

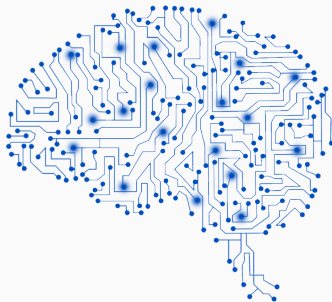
Apprentissage pour l'image

Machine learning for image processing

Course II – Introduction to Artificial Neural Networks: Backpropagation

Emile Pierret

Lundi 31 mars 2025



À la fin du cours :

- Comprendre ce qu'est un CNN (Réseau de Neurones Convolutifs)
- Implémenter l'entraînement d'un CNN pour la classification avec Pytorch

Pour cette session :

- Backpropagation pour calculer les gradients dans les Réseaux de Neurones

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ?

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

❸ $z = w$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

❸ $z = w$

Rappelons que pour $g, h : \mathbb{R} \rightarrow \mathbb{R}$,

$$(g \circ f)'(x) = g'(f(x))f'(x)$$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

❸ $z = w$

Rappelons que pour $g, h : \mathbb{R} \rightarrow \mathbb{R}$,

$$(g \circ f)'(x) = g'(f(x))f'(x)$$

Par conséquent,

$$\frac{\partial z}{\partial x}$$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

❸ $z = w$

Rappelons que pour $g, h : \mathbb{R} \rightarrow \mathbb{R}$,

$$(g \circ f)'(x) = g'(f(x))f'(x)$$

Par conséquent,

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial w} \frac{\partial w}{\partial x}$$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

❸ $z = w$

Rappelons que pour $g, h : \mathbb{R} \rightarrow \mathbb{R}$,

$$(g \circ f)'(x) = g'(f(x))f'(x)$$

Par conséquent,

$$\begin{aligned}\frac{\partial z}{\partial x} &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial x} \\ &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial v} \frac{\partial v}{\partial x}\end{aligned}$$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

❸ $z = w$

Rappelons que pour $g, h : \mathbb{R} \rightarrow \mathbb{R}$,

$$(g \circ f)'(x) = g'(f(x))f'(x)$$

Par conséquent,

$$\begin{aligned}\frac{\partial z}{\partial x} &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial x} \\ &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial v} \frac{\partial v}{\partial x} \\ &= 1 \times\end{aligned}$$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

❸ $z = w$

Rappelons que pour $g, h : \mathbb{R} \rightarrow \mathbb{R}$,

$$(g \circ f)'(x) = g'(f(x))f'(x)$$

Par conséquent,

$$\begin{aligned}\frac{\partial z}{\partial x} &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial x} \\ &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial v} \frac{\partial v}{\partial x} \\ &= 1 \times \frac{1}{v} \times\end{aligned}$$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

❸ $z = w$

Rappelons que pour $g, h : \mathbb{R} \rightarrow \mathbb{R}$,

$$(g \circ f)'(x) = g'(f(x))f'(x)$$

Par conséquent,

$$\begin{aligned}\frac{\partial z}{\partial x} &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial x} \\ &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial v} \frac{\partial v}{\partial x} \\ &= 1 \times \frac{1}{v} \times y\end{aligned}$$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

❸ $z = w$

Rappelons que pour $g, h : \mathbb{R} \rightarrow \mathbb{R}$,

$$(g \circ f)'(x) = g'(f(x))f'(x)$$

Par conséquent,

$$\begin{aligned}\frac{\partial z}{\partial x} &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial x} \\ &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial v} \frac{\partial v}{\partial x} \\ &= 1 \times \frac{1}{v} \times y \\ &= \frac{y}{v}\end{aligned}$$

Un exemple simple

Considérons la fonction $f : (x, y) \in \mathbb{R}^2 \mapsto \log(xy)$. Comment différencier f étape par étape par rapport à x ? Nous pouvons décomposer f ainsi :

❶ $v = xy$

❷ $w = \log(v) = \log(xy)$

❸ $z = w$

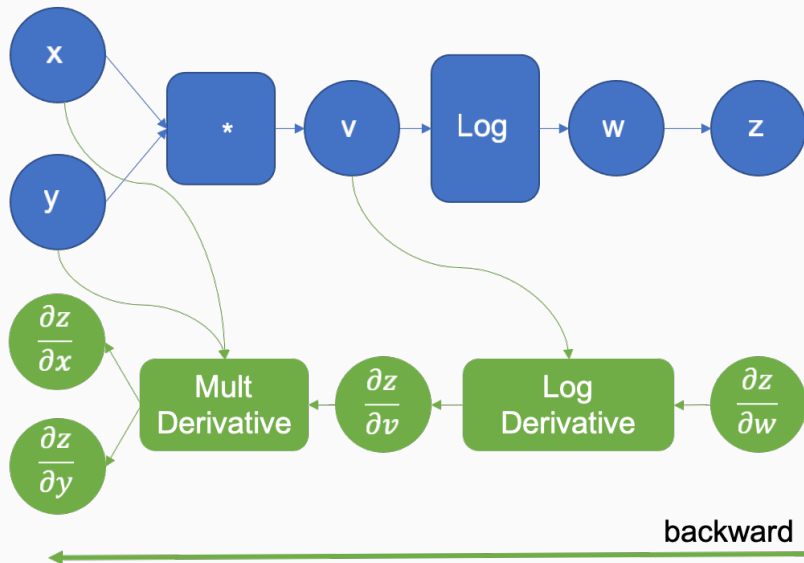
Rappelons que pour $g, h : \mathbb{R} \rightarrow \mathbb{R}$,

$$(g \circ f)'(x) = g'(f(x))f'(x)$$

Par conséquent,

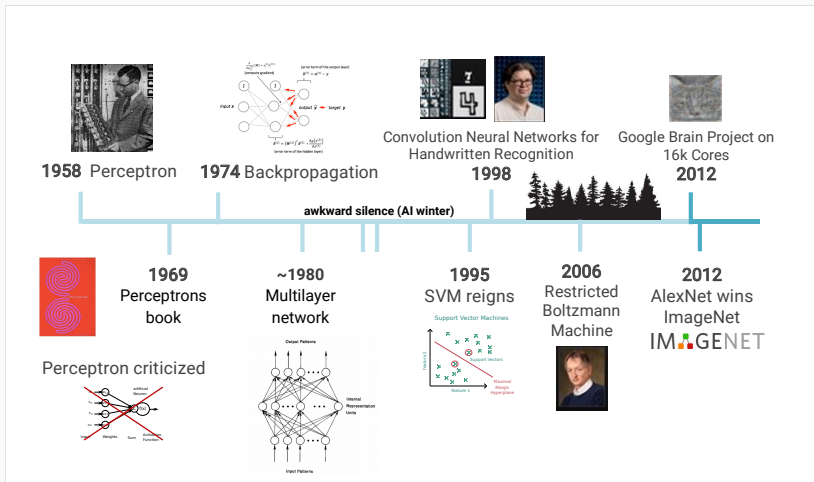
$$\begin{aligned}\frac{\partial z}{\partial x} &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial x} \\ &= \frac{\partial z}{\partial w} \frac{\partial w}{\partial v} \frac{\partial v}{\partial x} \\ &= 1 \times \frac{1}{v} \times y \\ &= \frac{y}{v} \\ &= \frac{1}{x}\end{aligned}$$

