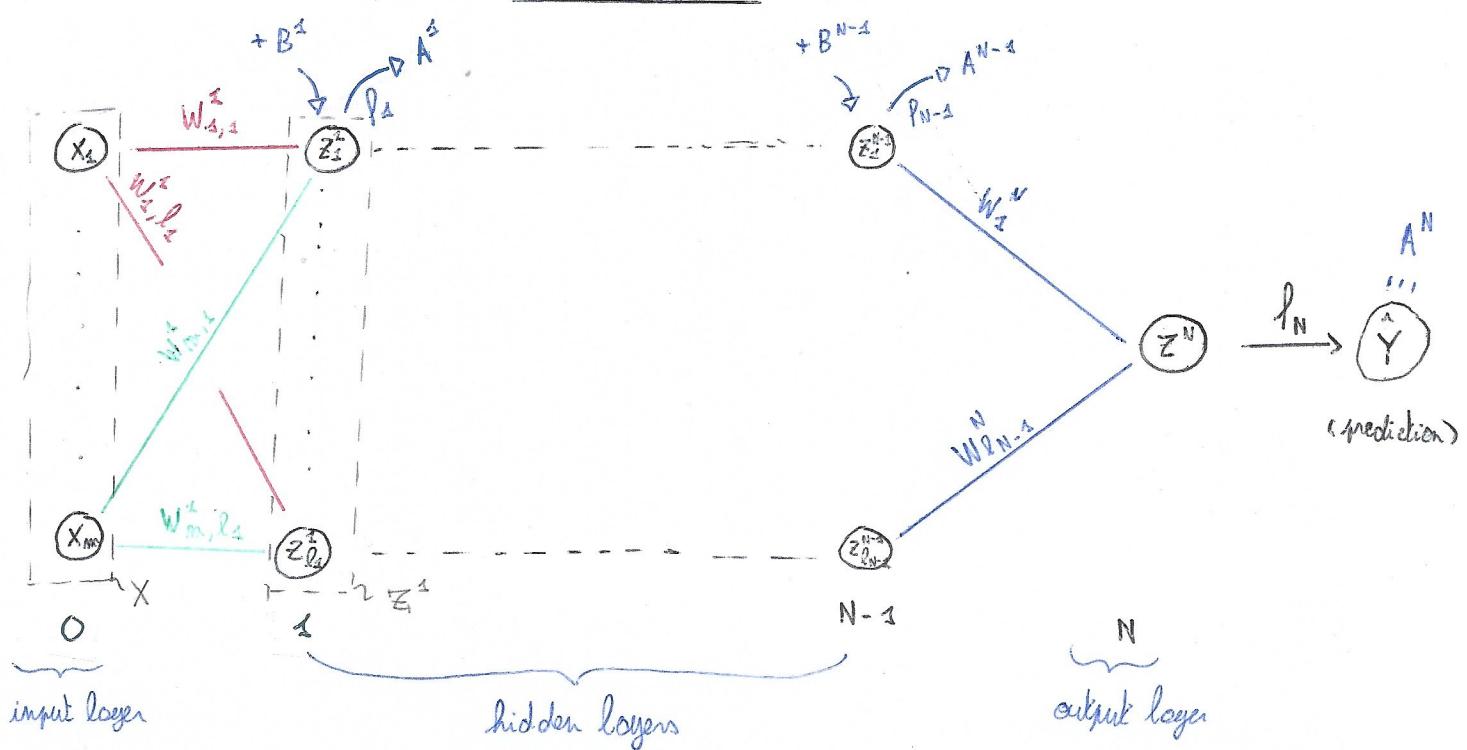


## Neural Network



### Notation

$X$  = training input (dataset without output)  $\equiv A_0$

$Y$  = training output (dataset with the output only).  
 $\hat{Y}$  = predicted output (activated output  $A^N$  associated with the  $N$ -th layer).

$W^i$  = weight matrix associated with the  $i$ -th layer.

$B^i$  = bias matrix associated with the  $i$ -th layer.

$$Z^i = (A^{i-1} \cdot W^i) + B^i$$

= output matrix associated with the  $i$ -th layer. (Linear)

$$A^i = f_i(Z^i)$$

= activated output associated with the  $i$ -th layer. (Non-linear)

$f_i$  = activation function associated with the  $i$ -th layer.

