

Mapi3 Machine Learning : Réseaux de Neurones

Contexte $(x, y) \sim \mathcal{P}$. Trouver $f \in \mathcal{F}$ tq $y \approx f(x)$.

Miseaux de neurones: Une classe particulière de \mathcal{F} 's !

$$f(x) = p_{\theta_N}^{(N)} \left(\dots \left(p_{\theta_1}^{(1)}(x) \right) \right) \quad (\theta_1, \dots, \theta_N) \in \mathbb{R}^{d_1} \times \dots \times \mathbb{R}^{d_N}.$$

$$\mathcal{F} = \left\{ p_{(\theta_1, \dots, \theta_N)} \ ; \ (\theta_1, \dots, \theta_N) \in \mathbb{R}^{d_1} \times \dots \times \mathbb{R}^{d_N} \right\}$$

Vocabulaire: • f = architecture

• $(\theta_1, \dots, \theta_N)$ = poids

• N = profondeur

• $(p^{(1)} \circ \dots \circ p^{(N)})$ = layer (ou niveau) n

• $(p^{(1)} \circ \dots \circ p^{(N)})(x)$ = activations au layer n .

I. Chaines architecturales classiques

Principe général des réseaux de neurones:

Chaque $p^{(i)}$ est en général très simple, la complexité et l'expressivité provient de la profondeur.

1) Les couches linéaires

Une classe très simple de fonctions est celle des fonctions linéaires.

$$p^{(i)}(x) = w^{(i)} \cdot x + b^{(i)}$$

$$\theta_i^{(1)}(x) = w_i x + b_i ; \quad \theta_i = (w_i, b_i).$$

Vocabulaire :

- w = poids
- b = biais

⚠ Le terme de "poids" a un double usage et le terme de fonction linéaire désigne en fait une fonction affine.

2) Les activations et couches non-linéaires

Simplement composer des applications affines n'est en général pas très intéressant car le résultat est aussi affine.

L'expressivité des réseaux de neurones provient de la composition d'applications linéaires et de rectifications non-linéaires.

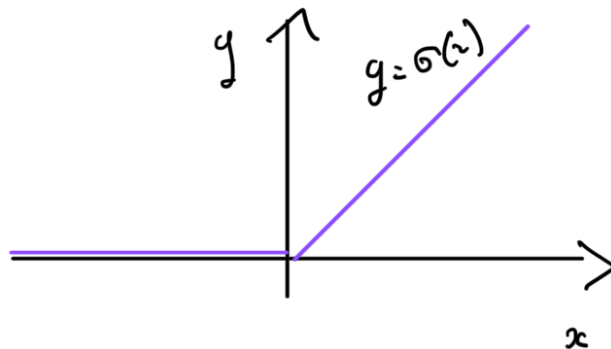
• Les activations

$$\text{Soient, } \rho_{\theta_i}^{(1)}(x) = \sigma(w_i x + b_i) ; \quad \theta_i = (w_i, b_i)$$

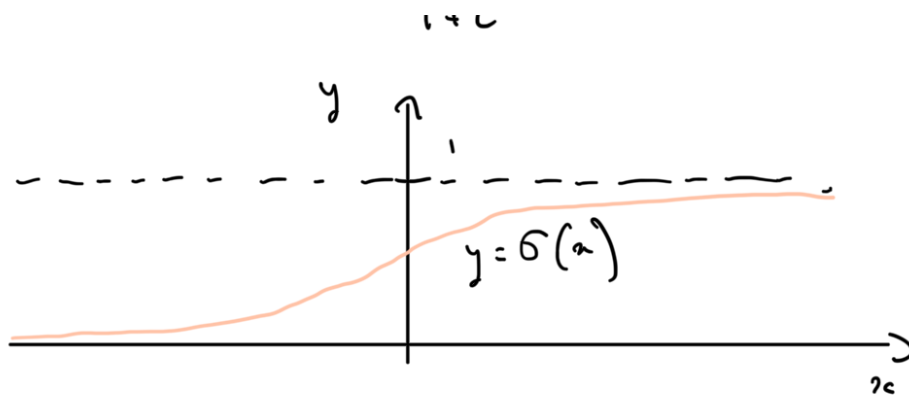
↑
fonction d'activation.

exemples classiques :

