

# Mise en œuvre de Neurones

①

Contexte : 
$$\hat{f} \in \arg \min_{f \in F} \frac{1}{n} \sum_{i=1}^n P(y_i, f(x_i)) + \Omega(f)$$

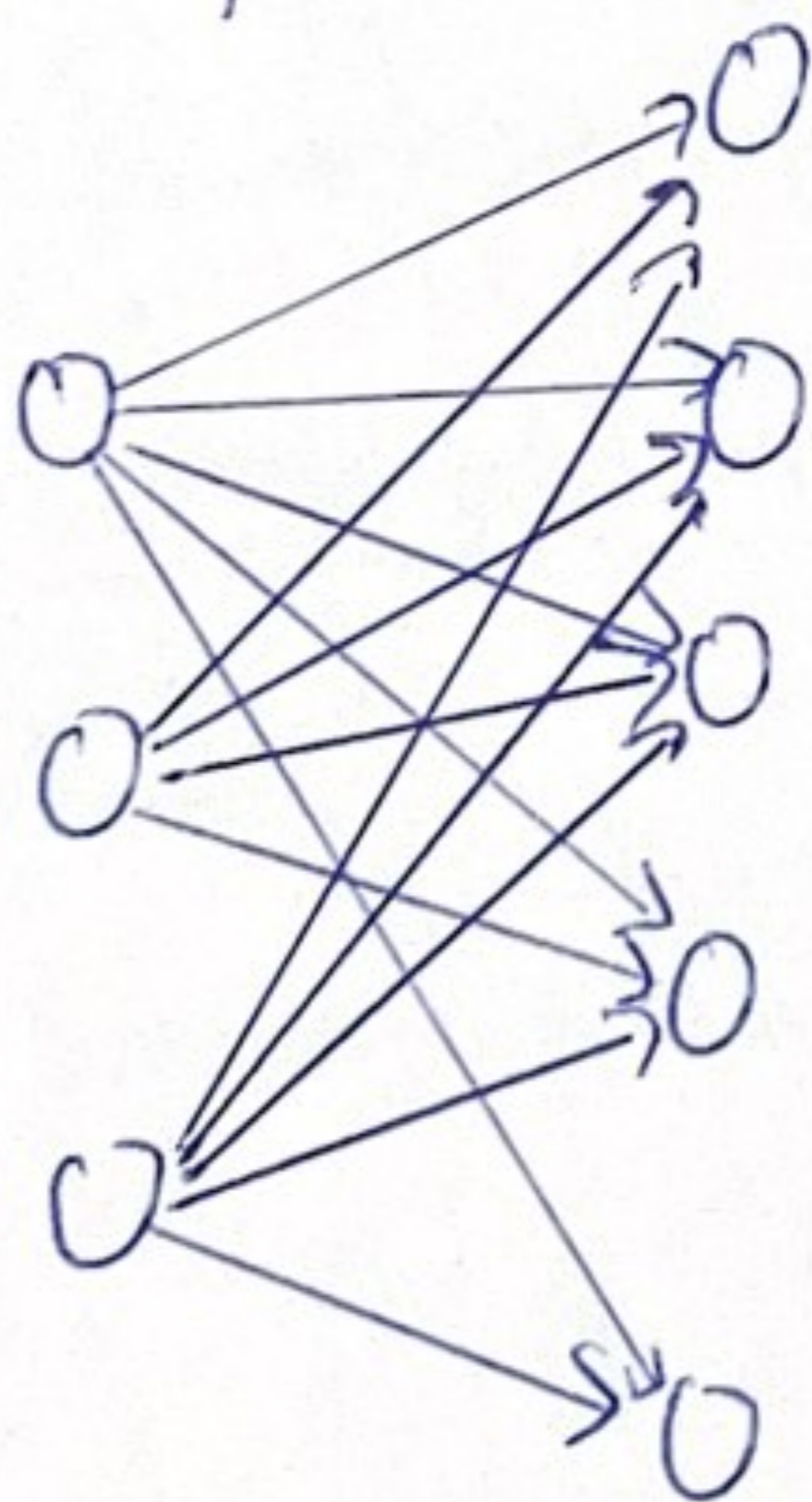
## Mise en œuvre de neurones :

Une classe de fonctions  $F$  qui est facilement représentable sur un ordinateur et pour laquelle il est facile de faire les calculs nécessaires à l'inférence et à l'optimisation.

