

ANNEX 2: Gradient descent in multi-layer networks

- We define an error function to minimize:

$$E(y) = \frac{1}{2} \cdot (r - y)^2$$

- For each weight w_{ij}^l , we characterize the impact of a change on E
 - Gradient descent on w_{ij}^l

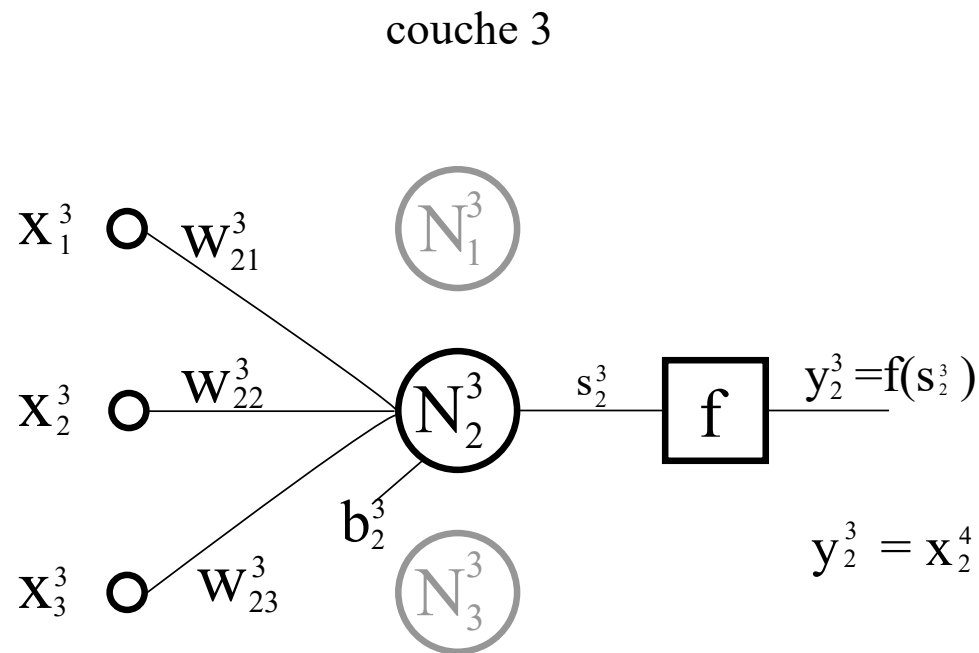
$$w_{ij}^l \leftarrow w_{ij}^l - \alpha \cdot \frac{\partial E(y)}{\partial w_{ij}^l}$$

- We have to find: $\frac{\partial E(y)}{\partial w_{ij}^l}$

ANNEX 2: Gradient descent in multi-layer networks

- **Gradient descent: notations**

- We note l a layer, i the index of a neuron, L the last (output) layer.
- A weight w_{ij}^l connects neuron N_i^l and neuron N_j^{l-1} ($N_j^{l-1} \rightarrow N_i^l$)
- We note $s_i^l = s(W_i^l, X^l, b_i^l) = \sum w_{ij}^l \cdot x_j^l + b_i^l$ the weighted sum of neuron N_i^l of layer l .
- We note $y_i^l = f(s_i^l)$ the output of neuron N_i^l after activation function f



ANNEX 2: Gradient descent in multi-layer networks

- **Gradient descent:
case of last layer L**

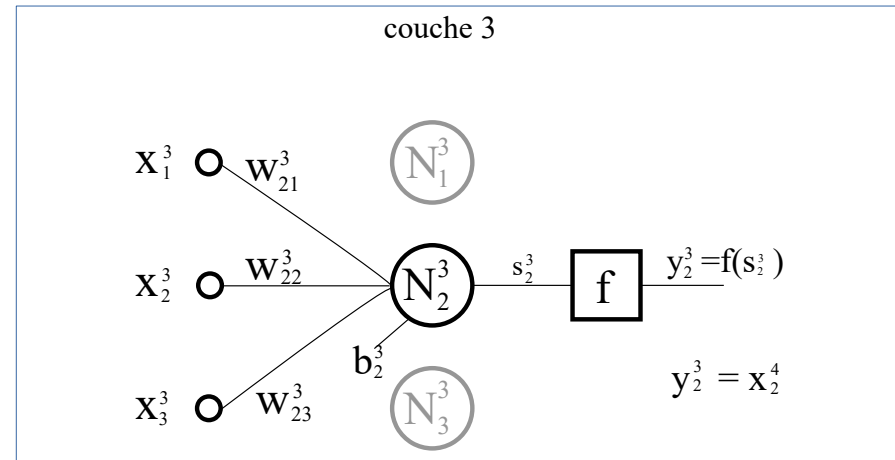
- We must find the term:

$$\frac{\partial E(y_i^L)}{\partial w_{ij}^L}$$

- We can decompose with intermediate values $w \rightarrow s \rightarrow y \rightarrow E$
 - Theorem of composed derivative functions

$$\frac{\partial E(y_i^L)}{\partial w_{ij}^L} = \frac{\partial E(y_i^L)}{\partial y_i^L} \cdot \frac{\partial y_i^L}{\partial s_i^L} \cdot \frac{\partial s_i^L}{\partial w_{ij}^L}$$