forward

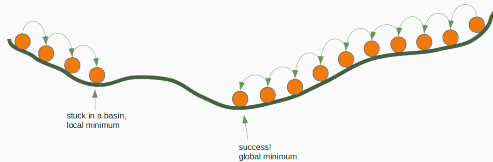


backward

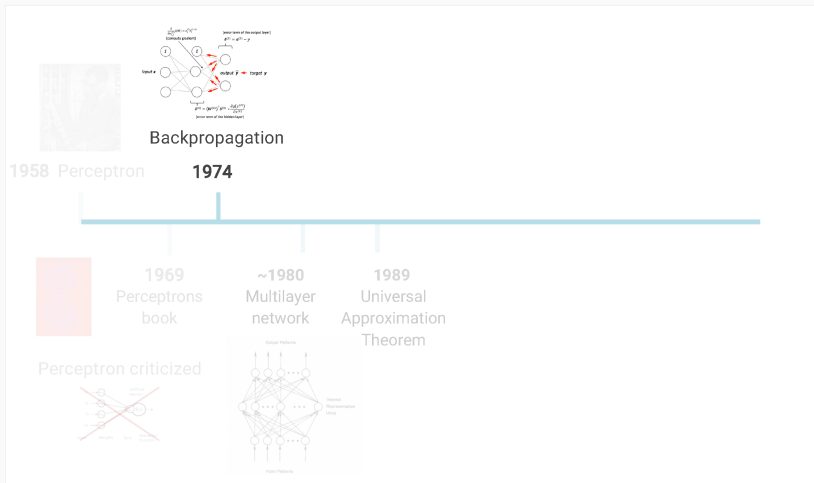
Timeline of (deep) learning



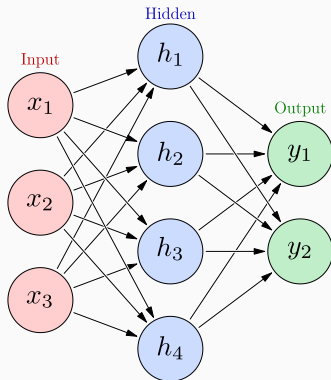
Backpropagation



Learning with backpropagation



Artificial neural network / Multilayer perceptron / NeuralNet



$$h_1 = g_1 (w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1)$$

$$h_2 = g_1 (w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1)$$

$$h_3 = g_1 (w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1)$$

$$h_4 = g_1 (w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1)$$

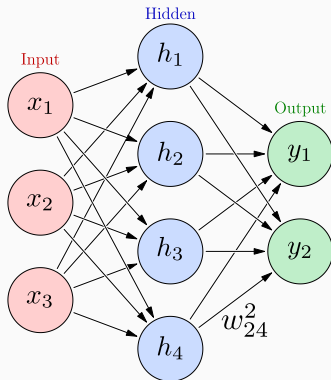
$$y_1 = g_2 (w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2)$$

$$y_2 = g_2 (w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2)$$

w_{ij}^k synaptic weight between previous node j and next node i at layer k .

g_k are any activation function applied to each coefficient of its input vector.

Artificial neural network / Multilayer perceptron / NeuralNet



$$h_1 = g_1 (w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1)$$

$$h_2 = g_1 (w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1)$$

$$h_3 = g_1 (w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1)$$

$$h_4 = g_1 (w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1)$$

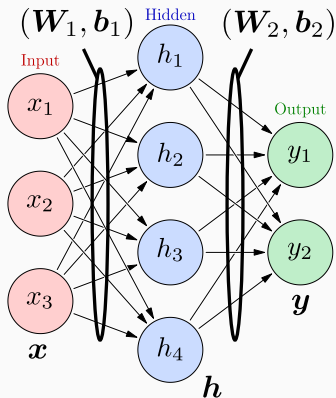
$$y_1 = g_2 (w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2)$$

$$y_2 = g_2 (w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2)$$

w_{ij}^k synaptic weight between previous node j and next node i at layer k .

g_k are any activation function applied to each coefficient of its input vector.

Artificial neural network / Multilayer perceptron / NeuralNet



$$h_1 = g_1 (w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1)$$

$$h_2 = g_1 (w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1)$$

$$h_3 = g_1 (w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1)$$

$$h_4 = g_1 (w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1)$$

$$h = g_1 (W_1 x + b_1)$$

$$y_1 = g_2 (w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2)$$

$$y_2 = g_2 (w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2)$$

$$y = g_2 (W_2 h + b_2)$$

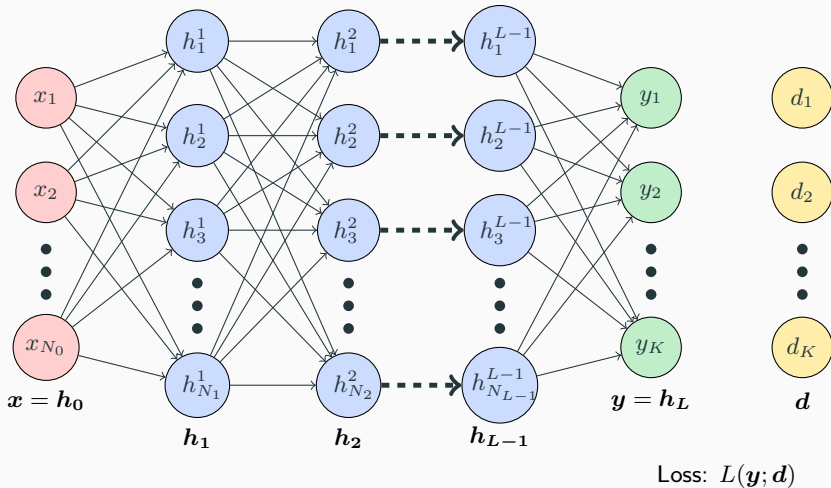
w_{ij}^k synaptic weight between previous node j and next node i at layer k .

g_k are any activation function applied to each coefficient of its input vector.

The matrices W_k and biases b_k are learned from labeled training data.

Feedforward Artificial Neural Network

Recall the feedforward structure



Input Layer

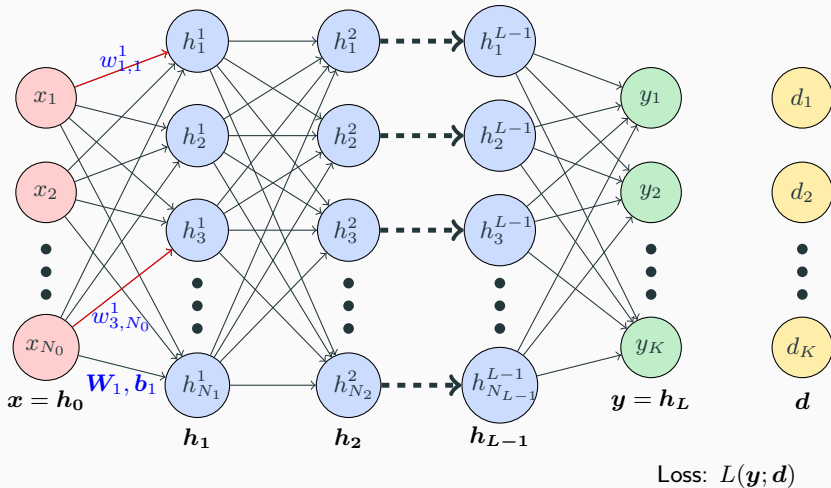
Hidden Layers

Output Layer

Label

Feedforward Artificial Neural Network

Recall the feedforward structure



Input Layer

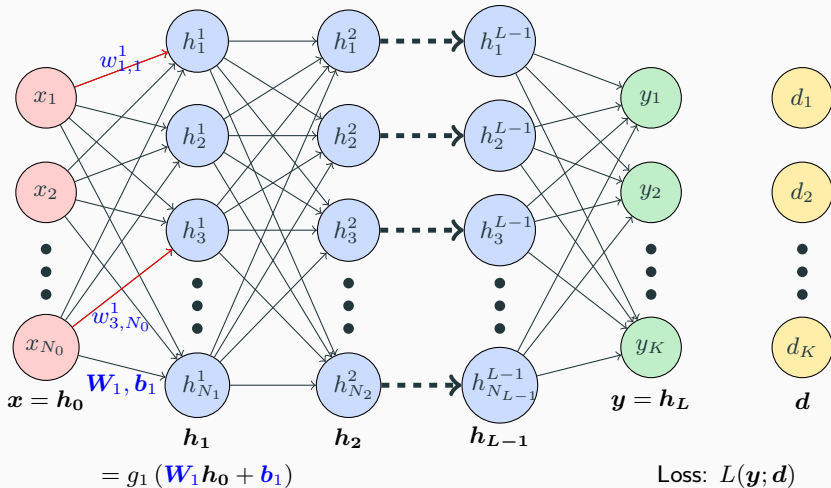
Hidden Layers

Output Layer

Label

Feedforward Artificial Neural Network

Recall the feedforward structure



Input Layer

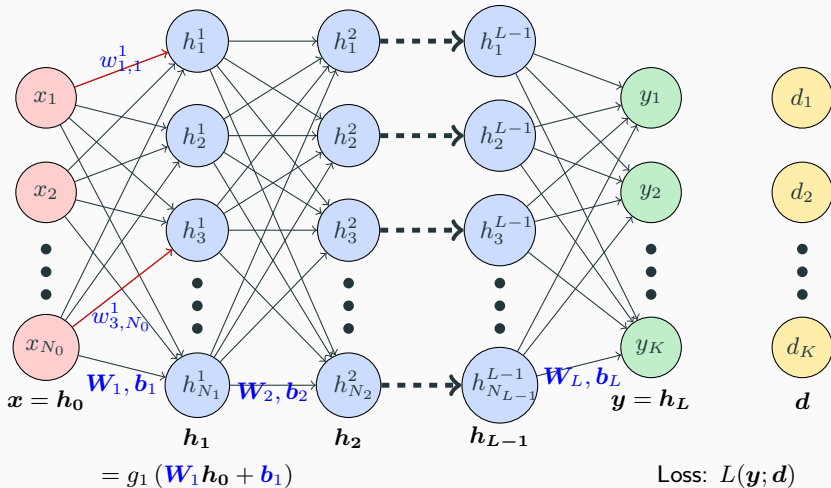
Hidden Layers

Output Layer

Label

Feedforward Artificial Neural Network

Recall the feedforward structure



Input Layer

Hidden Layers

Output Layer

Label

- Les hyperparamètres du réseau de neurones sont
le nombre de couches, la taille des couches, les fonctions d'activation

- Les hyperparamètres du réseau de neurones sont
le nombre de couches, la taille des couches, les fonctions d'activation
N'évoluent pas au cours de l'entraînement !

