

# Backpropagat° algo & auto-encoders

"AI" = perf. tirées par 1 tech: DNN + algo apprentissage supervisé (deep learning)  
↳ reco patterns (image, texte, sons, ...)

avatars : → deep reinforce<sup>t</sup> learning (Alpha go)  
→ generative models (GAN, generative adversarial netw.)

• Sources - livre "deep-learning" GoodFellow, Bengio, Courville  
MIT Press online

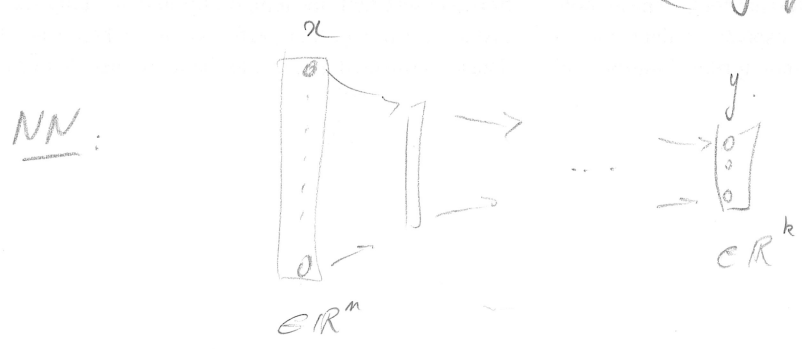
- Cours Stanford en ligne (Youtube)
- cours Berkeley : Ali Ghahsi (Ytubo)
- Wikipédia
- (Lebesgue Days)

↳ pas toujours bien formalisé / mathématique;  
bcp de "recettes" ... et de résultats impressionnants

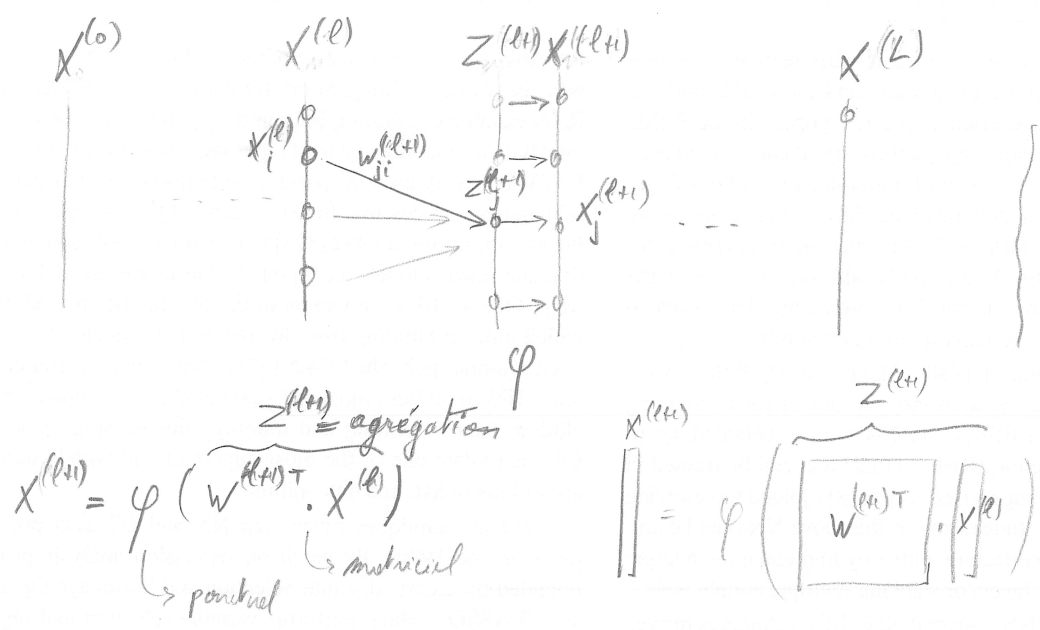
• apprent. supervise :

essayer d'approcher / d'apprendre  $f: \mathbb{R}^m \rightarrow \mathbb{R}^k$   
 $x \mapsto y = f(x)$

à partir de (bcp) d'exemples  $(x^i, y^i = f(x^i)) \quad i=1, \dots, N$



- $\min_w \mathbb{E} (\|y - f(x)\|^2)$   $\rightarrow$  erreur moyenne. où  $(x, y) \sim P_{x, y}$  loi des données
- on a des tirages iid  $(x^i, y^i) \quad i=1 \dots n$
- Remplace  $\mathbb{E}(\dots)$  par  $\frac{1}{n} \sum_i \|y^i - f_w(x^i)\|^2 = E(w)$   
 = erreur.
- Puis gradient descent:  $\nabla_w E$  et  $w^{(t+1)} \leftarrow w^{(t)} - \alpha \nabla_w E(w^{(t)})$
- en fait: procéder par mini-batches, puis reshuffling.  
 jusqu'à CV sur le dataset.
  - on va voir calcul grad. pour 1 donnée  $(x, y)$ .  
 or off