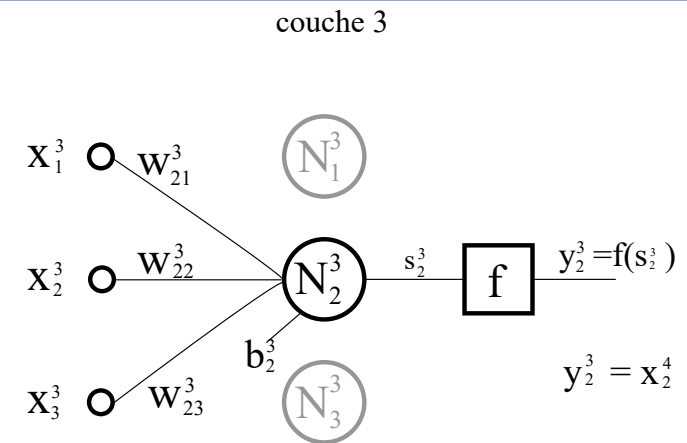
- Three derivatives to find



ANNEX 2: Gradient descent in multi-layer networks

- Gradient descent:
case of last layer L

$$\frac{\partial E(y_i^L)}{\partial w_{ij}^L} = \frac{\partial E(y_i^L)}{\partial y_i^L} \cdot \frac{\partial y_i^L}{\partial s_i^L} \cdot \frac{\partial s_i^L}{\partial w_{ij}^L}$$



- First term: $\frac{\partial E(y_i^L)}{\partial y_i^L}$

Derivation of composed functions:
 $(f \circ g)' = (f' \circ g) \cdot g'$

- Error function: $E(y) = \frac{1}{2} \cdot (r - y)^2$

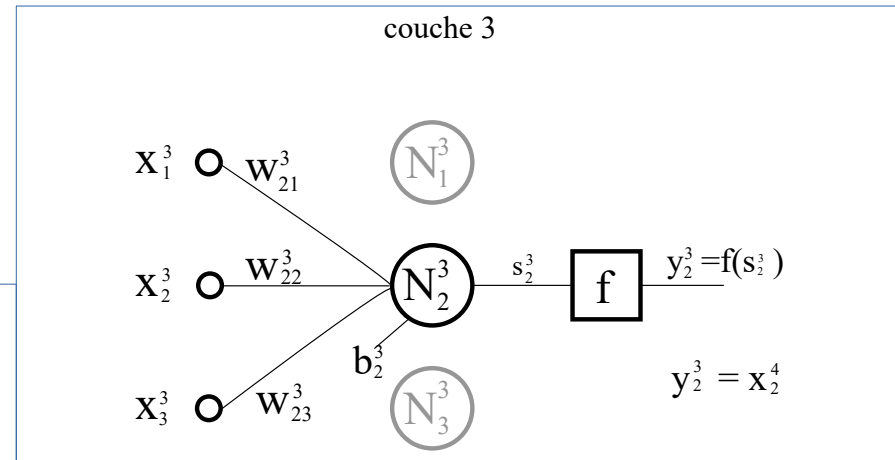
$$\frac{\partial E(y)}{\partial y} = \frac{1}{2} \cdot \frac{\partial (r - y)^2}{\partial y} = \cancel{\frac{2}{2}} \cdot (r - y) \cdot \frac{\partial (r - y)}{\partial y} \rightarrow -1$$

$$\frac{\partial E(y)}{\partial y} = -(r - y)$$

ANNEX 2: Gradient descent in multi-layer networks

- Gradient descent:
case of last layer L

$$\frac{\partial E(y_i^L)}{\partial w_{ij}^L} = \frac{\partial E(y_i^L)}{\partial y_i^L} \cdot \frac{\partial y_i^L}{\partial s_i^L} \cdot \frac{\partial s_i^L}{\partial w_{ij}^L}$$



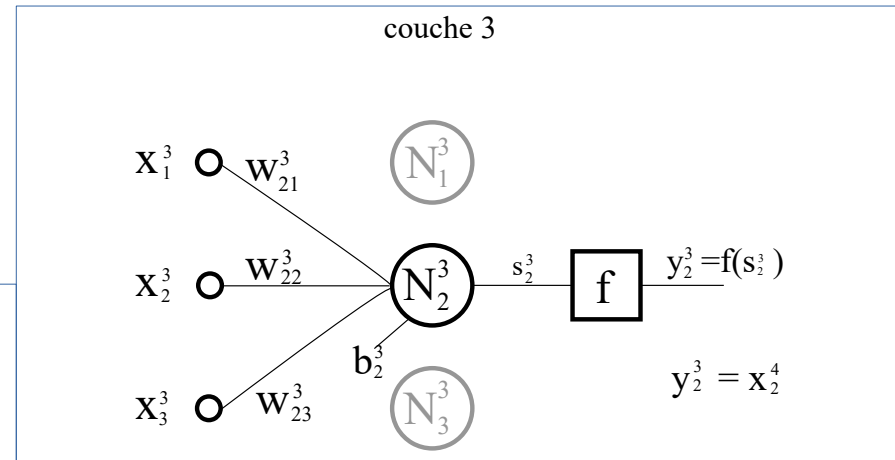
- second term: $\frac{\partial y_i^L}{\partial s_i^L} = \frac{\partial f(s_i^L)}{\partial s_i^L}$
- It is just the derivative of activation function !