- Les hyperparamètres du réseau de neurones sont
le nombre de couches, la taille des couches, les fonctions d'activation
N'évoluent pas au cours de l'entraînement !
- Les paramètres du réseau de neurones sont

$$W = (W_1, b_1, W_2, b_2, \dots, W_L, b_L)$$

- Les hyperparamètres du réseau de neurones sont
le nombre de couches, la taille des couches, les fonctions d'activation

N'évoluent pas au cours de l'entraînement !

- Les paramètres du réseau de neurones sont

$$\mathbf{W} = (\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \dots, \mathbf{W}_L, \mathbf{b}_L)$$

- Entraîner le réseau = minimiser la fonction de perte d'entraînement $E(\mathbf{W})$

Objectif: $\min_{\mathbf{W}} E(\mathbf{W})$ où $E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{y}^i; \mathbf{d}^i)$

- **Solution:** pas de solution sous forme analytique

- Les hyperparamètres du réseau de neurones sont
le nombre de couches, la taille des couches, les fonctions d'activation

N'évoluent pas au cours de l'entraînement !

- Les paramètres du réseau de neurones sont

$$\mathbf{W} = (\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \dots, \mathbf{W}_L, \mathbf{b}_L)$$

- Entraîner le réseau = minimiser la fonction de perte d'entraînement $E(\mathbf{W})$

Objectif: $\min_{\mathbf{W}} E(\mathbf{W})$ où $E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{y}^i; \mathbf{d}^i)$

- **Solution:** pas de solution sous forme analytique \Rightarrow utilisation de la descente de gradient stochastique (SGD).

$$\Rightarrow \nabla E(\mathbf{W}) = \left(\nabla_{\mathbf{W}_1} E(\mathbf{W}) \quad \nabla_{\mathbf{b}_1} E(\mathbf{W}) \quad \dots \quad \nabla_{\mathbf{W}_L} E(\mathbf{W}) \quad \nabla_{\mathbf{b}_L} E(\mathbf{W}) \right)$$

Minimisation de la perte d'entraînement

Pour les réseaux de neurones multicouches, $\mathbf{W} \mapsto E(\mathbf{W})$ est non convexe
 \Rightarrow Aucune garantie de convergence.

Même si la convergence a lieu, la solution dépend de l'initialisation et du pas d'apprentissage γ .

Néanmoins, de très bons minima ou points selles sont atteints en pratique grâce à

$$\mathbf{W}^{t+1} \leftarrow \mathbf{W}^t - \gamma \nabla E(\mathbf{W}^t), \quad \gamma > 0$$

La descente de gradient peut être exprimée coordonnée par coordonnée comme suit :

$$w_{i,j}^{k,t+1} \leftarrow w_{i,j}^{k,t} - \gamma \frac{\partial E(\mathbf{W}^t)}{\partial w_{i,j}^k}$$

pour tous les poids $w_{i,j}^k$ reliant un nœud j à un nœud i dans la couche suivante k .

Fonctions de perte : Les fonctions de perte classiques sont :

Fonctions de perte : Les fonctions de perte classiques sont :

Pour divers applications : $\mathbf{d}^i \in \mathbb{R}^K$

- Erreur quadratique

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \frac{1}{2} \|\mathbf{y}^i - \mathbf{d}^i\|_2^2$$

Fonctions de perte : Les fonctions de perte classiques sont :

Pour divers applications : $\mathbf{d}^i \in \mathbb{R}^K$

- Erreur quadratique

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \frac{1}{2} \|\mathbf{y}^i - \mathbf{d}^i\|_2^2$$

Pour la classification multi-classes : En notant $d^i \in \{1, \dots, K\}$ les labels,

- Cross-entropy avec softmax comme dernière couche

$$E(\mathbf{W}) = - \sum_{(\mathbf{x}^i, \mathbf{d}^i)} \sum_{k=1}^K d_k^i \log y_k^i$$

$$\text{with } \mathbf{y}^i = f(\mathbf{x}^i; \mathbf{W}) = \text{softmax}(\mathbf{a}^i) \in (0, 1)^K.$$

- Les fonctions de perte sont de la forme

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i)} L(\mathbf{y}^i; \mathbf{d}^i)$$

- Par linéarité,

$$\nabla E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i)} \nabla L(\mathbf{y}^i; \mathbf{d}^i)$$

- Ici, la sortie du réseau de neurones $\mathbf{y}^i = f(\mathbf{x}^i; \mathbf{W})$ est une fonction des données d'entrée \mathbf{x}^i et des poids du réseau \mathbf{W} .
- Nous connaissons le gradient de $L(\mathbf{y}^i; \mathbf{d}^i)$ par rapport à la variable \mathbf{y}
 - Régression / Erreur quadratique :

$$L(\mathbf{y}; \mathbf{d}) = \frac{1}{2} \|\mathbf{y} - \mathbf{d}\|_2^2 \quad \Rightarrow \quad \nabla_{\mathbf{y}} L(\mathbf{y}; \mathbf{d}) = \mathbf{y} - \mathbf{d}$$

- Classification multi-classes / Entropie croisée :

$$L(\mathbf{y}; \mathbf{d}) = -y_d + \log \left(\sum_{k=1}^K \exp(y_k) \right) \Rightarrow (\nabla_{\mathbf{y}} L(\mathbf{y}; \mathbf{d}))_{\ell} = \text{softmax}(\mathbf{y})_{\ell} - \delta_{\ell, d}.$$

- Les fonctions de perte sont de la forme

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i)} L(\mathbf{y}^i; \mathbf{d}^i)$$

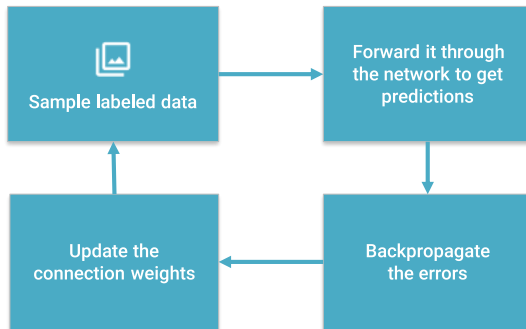
- Par linéarité,

$$\nabla E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i)} \nabla L(\mathbf{y}^i; \mathbf{d}^i)$$

- Ici, la sortie du réseau de neurones $\mathbf{y}^i = f(\mathbf{x}^i; \mathbf{W})$ est une fonction des données d'entrée \mathbf{x}^i et des poids du réseau \mathbf{W} .
- Nous connaissons le gradient de $L(\mathbf{y}^i; \mathbf{d}^i)$ par rapport à la variable \mathbf{y}
- Il reste à calculer

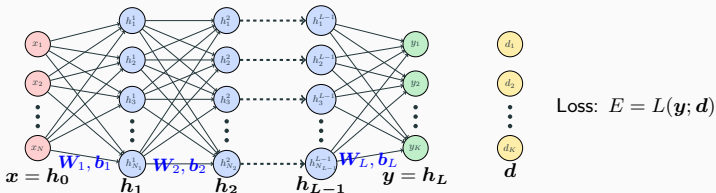
$$\nabla_{\mathbf{W}_k} L(\mathbf{y}; \mathbf{d}) \quad \text{and} \quad \nabla_{\mathbf{b}_k} L(\mathbf{y}; \mathbf{d}) \quad \text{pour } k = 0, \dots, L.$$

Training process



Learns by generating an error signal that measures the difference between the predictions of the network and the desired values and then **using this error signal to change the weights** (or parameters) so that predictions get more accurate.