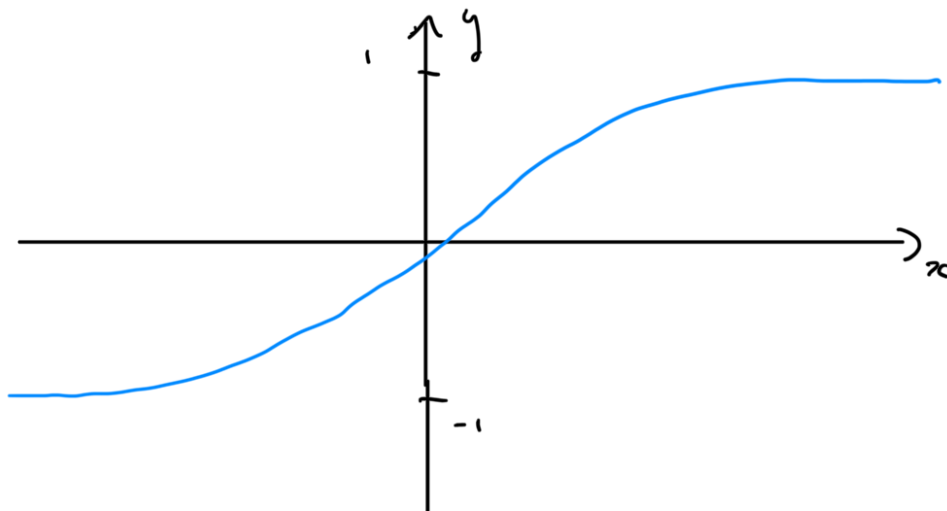
• ReLU : $\sigma(x) = (x)_+ = \max(0, x)$.



• Sigmoid : $\sigma(u) = \frac{1}{1 + e^{-u}}$



• Tanh $\sigma(v) = \frac{e^v - e^{-v}}{e^v + e^{-v}}$



• Le Softmax

En classification, on veut $f(x) \in \Delta_{d-1}$ avec

$$\Delta_{d-1} = \left\{ x \in \mathbb{R}^d \text{ st } x \geq 0 \text{ et } \mathbf{1}^T x = 1 \right\}.$$

L'interprétation est que $f(x)$ encode les probabilités d'être dans telle ou telle classe.

Il faut alors convertir un vecteur de \mathbb{R}^d en un point de Δ_{d-1} .
C'est réalisé grâce au softmax.

$$\text{Softmax} \left(\begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} \right) = \begin{pmatrix} \frac{e^{x_1}}{\sum_i e^{x_i}} \\ \vdots \\ \frac{e^{x_d}}{\sum_i e^{x_i}} \end{pmatrix}$$

$$\frac{e^{x_d}}{\sum_i e^{x_d}}$$

3) Les fonctions de perte

Un réseau est ensuite entraîné avec une fonction de perte $l(\cdot, \cdot)$

• En régression

La fonction de perte est souvent la fonction de perte quadratique

$$l(y, f(z)) = \|y - f(z)\|_2^2$$

• En classification

En classification, il faut pouvoir comparer deux vecteurs de probabilités de Δ_{d-1}

C'est effectué avec la cross-entropy loss

$$l\left(\begin{pmatrix} y_1 \\ \vdots \\ y_d \end{pmatrix}, \begin{pmatrix} p_1 \\ \vdots \\ p_d \end{pmatrix}\right) = l(z) = -\sum_{i=1}^d y_i \log(p_i)$$

Proposition $\forall p, q \in \Delta_{d-1}$, $l(p, q) \geq 0$ et $l(p, q) = 0$ si et seulement si $p = q$.

Preuve :

$$l(p, q) - l(p, p) = -\sum_i p_i \log(q_i) - \left(-\sum_i p_i \log(p_i)\right)$$

$$= -\sum_i p_i \log\left(\frac{q_i}{p_i}\right)$$

$$\underset{\text{Jensen}}{\geq} - \log \left(\sum_i p_i \frac{q_i}{p_i} \right)$$

$$= - \log \left(\sum_i q_i \right) = - \log(1) = 0.$$

Il y a égalité dans l'inégalité de Jensen ssi $\forall i, p_i = q_i$.

La suite au prochain cours...

II. Optimisation et règle de la chaîne

(5)

$$F(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_{\theta}(x_i)) + \Omega(\theta).$$

$$\text{d'où: } \nabla_{\theta} F(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \ell(y_i, f_{\theta}(x_i)) + \nabla_{\theta} \Omega(\theta).$$

Déscente de gradient:

- Commencer à θ_0

- Tant que (critère à vérifier pour continuer)

$$\theta_t \leftarrow \theta_{t-1} - \gamma_t \nabla_{\theta} F(\theta_{t-1})$$

↑
"Learning rate" à l'itération t .