$$\frac{\partial y_i^L}{\partial s_i^L} = f'_{(s_i^L)}$$

ANNEX 2: Gradient descent in multi-layer networks

- Gradient descent:
case of last layer L

$$\frac{\partial E(y_i^L)}{\partial w_{ij}^L} = \frac{\partial E(y_i^L)}{\partial y_i^L} \cdot \frac{\partial y_i^L}{\partial s_i^L} \cdot \frac{\partial s_i^L}{\partial w_{ij}^L}$$



- third term: $\frac{\partial s_i^L}{\partial w_{ij}^L}$

'constants'

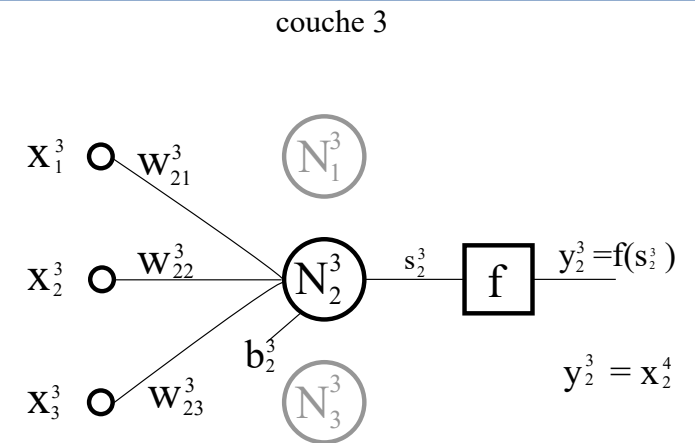
$$\frac{\partial s_i^L}{\partial w_{ij}^L} = \frac{\partial (w_{i1}^L \cdot x_1^L + w_{i2}^L \cdot x_2^L + \dots + w_{ij}^L \cdot x_i^L + \dots + b_i^L)}{\partial w_{ij}^L} = \frac{\partial (w_{ij}^L \cdot x_j^L)}{\partial w_{ij}^L}$$

$$\frac{\partial s_i^L}{\partial w_{ij}^L} = x_j^L$$

ANNEX 2: Gradient descent in multi-layer networks

- Gradient descent:
case of last layer L

$$\frac{\partial E(y_i^L)}{\partial w_{ij}^L} = \frac{\partial E(y_i^L)}{\partial y_i^L} \cdot \frac{\partial y_i^L}{\partial s_i^L} \cdot \frac{\partial s_i^L}{\partial w_{ij}^L}$$



- We summarize:

$$\frac{\partial E(y_i^L)}{\partial w_{ij}^L} = -(r - y_i^L) \cdot f'_{(s_i^L)} \cdot x_j^L$$

- Gradient descent: $w_{ij}^L \Leftarrow w_{ij}^L + \alpha \cdot x_j^L \cdot (r - y_i^L) \cdot f'_{(s_i^L)}$

- We set $\delta_i^L = (r - y_i^L) \cdot f'_{(s_i^L)}$

Note :

$$\delta_i^L = -\frac{\partial E(y_i^L)}{\partial s_i^L}$$

$$\rightarrow w_{ij}^L \Leftarrow w_{ij}^L + \alpha \cdot x_j^L \cdot \delta_i^L$$

ANNEX 2: Gradient descent in multi-layer networks

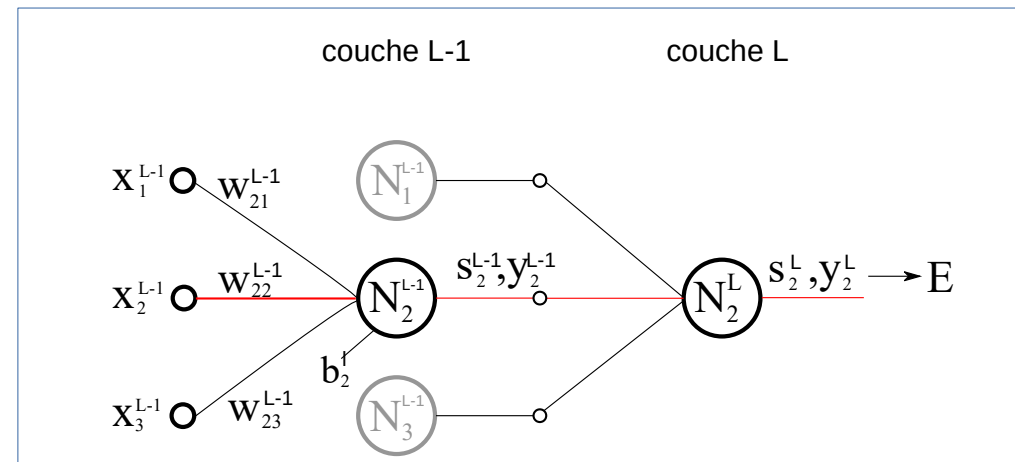
- Gradient descent:
case of hidden layer L-1**

- We must find

$$\frac{\partial E(y_m^L)}{\partial w_{ij}^{L-1}}$$

- We decompose again:

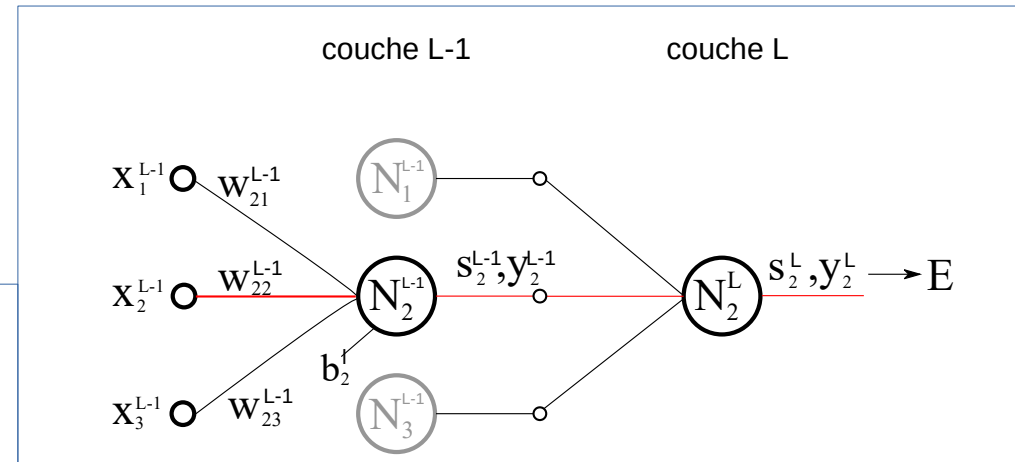
$$\frac{\partial E(y_m^L)}{\partial w_{ij}^{L-1}} = \frac{\partial E(y_m^L)}{\partial y_i^{L-1}} \cdot \frac{\partial y_i^{L-1}}{\partial s_i^{L-1}} \cdot \frac{\partial s_i^{L-1}}{\partial w_{ij}^{L-1}}$$



ANNEX 2: Gradient descent in multi-layer networks

- **Gradient descent:
case of hidden layer L-1**

$$\frac{\partial E(y_m^L)}{\partial w_{ij}^{L-1}} = \frac{\partial E(y_m^L)}{\partial y_i^{L-1}} \cdot \frac{\partial y_i^{L-1}}{\partial s_i^{L-1}} \cdot \frac{\partial s_i^{L-1}}{\partial w_{ij}^{L-1}}$$



- Second and third terms are the same:

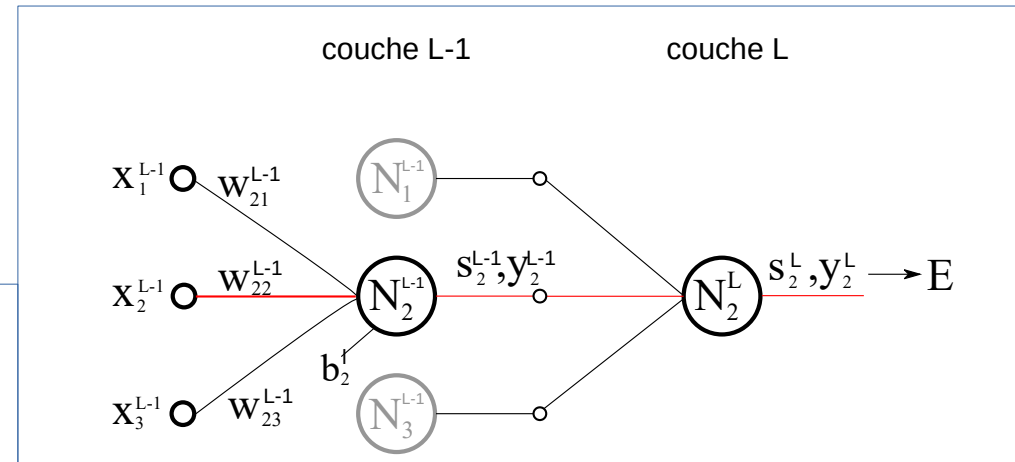
$$\frac{\partial y_i^{L-1}}{\partial s_i^{L-1}} = f'_{(s_i^{L-1})}$$

$$\frac{\partial s_i^{L-1}}{\partial w_{ij}^{L-1}} = x_j^{L-1}$$

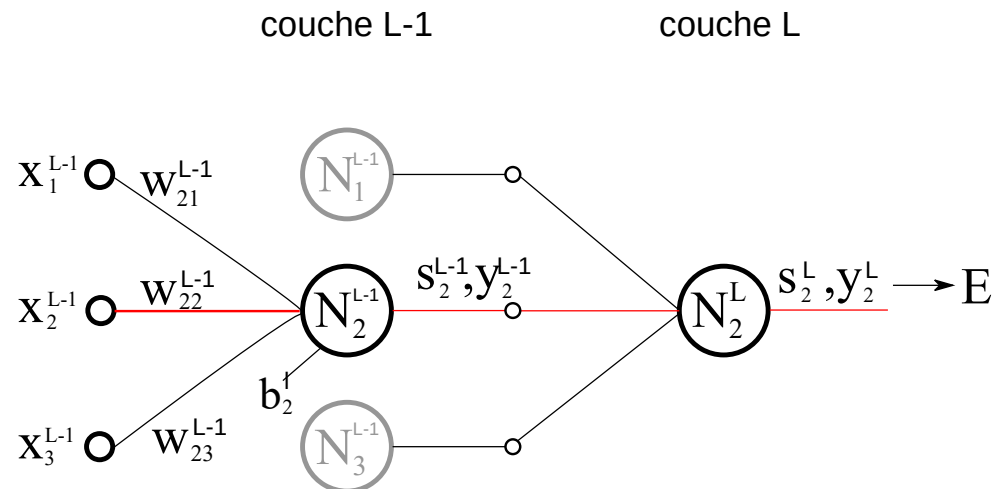
ANNEX 2: Gradient descent in multi-layer networks