RNA – Backpropagation



Passage 'forward'

Initialisation :

$$h_0 = x$$

for couche $k = 1$ **to** L **do**

$$a_k = W_k h_{k-1} + b_k$$

Activation :

$$h_k = g_k(a_k)$$

end

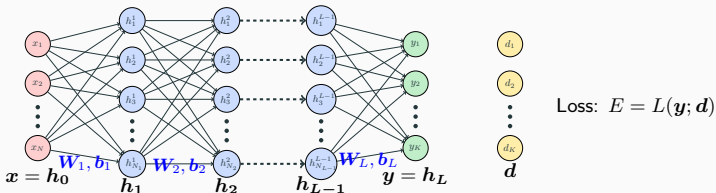
Output :

$$y = h_L$$

Calcul de la perte (loss) :

$$E = L(y; d)$$

RNA – Backpropagation



Passage 'forward'

Initialisation :

$$h_0 = x$$

for couche $k = 1$ **to** L **do**

$$a_k = W_k h_{k-1} + b_k$$

Activation :

$$h_k = g_k(a_k)$$

end

Output :

$$y = h_L$$

Calcul de la perte (loss) :

$$E = L(y; d)$$

Passage 'backward'

Objectif : Calculer le gradient par rapport à tous les paramètres

$$\frac{\partial E}{\partial w_{i,j}^k} = ? \quad \frac{\partial E}{\partial b_i^k} = ?$$

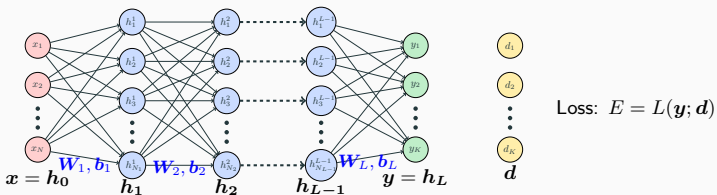
pour all

$$k \in \{1, \dots, L\},$$

$$i \in \{1, \dots, N_k\},$$

$$j \in \{1, \dots, N_{k-1}\}.$$

Backpropagation

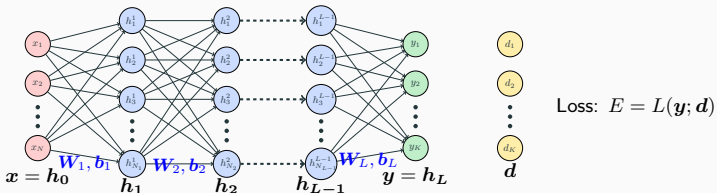


Backward

- Nous savons comment calculer la fonction de perte et son gradient :

$$\nabla_{h_L} E = \nabla L(y; d)$$

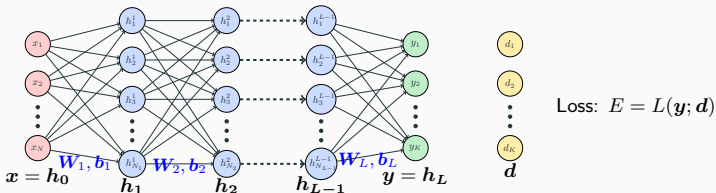
Backpropagation



Gradient par rapport à la sortie de la dernière unité linéaire a_L

$$h_L = g_L(a_L)$$

Backpropagation

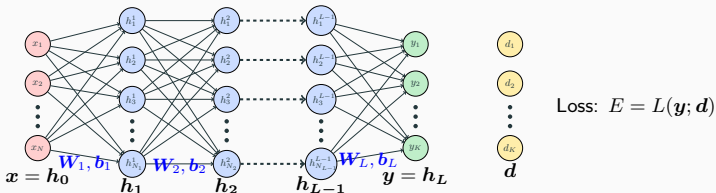


Gradient par rapport à la sortie de la dernière unité linéaire a_L

$$h_L = g_L(a_L)$$

C'est-à-dire que pour tout $i \in \{1, \dots, N_L\}$, $h_i^L = g_L(a_i^L)$.

Backpropagation



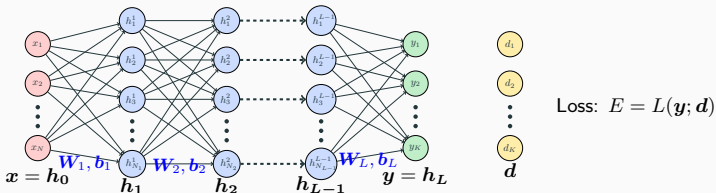
Gradient par rapport à la sortie de la dernière unité linéaire a_L

$$h_L = g_L(a_L)$$

C'est-à-dire que pour tout $i \in \{1, \dots, N_L\}$, $h_i^L = g_L(a_i^L)$. Par la règle de la chaîne,

$$\frac{\partial E}{\partial a_i^L} = \frac{\partial E}{\partial h_i^L} \frac{\partial h_i^L}{\partial a_i^L} = [\nabla_{h_L} E]_i g'_L(a_i^L)$$

Backpropagation



Gradient par rapport à la sortie de la dernière unité linéaire a_L

$$h_L = g_L(a_L)$$

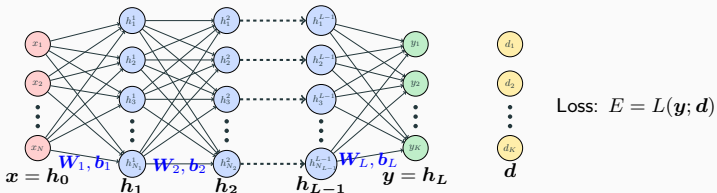
C'est-à-dire que pour tout $i \in \{1, \dots, N_L\}$, $h_i^L = g_L(a_i^L)$. Par la règle de la chaîne,

$$\frac{\partial E}{\partial a_i^L} = \frac{\partial E}{\partial h_i^L} \frac{\partial h_i^L}{\partial a_i^L} = [\nabla_{h_L} E]_i g'_L(a_i^L)$$

Formule vectorielle : $\nabla_{a_L} E = \nabla_{h_L} E \odot g'_L(a_L)$

où \odot est le produit composant par composantes entre les vecteurs, c'est-à-dire le produit de Hadamard.

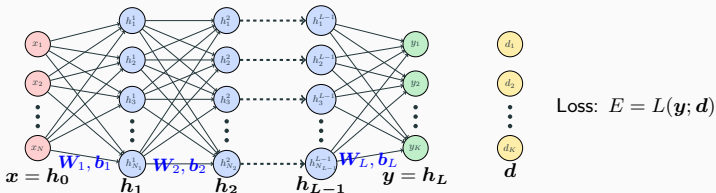
Backpropagation



Gradient par rapport au dernier biais et au dernier poids synaptique

$$a_L = W_L h_{L-1} + b_L$$

Backpropagation

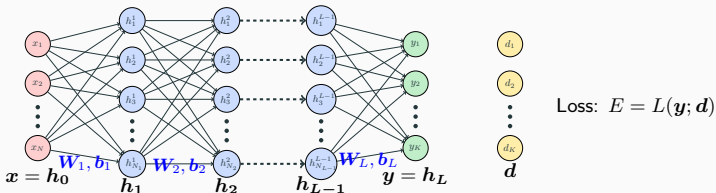


Gradient par rapport au dernier biais et au dernier poids synaptique

$$a_L = W_L h_{L-1} + b_L$$

C'est-à-dire que pour tout $i \in \{1, \dots, N_L\}$, $a_i^L = \sum_{j=1}^{N_{L-1}} w_{i,j}^L h_j^{L-1} + b_i^L$.

Backpropagation



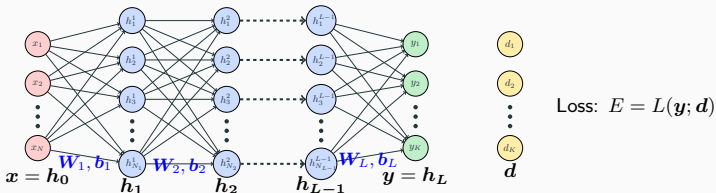
Gradient par rapport au dernier biais et au dernier poids synaptique

$$a_L = W_L h_{L-1} + b_L$$

C'est-à-dire que pour tout $i \in \{1, \dots, N_L\}$, $a_i^L = \sum_{j=1}^{N_L-1} w_{i,j}^L h_j^{L-1} + b_i^L$. Par la règle de la chaîne, pour tout $i \in \{1, \dots, N_L\}$,

$$\frac{\partial E}{\partial b_i^L} = \frac{\partial E}{\partial a_i^L} \underbrace{\frac{\partial a_i^L}{\partial b_i^L}}_{=1} = \frac{\partial E}{\partial a_i^L} = [\nabla_{a_L} E]_i$$

Backpropagation



Gradient par rapport au dernier biais et au dernier poids synaptique

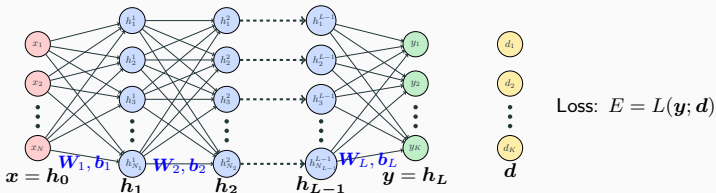
$$a_L = W_L h_{L-1} + b_L$$

C'est-à-dire que pour tout $i \in \{1, \dots, N_L\}$, $a_i^L = \sum_{j=1}^{N_L-1} w_{i,j}^L h_j^{L-1} + b_i^L$. Par la règle de la chaîne, pour tout $i \in \{1, \dots, N_L\}$,

$$\frac{\partial E}{\partial b_i^L} = \frac{\partial E}{\partial a_i^L} \underbrace{\frac{\partial a_i^L}{\partial b_i^L}}_{=1} = \frac{\partial E}{\partial a_i^L} = [\nabla_{a_L} E]_i$$

Formule vectorielle : $\nabla_{b_L} E = \nabla_{a_L} E$

Backpropagation

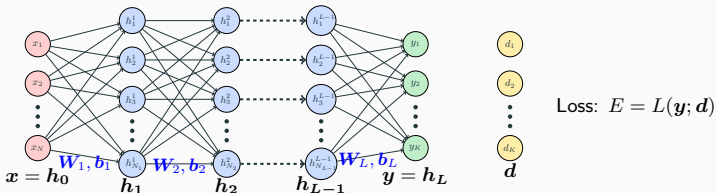


Gradient par rapport aux poids de la dernière unité linéaire W_L

$$a_L = W_L h_{L-1} + b_L$$

C'est-à-dire que pour tout $i \in \{1, \dots, N_L\}$, $a_i^L = \sum_{j=1}^{N_{L-1}} w_{i,j}^L h_j^{L-1} + b_i^L$.