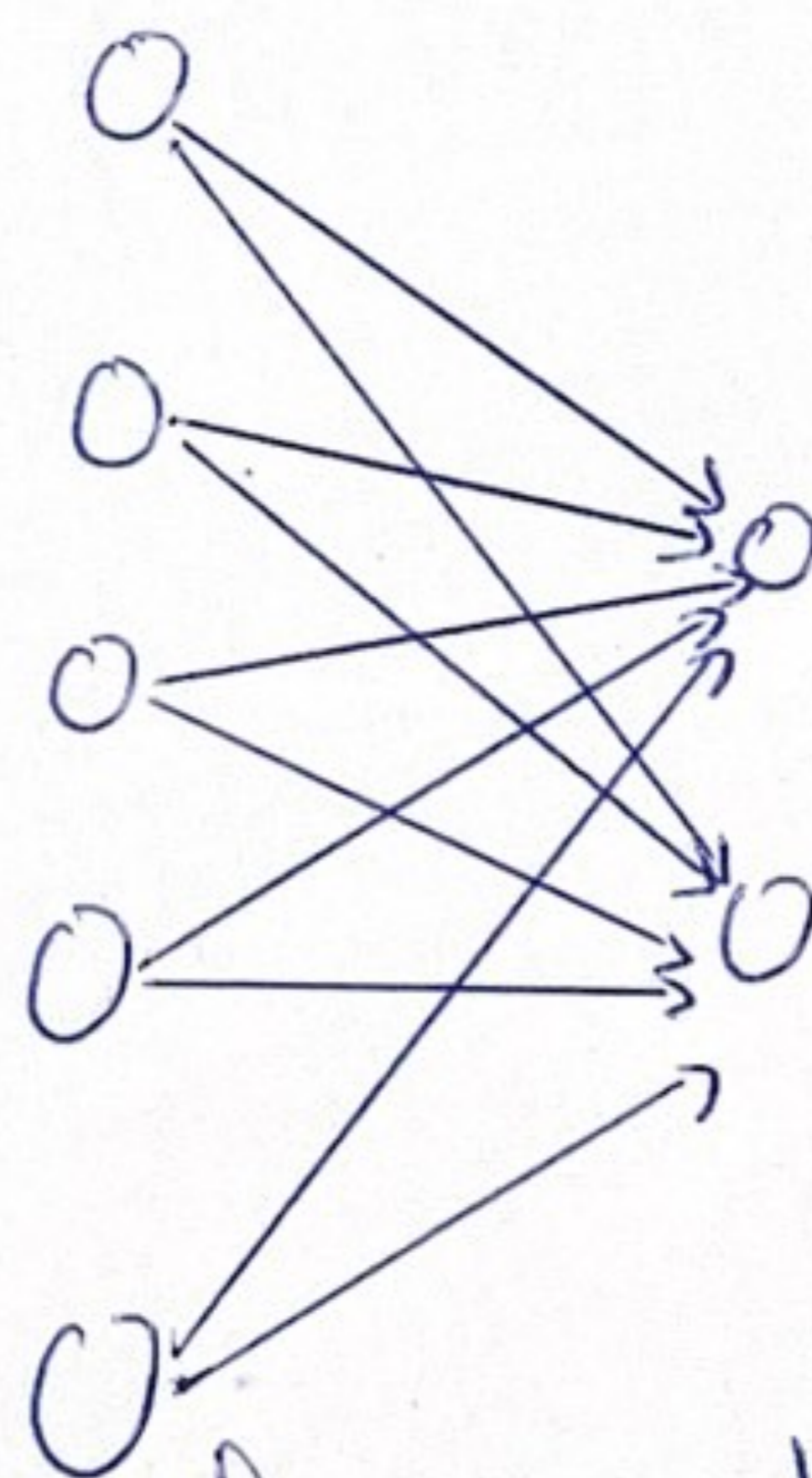
Le graphe des calculs d'une fonction  $f$  de  $F$  doit être acyclique.

## Exemples :

• Feedforward, fully connected



...



input layer hidden layer

hidden layer output layer

$$\text{output} = A_{\text{out}} \sigma \left( A_n \dots \sigma \left( A_2 \sigma \left( A_1 \text{input} + b_1 \right) + b_2 \right) \dots + b_n \right) + B_{\text{out}}$$



- (2)
- Plus généralement: Un graphe de calcul orienté sans cycle dont les nœuds sont paramétrés.