- **Gradient descent:**
case of hidden layer L-1

$$\frac{\partial E(y_m^L)}{\partial w_{ij}^{L-1}} = \frac{\partial E(y_m^L)}{\partial y_i^{L-1}} \cdot \frac{\partial y_i^{L-1}}{\partial s_i^{L-1}} \cdot \frac{\partial s_i^{L-1}}{\partial w_{ij}^{L-1}}$$



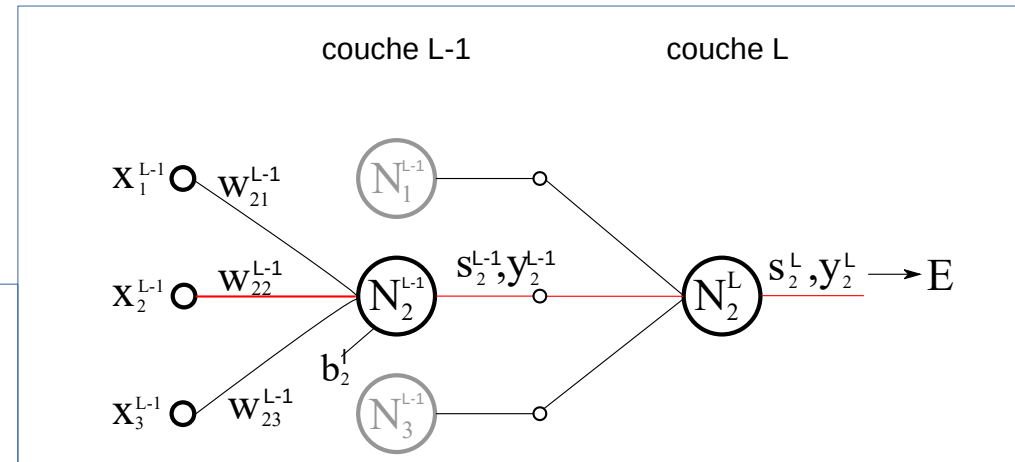
- First term: $\frac{\partial E(y_m^L)}{\partial y_i^{L-1}}$
- More complex: if y_i^{L-1} is modified, E is modified through s_i^{L-1} (and thus y_j^{L-1}) of layer L



ANNEX 2: Gradient descent in multi-layer networks

- **Gradient descent:**
case of hidden layer L-1

$$\frac{\partial E(y_m^L)}{\partial w_{ij}^{L-1}} = \frac{\partial E(y_m^L)}{\partial y_i^{L-1}} \cdot \frac{\partial y_i^{L-1}}{\partial s_i^{L-1}} \cdot \frac{\partial s_i^{L-1}}{\partial w_{ij}^{L-1}}$$



$$\frac{\partial E(y_m^L)}{\partial y_i^{L-1}} = \frac{\partial E(y_m^L)}{\partial s_m^L} \cdot \frac{\partial s_m^L}{\partial y_i^{L-1}}$$

It is $-\delta_m^L$!

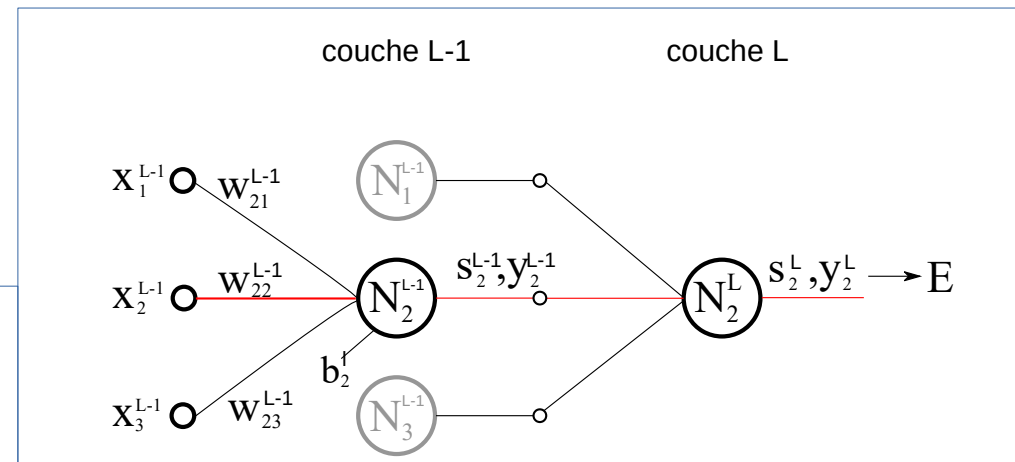
$$\frac{\partial s_m^L}{\partial y_i^{L-1}} = \frac{\partial (w_{m1}^L \cdot x_1^L + w_{m2}^L \cdot x_2^L + \dots + w_{mi}^L \cdot x_i^L + \dots + b_m^L)}{\partial x_i^L} = w_{mi}^L$$

$$\frac{\partial E(y_m^L)}{\partial y_i^{L-1}} = -\delta_m^L \cdot w_{mi}^L$$

ANNEX 2: Gradient descent in multi-layer networks

- **Gradient descent:
case of hidden layer L-1**

$$\frac{\partial E(y_m^L)}{\partial w_{ij}^{L-1}} = \frac{\partial E(y_m^L)}{\partial y_i^{L-1}} \cdot \frac{\partial y_i^{L-1}}{\partial s_i^{L-1}} \cdot \frac{\partial s_i^{L-1}}{\partial w_{ij}^{L-1}}$$