Lorsque n est grand, calculer $\nabla_{\theta} F (= \frac{1}{n} \sum_{i=1}^n \dots)$ peut être trop cher, il est alors possible de remplacer $\nabla_{\theta} F$ par un estimateur calculé sur un sous-ensemble du jeu de données uniquement

Soit $B \in \llbracket 1, n \rrbracket$

$$\nabla_{\theta}^{(B)}(\theta) = \frac{1}{|B|} \sum_{b \in B} \nabla_{\theta} \ell(y_b, f_{\theta}(x_b)) + \nabla_{\theta} \Omega(\theta)$$

Déscende de gradient stochastique:

- Commencer à Θ_0
- Tant que (critère à vérifier pour continuer)
- Sélectionner $B \subseteq \{1, \dots, n\}$ selon la règle choisie (aléa ou déterministe)
- $\Theta_t \leftarrow \Theta_{t-1} - \gamma_t \nabla_{\Theta}^{(B)} F(\Theta_{t-1})$

B s'appelle un "batch".

1) Algorithme de backpropagation

$$x \xrightarrow{A_1 + b_1} a_1 \xrightarrow{\sigma(\cdot)} b_1 \xrightarrow{A_2 + b_2} a_2 \dots a_n \xrightarrow{\sigma(\cdot)} b_n \xrightarrow{A_{out} + b_{out}} a_{out} \xrightarrow{P(y, \cdot)}$$

Comment calculer le gradient en les paramètres de $P(y, h(z))$?
 \uparrow
 Fonction encodée par le réseau.

Solution: Règle de la chaîne !

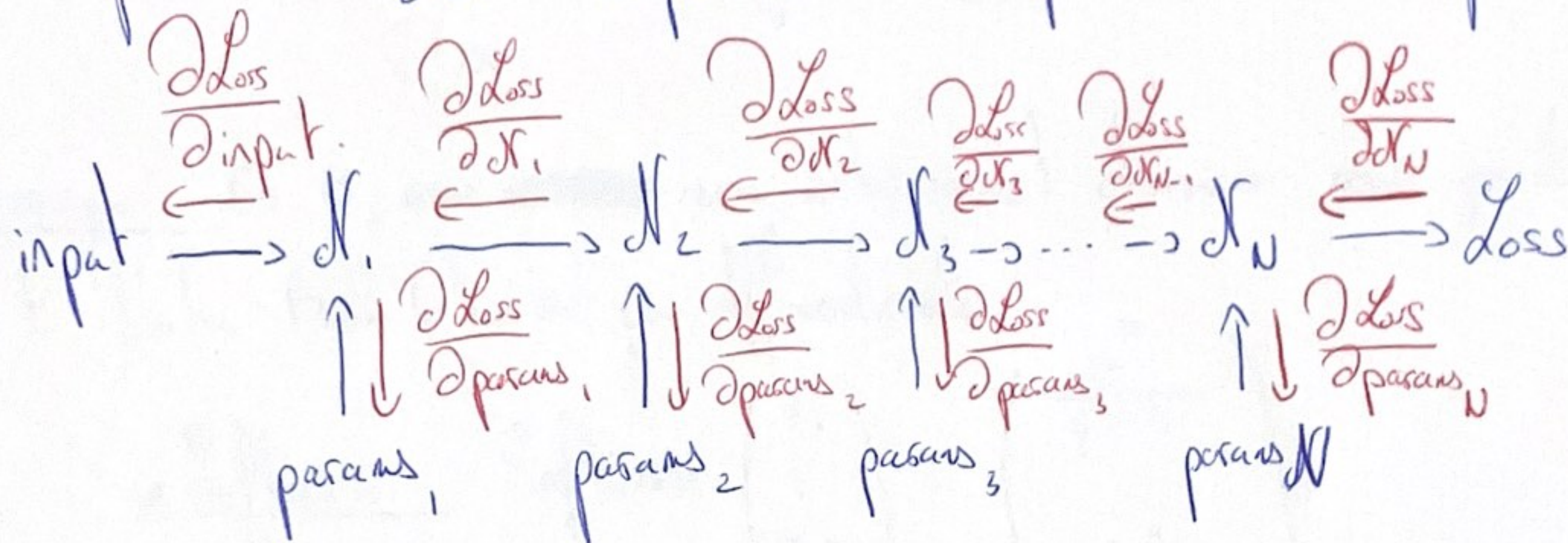
Si f et g sont différentiables, $d(g \circ f) = (dg \circ f) \circ df$

Ce qui donne, en notation informelle :

(4)