Backpropagation



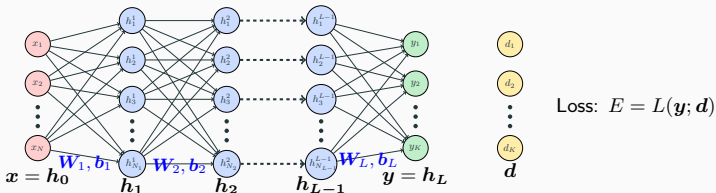
Gradient par rapport aux poids de la dernière unité linéaire W_L

$$a_L = W_L h_{L-1} + b_L$$

C'est-à-dire que pour tout $i \in \{1, \dots, N_L\}$, $a_i^L = \sum_{j=1}^{N_{L-1}} w_{i,j}^L h_j^{L-1} + b_i^L$. Par la règle de la chaîne, pour tout $i \in \{1, \dots, N_L\}$ et $j \in \{1, \dots, N_{L-1}\}$,

$$\frac{\partial E}{\partial w_{i,j}^L} = \frac{\partial E}{\partial a_i^L} \underbrace{\frac{\partial a_i^L}{\partial w_{i,j}^L}}_{=h_j^{L-1}} = \frac{\partial E}{\partial a_i^L} h_j^{L-1} = [\nabla_{a_L} E]_i [h_{L-1}]_j$$

Backpropagation



Gradient par rapport aux poids de la dernière unité linéaire W_L

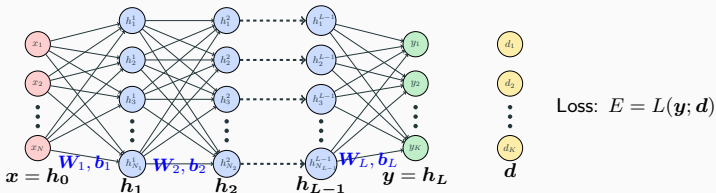
$$a_L = W_L h_{L-1} + b_L$$

C'est-à-dire que pour tout $i \in \{1, \dots, N_L\}$, $a_i^L = \sum_{j=1}^{N_{L-1}} w_{i,j}^L h_j^{L-1} + b_i^L$. Par la règle de la chaîne, pour tout $i \in \{1, \dots, N_L\}$ et $j \in \{1, \dots, N_{L-1}\}$,

$$\frac{\partial E}{\partial w_{i,j}^L} = \frac{\partial E}{\partial a_i^L} \underbrace{\frac{\partial a_i^L}{\partial w_{i,j}^L}}_{=h_j^{L-1}} = \frac{\partial E}{\partial a_i^L} h_j^{L-1} = [\nabla_{a_L} E]_i [h_{L-1}]_j$$

Formule matricielle : $\nabla_{W_L} E = \nabla_{a_L} E h_{L-1}^T$

Backpropagation

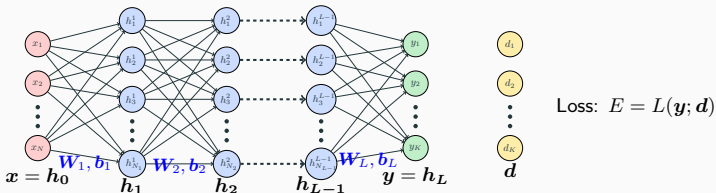


Gradients pour les paramètres de la dernière couche

Étant donné le gradient par rapport à la couche de sortie $\nabla_{h_L} E$, nous pouvons maintenant calculer :

- $\nabla_{a_L} E = \nabla_{h_L} E \odot g'_L(a_L)$
- $\nabla_{b_L} E = \nabla_{a_L} E$
- $\nabla_{W_L} E = \nabla_{a_L} E h_{L-1}^T$

Backpropagation



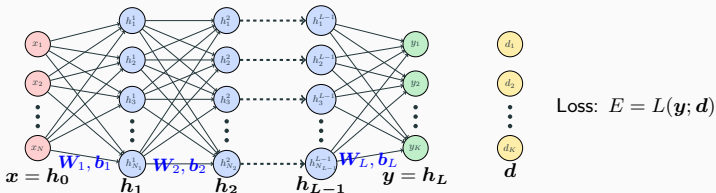
Gradients pour les paramètres de la dernière couche

Étant donné le gradient par rapport à la couche de sortie $\nabla_{h_L} E$, nous pouvons maintenant calculer :

- $\nabla_{a_L} E = \nabla_{h_L} E \odot g'_L(a_L)$
- $\nabla_{b_L} E = \nabla_{a_L} E$
- $\nabla_{W_L} E = \nabla_{a_L} E h_{L-1}^T$

Comment pouvons-nous calculer les gradients pour les paramètres de la couche $L - 1$?

Backpropagation



Gradients pour les paramètres de la dernière couche

Étant donné le gradient par rapport à la couche de sortie $\nabla_{h_L} E$, nous pouvons maintenant calculer :

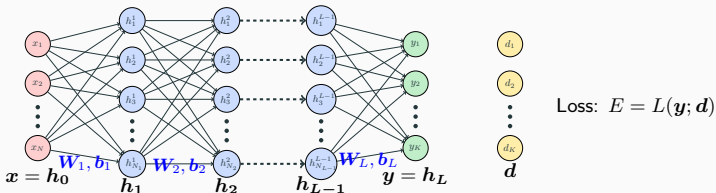
- $\nabla_{a_L} E = \nabla_{h_L} E \odot g'_L(a_L)$
- $\nabla_{b_L} E = \nabla_{a_L} E$
- $\nabla_{W_L} E = \nabla_{a_L} E h_{L-1}^T$

Comment pouvons-nous calculer les gradients pour les paramètres de la couche $L - 1$?

Nous avons besoin de l'expression du gradient par rapport à la couche cachée juste avant la dernière, h_{L-1} ... et ensuite, les mêmes formules s'appliquent !

$$\nabla_{h_{L-1}} E = ?$$

Backpropagation

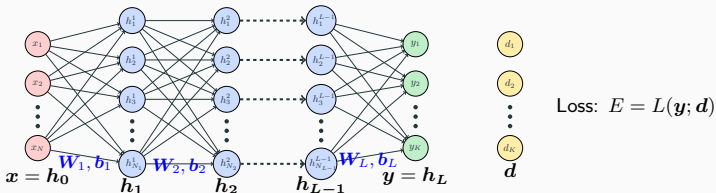


Gradient par rapport à la couche cachée juste avant la dernière h_{L-1}

Ici, même pour calculer la dérivée partielle $\frac{\partial E}{\partial h_j^{L-1}}$, nous devons utiliser le calcul différentiel pour les fonctions multidimensionnelle, car h_j^{L-1} apparaît dans chaque composant de \mathbf{a}_L :

$$\text{Pour tout } i \in \{1, \dots, N_L\}, a_i^L = \sum_{j=1}^{N_{L-1}} w_{i,j}^L h_j^{L-1} + b_i^L.$$

Backpropagation



Gradient par rapport à la couche cachée juste avant la dernière h_{L-1}

Rappelons la règle de dérivation pour la composition avec des applications affines :

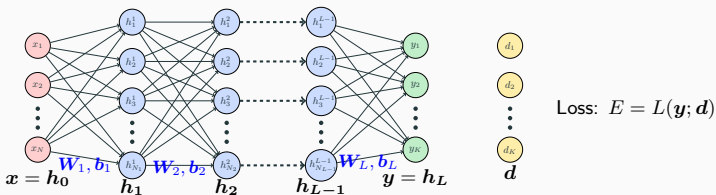
$$\text{Pour } \varphi(x) = f(\mathbf{A}x + \mathbf{b}) \text{ on a } \nabla \varphi(x) = \mathbf{A}^T \nabla f(\mathbf{A}x + \mathbf{b}).$$