## Exemples de nœuds:

### Fonctions d'activation

- Sigmoid  $\sigma(z) = \frac{1}{1+e^{-z}}$

- tanh  $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

- ReLU  $\text{relu}(z) = \max(z, 0)$

- Softmax  $\sigma_i(z) = \frac{e^{z_i}}{\sum_j e^{z_j}}$

### Fonctions linéaires

- Transformation affine  $f(z) = Az + b$

- Convolution  $f(z) = x * k$  (10, 20, ...)

### Fonctions arbitraires

- Pooling

- Normalisation

- Clipping.



# I. Optimisation et règle de la chaîne

## 1) SGD

Les réseaux de neurones sont entraînés avec SGD, dont la fonction objectif est

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

↑ Réseau de neurones paramétré par  $\theta$ .

Il est alors possible de construire un estimateur du gradient en tirant aléatoirement  $B \subset \{1, \dots, n\}$  (de taille fixe ou aléatoire) et d'utiliser l'estimateur

$$\nabla_B F(\theta) = \frac{1}{|B|} \sum_{b \in B} \mathbb{P}(y_b, f_{\theta}(x_b)) + \Omega(\theta)$$

$B$  est appelé le "batch".

L'algorithme d'entraînement est alors :

Commence à $\theta_0$
Tant que (critère à vérifier pour continuer)
$\theta_t \leftarrow \theta_{t-1} - \gamma_t \nabla_B F(\theta_t)$ ( $B$ est tiré à chaque étape)



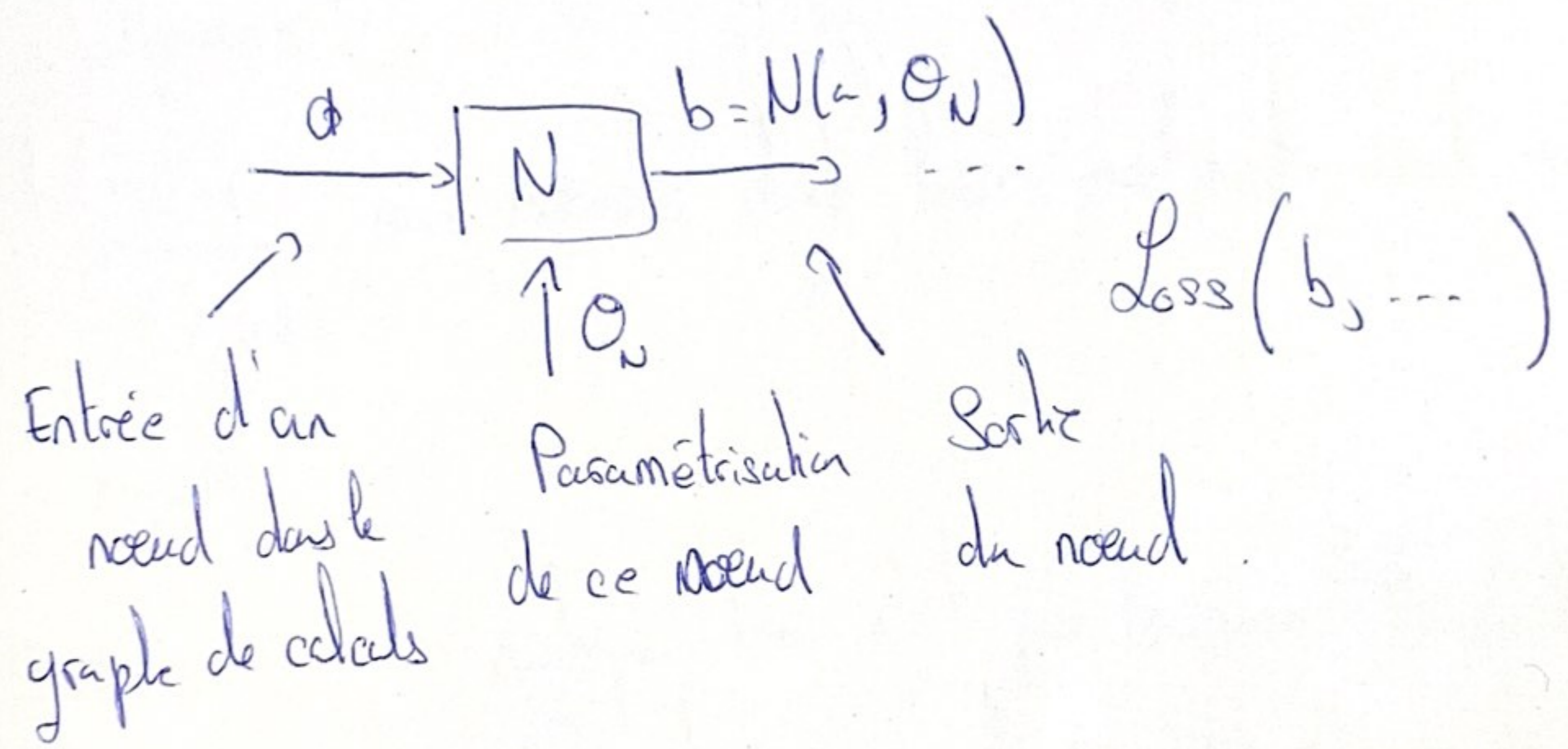
Remarque 1:

En pratique, B n'est pas toujours aléatoire et il est possible de voir les données dans un ordre déterminé.

Remarque 2:

En pratique, il existe des algorithmes plus avancés (Adagrad, RMS Prop, Adam, ...), mais ils se basent tous sur la construction de gradients stochastiques.

2) Back propagation: Comment calculer les gradients?



Remarque: Si tout est différentiable, alors la règle de la chaîne implique

$$\frac{\partial \text{Loss}}{\partial a}(a) = \frac{\partial \text{Loss}}{\partial b}(b) \frac{\partial b}{\partial a}(a) \quad \frac{\partial \text{Loss}}{\partial \theta_N}(a) = \frac{\partial \text{Loss}}{\partial b}(b) \frac{\partial b}{\partial \theta_N}(a)$$



Algorithme de Backpropagation:

- (Forward): Calculer Loss normalement en gardant la valeur de  $N(u)$  pour chaque nœud (Les activations)

- (Backward): Memoriser le graphe de calculs en partant de la fin.

Pour chaque nœud, on calcule  $\frac{\partial Loss}{\partial inputs}$  et

$\frac{\partial Loss}{\partial outputs}$  avec les formules précédentes, puis on passe

$\frac{\partial Loss}{\partial inputs}$  aux nœuds dont leurs outputs sont nos inputs.

Remarque: En pratique, tout n'est pas toujours différentiable, il est donc possible d'utiliser des "proxys" de gradients (comme les sous-gradients).



Mémoire: Le problème est non-convexe, SGD n'a pas usant <sup>(6)</sup>  
de garanties fortes de convergence. Cependant, il a été observé  
que le réseau converge vers l'optimum avec beaucoup de paramètres,  
et il est possible d'obtenir des garanties dans les régimes  
avec une infinité de neurones.

## II. Réseau Meta à une couche cachée

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

Modèle déjà très riche qui peut être étudié  
"simplement".

### 1) Propriété d'homogénéité

Si  $\alpha_j > 0$ , le réseau est invariant par la transformation

$$\begin{cases} \eta_j \leftarrow \eta_j \alpha_j \\ w_j \leftarrow w_j / \alpha_j \\ b_j \leftarrow b_j / \alpha_j \end{cases}$$

Il s'agit de la propriété d'homogénéité des



réseaux MLP.

Remarque: Il existe donc plusieurs paramétrisations donnant les mêmes fonctions

Remarque: Il est donc possible d'envisager l'étude théorique comme restreinte à un sous-ensemble de poids.

## 2) Régularisation

En ajoutant pour chaque  $j$  un terme de pénalisation de la forme  $\alpha \eta_j^2 + \|w_j\|_2^2 + b_j^2 / M^2$  pour homogénéiser,

ce qui est la même chose que de pénaliser par

$$\alpha \alpha_j^2 \eta_j^2 + (\|w_j\|_2^2 + b_j^2 / M^2) / \alpha_j^2, \forall \alpha_j > 0$$

par homogénéité, on voit que c'est équivalent (à optimiser en  $\alpha_j$ ) à pénaliser par

$$2 |\eta_j| (\|w_j\|_2^2 + b_j^2 / M^2)^{1/2}$$

Ainsi, par homogénéité, on peut se restreindre

$$\text{avec ces } \|w_j\|_2^2 + b_j^2 / M^2 = 1 \text{ avec une pénalisation } \alpha |\eta_j|.$$



27  
⑧

Pour simplicité, nous allons nous restreindre au cas où les paramètres sont ~~restreints~~. Pour le cas général, merci de lire (Bach 2024) <sup>contraints</sup>.