- We summarize:

$$\frac{\partial E(y_i^L)}{\partial w_{ij}^{L-1}} = - f'_{(s_i^{L-1})} \cdot x_j^{L-1} \cdot \delta_m^L \cdot w_{mi}^L$$

- Gradient descent: $w_{ij}^{L-1} \Leftarrow w_{ij}^{L-1} + \alpha \cdot x_j^{L-1} \cdot f'_{s_i^{L-1}} \cdot \delta_m^L \cdot w_{mi}^L$

- We set: $\delta_i^{L-1} = f'_{s_i^{L-1}} \cdot \delta_m^L \cdot w_{mi}^L$

Note:

$$\delta_i^{L-1} = - \frac{\partial E(y_m^L)}{\partial s_i^{L-1}}$$

And again:

$$w_{ij}^{L-1} \Leftarrow w_{ij}^{L-1} + \alpha \cdot x_j^{L-1} \cdot \delta_i^{L-1}$$

ANNEX 2: Gradient descent in multi-layer networks

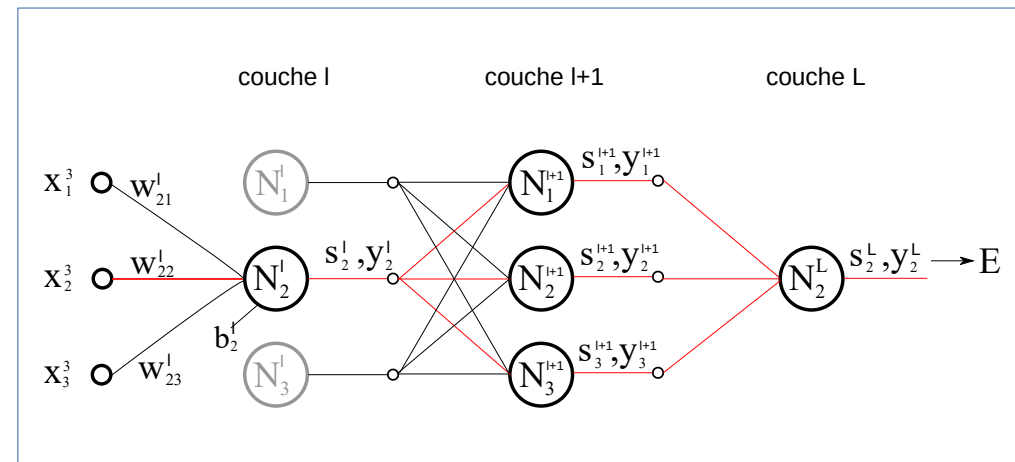
- Gradient descent:**
Generalization for layers l

- We must find:

$$\frac{\partial E(y_m^L)}{\partial w_{ij}^l}$$

- We decompose again:

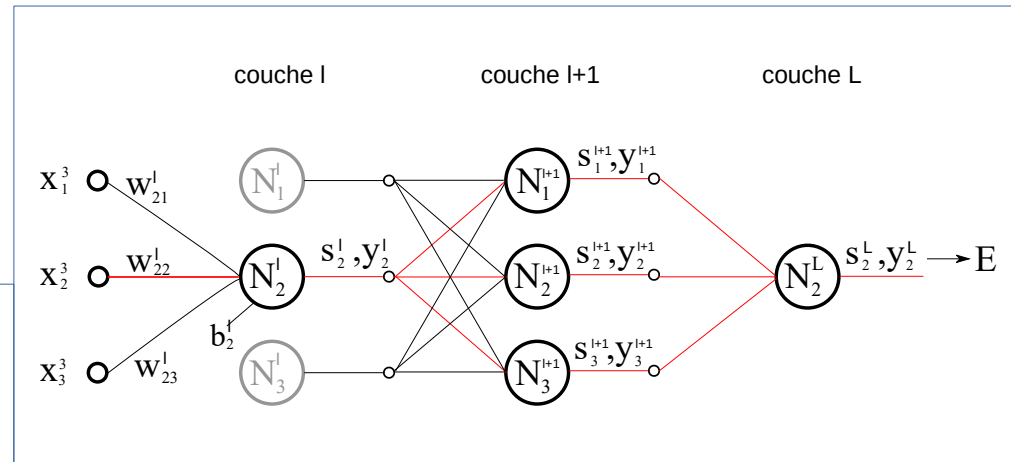
$$\frac{\partial E(y_m^L)}{\partial w_{ij}^l} = \frac{\partial E(y_m^L)}{\partial y_i^l} \cdot \frac{\partial y_i^l}{\partial s_i^l} \cdot \frac{\partial s_i^l}{\partial w_{ij}^l}$$



