

Part I

Théorie des réseaux de neurones

1 Backward propagation à travers le réseau

En utilisant la règle de la chaîne :

$$\frac{\partial E}{\partial w} = \sum_j \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial w}$$

Connaitre $\frac{\partial E}{\partial Y}$ et $\frac{\partial Y}{\partial w}$ nous assure donc de connaitre $\frac{\partial E}{\partial w}$.
Or on a, toujours d'après la règle de la chaîne :

$$\frac{\partial E}{\partial x_i} = \sum_j \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

Connaitre $\frac{\partial E}{\partial Y}$ et $\frac{\partial Y}{\partial X}$ nous assure donc de connaitre $\frac{\partial E}{\partial X}$.
Et si $n \in \mathbf{N}$ alors $\frac{\partial E}{\partial X_{n+1}}$ correspond à $\frac{\partial E}{\partial Y_n}$.

2 Fully connected

Forward propagation :

$$Y = XW + B$$

$$\text{Ici : } X = [x_1 \quad x_2 \quad \cdots \quad x_i], W = \begin{bmatrix} w_{11} & w_{12} & \cdots & w_{1j} \\ \vdots & \vdots & \ddots & \vdots \\ w_{i1} & w_{i2} & \cdots & w_{ij} \end{bmatrix}, B = [b_1 \quad b_2 \quad \cdots \quad b_i]$$

$$\text{D'où } Y = \left[\sum_i x_i \times w_{i1} + b_1 \quad \cdots \quad \sum_i x_i \times w_{ij} + b_j \right].$$

Backward propagation :

$$\frac{\partial E}{\partial W} = \begin{bmatrix} \frac{\partial E}{\partial w_{1,1}} & \cdots & \frac{\partial E}{\partial w_{1,j}} \\ \vdots & \ddots & \vdots \\ \frac{\partial E}{\partial w_{i,1}} & \cdots & \frac{\partial E}{\partial w_{i,j}} \end{bmatrix}$$

$$\frac{\partial E}{\partial X} = \frac{\partial E}{\partial Y} W^t$$

$$\frac{\partial E}{\partial W} = X^t \cdot \frac{\partial E}{\partial Y}$$

$$\frac{\partial E}{\partial B} = \frac{\partial E}{\partial Y}$$

Démonstration :

Règle de la chaîne :

$$\frac{\partial E}{\partial w_{i,j}} = \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial w_{i,j}} + \dots + \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial w_{i,j}} = \frac{\partial E}{\partial y_j} x_i$$

D'où :

$$\frac{\partial E}{\partial W} = \begin{bmatrix} \frac{\partial E}{\partial y_1} x_1 & \dots & \frac{\partial E}{\partial y_j} x_1 \\ \vdots & \ddots & \vdots \\ \frac{\partial E}{\partial y_1} x_i & \dots & \frac{\partial E}{\partial y_j} x_i \end{bmatrix} = X^t \frac{\partial E}{\partial Y}$$

$$\frac{\partial E}{\partial B} = \left[\frac{\partial E}{\partial b_1} \quad \dots \quad \frac{\partial E}{\partial b_j} \right]$$

Règle de la chaîne :

$$\frac{\partial E}{\partial b_i} = \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial b_i} + \dots + \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial b_i} = \frac{\partial E}{\partial y_i}$$

D'où

$$\frac{\partial E}{\partial B} = \left[\frac{\partial E}{\partial y_1} \quad \dots \quad \frac{\partial E}{\partial y_j} = \frac{\partial E}{\partial Y} \right]$$

$$\frac{\partial E}{\partial X} = \left[\frac{\partial E}{\partial x_1} \quad \dots \quad \frac{\partial E}{\partial x_i} \right]$$

Règle de la chaîne :

$$\frac{\partial E}{\partial x_i} = \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial x_i} + \dots + \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_i} = \frac{\partial E}{\partial y_1} w_{i,1} + \dots + \frac{\partial E}{\partial y_j} w_{i,j}$$

D'où :

$$\frac{\partial E}{\partial X} = \left[\frac{\partial E}{\partial y_1} w_{1,1} + \dots + \frac{\partial E}{\partial y_j} w_{1,j} \quad \dots \quad \frac{\partial E}{\partial y_1} w_{i,1} + \dots + \frac{\partial E}{\partial y_j} w_{i,j} \right] = \frac{\partial E}{\partial Y} W^t$$

3 Activation Layer

Forward propagation :

$$Y = [f(x_1) \quad \cdots \quad f(x_i)]$$

Backward propagation :

$$\frac{\partial E}{\partial X} = \frac{\partial E}{\partial Y} \odot f'(X) = \left[\frac{\partial E}{\partial y_1} \times f'(x_1) \quad \cdots \quad \frac{\partial E}{\partial y_i} \times f'(x_i) \right]$$

Où \odot est le produit d'Hadamard.

