

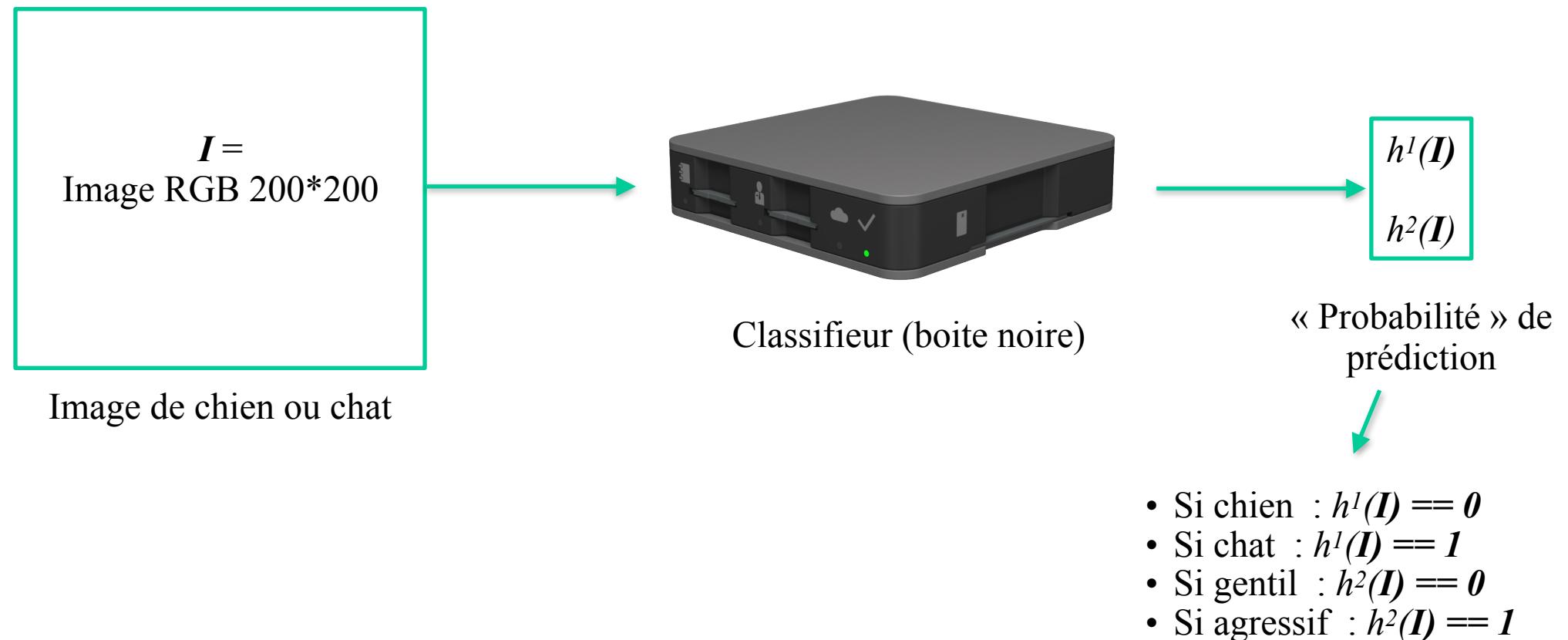
Introduction au calcul GP-GPU

Partie 3 : Intérêt en deep-learning

Laurent Risser

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

Prédiction



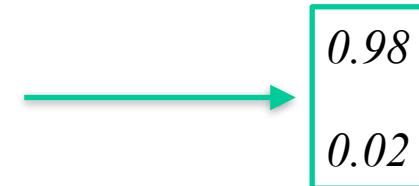
Prédiction



Image de chien ou chat



Classifieur (boite noire)



0.98
0.02

- Si chien : $h^1(\mathbf{I}) == 0$
- Si chat : $h^1(\mathbf{I}) == 1$
- Si gentil : $h^2(\mathbf{I}) == 0$
- Si agressif : $h^2(\mathbf{I}) == 1$

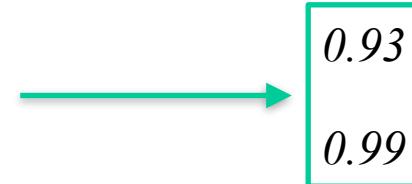
Prédiction



Image de chien ou chat



Classifieur (boite noire)



0.93
0.99

- Si chien : $h^1(\mathbf{I}) == 0$
- Si chat : $h^1(\mathbf{I}) == 1$
- Si gentil : $h^2(\mathbf{I}) == 0$
- Si agressif : $h^2(\mathbf{I}) == 1$

Prédiction

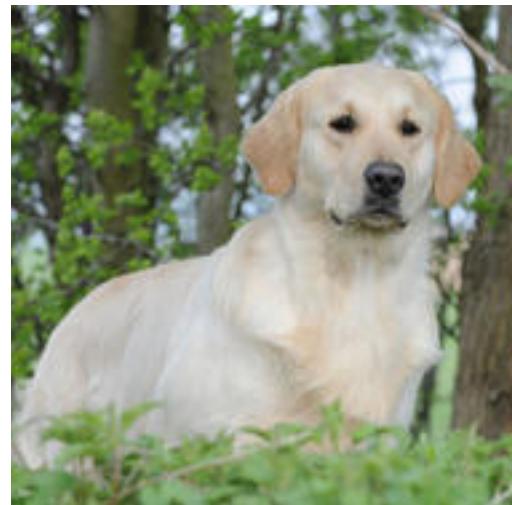
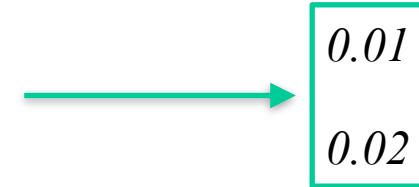


Image de chien ou chat



Classifieur (boite noire)

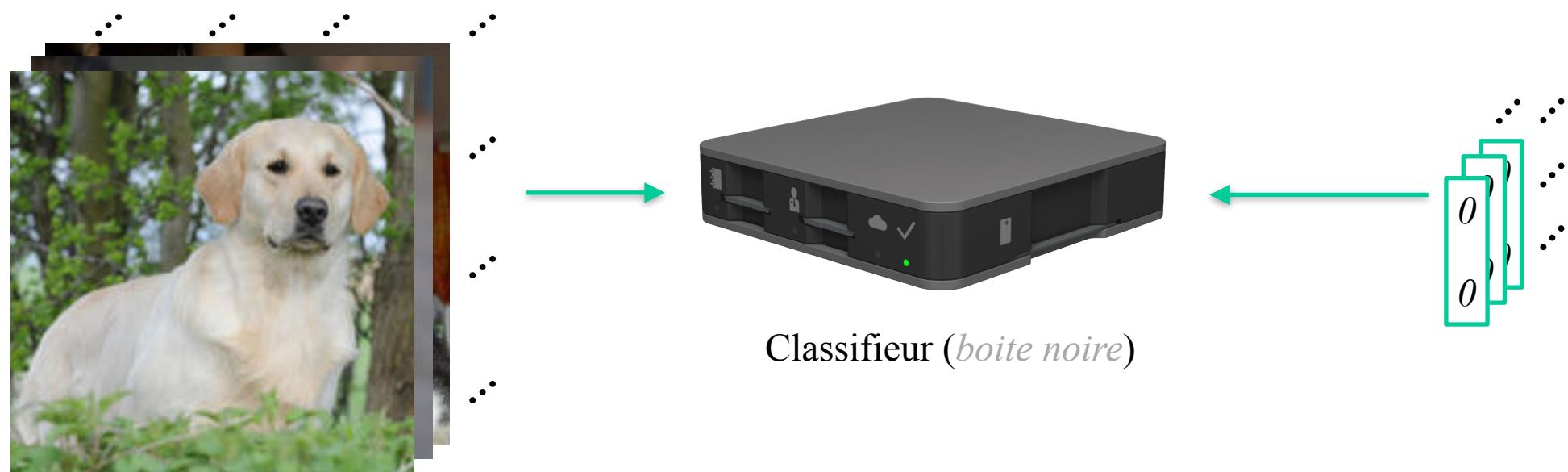


0.01
0.02

- Si chien : $h^1(\mathbf{I}) == 0$
- Si chat : $h^1(\mathbf{I}) == 1$
- Si gentil : $h^2(\mathbf{I}) == 0$
- Si agressif : $h^2(\mathbf{I}) == 1$

