

Autograd — Rétropropagation depuis Zéro

Jour 3 — Matin

Julien Rolland

Formation M2 Développement Fullstack

Jour 3

- 1 Limite du Linéaire
- 2 La Non-Linearité
- 3 Credit Assignment
- 4 Graphe de Calcul
- 5 Forward & Backward Pass
- 6 Micro-Autograd en Python
- 7 Vers PyTorch

Rappel J2 : Un neurone = une droite

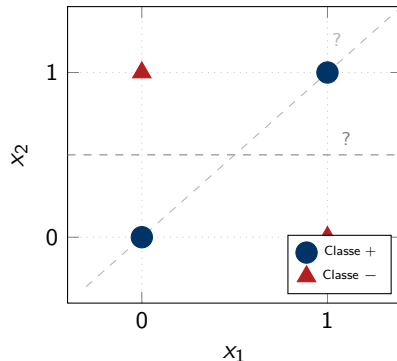
- $y = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$
- Sépare l'espace en deux **demi-plans**

XOR : Non linéairement séparable

- Classe \bullet : (0, 0) et (1, 1)
- Classe \blacktriangle : (0, 1) et (1, 0)
- **Aucune droite** ne les sépare

Solution

« Tordre l'espace » \Rightarrow **Couches cachées**
(*Hidden Layers*)



Sans activation : réseau inutile

$$W_2(W_1 \cdot \mathbf{X}) = \underbrace{(W_2 W_1)}_{W_{\text{comb}}} \cdot \mathbf{X}$$

1000 couches = 1 seule couche linéaire

Avec activation σ : expressivité

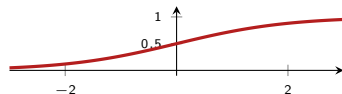
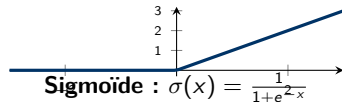
$h = \sigma(W_1 \cdot \mathbf{X})$ puis $y = W_2 h$

⇒ **Non réductible** à une seule couche

ReLU : La Norme Industrielle

$f(x) = \max(0, x)$ Dérivée : 0 ou 1

ReLU : $f(x) = \max(0, x)$



Le « Credit Assignment Problem »

Dans un réseau profond

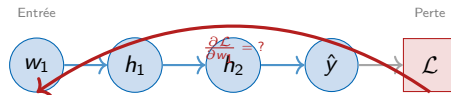
- On calcule \mathcal{L} à la **sortie**
- Les poids w_1, w_2, \dots sont au **début**
- **Question** : de combien modifier w_i pour réduire \mathcal{L} ?

Le problème fondamental

L'influence d'un poids précoce passe par **toutes les couches suivantes**

Solution : Rétropropagation

Propager $\frac{\partial \mathcal{L}}{\partial w_i}$ à l'**envers** le long du graphe



Modélisation

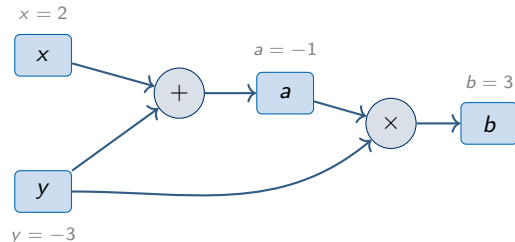
- Toute fonction \Rightarrow graphe **orienté acyclique**
- **Nœuds** : opérations élémentaires (+, \times , exp...)
- **Arêtes** : variables / tenseurs

Intérêt

Chaque nœud calcule sa **dérivée locale** \Rightarrow les gradients se chaînent vers l'entrée

Notre exemple

$$f(x, y) = (x + y) \times y$$
$$a = x + y, \quad b = a \times y$$



Mécanique

- 1 Parcourir le graphe de **gauche à droite**
- 2 Calculer chaque nœud
- 3 **Stocker** les valeurs intermédiaires

Pourquoi stocker ?

