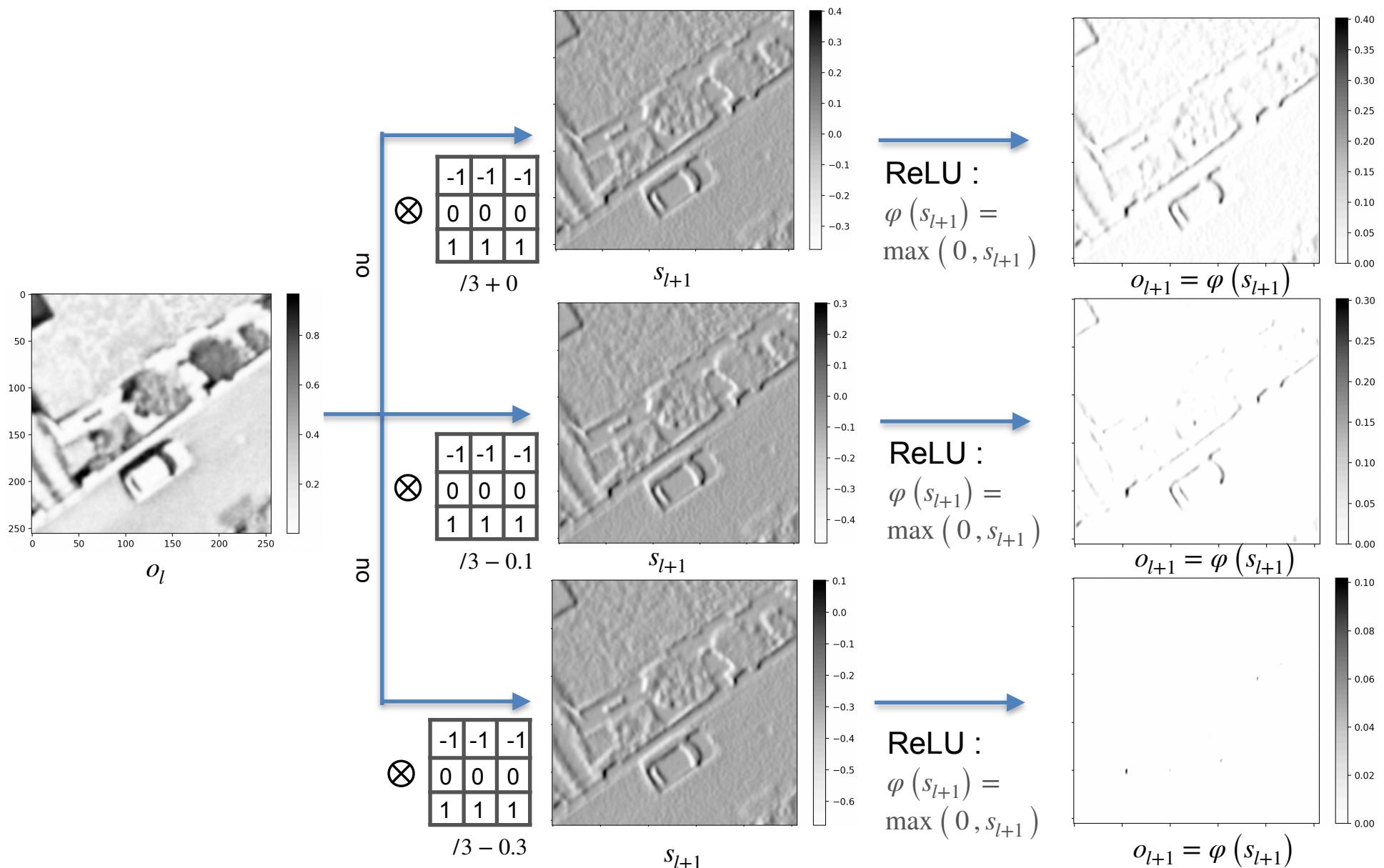
2) Réseaux convolutifs

Briques de base : Convolution et biais

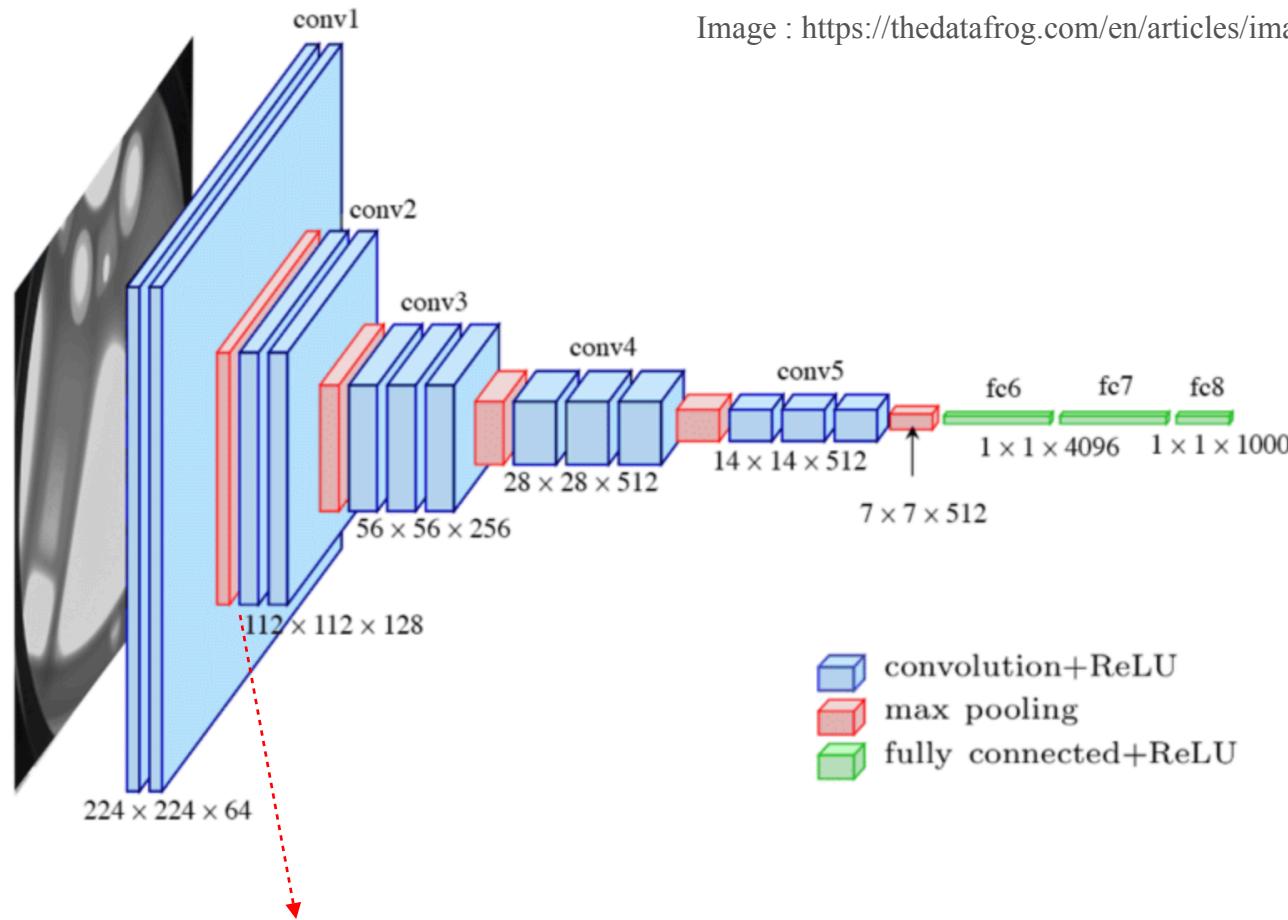


2) Réseaux convolutifs

Briques de base : Convolution et biais suivi de ReLU



2) Réseaux convolutifs



Convolution de [64 canaux de dimension 112×112] vers [128 canaux de dimension 112×112] pour l'observation b du mini-batch

$$\left(o_{c,f,g}^l \right)_b^b \quad c = 1, \dots, 64 \\ f = 1, \dots, 112 \\ g = 1, \dots, 112 \\ b = 1, \dots, B$$

$$\left(o_{c,f,g}^{l+1} \right)_b^b \quad c = 1, \dots, 128 \\ f = 1, \dots, 112 \\ g = 1, \dots, 112 \\ b = 1, \dots, B$$