Deep learning – Principe de l'apprentissage supervisé

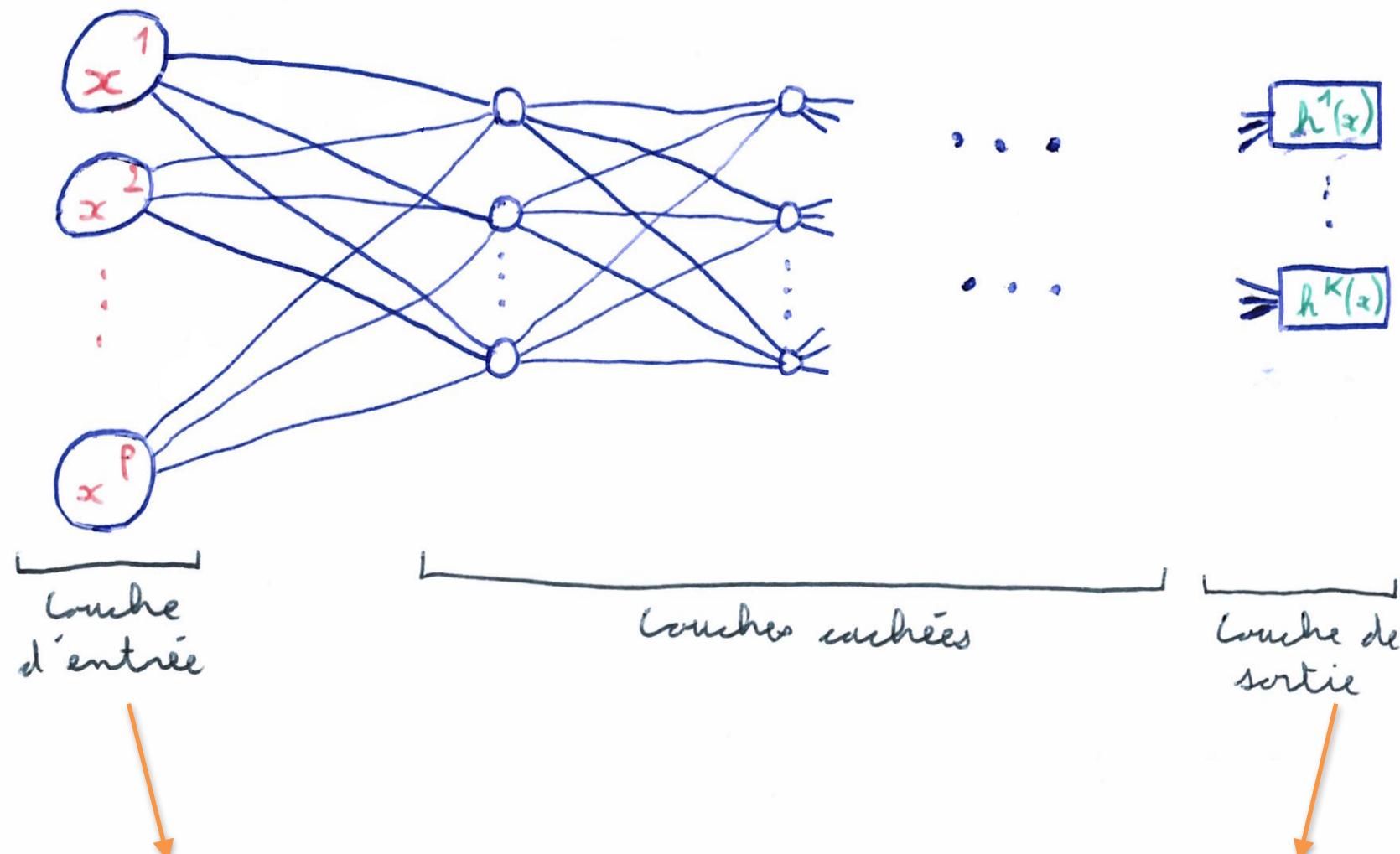
Apprentissage : Optimisation des paramètres pour obtenir les meilleures prédictions en moyenne



Données d'apprentissage (1/2)
Beaucoup d'images de chien ou chat

Données d'apprentissage (2/2)
Labels de chaque image

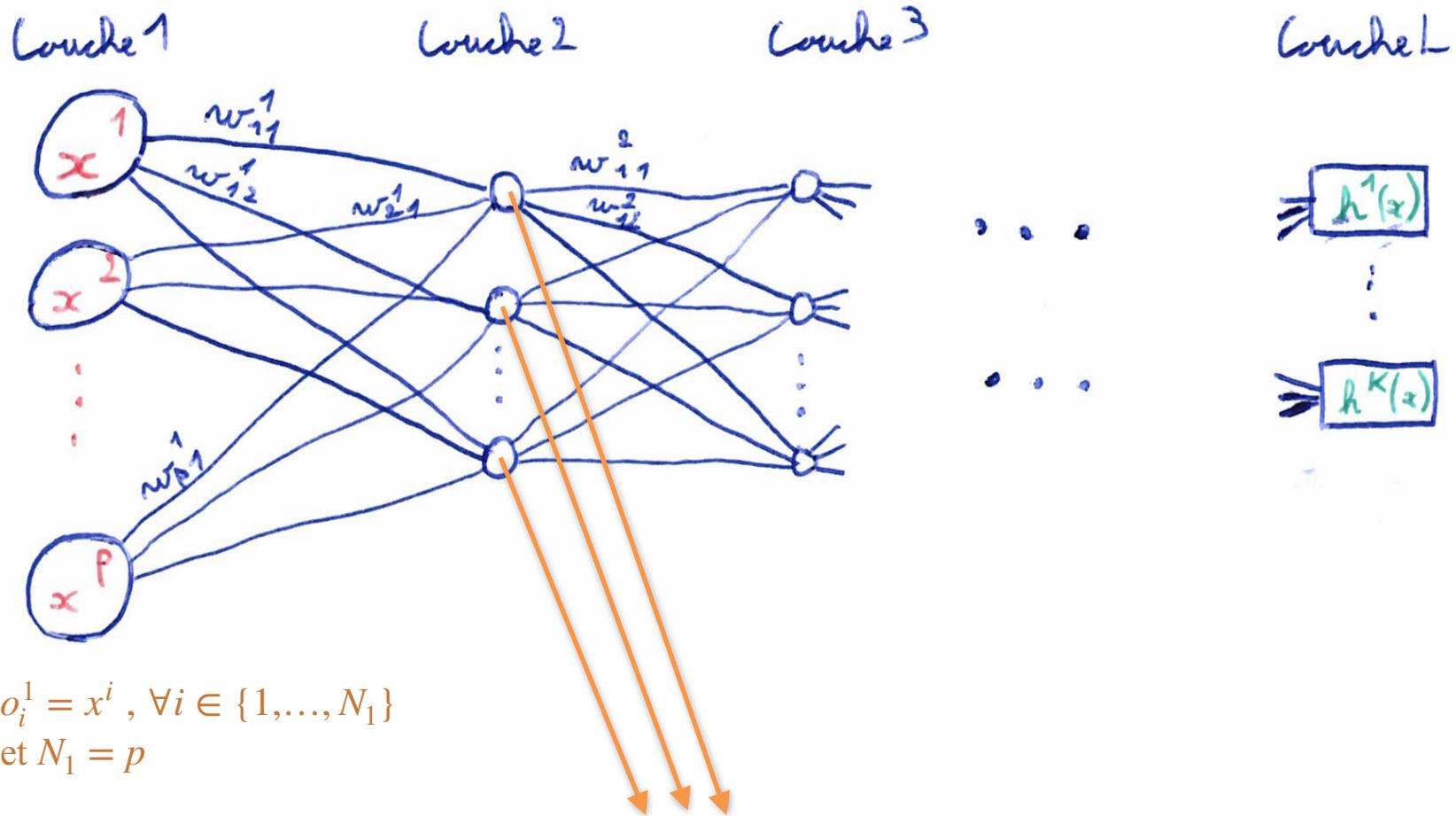
Deep learning – Apprentissage d'un perceptron multi-couches