ANNEX 2: Gradient descent in multi-layer networks

- **Gradient descent:**
Generalization for layers l

$$\frac{\partial E(y_m^L)}{\partial w_{ij}^l} = \frac{\partial E(y_m^L)}{\partial y_i^l} \cdot \frac{\partial y_i^l}{\partial s_i^l} \cdot \frac{\partial s_i^l}{\partial w_{ij}^l}$$



- Second and third terms are the same (again):

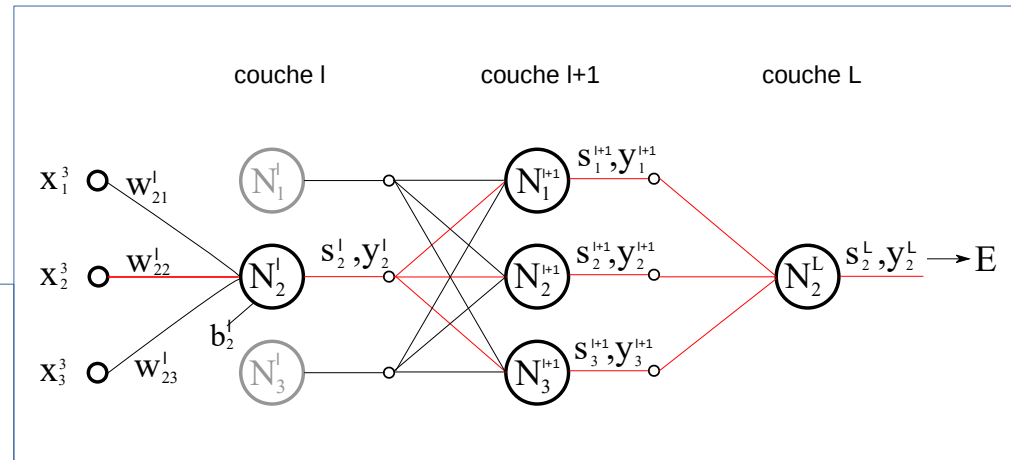
$$\frac{\partial y_i^l}{\partial s_i^l} = f'(s_i^l)$$

$$\frac{\partial s_i^l}{\partial w_{ij}^l} = x_j^l$$

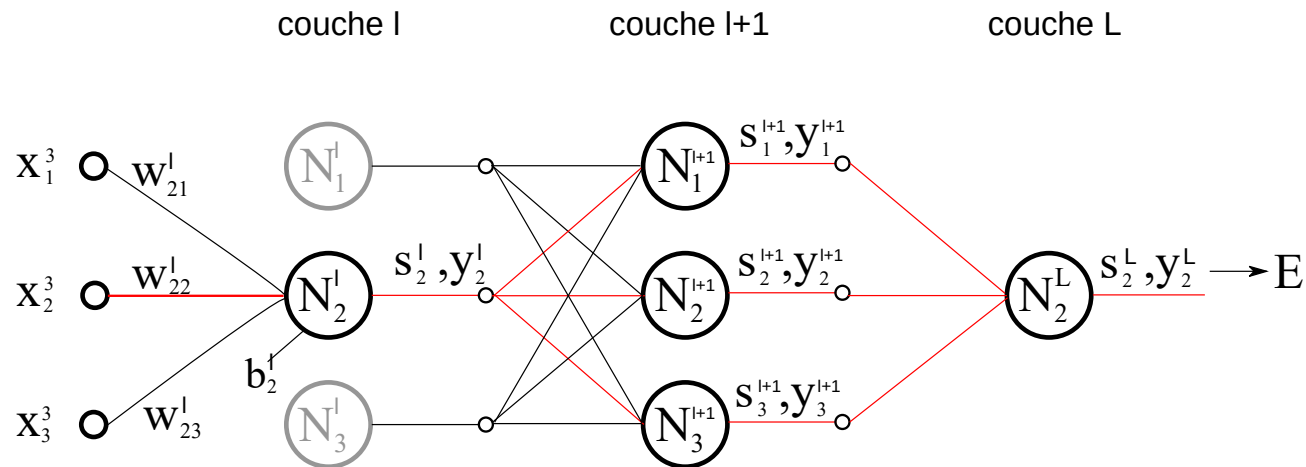
ANNEX 2: Gradient descent in multi-layer networks

- **Gradient descent:**
Generalization for layers l

$$\frac{\partial E(y_m^L)}{\partial w_{ij}^l} = \frac{\partial E(y_m^L)}{\partial y_i^l} \cdot \frac{\partial y_i^l}{\partial s_i^L} \cdot \frac{\partial s_i^L}{\partial w_{ij}^l}$$



