

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.

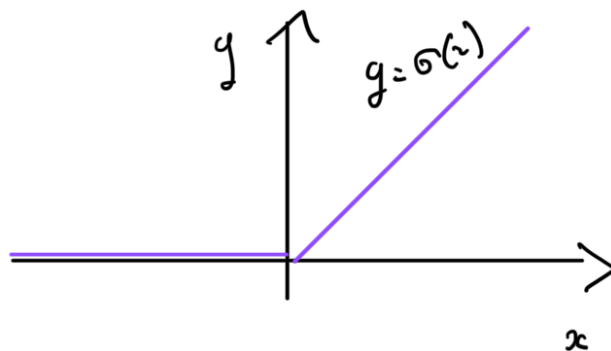
• Les activations

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

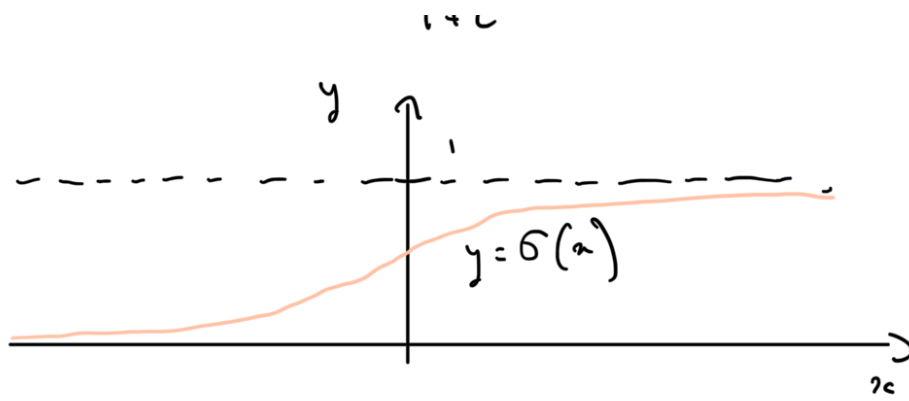
↑
fonction d'activation.

exemples classiques :

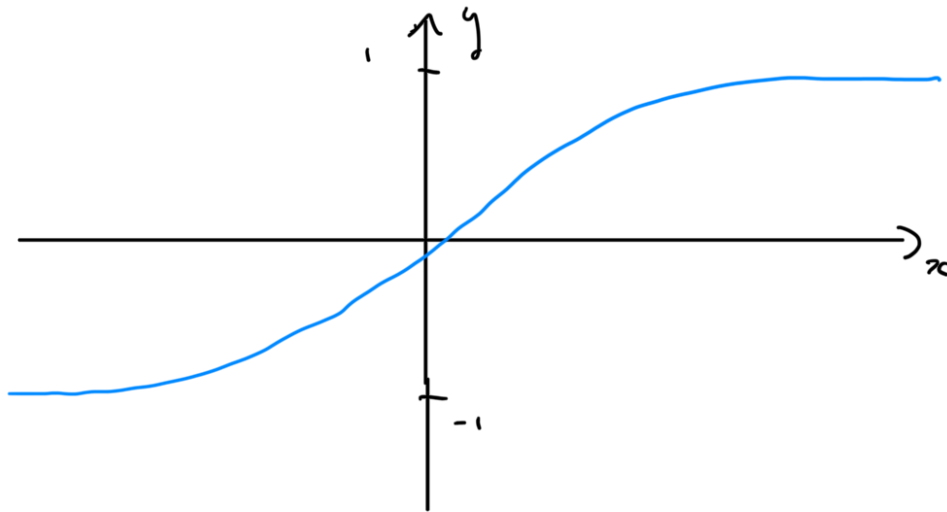
• ReLU : $\sigma(x) = (x)_+ = \max(0, x)$.



• 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

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

• En classification

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

C'est effectué avec la cross-entropy loss

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

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...