Les x^i peuvent être les intensités de l'image I sur chaque canal RGB

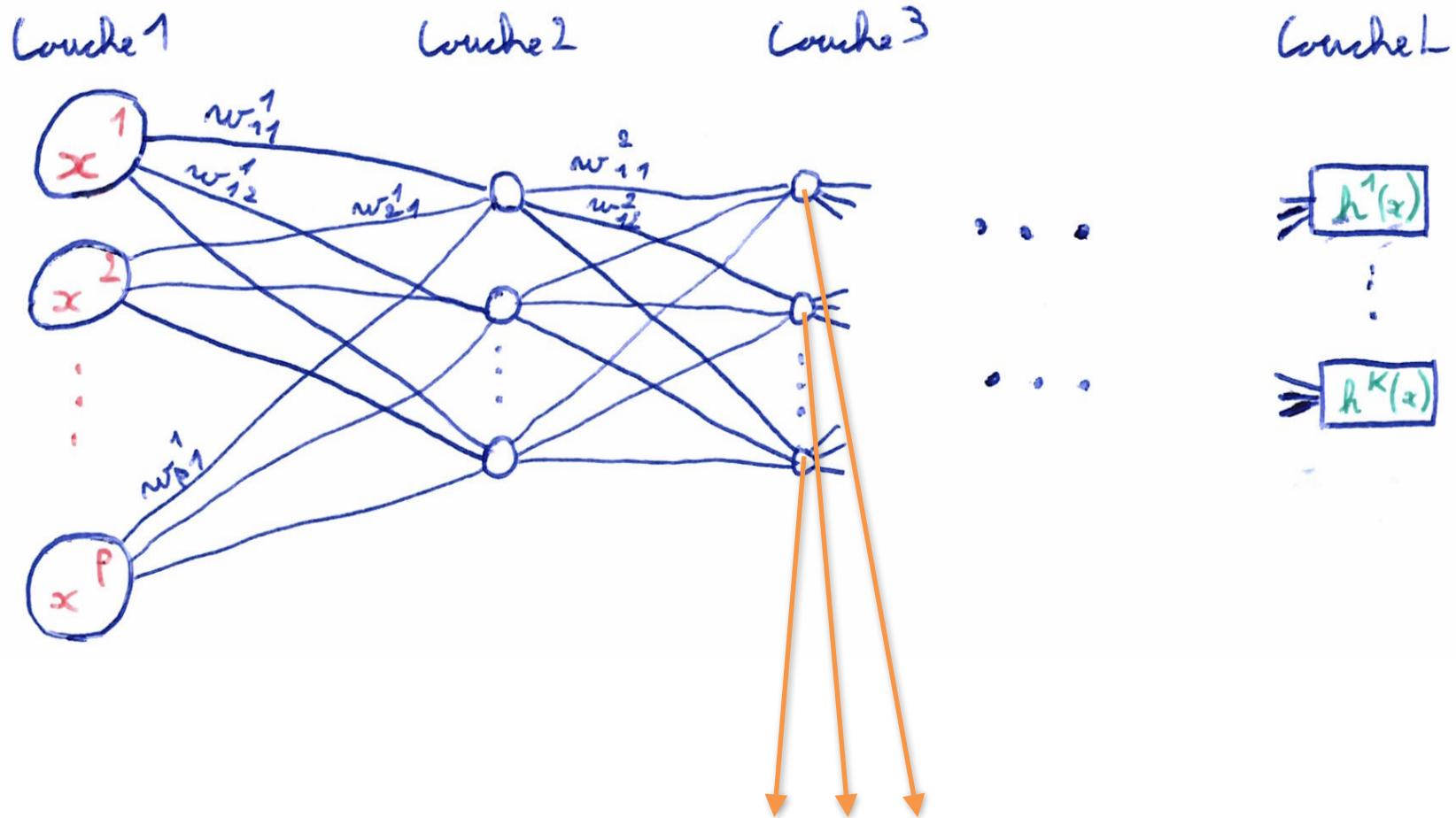
Labels prédits

Deep learning – Apprentissage d'un perceptron multi-couches – Prédiction



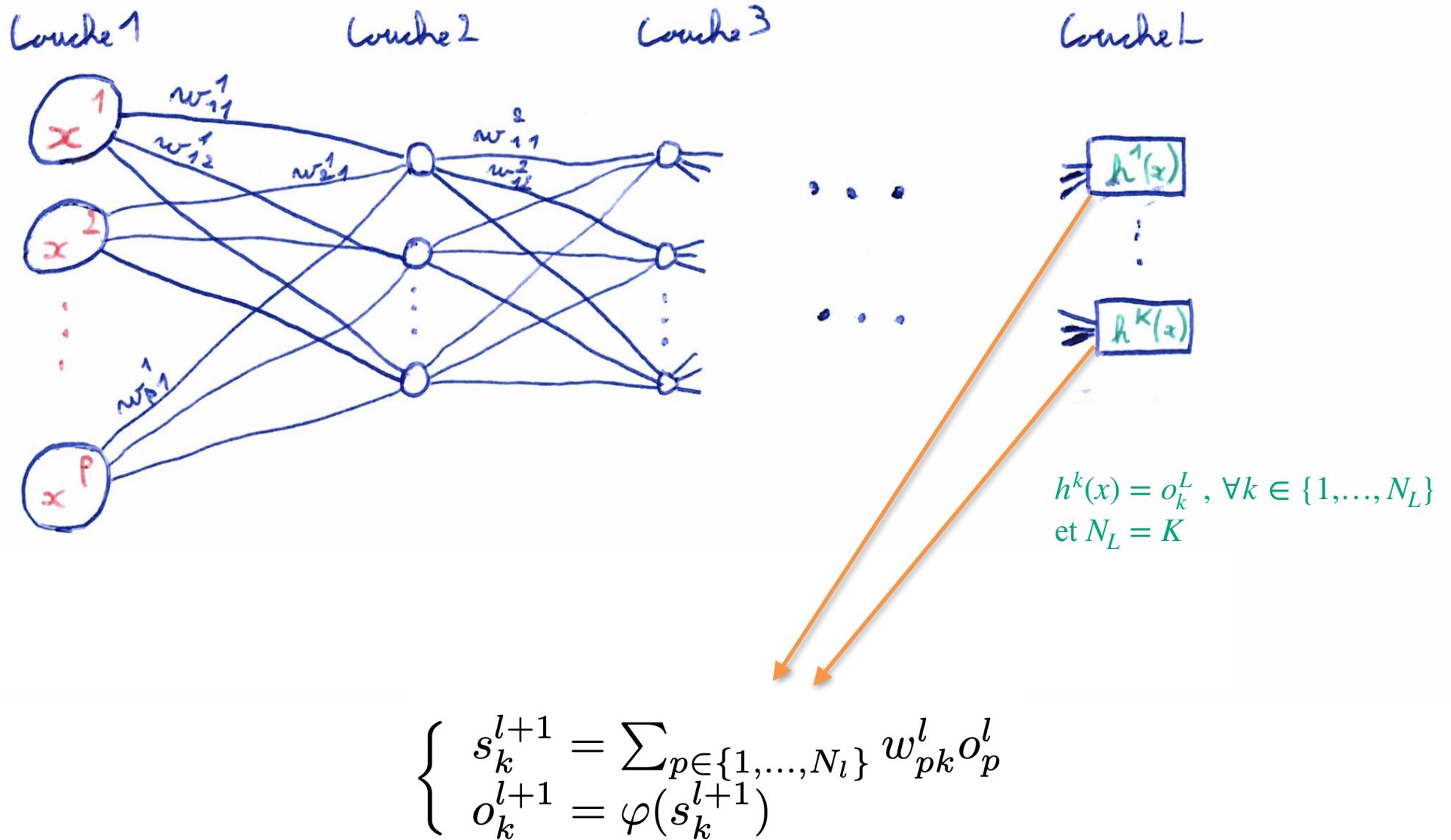
$$\begin{cases} s_k^{l+1} = \sum_{p \in \{1, \dots, N_l\}} w_{pk}^l o_p^l \\ o_k^{l+1} = \varphi(s_k^{l+1}) \end{cases}$$