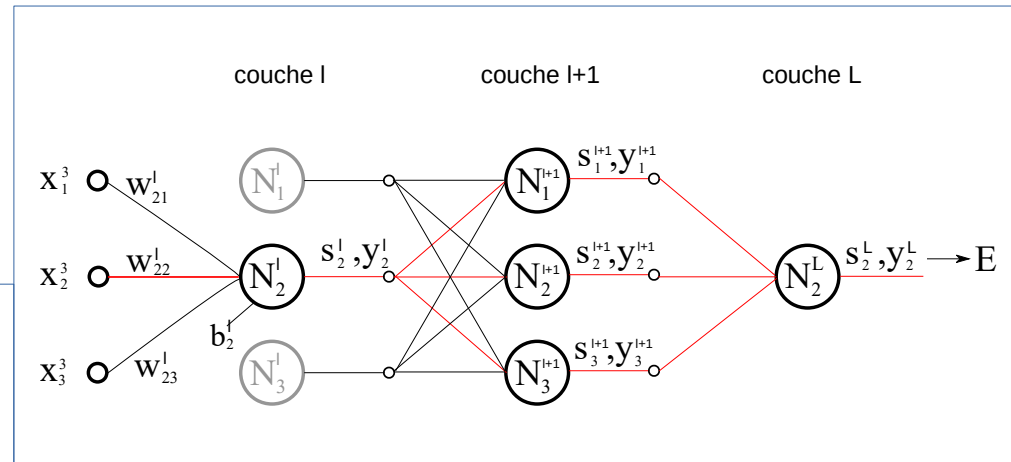
- First term: $\frac{\partial E(y_m^L)}{\partial y_i^l}$
- This time, if y_i^l is modified, E is modified through outputs s (and thus y_j^{l+1}) of next layer $l+1$



ANNEX 2: Gradient descent in multi-layer networks

- **Gradient descent:**
Generalization for layers l

$$\frac{\partial E(y_m^L)}{\partial w_{ij}^l} = \frac{\partial E(y_m^L)}{\partial y_i^l} \cdot \frac{\partial y_i^l}{\partial s_i^L} \cdot \frac{\partial s_i^L}{\partial w_{ij}^l}$$



$$\frac{\partial E(y_m^L)}{\partial y_i^l} = \frac{\partial E(y_m^L)}{\partial s_1^{l+1}} \cdot \frac{\partial s_1^{l+1}}{\partial y_i^l} + \frac{\partial E(y_m^L)}{\partial s_2^{l+1}} \cdot \frac{\partial s_2^{l+1}}{\partial y_i^l} + \dots + \frac{\partial E(y_m^L)}{\partial s_n^{l+1}} \cdot \frac{\partial s_n^{l+1}}{\partial y_i^l}$$

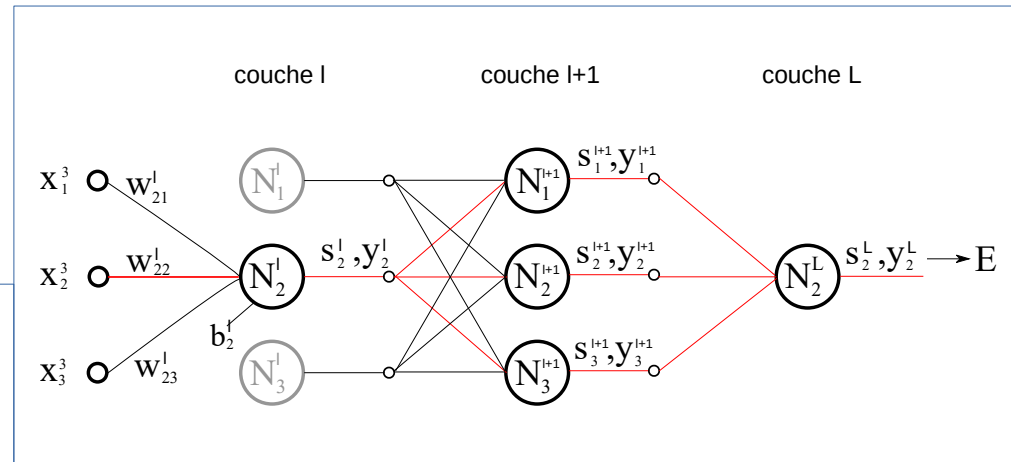
$$\frac{\partial E(y_m^L)}{\partial y_i^l} = \sum_{k \in [1, n]} \frac{\partial E(y_m^L)}{\partial s_k^{l+1}} \cdot \frac{\partial s_k^{l+1}}{\partial y_i^l} \longrightarrow w_{ki}^{l+1}$$

These are $-\delta_k^{l+1}$!

ANNEX 2: Gradient descent in multi-layer networks

- **Gradient descent:**
Generalization for layers l

$$\frac{\partial E(y_m^L)}{\partial w_{ij}^l} = \frac{\partial E(y_m^L)}{\partial y_i^l} \cdot \frac{\partial y_i^l}{\partial s_i^L} \cdot \frac{\partial s_i^L}{\partial w_{ij}^l}$$



- We summarize (a last time):

$$\frac{\partial E(y_i^L)}{\partial w_{ij}^l} = -f'_{(s_i^l)} \cdot x_j^l \cdot \sum_{k \in [1, n]} \delta_k^{l+1} \cdot w_{ki}^{l+1}$$

- Gradient descent: $w_{ij}^l \leftarrow w_{ij}^l + \alpha \cdot x_j^l \cdot f'_{(s_i^l)} \cdot \sum_{k \in [1, n]} \delta_k^{l+1} \cdot w_{ki}^{l+1}$

- We set: $\delta_i^l = f'_{s_i^l} \cdot \sum_{k \in [1, n]} \delta_k^{l+1} \cdot w_{ki}^{l+1}$

And again:

$$w_{ij}^l \leftarrow w_{ij}^l + \alpha \cdot x_j^l \cdot \delta_i^l$$