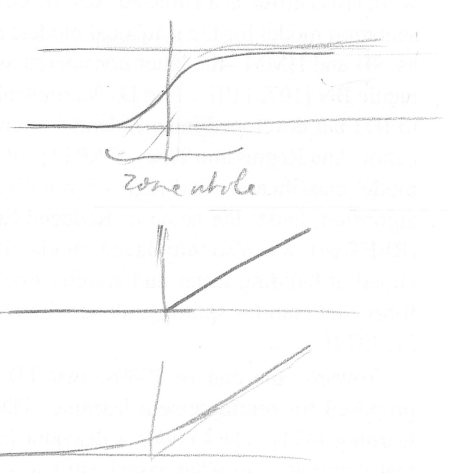


noter intensité  $w_{ji}$  va de  $i$  vers  $j$   
 d'où terme  $w^T$  dans les équations.

$\varphi =$  fonction d'activation (non linéaire, saillie)

- sigmoïde (fd. logistique)
- Relu
- Soft. +

$\varphi(x) = \frac{1}{1+e^{-x}}$   
 $\varphi(x) = \varphi(x)(1-\varphi(x))$   
 $\varphi(x) = \max(0, x)$   
 $\varphi(x) = \log(1+e^x)$



Erreur sur 1 donnée pour l'approximation de  $f$ .

soit  $f_w: x^{(0)} \mapsto x^{(L)}$   
 on a la donnée  $x$  et sortie désirée  $y$   
 mesurer écart de  $f(x)$  à  $y$ :  $E_v = \frac{1}{2} \|f_w(x) - y\|^2$  quadratique  
 $\hookrightarrow$  ajuster les poids  $w^{(n)}$  pour réduire l'erreur,  
 pb optimisat<sup>o</sup> numérique:  $\min_w E_v$   
 résolu par descente de gradient:

$$w^{(t+1)} = w^{(t)} - \alpha \underbrace{\nabla_w E_v(t)}_{\text{à calculer}}$$

(Retro) Propagation du gradient.

- On voudrait calculer les valeurs pour tout  $l=0 \dots L$

$$\nabla_{z^{(l)}} E = \begin{bmatrix} \frac{\partial E}{\partial z_1^{(l)}} \\ \vdots \\ \frac{\partial E}{\partial z_j^{(l)}} \end{bmatrix}$$

pour obtenir

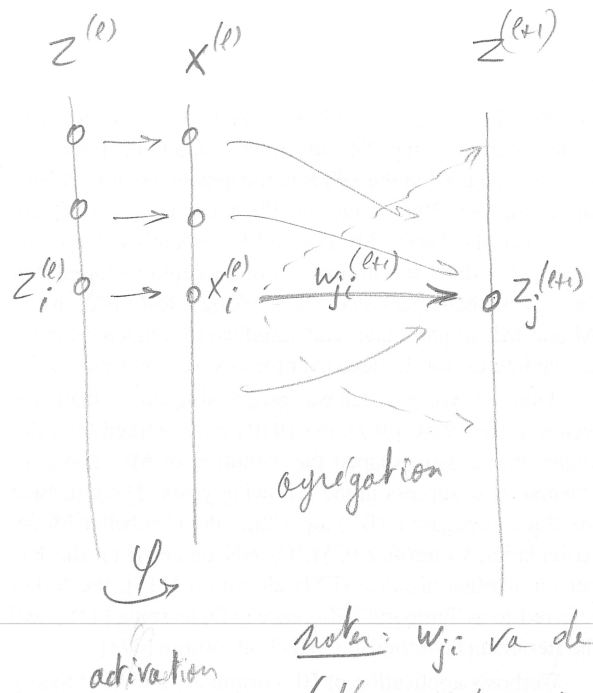
$$\frac{\partial E}{\partial w_{ji}^{(l+1)}} = \frac{\partial E}{\partial z_j^{(l+1)}} \cdot \frac{\partial z_j^{(l+1)}}{\partial w_{ji}^{(l+1)}}$$

$$= \frac{\partial E}{\partial z_j^{(l+1)}} \cdot x_i^{(l)}$$

(matricielle):

$$\frac{\partial E}{\partial w^{(l+1)}} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \cdot \overline{x^{(l)T}}$$

$$\nabla_{z^{(l+1)}} E$$



note:  $w_{ji}$  va de  $i$  vers  $j$   
(d'où  $w^T$  dans les eq°)