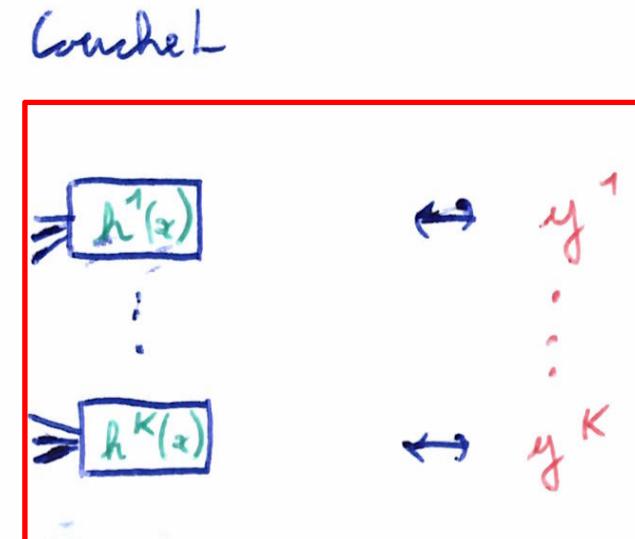
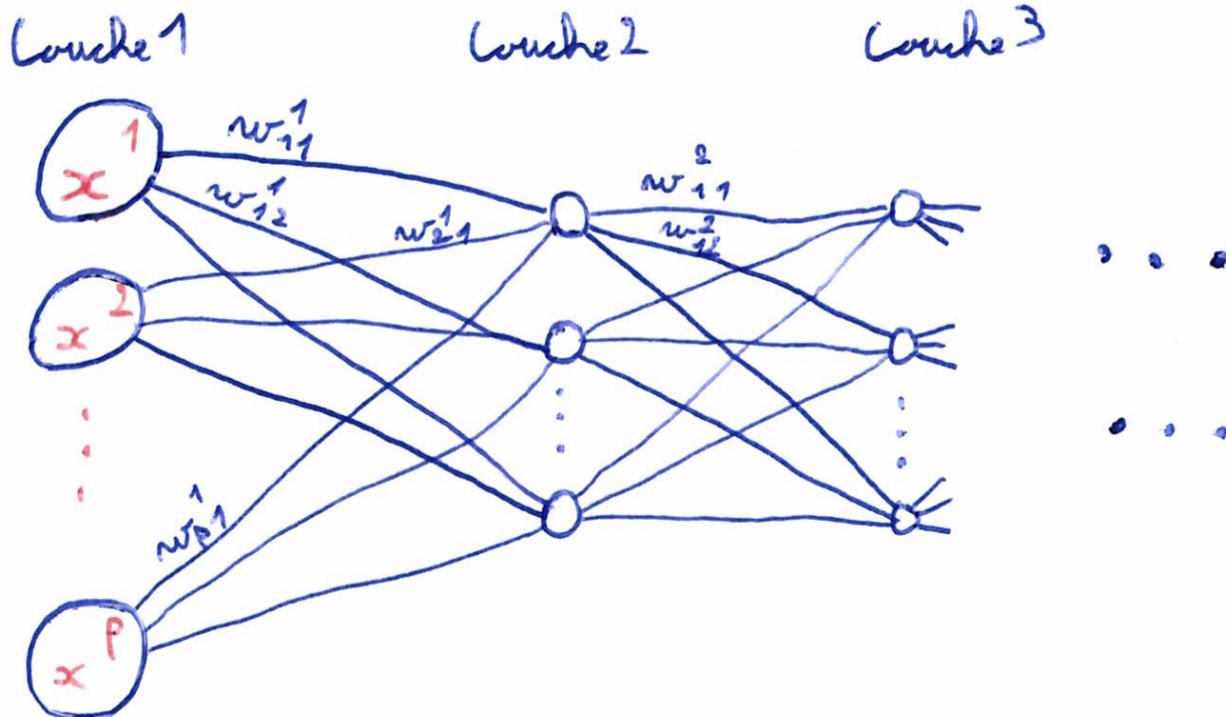
Deep learning – Apprentissage d'un perceptron multi-couches – Prédiction



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

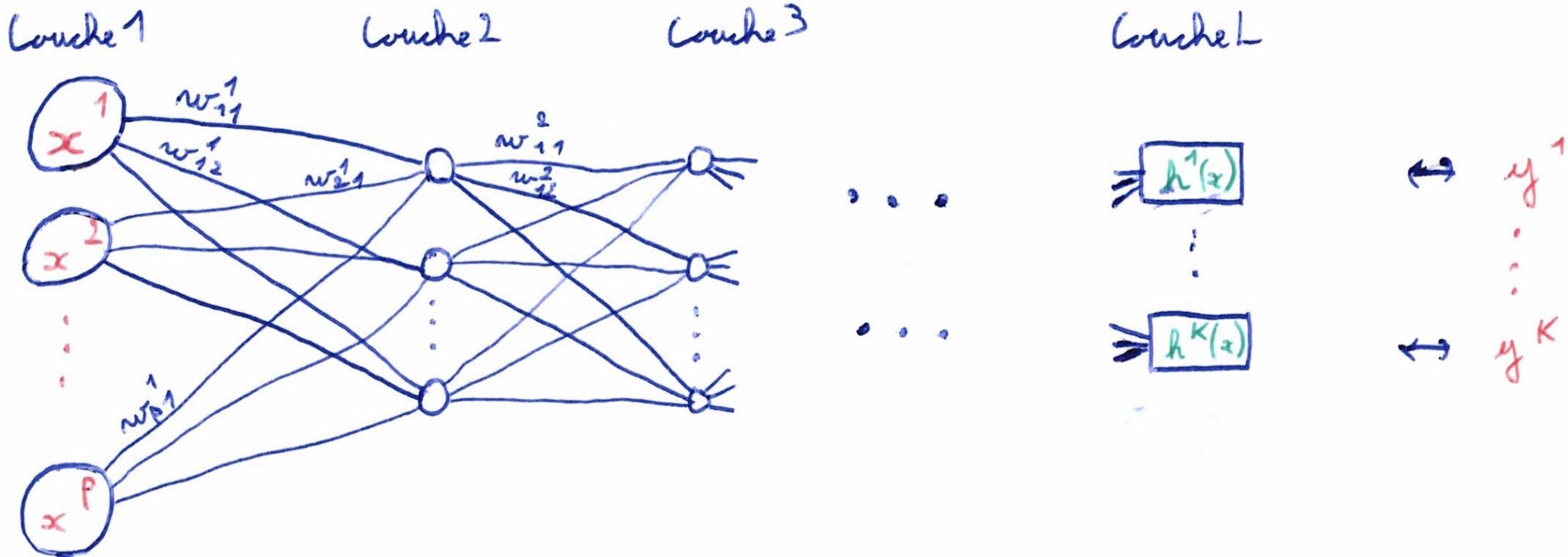
Deep learning – Apprentissage d'un perceptron multi-couches – Prédiction





Erreur d'estimation dans un mini-batch :

$$E \left(\underbrace{\mathbf{x}^1, \dots, \mathbf{x}^P}_{\mathbf{x}} \right) = \sum_{k=1}^K (h^k(\mathbf{x}) - y^k)^2$$