### III. Erreur de généralisation - Cas contraint

Le paragraphe précédent a servi d'heuristique pour justifier pourquoi nous pouvons nous restreindre au cas

$$\begin{cases} \forall j, \|w_j\|_2^2 + b_j^2 / M^2 = 1 \\ \|A\|_1 \leq D \end{cases}$$

~~3. de plus~~

Notes (H) L'ensemble des paramètres vérifiant ces contraintes.

De plus, faisons les hypothèses suivantes :

- $\|x\| \leq M$  p.s.

- $f$  est  $G$ -Lipschitz en son second argument.



Analyse par Mademacher:

D'après le principe de symétrisation:

$$E(M(\hat{f}) - \inf_{\mathcal{F}} M(f_0)) \leq 4M_n(\mathcal{H})$$

$$\text{où } M_n(\mathcal{H}) = E_{\substack{x_i, y_i \\ \varepsilon_1, \dots, \varepsilon_n \text{ iid } N(0, 1/2)}} \left( \sup_{\mathcal{H}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f(y_i, f_0(x_i)) \right)$$

Par principe de contraction Lipschitz,

$$M_n(\mathcal{H}) \leq G E \left( \sup_{\mathcal{H}} \frac{1}{n} \sum_{i=1}^n \varepsilon_i f_0(x_i) \right)$$

$$= G E \left( \sup_{\mathcal{H}} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \gamma_j (\omega_j^T x_i + b_j)_+ \right)$$

En réinjectant les contraintes et en utilisant le résultat de dualité  $\sup_{\|\gamma\|_1 \leq D} \sum^T \gamma = D \|\sum\|_\infty$ ,

$$M_n(\mathcal{H}) \leq G E \left( \sup_i \sup_{\|w_j\|_2 + b_j^2 / m = 1} D \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (w_j^T x_i + b_j)_+ \right| \right)$$



De plus, comme l'information est redondante en  $j$ ,

$$R_n(\Theta) \leq \mathbb{E} \left( \sup_{\|w\|_2 + |b| \leq R} \left| \frac{1}{n} \sum_{i=1}^n \varepsilon_i (w^T x_i + b)_+ \right| \right)$$

Nous aimerions utiliser le fait que  $(\cdot)_+$  est 1-Lipschitz, mais le 1.1 fait que nous ne sommes pas dans le cas vu en cours.

À la place, nous avons besoin de la version plus générale suivante:

Proposition: Si  $\psi_i$ ,  $\phi_i$  est L-Lipschitz, et  $\phi_i(0) = 0$ , alors

$$\mathbb{E}_\varepsilon \left( \sup_{\theta} \left| \sum_{i=1}^n \varepsilon_i \psi_i(a_i(\theta)) \right| \right) \leq 2 \mathbb{E}_\varepsilon \left( \sup_{\theta} \left| \sum_{i=1}^n \varepsilon_i a_i(\theta) \right| \right)$$

démonstration: Admis.

$$\text{Alors } R_n(\Theta) \leq 2 \mathbb{E} \left( \sup_{\|w\|_2 + |b| \leq R} \left| w^T \left( \frac{1}{n} \sum_{i=1}^n \varepsilon_i x_i \right) + b \left( \frac{1}{n} \sum_{i=1}^n \varepsilon_i \right) \right| \right)$$

$$\stackrel{\text{(dualité)}}{=} 2 \mathbb{E} \left( \sqrt{\left\| \frac{1}{n} \sum_{i=1}^n \varepsilon_i x_i \right\|_2^2 + R^2 \left( \frac{1}{n} \sum_{i=1}^n \varepsilon_i \right)^2} \right)$$



$$\begin{aligned}
 &\leq 2GD \sqrt{E(\dots)} \\
 (\text{Jensen}) & \\
 &= 2GD \sqrt{\left( \frac{1}{n} E(\|z\|_2^2) + \frac{M^2}{n} \right)} \\
 &\leq \frac{2GDM\sqrt{2}}{\sqrt{n}}
 \end{aligned}$$