Si $X = X(\text{input}, \text{params})$, et $\text{Loss} = \text{Loss}(N)$,

$$(*) \quad \frac{\partial \text{Loss}}{\partial \text{params}} = \frac{\partial \text{Loss}}{\partial X} \frac{\partial X}{\partial \text{params}} ; \quad \frac{\partial \text{Loss}}{\partial \text{input}} = \frac{\partial \text{Loss}}{\partial X} \frac{\partial X}{\partial \text{input}} (**)$$



Algorithme Forward Backward

• Forward: Calculer $N(x, \theta)$ en gardant les activations dans chaque nœud

• Backward: En remontant à partir de la Loss, pour chaque nœud X

• Calculer $\frac{\partial \text{Loss}}{\partial \text{params}(X)}$ (avec $(*)$)

• Calculer $\frac{\partial \text{Loss}}{\partial \text{input}(X)}$ (avec $(**)$)

2) Calculs des gradients à la main

(8)

• ReLU $\text{ReLU} \begin{pmatrix} x_1 \\ \vdots \\ x_c \end{pmatrix} = \begin{pmatrix} (x_1)_+ \\ \vdots \\ (x_c)_+ \end{pmatrix}$

$$\text{donc } \frac{\partial \text{ReLU}(x)_j}{\partial x_i} = \delta_{ij} \mathbb{1}_{x_i > 0}$$

Remarque: En 0, nous ~~avons~~ avons arbitrairement choisi un sous-gradient car la fonction $\text{ReLU}(\cdot)$ n'est pas différentiable.

• Softmax $\text{Softmax} \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} = \begin{pmatrix} \frac{e^{x_1}}{\sum_j e^{x_j}} \\ \vdots \\ \frac{e^{x_p}}{\sum_j e^{x_j}} \end{pmatrix}$

$$\frac{\partial \text{Softmax}(x)_j}{\partial x_i} = \text{Softmax}(x)_j (\delta_{ij} - \text{Softmax}(x)_j)$$

• Affine: $\text{aff}(x) = Ax + b$

$$\frac{\partial \text{aff}(x)_j}{\partial x_i} = A_{ji}$$

$$\frac{\partial \text{aff}(x)_j}{\partial b_i} = 1$$

$$\frac{\partial \text{aff}(x)_b}{\partial A_{ij}} = \delta_{ij} x_j$$

III Approximation Universelle

(9)

Nous allons prouver que les réseaux de neurones ont la capacité d'approximer les fonctions continues aussi proche que l'on veut.

Mécanisme à une couche cachée

$$f(x) = \sum_{j=1}^m \eta_j (\omega_j^T x + b_j)_+$$

Étape 1: Les fonctions continues sont approchables par des fonctions continues affines par morceaux en 1.1.2.

Soit $g \in C^0([a, b], \mathbb{R})$.

Alors d'après le théorème de Heine, g est absolument continue.

en effet, sinon, $\exists \varepsilon > 0 \forall \eta, \exists a_\eta, b_\eta \in [a, b]$ t.q. $|a_\eta - b_\eta| < \eta$ et $|g(a_\eta) - g(b_\eta)| \geq \varepsilon$.

$\eta_n = \frac{1}{n}$ donne deux suites $(a_n)_{n \geq 1}$ et $(b_n)_{n \geq 1}$

Comme $[a, b]$ est compact, on extrait $a_{\varphi(n)} \rightarrow a_\infty \in [a, b]$

De même, on extrait $b_{\psi(n)} \rightarrow b_\infty \in [a, b]$

Alors $a_{\varphi(n)} \rightarrow a_\infty$; $b_{\psi(n)} \rightarrow b_\infty$

Comme $\forall n, |a_{\varphi(n)} - b_{\psi(n)}| \leq \frac{1}{\varphi(n)} \leq \frac{1}{n}$, $a_\infty = b_\infty$.

Par continuité, $|g(a_{\varphi(n)}) - g(b_{\psi(n)})| \xrightarrow{n \rightarrow +\infty} 0$ ce qui est absurde.