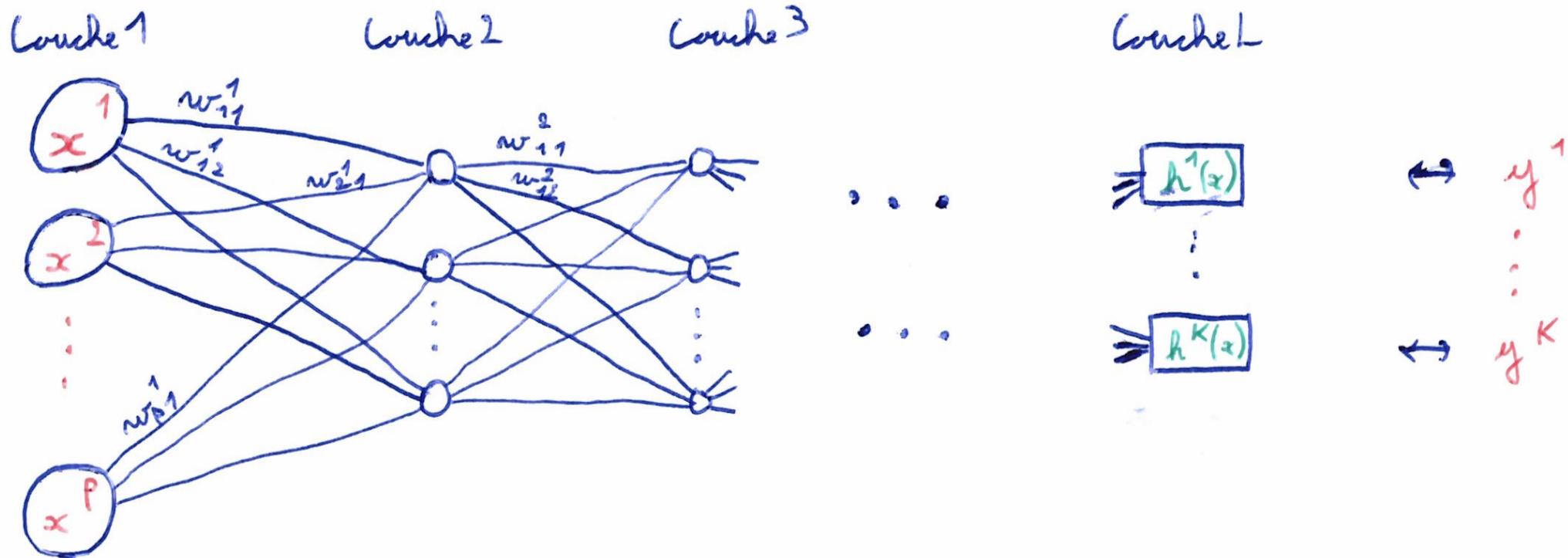


Erreur d'estimation dans un mini-batch :

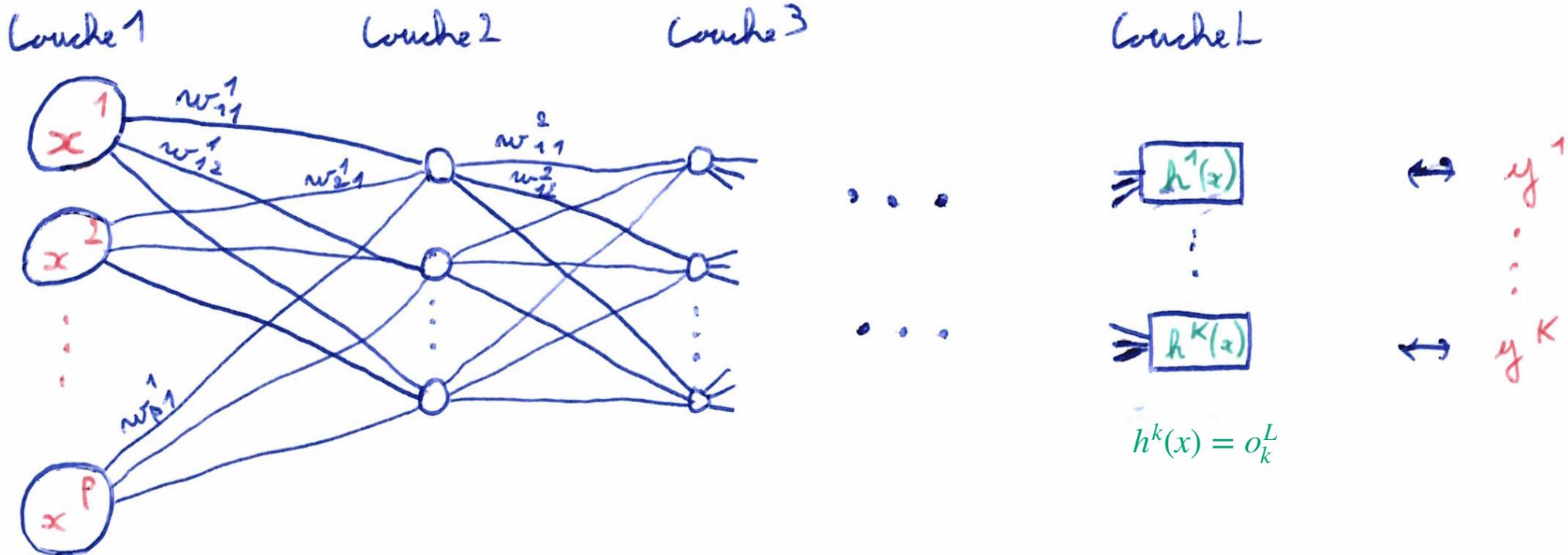
$$E \left(\underbrace{x^1, \dots, x^p}_{x} \right) = \sum_{k=1}^K (h^k(x) - y^k)^2$$

Apprentissage par descente de gradient stochastique \approx

- Descente de gradient sur l'erreur moyenne (*risque empirique*) par rapport aux paramètres w_{pk}^l du réseau
- A chaque itération, l'erreur moyenne est approchée avec un sous-ensemble des observations (*mini-batch*)



Itération de la descente de gradient stochastique : $w_{pk}^l = w_{pk}^l - \eta \frac{\partial E}{\partial w_{pk}^l}, \forall \{p, k, l\}$

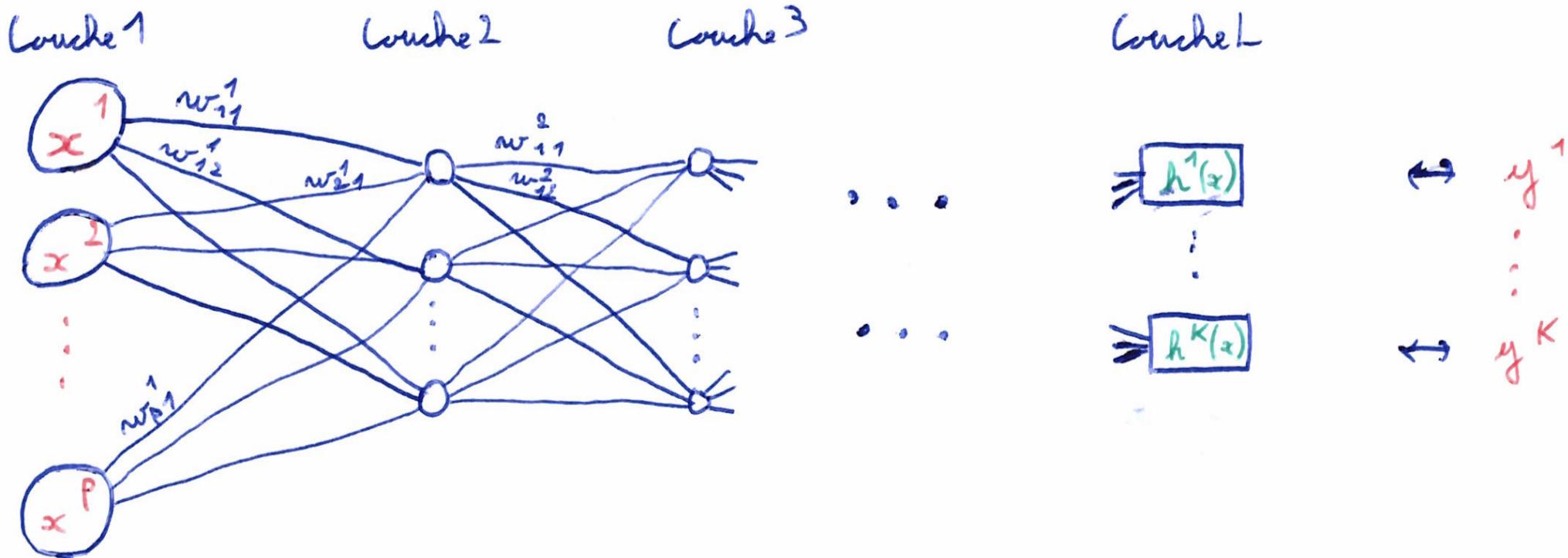


Calcul des dérivées $\frac{\partial E}{\partial w_{pk}^l}$ du risque empirique :

(1) Calcul analytique de $\frac{\partial E}{\partial o_k^L}$ à la sortie du réseau ; par exemple avec la MSE :

$$\frac{\partial E}{\partial o_k^L} = \frac{\partial}{\partial o_k^L} \frac{1}{2n} \sum_{i=1}^n (\underbrace{h^k(x_i)}_{=o_k^L \text{ for obs. } i} - y_i^k)^2 = \frac{1}{n} \sum_{i=1}^n (h^k(x_i) - y_i^k)$$