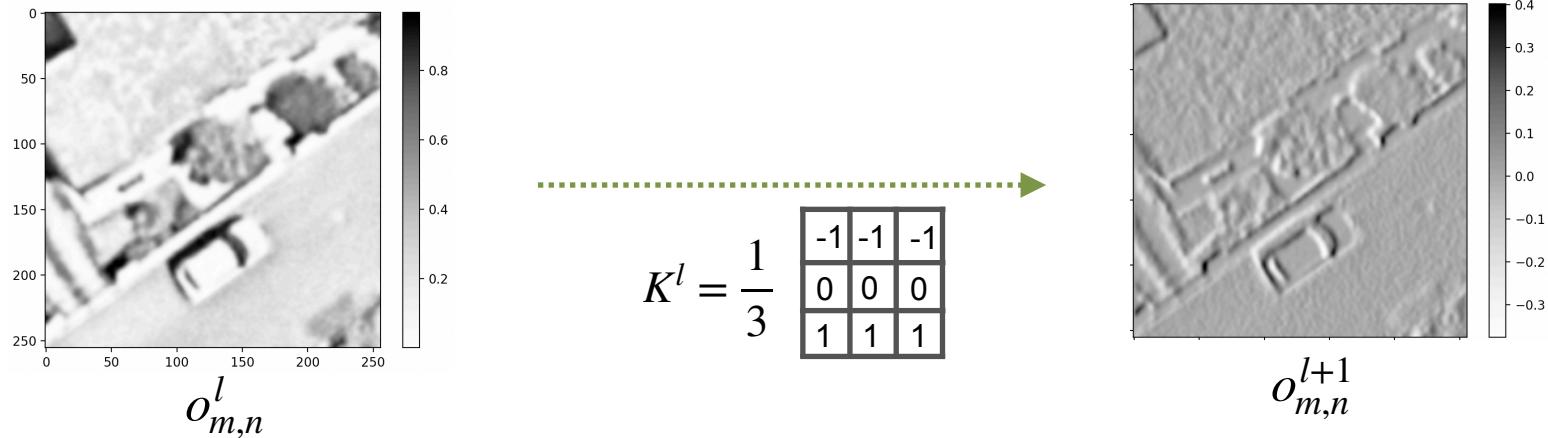
$$\left(s_{c,f,g}^{l+1} \right)_b^b = bias_c^l + \sum_{c'=1}^{64} \sum_{v=-V}^V \sum_{w=-W}^W K_{c',c,v,w}^l \left(o_{c',f+v,g+w}^l \right)_b^b$$

$$\left(o_{c,f,g}^{l+1} \right)_b^b = \varphi \left(\left(s_{c,f,g}^{l+1} \right)_b^b \right) \quad \forall c = 1, \dots, 128$$

2) Calcul d'un gradient dans un réseau de neurones

TP2

On se restreint à la convolution d'un canal à un canal...
... mais on va optimiser en CUDA le code la backpropagation d'une couche convulsive !



Modèle forward :

$$\begin{cases} s_{f,g}^{l+1} = \sum_{v=-V}^V \sum_{w=-W}^W K_{v,w}^l o_{f+v,g+w}^l \\ o_{f,g}^{l+1} = \varphi(s_{f,g}^{l+1}) \end{cases}$$

Equations backward :

$$\frac{\partial R_\theta^i}{\partial o_{m,n}^l} = \sum_{v=-V}^V \sum_{w=-W}^W \frac{\partial R_\theta^i}{\partial o_{m+v,n+w}^{l+1}} \frac{\partial o_{m+v,n+w}^{l+1}}{\partial s_{m+v,n+w}^{l+1}} \frac{\partial s_{m+v,n+w}^{l+1}}{\partial o_{m,n}^l} = \frac{\partial R_\theta^i}{\partial o_{m+v,n+w}^{l+1}} \varphi'(s_{m+v,n+w}^{l+1}) K_{-v,-w}^l$$

$$\frac{\partial R_\theta^i}{\partial K_{v,w}^l} = \sum_{m=1}^M \sum_{n=1}^N \frac{\partial R_\theta^i}{\partial o_{m,n}^{l+1}} \frac{\partial o_{m,n}^{l+1}}{\partial s_{m,n}^{l+1}} \frac{\partial s_{m,n}^{l+1}}{\partial K_{v,w}^l} = \frac{\partial R_\theta^i}{\partial o_{m,n}^{l+1}} \varphi'(s_{m,n}^{l+1}) o_{m+v,n+w}^l$$