$L$  = loss function (used to optimize parameters  $W^i$  and  $B^i$ )

Examples of activation function:

$$\text{• Sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

$$\text{• ReLU}(x) = \begin{cases} x, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

Examples of loss function:

$$\text{• Mean Squared Error : } L(Y, \hat{Y}) = \frac{1}{2} \sum_{i=1}^m (Y_i - \hat{Y}_i)^2 \quad (\text{MSE})$$

$$\text{• Cross-Entropy : } L(Y, \hat{Y}) = - \sum_{i=1}^m Y_i \ln(\hat{Y}_i)$$

Dimensions:

$0, \dots, N$  : number of layers.

$m$  : number of entries in the training dataset.

$n$  : number of categories in the training dataset. (excepted "output category")  
= number of inputs

$l_i$  = number of nodes for the  $i$ -th layer ;  $i = 1, \dots, N$

1st layer:

$$\dim(X) \equiv \dim(A^0) = m \times n$$

$$\dim(W^1) = n \times l_1$$

$$\dim(X \cdot W^1) = m \times l_1 \Rightarrow \dim(B^1) = m \times l_1$$

$$\dim(Z^1) = \dim(m \times l_1) \Rightarrow \dim(A^1) = m \times l_1$$

Any  $i$ -th layer

$$\dim(A^{i-1}) = m \times l_{i-1} \quad \dim(B^i) = m \times l_i$$

$$\dim(W^i) = l_{i-1} \times l_i \quad \Rightarrow \dim(Z^i) = m \times l_i \Rightarrow \dim(A^i) = m \times l_i$$

# Neural Network

Feed Forward:  $x \rightarrow z_1 \rightarrow z_2 \dots \rightarrow \hat{y}$

$$z^i = (A^{i-1} \cdot W^i) + B^i$$

$$A^i = f_i(z^i)$$

with:  $A^0 = X$ ;  $A^N = \hat{Y}$

Backpropagation:  $x \dots \leftarrow \hat{Y}$

$$L(Y, \hat{Y}) = L(Y, A^N) \in \mathbb{R}$$

intuitively:  $\begin{cases} \text{if } L(Y, A^N) \text{ is small} \Rightarrow \text{good perf.} \\ \text{if } L(Y, A^N) \text{ is big} \Rightarrow \text{bad perf.} \end{cases}$

To optimize  $L(Y, A^N)$  we have to compute  $dL(Y, A^N) = 0$ .

Error on layer

$$\frac{\partial L}{\partial z^i} \approx \frac{\partial L}{\partial z^i}$$

$$= \frac{\partial L}{\partial A^N} \cdot \frac{\partial A^N}{\partial z^i} = \frac{\partial L}{\partial A^N} \cdot \rho'(z^i)$$

Error on bias and weight

$$\frac{\partial L}{\partial B^i} = \frac{\partial L}{\partial z^i}$$

$$\frac{\partial L}{\partial W^i} = \frac{\partial L}{\partial z^i} \cdot \frac{\partial z^i}{\partial W^i}$$

$$= \frac{\partial L}{\partial z^i} \cdot [A^{i-1}]^T$$

$$\frac{\partial L}{\partial W^N} = \frac{\partial L}{\partial A^N} \cdot \frac{\partial A^N}{\partial Z^N} \cdot \frac{\partial Z^N}{\partial W^N}$$

$$= \frac{\partial L}{\partial A^N} \cdot f'_N(Z^N) \cdot [A^{N-1}]^T$$

$$\frac{\partial L}{\partial W^{N-1}} = \left( \frac{\partial L}{\partial A^N} \cdot \frac{\partial A^N}{\partial Z^N} \right) \cdot \left( \frac{\partial Z^N}{\partial A^{N-1}} \cdot \frac{\partial A^{N-1}}{\partial Z^{N-1}} \right) \cdot \frac{\partial Z^{N-1}}{\partial W^{N-1}}$$

$$= \left( \frac{\partial L}{\partial A^N} f'_N(Z^N) \right) \cdot \left( [W^N]^T \cdot f'_{N-1}(Z^{N-1}) \right) \cdot [A^{N-2}]^T$$