- récurrence arrière:  $E$  dépend de  $z^{(l)}$  via  $z^{(l+1)}$

$$\frac{\partial E}{\partial z_i^{(l)}} = \sum_j \frac{\partial E}{\partial z_j^{(l+1)}} \cdot \frac{\partial z_j^{(l+1)}}{\partial z_i^{(l)}} = \sum_j \frac{\partial E}{\partial z_j^{(l+1)}} \cdot \underbrace{\frac{\partial z_j^{(l+1)}}{\partial x_i^{(l)}}}_{w_{ji}^{(l+1)}} \cdot \underbrace{\frac{\partial x_i^{(l)}}{\partial z_i^{(l)}}}_{\psi'(z_i^{(l)})}$$

$$= \overline{w_{j,i}^{(l+1)}} \cdot \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \cdot \psi'(z_i^{(l)})$$

$$\nabla_{z^{(l+1)}} E$$

! vient atténuer le gradient si on prend la sigmoïde, hors zone utile.

(matriciellement):

$$\nabla_{z^{(l)}} E = \left( \overline{w^{(l+1)T}} \cdot \nabla_{z^{(l+1)}} E \right) * \psi'(z^{(l)})$$

$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} * \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

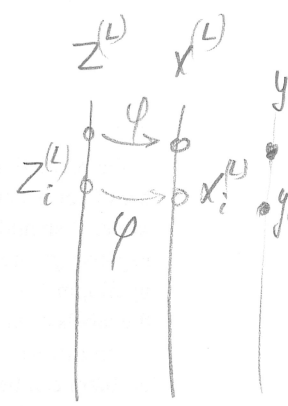
(produit terme à terme)  
non linéaire

Passer de  $\sigma$  à RELU a permis de propager le gradient + profond  $\Rightarrow$  deep neural nets.

• Initialisation

$$E = \frac{1}{2} \|f_w(x) - y\|^2 = \frac{1}{2} \|X^{(L)}(x) - y\|^2$$

$$\nabla_{z_i^{(L)}} = \frac{1}{2} \sum_i (x_i^{(L)}(x) - y_i)^2$$



$$\frac{\partial E}{\partial z_i^{(L)}} = \frac{\partial E}{\partial x_i^{(L)}} \cdot \frac{\partial x_i^{(L)}}{\partial z_i^{(L)}} = (x_i^{(L)} - y_i) \cdot \phi'(z_i^{(L)})$$

matriciellement:  $\nabla_{z^{(L)}} E = (X^{(L)} - y) \cdot \phi'(z^{(L)})$   
 ↳ produit terme à terme

Remarques

1) Algo: forward = calculer les \$x^{(l)}\$ \$z^{(l)}\$, backward = repasser gradient en arrière.

1) on prend souvent  $x_i^{(l+1)} = \phi(\sum_j w_{ji} x_j^{(l)} + b_i^{(l+1)})$   
 ↳ biais, change sens de réact° du neurone.  
 ⇔ rajouter ds \$x\_i^{(l)}\$ un "neurone" toujours à 1

2) on ne procède pas exemple par exemple par règle les poids, mais par "mini batches": paires \$(x^{(k)} y^{(k)})\$ d'exemples, \$k=1 \dots K\$

On prend alors  $E = \frac{1}{2n} \sum_k \|f_w(x^{(k)}) - y^{(k)}\|^2$

3) Le signal d'erreur incorpore souvent un terme de régularisation pour garder les poids petits.

$$E = \frac{1}{2} \|X^{(L)}(x) - y\|^2 + \frac{\lambda}{2} \sum_{l=1}^L \|W^{(l)}\|^2$$

↳ rajoute 1 terme linéaire  $\frac{\partial E}{\partial w^{(l+1)}} = \dots + \lambda \cdot W^{(l)}$

L'évolution des poids devient

$$W(t+1) = W(t) + \beta \left[ \underbrace{-\alpha \nabla_w E_w(t)}_{\text{gradient}} + \underbrace{+\alpha \lambda W(t)}_{\text{régul.}} \right] + \underbrace{(1-\beta) \Delta W(t)}_{\text{inertie dans la descente - évite de rester piégé ds min locaux.}}$$

inertie dans la descente - évite de rester piégé ds min locaux.

4) La profondeur de propag. repose sur RELLU  
↳ n pb linéaire + seuils d'activite.

- Bes. de réglages "à la main", expérimentaux.
- On connaît des struct. de rse qui marchent bien pour la reco de patterns en image (rse conv.) et en parole (short-termalng km).

- Les tête couches sont "standard", multi-ops, et pré-entraînées → on n'a plus qu'à entraîner les couches sups p l'appli donnée.

ex image = convolution locale + décimation.  
(conv. locale, invariance transla°, rotat°)

- Les dern. couches sont totalement connectées.  
↳ elles portent une descript° compacte des images en "features".

- La mise en corresp des features de haut niveau permet d'associer texte/son/image dans un m<sup>e</sup> domaine de représentation → interrogat° d'images.