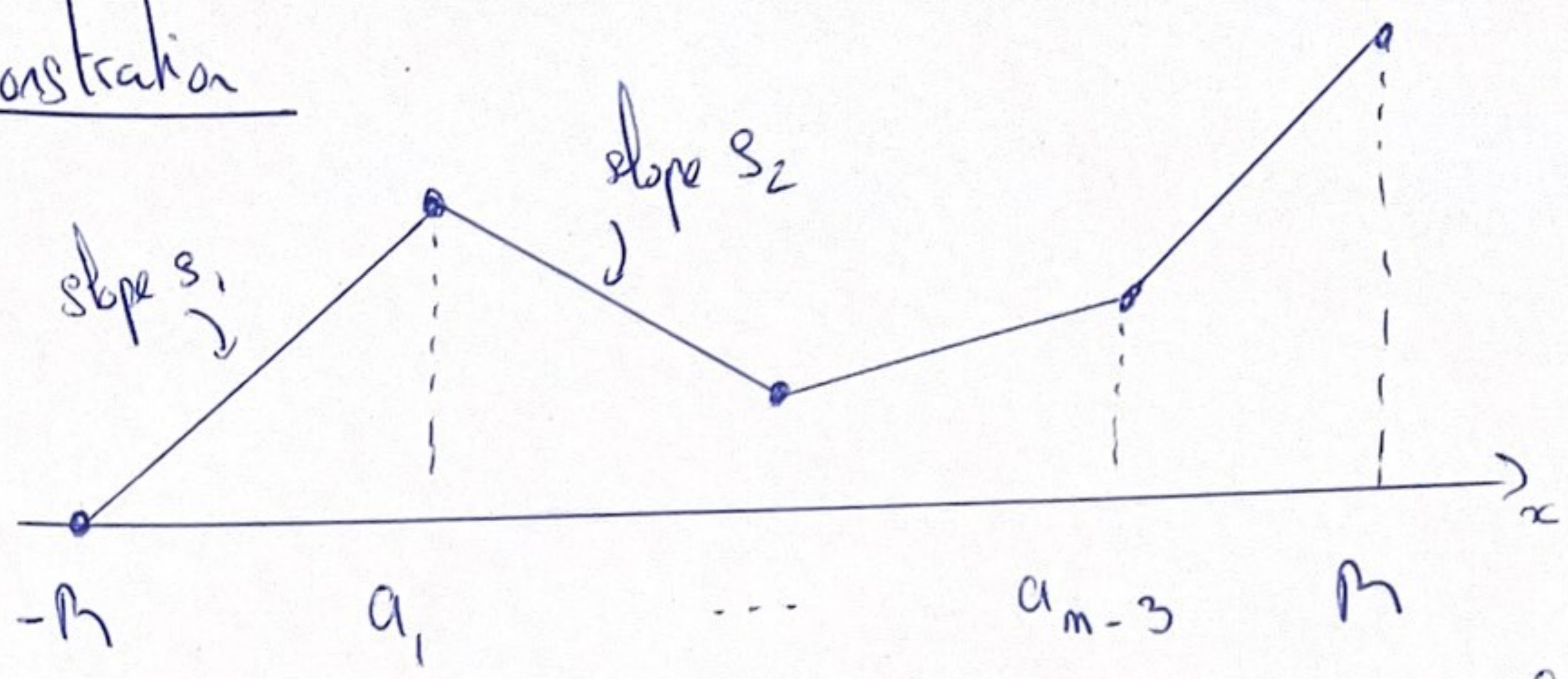
### IV. Approximation Universelle

Analyser l'erreur d'approximation est ~~très~~ en dehors des objectifs de ce cours. À la place, nous allons regarder la propriété d'approximation universelle des réseaux à une couche cachée. On se place dans le cas  $d=1$ .

Théorème: Toute fonction affine par morceaux sur  $(-M, M)$ , continue, et avec au plus  $m-2$  morceaux sur  $(-M, M)$  peut être approximée par un réseau Mélu à une couche cachée avec au plus  $m$  neurones.



démonstration



Commençons par le cas  $f(-M) = 0$ . Le cas général est obtenu en rajoutant un neurone et observant que

$$(x+M)_+ + (-x+M)_+ = 2M \text{ sur } [-M, M]$$

Alors  $f(x) = s_1 (x - (-M))_+$  sur  $[-M, a_1]$

$$f(x) = s_1 (x - (-M))_+ + (s_2 - s_1) (x - (a_1))_+$$

sur  $[a_1, a_2]$

...

□