$= \frac{\partial L}{\partial Z^{i+1}}$

$$\frac{\partial L}{\partial W^i} = \underbrace{\left( \frac{\partial L}{\partial A^N} f'_N(Z^N) \right) \cdot \left( [W^N]^T \cdot f'_{N-1}(Z^{N-1}) \right) \cdots \left( [W^{i+1}] f'_i(Z^i) \right)}_{= \frac{\partial L}{\partial Z^i}} [A^{i-1}]^T$$

$$\frac{\partial L}{\partial B^i} = \frac{\partial L}{\partial Z^i}$$

Summary

① Hadamard product

$$\text{Error: } \Delta Z = \frac{\partial L}{\partial Z^i}$$

$$= \frac{\partial L}{\partial A^N} \odot \frac{\partial A^N}{\partial Z^i}$$

$$= \frac{\partial L}{\partial A^N} \odot f'_i(Z^i)$$

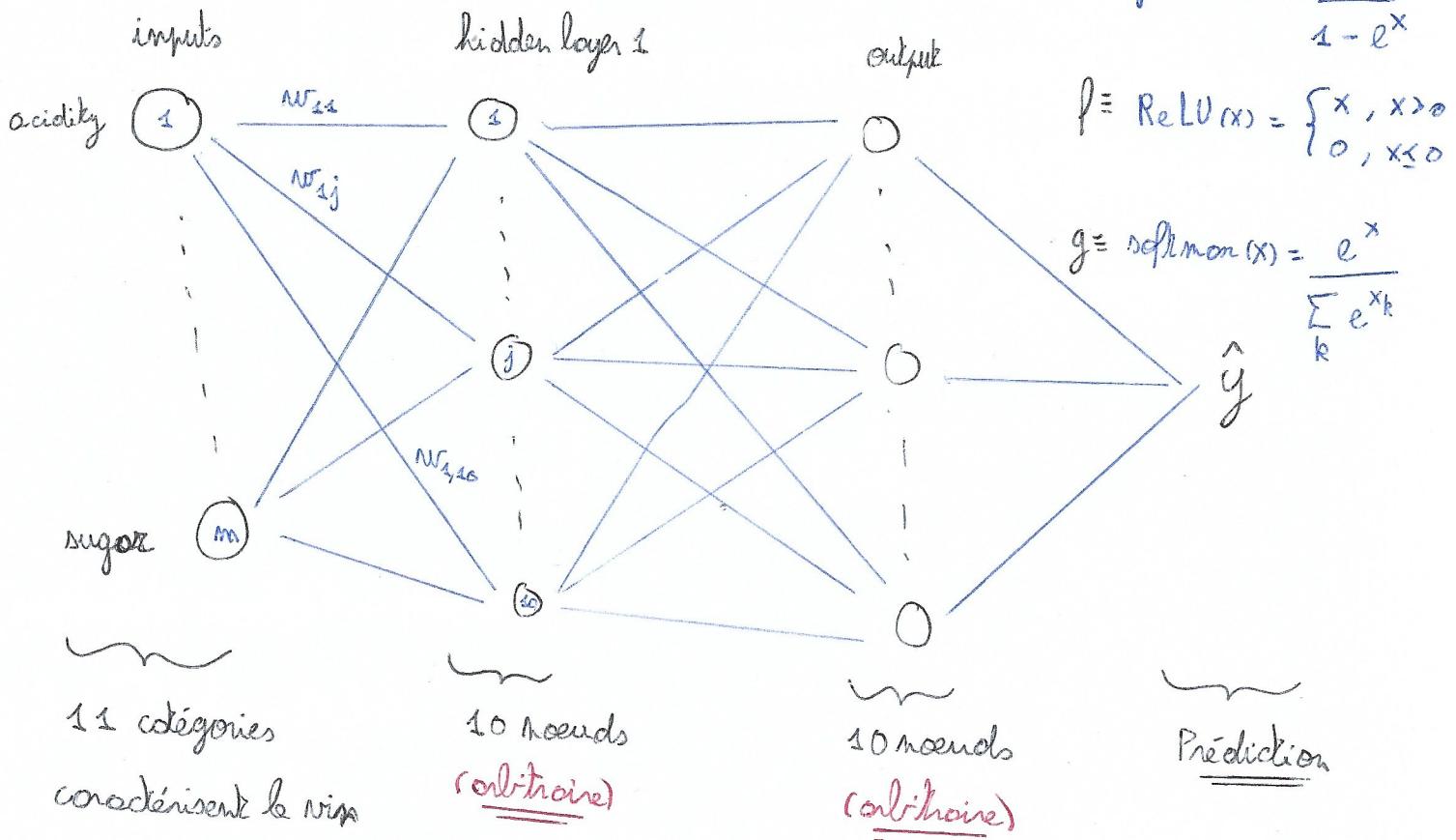
Error on parameters:

$$\frac{\partial L}{\partial W^i} = \frac{\partial L}{\partial Z^i} \cdot \frac{\partial Z^i}{\partial W^i}$$

$$= \frac{\partial L}{\partial Z^i} \cdot [A^{i-1}]^T = \frac{\partial L}{\partial Z^{i+1}} \cdot [W^{i+1}] \odot f'_i(Z^i) [A^{i-1}]^T$$

$$\frac{\partial L}{\partial B^i} = \frac{\partial L}{\partial Z^i}$$

# Neuronine



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

$$f = \text{ReLU}(x) = \begin{cases} x, & x \geq 0 \\ 0, & x \leq 0 \end{cases}$$

$$g = \text{softmax}(x) = \frac{e^x}{\sum_k e^{x_k}}$$

$$\hat{y}$$

training - dataset

testing - dataset

$\Rightarrow n =$  nombre de lignes du dataset (équitablement divisé en 2)

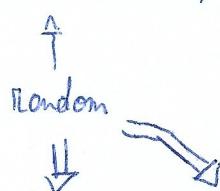
$\Rightarrow n =$  nombre de catégories = nombre d'inputs

$\Rightarrow l =$  nombre de noeuds arbitraires

Feed Forward :

Chaque couche est calculée comme suit,

W. random. uniform.



$$\text{NEXT\_LAYER} = \text{PREVIOUS\_LAYER} \cdot \text{WEIGHTS} + \text{BIAS}$$

$$\text{hidden-layer-1} = \underbrace{\text{inputs} \cdot \text{weights} + \text{bias}}_{R^{n \times m} \cdot R^{m \times l} \underbrace{R^{m \times l}}_{R^{n \times l}}} \quad \text{linéaire} \rightarrow \text{non-linéaire}$$

$$\Rightarrow f(\text{hidden-layer-1}) = f(Z^{(1)})$$

↳ sigmoid for hidden layers.  
↳ ReLU

$$\text{output-layer} = f(\text{hidden-layer-1}) \cdot \text{weights} + \text{bias} \quad \Rightarrow g(\text{output-layer}) = g(Z^{(2)})$$

↳ softmax

$$\hat{y}$$

## Backpropagation:

on cherche à minimiser l'erreur de la valeur prédictive ( $\hat{y}$ ) par rapport à la vraie valeur ( $y$ ).

On définit pour cela une "loss-function",

$$L(y, \hat{y}) := - \sum_i y_i \cdot \ln(\hat{y}_i)$$

ce qui ressemble à l'entropie statistique,

$$S[\rho] = -k_B \sum_i p_i \ln(p_i)$$

Elle permet de calculer l'erreur de la prédiction ( $\hat{y}$ ) par rapport à la vrai valeur ( $y$ ).

Minimiser cette erreur consiste en un problème d'optimisation,

$$\frac{\partial L}{\partial W} = \left( \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial W} \right) = 0$$

On sait que  $\hat{y} = X \cdot W + B$

$$\begin{cases} \frac{\partial L}{\partial W} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial W} = - \sum_i y_i \cdot \frac{1}{\hat{y}_i} \cdot X_i = 0 \\ \frac{\partial L}{\partial B} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial B} = - \sum_i y_i \cdot \frac{1}{\hat{y}_i} \cdot 1 = 0 \end{cases}$$