Les valeurs intermédiaires sont **nécessaires** lors du backward pass pour calculer les dérivées locales

```
1 x = 2.0
2 y = -3.0
3
4 # Noeud 1 : addition
5 a = x + y    # a = -1.0
6
7 # Noeud 2 : multiplication
8 b = a * y    # b = 3.0  <- Sortie
9
10 print(f"Sortie b = {b}")
11 # Sortie b = 3.0
```

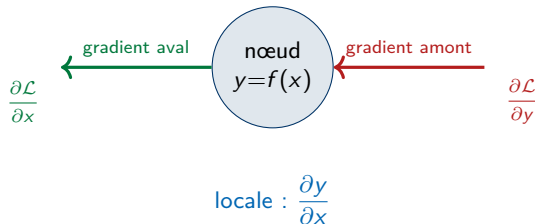
Formule

Si $z = f(y)$ et $y = g(x)$, alors :

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial x}$$

Pour les graphes de calcul

- **Gradient amont** $\frac{\partial \mathcal{L}}{\partial y}$ arrive de la sortie
- Nœud calcule sa **dérivée locale** $\frac{\partial y}{\partial x}$
- **Gradient aval** = amont \times local



$$\frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial y} \cdot \frac{\partial y}{\partial x}$$

Algorithme

- 1 Initialiser $\frac{\partial \mathcal{L}}{\partial \mathcal{L}} = 1$
- 2 Parcourir **de droite à gauche**
- 3 Chaque nœud : grad aval \times dérivée locale

Résultat clé

Un seul backward

\Rightarrow gradients de **tous** les paramètres

```
1 # b = a * y  =>  db/da = y,   db/dy = a
2 # a = x + y  =>  da/dx = 1,   da/dy = 1
3
4 # Init : gradient de la sortie
5 db = 1.0
6
7 # Noeud * : b = a * y
8 da  = db * y      # 1.0 * (-3.0) = -3.0
9 dy_2 = db * a      # 1.0 * (-1.0) = -1.0
10
11 # Noeud + : a = x + y
12 dx  = da * 1.0    # -3.0
13 dy_1 = da * 1.0    # -3.0
14
15 # Gradients finaux (accumulation)
16 grad_x = dx        # -3.0
17 grad_y = dy_2 + dy_1 # -4.0
```

Encapsuler le float

- data : valeur numérique
- grad : gradient (init. à 0)
- _prev : parents dans le graphe
- _backward : gradient local

Principe

Code mathématique **normal** +
construction du graphe
automatique

```
1 class Value:
2     def __init__(self, data, _children=()):
3         self.data = data
4         self.grad = 0.0
5         self._prev = set(_children)
6         self._backward = lambda: None
7
8     def __repr__(self):
9         return (f"Value(data={self.data:.3f}, "
10                f"grad={self.grad:.3f})")
11
12 # Utilisation
13 x = Value(2.0)
14 y = Value(-3.0)
15 print(x)  # Value(data=2.000, grad=0.000)
```

Dunder methods

- `__add__` → opérateur +
- `__mul__` → opérateur ×
- Graphe construit **implicitement**

Règle pour +

$$\frac{\partial(a+b)}{\partial a} = 1, \quad \frac{\partial(a+b)}{\partial b} = 1$$

Règle pour ×

$$\frac{\partial(a \cdot b)}{\partial a} = b, \quad \frac{\partial(a \cdot b)}{\partial b} = a$$

```
1 def __add__(self, other):
2     out = Value(self.data + other.data,
3                 (self, other))
4
5     def _backward():
6         self.grad += 1.0 * out.grad
7         other.grad += 1.0 * out.grad
8     out._backward = _backward
9     return out
10
11 def __mul__(self, other):
12     out = Value(self.data * other.data,
13                 (self, other))
14
15     def _backward():
16         self.grad += other.data * out.grad
17         other.grad += self.data * out.grad
18     out._backward = _backward
19     return out
```

Ce matin : Micro-Autograd scalaire

- Autograd sur des **scalaires** ($x \in \mathbb{R}$)
- Graphe construit manuellement
- Gradients : nombres simples

PyTorch : Même logique, Tenseurs

- Variables : **tenseurs** $X \in \mathbb{R}^{m \times n}$
- Graphe construit **automatiquement**
- Gradients : matrices (Jacobiennes)

	Scalaire	Tenseur
Valeur	$x \in \mathbb{R}$	$X \in \mathbb{R}^{m \times n}$
Gradient	$\partial \mathcal{L} / \partial x$	$\partial \mathcal{L} / \partial X$
Type	scalaire	Jacobienne
Règle	Chaîne	Chaîne vectorielle
Classe	Value	<code>torch.Tensor</code>

Cet après-midi : PyTorch

```
x = torch.tensor(..., requires_grad=True)
```