En utilisant la décomposition

$$\begin{aligned} \mathbb{R}^{N_{L-1}} &\rightarrow \mathbb{R}^{N_L} \rightarrow \mathbb{R} \\ \mathbf{h}_{L-1} &\mapsto \mathbf{a}_L = \mathbf{W}_L \mathbf{h}_{L-1} + \mathbf{b}_L \mapsto E \end{aligned}$$

$$\text{Formule vectorielle : } \nabla_{\mathbf{h}_{L-1}} E = \mathbf{W}_L^T \nabla_{\mathbf{a}_L} E$$

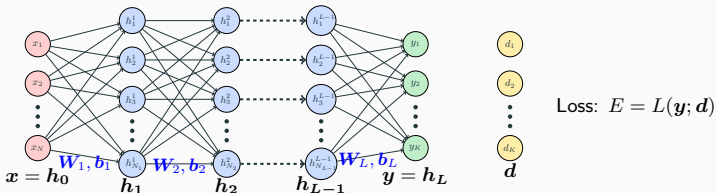
Backpropagation



On connaît $\nabla_{h_{L-1}} E$ et :

- $h_{L-1} = g_{L-1}(a_{L-1}) \Rightarrow \nabla_{a_{L-1}} E = \nabla_{h_{L-1}} E \odot g'_{L-1}(a_{L-1})$
- $a_{L-1} = W_{L-1} h_{L-2} + b_{L-1}$
- Comme précédemment,
- $\nabla_{b_{L-1}} E = \nabla_{a_{L-1}} E$
- $\nabla_{W_{L-1}} E = \nabla_{a_{L-1}} E h_{L-2}^T$

Backpropagation



Passage 'Forward'

Initialization:

$$\mathbf{h}_0 = \mathbf{x}$$

for layer $k = 1$ **to** L **do**

Linear unit:

$$\mathbf{a}_k = \mathbf{W}_k \mathbf{h}_{k-1} + \mathbf{b}_k$$

Componentwise non-linear activation:

$$\mathbf{h}_k = g_k(\mathbf{a}_k)$$

end

Output layer:

$$\mathbf{y} = \mathbf{h}_L$$

Compute loss:

$$E = L(\mathbf{y}; \mathbf{d})$$

Passage 'Backward'

Initialization: Gradient of output layer:

$$\nabla_{\mathbf{h}_L} E = \nabla L(\mathbf{y}; \mathbf{d})$$

for layer $k = L$ **to** 1 **do**

Componentwise gain of error:

$$\delta_k = \nabla_{\mathbf{a}_k} E = \nabla_{\mathbf{h}_k} E \odot g'_k(\mathbf{a}_k)$$

Gradient of layer bias:

$$\nabla_{\mathbf{b}_k} E = \delta_k$$

Gradient of weights:

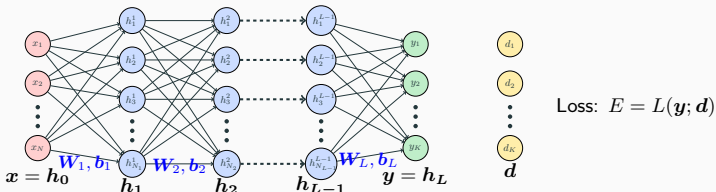
$$\nabla_{\mathbf{W}_k} E = \delta_k \mathbf{h}_{k-1}^T$$

Gradient of previous hidden layer:

$$\nabla_{\mathbf{h}_{k-1}} E = \mathbf{W}_k^T \delta_k$$

end

Backpropagation



Passage 'Forward'

Initialization:

$$\mathbf{h}_0 = \mathbf{x}$$

for layer $k = 1$ **to** L **do**

Linear unit:

$$\mathbf{a}_k = \mathbf{W}_k \mathbf{h}_{k-1} + \mathbf{b}_k \text{ (stored)}$$

Componentwise non-linear activation:

$$\mathbf{h}_k = g_k(\mathbf{a}_k) \text{ (stored)}$$

end

Output layer:

$$\mathbf{y} = \mathbf{h}_L$$

Compute loss:

$$E = L(\mathbf{y}; \mathbf{d})$$

Passage 'Backward'

Initialization: Gradient of output layer:

$$\nabla_{\mathbf{h}_L} E = \nabla L(\mathbf{y}; \mathbf{d})$$

for layer $k = L$ **to** 1 **do**

Componentwise gain of error:

$$\delta_k = \nabla_{\mathbf{a}_k} E = \nabla_{\mathbf{h}_k} E \odot g'_k(\mathbf{a}_k)$$

Gradient of layer bias:

$$\nabla_{\mathbf{b}_k} E = \delta_k$$

Gradient of weights:

$$\nabla_{\mathbf{W}_k} E = \delta_k \mathbf{h}_{k-1}^T$$

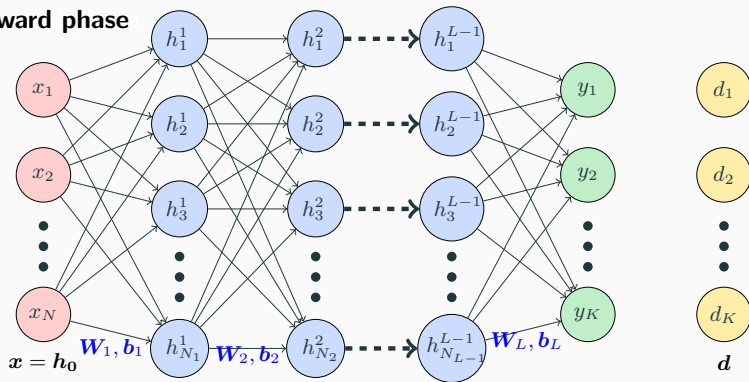
Gradient of previous hidden layer:

$$\nabla_{\mathbf{h}_{k-1}} E = \mathbf{W}_k^T \delta_k$$

end

Error backpropagation

Forward phase



Input Layer

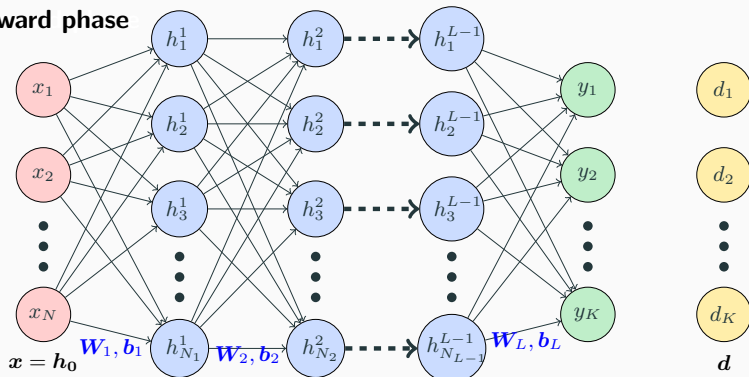
Hidden Layers

Output Layer

Label

Error backpropagation

Forward phase



$$a_1 = W_1 h_0 + b_1$$

$$h_1 = g_1(a_1)$$

$$e = \sum_k d_k y_k$$

Input Layer

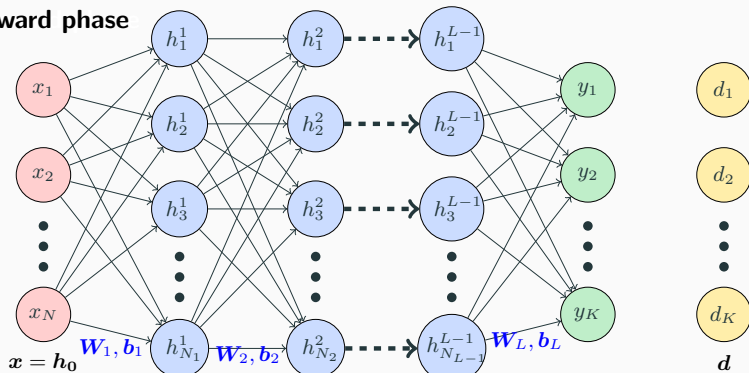
Hidden Layers

Output Layer

Label

Error backpropagation

Forward phase



Input Layer

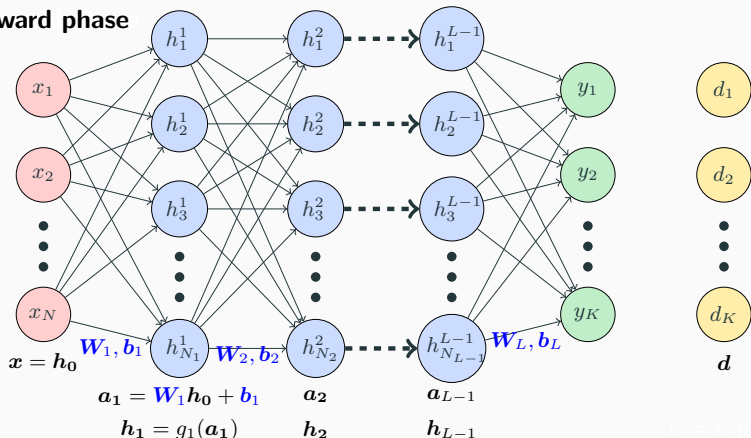
Hidden Layers

Output Layer

Label

Error backpropagation

Forward phase



Input Layer

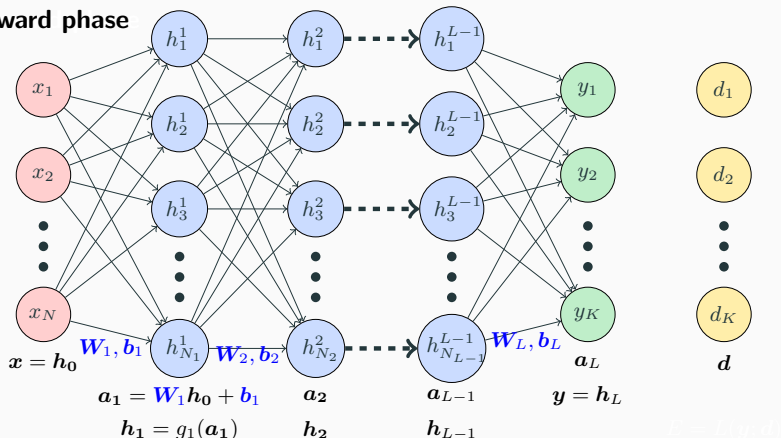
Hidden Layers

Output Layer

Label

Error backpropagation

Forward phase



Input Layer

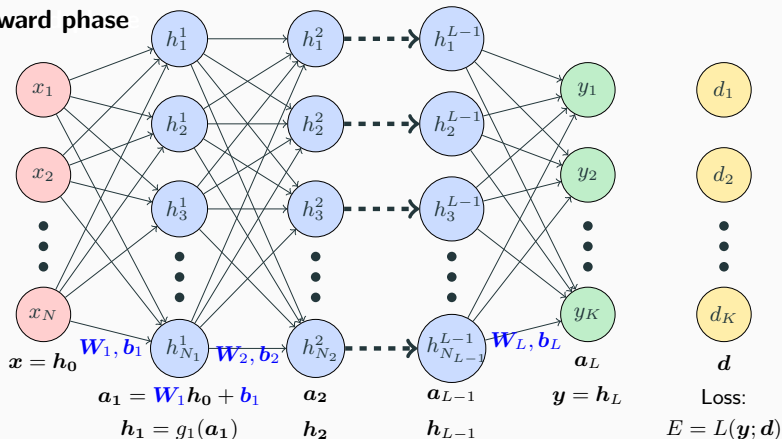
Hidden Layers

Output Layer

Label

Error backpropagation

Forward phase



Input Layer

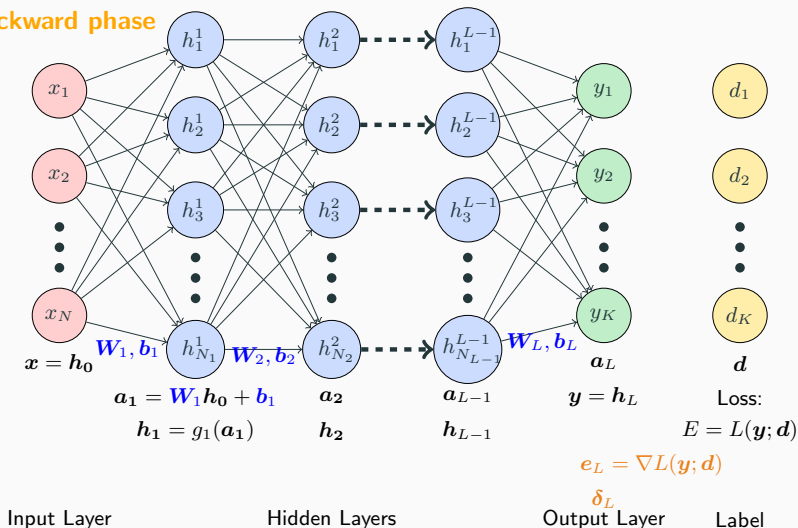
Hidden Layers

Output Layer

Label

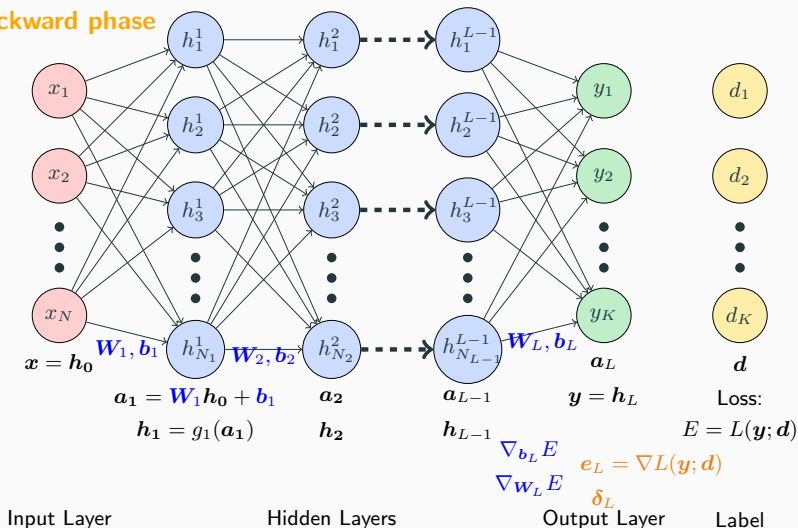
Error backpropagation

Backward phase



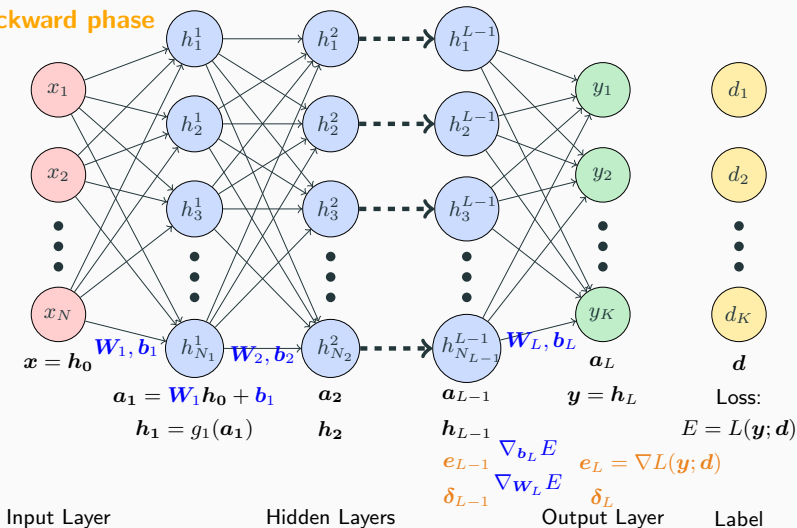
Error backpropagation

Backward phase



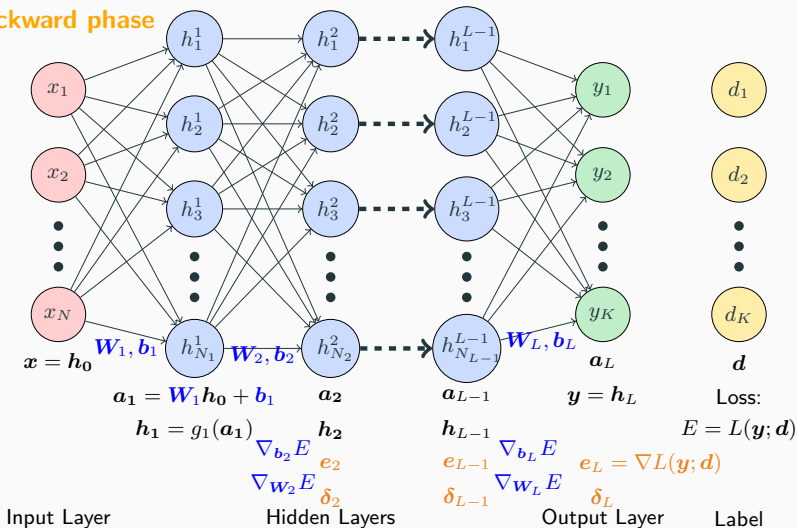
Error backpropagation

Backward phase



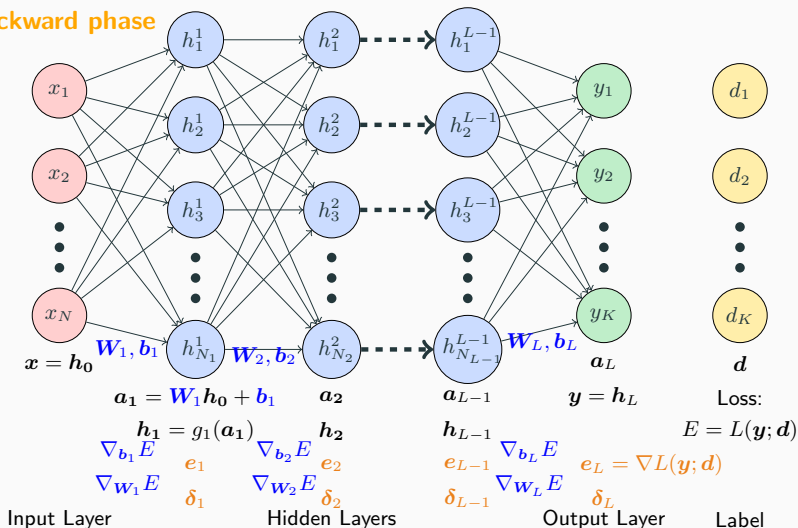
Error backpropagation

Backward phase



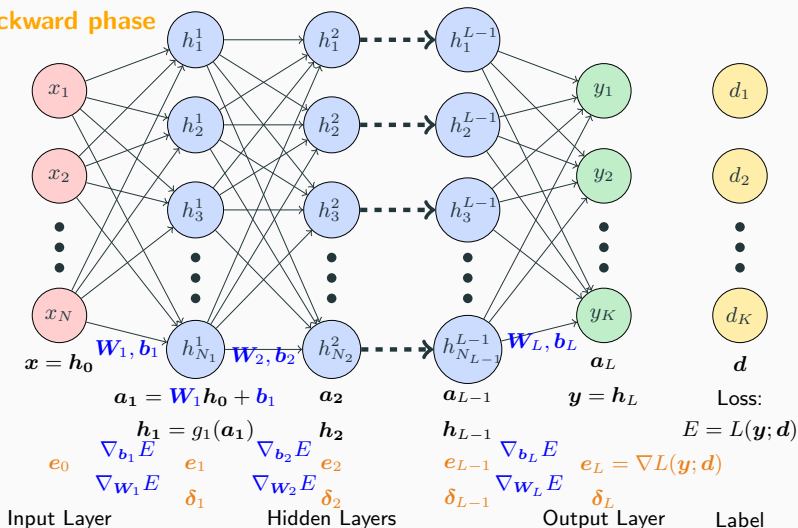
Error backpropagation

Backward phase



Error backpropagation

Backward phase



Pourquoi la backpropagation est-elle si efficace ?

- Évite de refaire de nombreux calculs
- Aucun paramètre à ajuster
- Facile à implémenter
- Cette méthode peut être vue comme une multiplication Jacobienne "dans la bonne direction"

Quelques limitations :

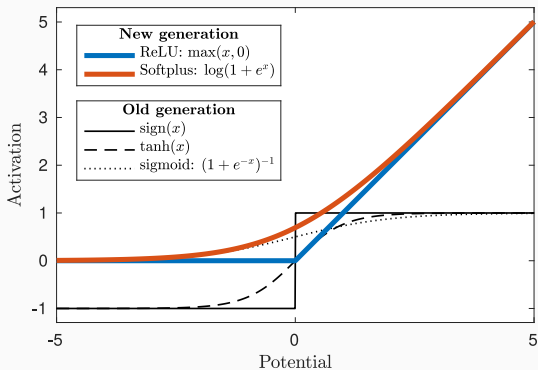
- Gourmande en mémoire
- L'apprentissage par backpropagation ne nécessite pas la normalisation des vecteurs d'entrée ; cependant, la normalisation peut améliorer les performances.

Retour sur les fonctions d'activation

En considérant

$$\mathbf{h}_{L-1} = g_{L-1}(\mathbf{a}_{L-1}) \Rightarrow \nabla_{\mathbf{a}_{L-1}} E = \nabla_{\mathbf{h}_{L-1}} E \odot g'_{L-1}(\mathbf{a}_{L-1})$$

Pourquoi ReLU est-il préférable à la fonction signe ?



Dans le TP1:

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

Algorithme utilisé:

- Initialiser \mathbf{W} aléatoirement
- Jusqu'à convergence
 - $\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} E(\mathbf{W})$
(où $\nabla_{\mathbf{W}} E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$)

Dans le TP1:

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

Algorithme utilisé:

- Initialiser \mathbf{W} aléatoirement
- Jusqu'à convergence
 - $\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} E(\mathbf{W})$
(où $\nabla_{\mathbf{W}} E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$)

Rien de stochastique ! Il s'agit d'une descente de gradient à pas constant.

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

$$\nabla_{\mathbf{W}} E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

- Initialiser \mathbf{W} aléatoirement
- Jusqu'à convergence
 - $\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} E(\mathbf{W})$

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

$$\nabla_{\mathbf{W}} E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

- | | |
|---|---|
| <ul style="list-style-type: none">• Initialiser \mathbf{W} aléatoirement• Jusqu'à convergence<ul style="list-style-type: none">• $\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} E(\mathbf{W})$ | <ul style="list-style-type: none">• Initialiser \mathbf{W} aléatoirement• Jusqu'à convergence<ul style="list-style-type: none">• Pour tout $(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}$