Démonstration :

On remarque $i=j$

$$\begin{aligned} \frac{\partial E}{\partial X} &= \left[\frac{\partial E}{\partial x_1} \quad \cdots \quad \frac{\partial E}{\partial x_i} \right] \\ &= \left[\frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial x_1} \quad \cdots \quad \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial x_i} \right] \\ &= \left[\frac{\partial E}{\partial y_1} f'(x_1) \quad \cdots \quad \frac{\partial E}{\partial y_i} f'(x_i) \right] \\ &= \left[\frac{\partial E}{\partial y_1} \quad \cdots \quad \frac{\partial E}{\partial y_i} \right] \odot [f'(x_1) \quad \cdots \quad f'(x_i)] \\ &= \left[\frac{\partial E}{\partial y_1} \quad \cdots \quad \frac{\partial E}{\partial y_i} \right] \odot f'(X) \end{aligned}$$

4 Théorie des GAN

4.1 Définitions

Définition 4.1.0.1 Tribu :

Soit E un ensemble. On appelle tribu tout sous-ensemble \mathcal{A} de $\mathcal{P}(E)$ tel que :

- $\emptyset \in \mathcal{A}$
- si $A \in \mathcal{A}$, alors $\bar{A} \in \mathcal{A}$
- si $(A_n)_{n \in \mathbb{N}}$ est une suite de \mathcal{A} , alors $\cup_{n \in \mathbb{N}} A_n \in \mathcal{A}$

Un espace mesurable est un ensemble muni d'une tribu.

Définition 4.1.0.2 Mesure :

Soit (E, \mathcal{A}) un espace mesurable.

On appelle mesure sur E toute application $\mu : \mathcal{A} \rightarrow \mathbb{R}^+$ vérifiant :

- $\mu(\emptyset) = 0$
- Si $(A_n)_{n \in \mathbb{N}}$ est une suite d'éléments de \mathcal{A} deux à deux disjoints, alors $\mu(\sqcup_{n \in \mathbb{N}} A_n) = \sum_{n \in \mathbb{N}} \mu(A_n)$

Définition 4.1.0.3 Fonction mesurable : Soit (E, \mathcal{T}) et (F, \mathcal{A}) deux espaces mesurables.

$f : E \rightarrow F$ est dite mesurable ssi $f^{-1}(A) \in \mathcal{T}$

Définition 4.1.0.4 Variable aléatoire continue :

$X : \Omega \rightarrow \mathbb{R}$ mesurable.

Remarque 1 Les densité de probabilités sont des mesures dont l'intégrale sur l'espace vaut 1.

Théorème 4.1.0.1 (Ω, \mathcal{T}, P) un espace probabilisé, $X : \Omega \rightarrow \mathbb{R}$ de loi notée P_X .

Soit $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ mesurable.

$$\text{Alors } \int_{\Omega} \phi(X(\omega)) dP(\omega) = \int_{\mathbb{R}} \phi(x) dP_X(x)$$

Corollaire 4.1.0.1.1 $E(\phi(X)) = \int_{\mathbb{R}} \phi(x) f(x) dx$

4.2 Discriminant optimal

Démonstration :

On note p_g la densité créée par le générateur, p_r la densité de données réelles.

A G fixé, on cherche D_G^* le discriminant optimal (qui maximise la fonction $D \mapsto V(G, D)$)

$$\begin{aligned} V(G, D) &= E_{x \sim p_r}[\log(D(x))] + E_{z \sim p_z}[\log(1 - D(G(z)))] \\ &= \int_x p_r(x) \log(D(x)) dx + \int_z p_z(z) \log(1 - D(G(z))) dz \end{aligned}$$

changement de variable : $x = G(z)$ (par Théorème de transfert)

$$\int_z p_z(z) \log(1 - D(G(z))) dz = \int_x p_g(x) \log(1 - D(x)) dx$$

Alors :

$$\begin{aligned} V(G, D) &= \int_x p_g(x) \log(1 - D(x)) + p_r(x) \log(D(x)) dx \\ \text{Soit } f : x &\mapsto a \cdot \log(x) + b \cdot \log(1 - x), \quad a, b > 0 \end{aligned}$$

f est maximale en $\frac{a}{a+b}$ sur $[0, 1]$

$$\text{Alors } V(G, D) \leq \int_x p_g(x) \log\left(1 - \frac{p_r(x)}{p_g(x) + p_r(x)}\right) + p_r(x) \log\left(\frac{p_r(x)}{p_g(x) + p_r(x)}\right) dx$$