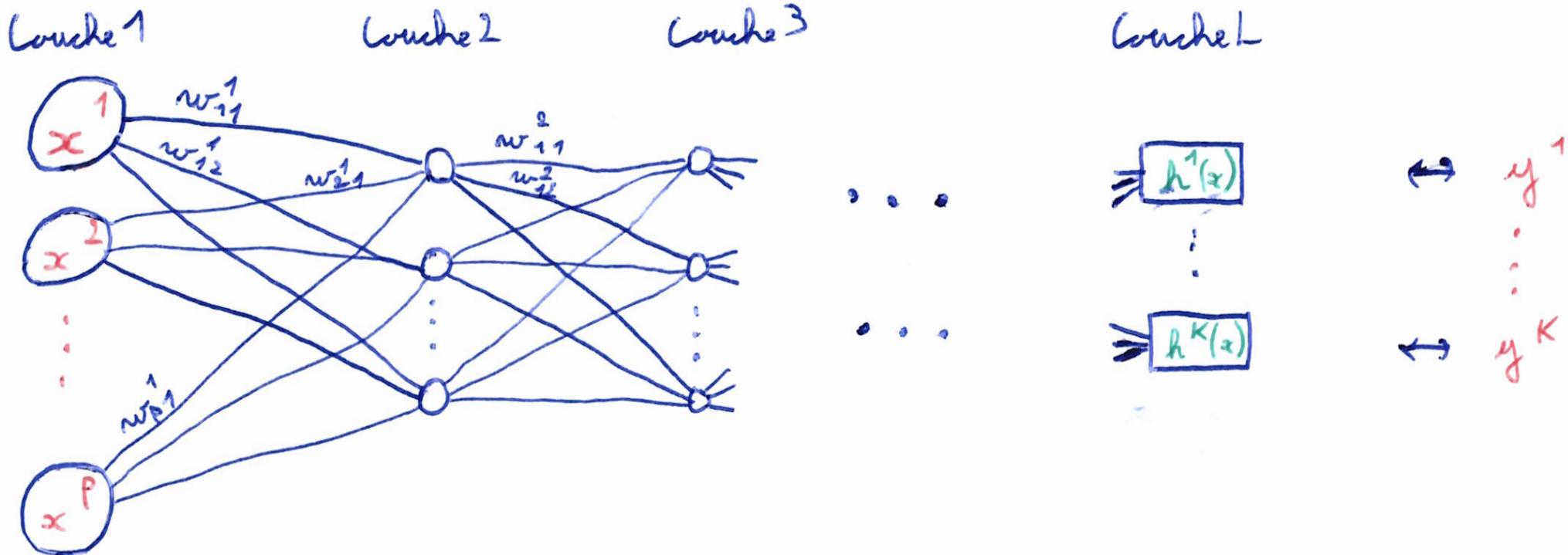
Deep learning – Apprentissage d'un perceptron multi-couches – Apprentissage



Calcul des dérivées $\frac{\partial E}{\partial w_{pk}^l}$ du risque empirique :

$$(2) \text{ Backpropagation des } \frac{\partial E}{\partial o_k^l} : \quad \frac{\partial E}{\partial o_k^l} = \sum_{p=1}^{N_{l+1}} \underbrace{\frac{\partial E}{\partial o_k^{l+1}}}_{\substack{\text{Connu à la couche } l+1 \\ \vdots}} \underbrace{\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}}}_{\substack{\text{Dérivée de } \varphi \\ \vdots}} \underbrace{\frac{\partial s_k^{l+1}}{\partial o_k^l}}_{=w_{pk}^l}$$

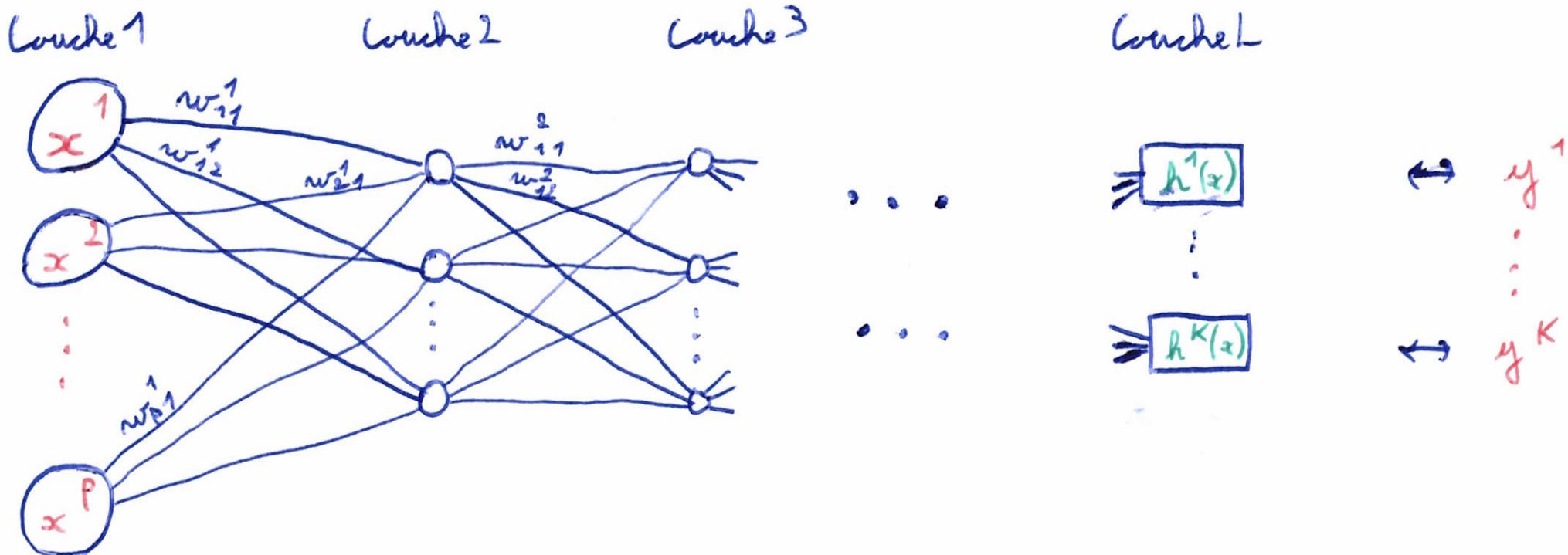
Deep learning – Apprentissage d'un perceptron multi-couches – Apprentissage



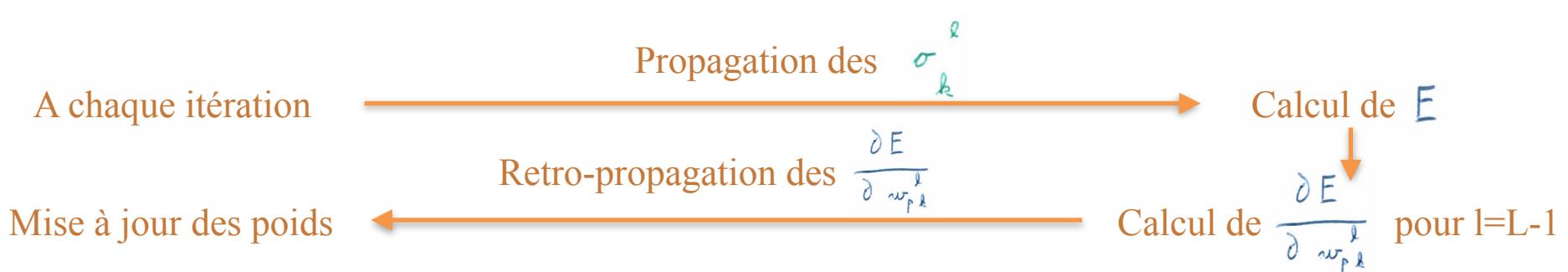
Calcul des dérivées $\frac{\partial E}{\partial w_{pk}^l}$ du risque empirique :