$\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$ |
|---|---|

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

$$\nabla_{\mathbf{W}} E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

- Initialiser \mathbf{W} aléatoirement
- Jusqu'à convergence
 - Pour tout $(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}$
 $\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

$$\nabla_{\mathbf{W}} E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

- Initialiser \mathbf{W} aléatoirement

- Jusqu'à convergence

- Pour tout $(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}$

$$\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

- Initialiser \mathbf{W} aléatoirement

- Jusqu'à convergence

- Échantillonner $(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}$

$$\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

$$\nabla_{\mathbf{W}} E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

- Initialiser \mathbf{W} aléatoirement

- Jusqu'à convergence

- Pour tout $(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}$

$$\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

- Initialiser \mathbf{W} aléatoirement

- Jusqu'à convergence

- Échantillonner $(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}$

$$\mathbf{W} \leftarrow \mathbf{W} - \gamma \nabla_{\mathbf{W}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

En fait, on a:

$$E(\mathbf{W}) \propto \mathbb{E}_{(\mathbf{x}^i, \mathbf{d}^i) \sim p_{\mathcal{T}}} \left[L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i) \right]$$

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

On découpe l'ensemble de données \mathcal{T} en "batchs" \mathcal{B} tels que $\mathcal{T} = \bigsqcup \mathcal{B}$

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

On découpe l'ensemble de données \mathcal{T} en "batches" \mathcal{B} tels que $\mathcal{T} = \bigsqcup \mathcal{B}$

- Pour chaque époque :
 - Pour chaque batch \mathcal{B} de l'ensemble de données d'entraînement, pris dans un ordre aléatoire
 - Calculer $L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$ pour $(\mathbf{y}^i; \mathbf{d}^i)$ dans le batch
 - Calculer $\mathbf{g} = \sum_{\mathbf{x} \in \mathcal{B}} \nabla_{\mathbf{W}} \ell(\mathbf{W}, \mathbf{x})$
 - Faire un pas de SGD :

$$\mathbf{W} \leftarrow \mathbf{W} - \gamma \mathbf{g}$$

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$$

On découpe l'ensemble de données \mathcal{T} en "batches" \mathcal{B} tels que $\mathcal{T} = \bigsqcup \mathcal{B}$

- Pour chaque époque :
 - Pour chaque batch \mathcal{B} de l'ensemble de données d'entraînement, pris dans un ordre aléatoire
 - Calculer $L(\mathbf{W}; \mathbf{y}^i; \mathbf{d}^i)$ pour $(\mathbf{y}^i; \mathbf{d}^i)$ dans le batch
 - Calculer $\mathbf{g} = \sum_{\mathbf{x} \in \mathcal{B}} \nabla_{\mathbf{W}} \ell(\mathbf{W}, \mathbf{x})$
 - Faire un pas de SGD :

$$\mathbf{W} \leftarrow \mathbf{W} - \gamma \mathbf{g}$$

-
- époque = désigne un passage complet du jeu de données d'entraînement par l'algorithme.
 - batch = morceau du jeu de données
 - Pour cette raison, input d'un réseau de la forme $[b, c, M, N]$

Questions?

Sources, images courtesy and acknowledgment

Charles Deledalle

V. Lepetit

L. Masuch