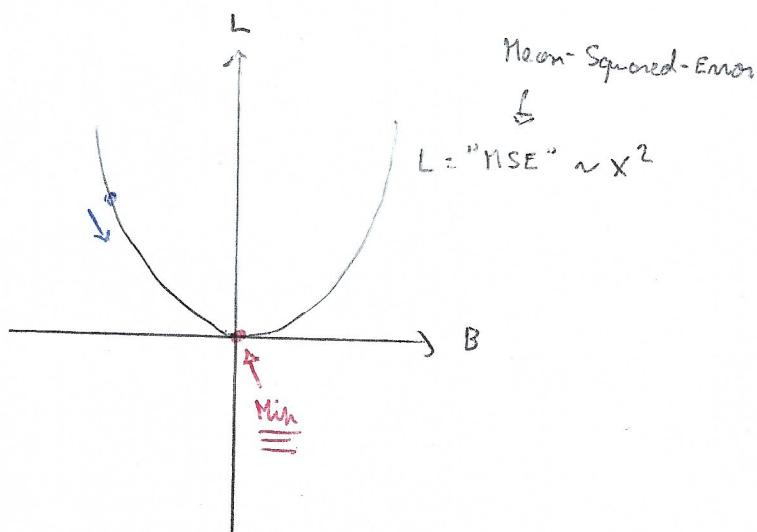
On peut aussi utiliser MSE:

$$\begin{aligned} L(y_i, \hat{y}_i) &= \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2 \\ &= \frac{1}{n} \end{aligned}$$

Gradient descent

Afin de trouver la valeur minimale de la "fonction de Loss", de manière non-analytique. On utilise un algorithme appelé le descente de gradient.

exemple :  $L = \text{"MSE"}$



On va petit à petit tweaker les paramètres  $B$  et  $W$  afin d'atteindre le minimum.

$$\left\{ \begin{array}{l} W = W - \frac{\partial L}{\partial W} . \text{"learning-rate"} \\ B = B - \frac{\partial L}{\partial B} . \text{"learning-rate"} \end{array} \right. \text{LD "vitesse / précision de l'algorithme"}$$

$$\triangleright \frac{\partial L}{\partial W^{(2)}} = - \sum_i \frac{\partial}{\partial W^{(2)}} \left( y_i \ln(g_0 z^{(1)}) \right) \quad \text{avec } g_0 z^{(1)} = g_0((f_0 z^{(1)}) \cdot W^{(2)} + b_0)$$

$$= - \sum_i \frac{\partial}{\partial W^{(2)}} \left( y_i \ln(g_0 ((f_0 z^{(1)} \cdot W^{(2)}) + b_0)) \right)$$

$$= - \sum_i y_i \cdot \frac{1}{g_0 z^{(1)}} \cdot f_0 z^{(1)} = - Y \cdot \frac{1}{g_0 z^{(1)}} \cdot f(z^{(1)})$$

$$\triangleright \frac{\partial L}{\partial B^{(2)}} = -Y \cdot \frac{1}{g(Z^{(2)})}$$

$$\triangleright \frac{\partial L}{\partial W^{(2)}} = -Y \cdot \frac{1}{g(Z^{(2)})} \cdot \frac{\partial(g \circ Z^{(2)})}{\partial W^{(2)}}$$

$$\begin{aligned} \text{or, } \frac{\partial(g \circ Z^{(2)})}{\partial W^{(2)}} &= \frac{\partial}{\partial W^{(2)}} (g(f(Z^{(2)}), W^{(2)} + B^{(2)})) \\ &= g'(f(Z^{(2)}), W^{(2)} + B^{(2)}) \cdot W^{(2)} f'(Z^{(2)}) \cdot \frac{\partial Z^{(2)}}{\partial W^{(2)}} \\ &= g'(Z^{(2)}), W^{(2)} f'(Z^{(2)}), X \end{aligned}$$

$$\Rightarrow \frac{\partial L}{\partial W^{(2)}} = -Y \cdot \frac{1}{g(Z^{(2)})} \cdot g'(Z^{(2)}) \cdot W^{(2)} f'(Z^{(2)}) \cdot X$$

$$\triangleright \frac{\partial L}{\partial \theta^{(1)}} = -Y \cdot \frac{1}{g(Z^{(2)})} \cdot g'(Z^{(2)}) \cdot W^{(2)} f'(Z^{(2)}) \cdot 1$$

$$g'(x) = e^x \cdot \sum_k e^{x_k} - e^x \cdot \sum_k e^{x_k} \cdot d(x_k)/dx$$

$$(\sum_k e^{x_k})^2$$

$$f'(x) = \begin{cases} 1, & x > 0 \\ 0, & x \leq 0 \end{cases}$$