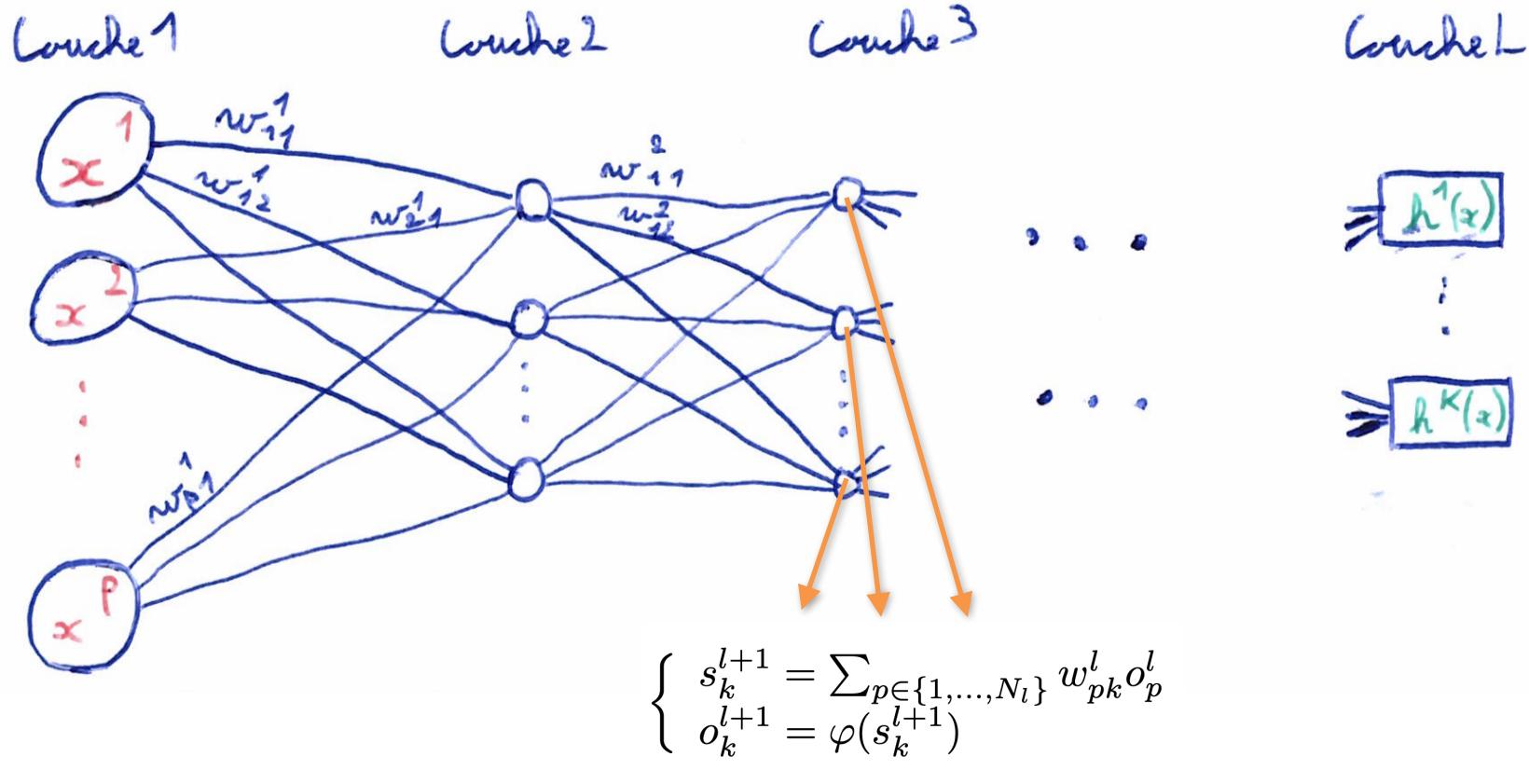
$$(3) \text{ Calcul en parallèle des } \frac{\partial E}{\partial w_{pk}^l} : \quad \frac{\partial E}{\partial o_k^l} = \sum_{p=1}^{N_{l+1}} \underbrace{\frac{\partial E}{\partial o_k^{l+1}}}_{\substack{\text{Connu à la couche } l+1 \\ \vdots}} \quad \underbrace{\frac{\partial o_k^{l+1}}{\partial s_k^{l+1}}}_{\substack{\text{Dérivée de } \varphi \\ = w_{pk}^l}} \quad \underbrace{\frac{\partial s_k^{l+1}}{\partial o_k^l}}_{\substack{\text{Connu à la couche } l+1 \\ \vdots}}$$

$$\boxed{\frac{\partial E}{\partial w_{pk}^l} = \frac{\partial E}{\partial o_k^{l+1}} \frac{\partial o_k^{l+1}}{\partial s_k^{l+1}} \frac{\partial s_k^{l+1}}{\partial w_{pk}^l}}$$

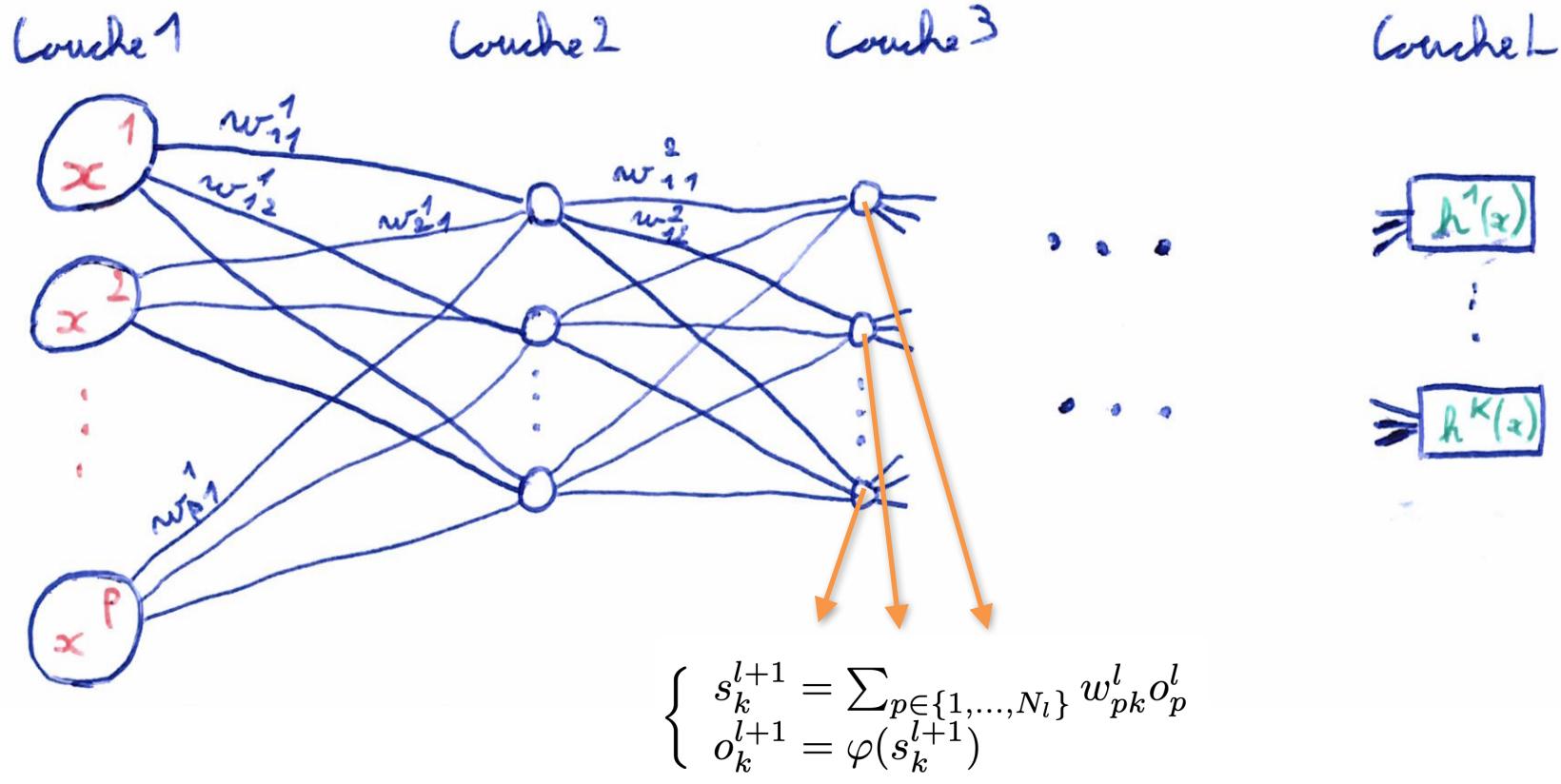


Descente de gradient stochastique (principe) :





Forme matricielle pour le calcul des o_k^l de la couche l ?