$D_G^*(x) \in [0, 1]$ (C'est bien une probabilité, cohérent)

On pose $C(G) = V(G, D_G^*)$

$$D_G^*(x) = \frac{p_r(x)}{p_g(x) + p_r(x)}$$

4.3 Divergence de Kullback-Leibler

Caractère positif :

$$\begin{aligned} D(p \parallel q) &= \int_x p(x) \cdot \log\left(\frac{p(x)}{q(x)}\right) dx = - \int_x p(x) \cdot \log\left(\frac{q(x)}{p(x)}\right) dx \\ &\geq - \log\left(\int_x p(x) \cdot \frac{q(x)}{p(x)} dx\right) \quad (\text{Inégalité de Jensen}) \\ &= - \log\left(\int_x q(x) dx\right) \\ &= 0 \end{aligned}$$

Cas d'égalité:

La divergence de Kullback-Leibler est positive.

De plus le log est strictement concave : il ya donc égalité dans l'égalité précédente si et seulement si $x \mapsto \frac{p(x)}{q(x)}$ est constante p-presque partout.

Alors il existe $\alpha \in \mathbb{R}_+^*$ tel que $p = \alpha \cdot q$.

De plus, $\int_x p(x) dx = \int_x q(x) dx = 1$, donc $\alpha = 1$ i.e. $p = q$ p-presque partout.

4.4 Divergence de Jensen-Shannon

Caractère fini et distance :

$\forall p, q, JS(p \parallel q)$ est fini :

$$\begin{aligned} D(p \parallel \frac{p+q}{2}) &= \int_x p(x) \cdot \log\left(\frac{2 \cdot p(x)}{p(x)+q(x)}\right) dx \\ &\leq \int_x p(x) \cdot \log\left(\frac{2 \cdot p(x)}{p(x)}\right) dx \\ &\leq \log(2) \end{aligned}$$

La divergence de Jensen-Shannon est bien une distance, car symétrique, positive, et nulle si et seulement si q et p sont égales p -presque partout.

Calcul de $C(G)$:

$$\begin{aligned} \int_x p_r(x) \log\left(\frac{p_r(x)}{p_g(x)+p_r(x)}\right) dx &= \int_x p_r(x) [\log\left(\frac{p_r(x)}{\frac{p_g(x)+p_r(x)}{2}}\right) - \log(2)] dx \\ &= \int_x p_r(x) \log\left(\frac{p_r(x)}{\frac{p_g(x)+p_r(x)}{2}}\right) dx - \log(2) \\ &= D(p_r \parallel \frac{p_r+p_g}{2}) - \log(2) \end{aligned}$$

$$\begin{aligned} \int_x p_g(x) \log\left(\frac{p_r(x)}{p_g(x)+p_r(x)}\right) dx &= \int_x p_g(x) [\log\left(\frac{p_r(x)}{\frac{p_g(x)+p_r(x)}{2}}\right) - \log(2)] dx \\ &= \int_x p_g(x) \log\left(\frac{p_r(x)}{\frac{p_g(x)+p_r(x)}{2}}\right) dx - \log(2) \\ &= D(p_g \parallel \frac{p_r+p_g}{2}) - \log(2) \end{aligned}$$

Alors $C(G) = 2 \cdot JS(p_r \parallel p_g) - \log(4)$

5 Distance de Wassertein

Soit $\theta \in [-1; 1]$ Si $\forall(x, y) \in P, x = 0, y \sim U([0; 1])$ et $\forall(x, y) \in Q, x = \theta, y \sim U([0; 1])$

$|\theta| \neq 0 :$

$$\begin{aligned} D_{KL}(P|Q) &= \sum_{x=0, y \sim U([0;1])} (1 \cdot \log \frac{1}{0}) = +\infty \\ D_{KL}(Q|P) &= \sum_{x=\theta, y \sim U([0;1])} (1 \cdot \log \frac{1}{0}) = +\infty \\ D_{JS}(Q|P) &= \frac{1}{2} \cdot \left(\sum_{x=0, y \sim U([0;1])} 1 \cdot \log \frac{1}{1/2} + \sum_{x=\theta, y \sim U([0;1])} 1 \cdot \log \frac{1}{1/2} \right) = \log(2) \\ W(P, Q) &= |\theta| \end{aligned}$$

$|\theta| = 0 :$

$$W(P, Q) = D_{KL}(P|Q) = D_{KL}(Q|P) = D_{JS}(Q|P) = |\theta| = 0$$