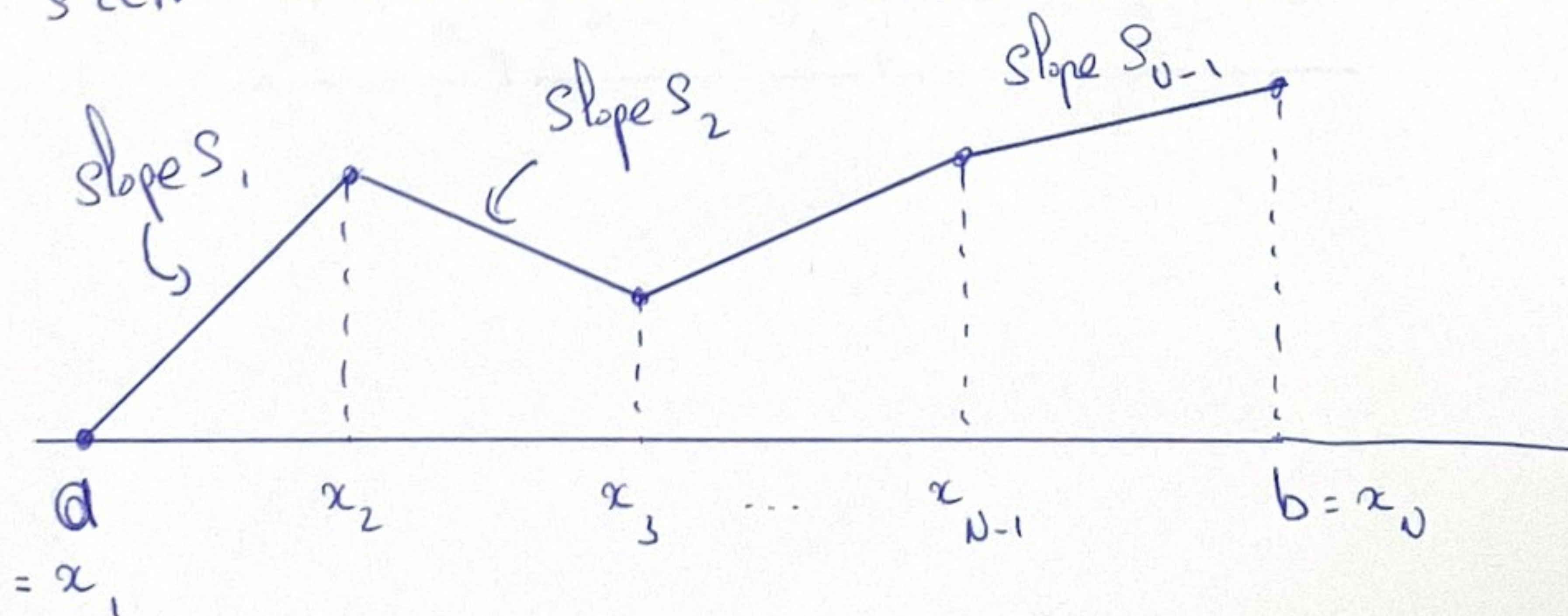
Comme g est absolument continue, $\forall \varepsilon > 0$, il existe $\eta > 0$ tel que $|x - y| < \eta \Rightarrow |g(x) - g(y)| \leq \varepsilon$.

Considérons la fonction (a=0 pour simplifier)

$$g_\eta(x) = \begin{cases} g(k_\eta) + (x - k_\eta) \frac{g((k+1)\eta) - g(k_\eta)}{(k+1)\eta - k_\eta} & \text{si } x \in [k_\eta, (k+1)\eta] \\ g(k_\eta) + (x - k_\eta) \frac{g(b) - g(k_\eta)}{b - k_\eta} & \text{si } k_\eta \in [a, b] \\ & \text{et } (k+1)\eta > b \\ & \text{et } x \in [k_\eta, b] \\ g(b) & \text{si } x = b \end{cases}$$

Alors g_η est continue et satisfait $\forall x \in [a, b] |g(x) - g_\eta(x)| \leq 2\varepsilon$

Étape 2: Toutes les fonctions continues et affines par morceaux sur $[a, b]$ peuvent s'écrire comme un réseau de neurones à une couche cachée



Cas n° 1: $f(a) = 0$

$$f(x) = S_1(x - x_1)_+ \text{ sur } [x_1, x_2]$$

$$f(x) = S_1(x - x_1)_+ + (S_2 - S_1)(x - x_2)_+ \text{ sur } [x_2, x_3].$$

...

$$\underline{f(x) = \sum_{i \geq 1} (S_i - S_{i-1})(x - x_i)_+ \text{ sur } [x_i, x_{i+1}]}$$

avec la convention $S_0 = 0$.

Cas général:

On voit f sur $[a, b]$ comme la restriction d'une fonction f' affine par morceaux et continue sur $[a-1, b]$.

$$\underline{f(x) = \sum_{i \geq 0} (S_i - S_{i-1})(x - x_i)_+}$$

$$\underline{\text{où } x_0 = a-1, S_0 = f(a) \text{ et } S_{-1} = 0}$$