Forme matricielle pour le calcul des o_k^l de la couche l ?

$$S^l = W^{l-1} O^{l-1} \text{ puis } O^l = \varphi(S^l)$$

$$\in \mathbb{R}^{(\#N_l) \times (\#N_{l+1})}$$

$$\in \mathbb{R}^{\#N_{l+1}}$$

Deep learning – Apprentissage d'un perceptron multi-couches - Lien avec le GPU

```
__kernel void mv_mul(const int N, __global float *W_lm1, __global float *O_lm1, __global float *S_1)
{
    int i;
    int k = get_global_id(0);
    float tmp; //in private memory

    tmp = 0.0f;
    for (i = 0; i < N; i++)
        tmp += W_lm1[k*N+i]*O_lm1[i];
    S_1[k] = tmp;
}
```



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

Forme matricielle pour le calcul des o_k^l de la couche l ?

$$S^l = W^{l-1} O^{l-1} \text{ puis } O^l = \varphi(S^l)$$

Deep learning – Apprentissage d'un perceptron multi-couches - Lien avec le GPU

```
__kernel void mv_mul(const int N, __global float *W_lml, __global float *O_lml, __global float *S_l)  
{  
    int i;  
    int k = get_global_id(0);  
    float tmp; //in private memory  
  
    tmp = 0.0f;  
    for (i = 0; i < N; i++)  
        tmp += W_lml[k*N+i]*O_lml[i];  
    S_l[k] = tmp;  
}
```

OpenCL

```
import tensorflow as tf  
N=784  
O_lml = tf.placeholder(tf.float32, [None, N])  
W_lml = tf.Variable(tf.zeros([N, 128]))  
o_l = tf.matmul(O_lml, W_lml)
```

Tensorflow/Python

Forme matricielle pour le calcul des o_k^l de la couche l ?

$$S^l = \boxed{W^{l-1} O^{l-1}} \text{ puis } O^l = \varphi(S^l)$$

Parallélisation encore plus intéressantes dans des réseaux convolutionnels :

- Modèle *forward* :

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

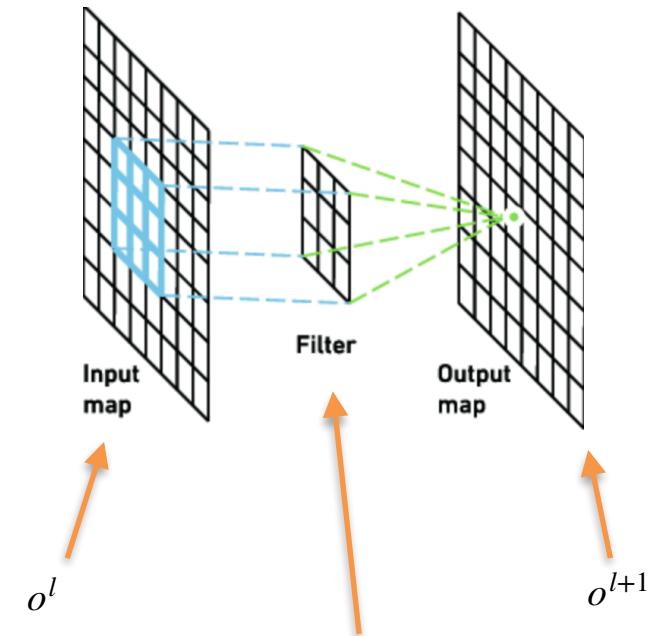
- Modèle *backward* :

$$\frac{\partial E}{\partial o_{m,n}^l} = \sum_{v=-V}^V \sum_{w=-W}^W \frac{\partial E}{\partial o_{m+v, n+w}^{l+1}} \frac{\partial o_{m+v, n+w}^{l+1}}{\partial s_{m+v, n+w}^{l+1}} \frac{\partial s_{m+v, n+w}^{l+1}}{\partial o_{m,n}^l}$$

$$= \frac{\partial E}{\partial o_{m+v, n+w}^{l+1}} \varphi'(s_{m+v, n+w}^{l+1}) K_{-v, -w}^l$$

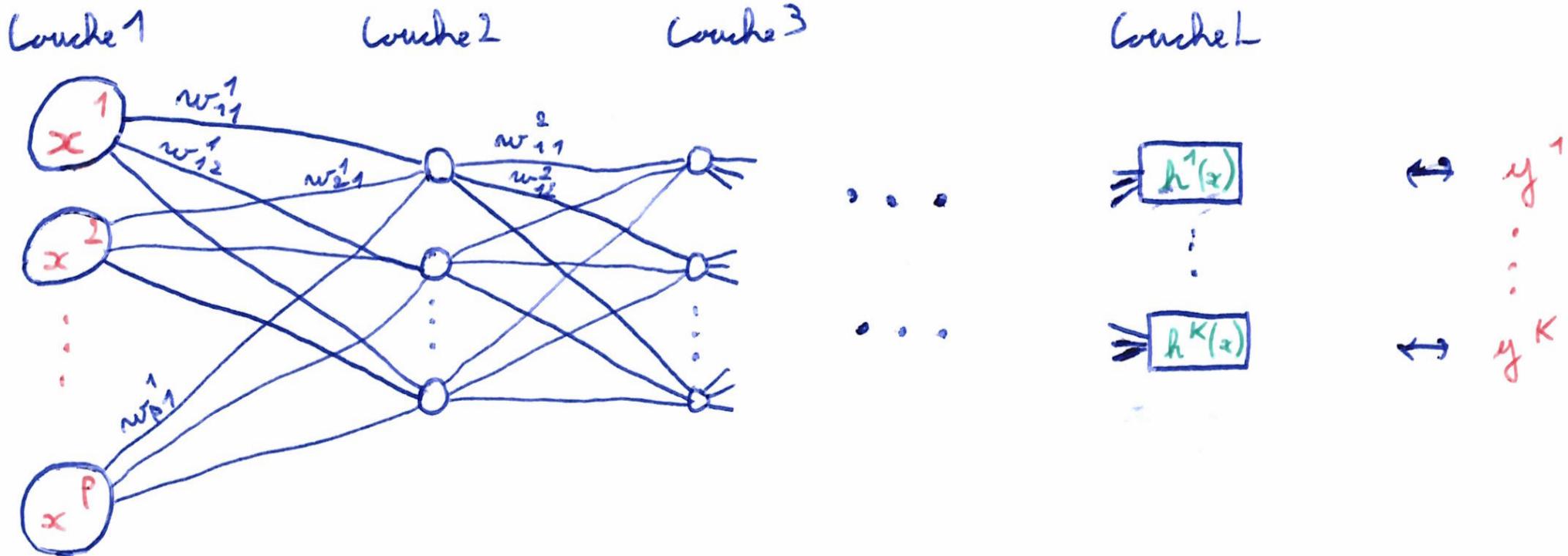
$$\frac{\partial E}{\partial K_{v,w}^l} = \sum_{m=1}^M \sum_{n=1}^N \frac{\partial E}{\partial o_{m,n}^{l+1}} \frac{\partial o_{m,n}^{l+1}}{\partial s_{m,n}^{l+1}} \frac{\partial s_{m,n}^{l+1}}{\partial K_{v,w}^l}$$

$$= \frac{\partial E}{\partial o_{m,n}^{l+1}} \varphi'(s_{m,n}^{l+1}) o_{m+v, n+w}^l$$



$$K_{v,w}^l, \begin{cases} v \in \{-V, \dots, V\} \\ w \in \{-W, \dots, W\} \end{cases}$$

Deep learning – Pour conclure



Dans le cas général :

- Prédiction et apprentissage particulièrement simples à paralléliser sur GPU
- Librairie Nvidia cuDNN massivement utilisée et sous jacente à Caffe, Keras, TensorFlow, Theano...