



$$\text{ReLU}(x) = \max(0, x)$$

$$\text{softmax}(X) = \left(\frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \right)_{1 \leq i \leq n}$$

↓
vecteur taille n

fonction de décision : prend $S(X) = (x_i = 1 \text{ si } x_i = \max_{1 \leq j \leq n} x_j ; 0 \text{ sinon})$

Fonction de perte : $L = -\frac{1}{m} \sum_{k=1}^m (y_k \log(\hat{y}_k) + (1-y_k) \log(1-\hat{y}_k))$: (produit scalaire np. vdot.)

m le nombre d'exemples, (y_k) les labels, colonnes de Y
 (\hat{y}_k) les prédictions, colonnes de \hat{Y}

On note $z^1 = W_1 X + b_1$

$a_1 = \text{relu}(z_1)$; $a_0 = x$

$\forall l \in \llbracket 2, L \rrbracket : z_l = W_l a_{l-1} + b_l$

$a_l = \text{relu}(z_l)$

$\hat{y} = \text{softmax}(a_L)$

Objectif : calculer $\frac{\partial L}{\partial W_l}$ et $\frac{\partial L}{\partial b_l}$, pour $l \in \llbracket 1, L \rrbracket$

étapes : calculer $\frac{\partial L}{\partial \hat{y}}$, $\frac{\partial L}{\partial a_L}$, $\frac{\partial L}{\partial z_L}$, $\frac{\partial L}{\partial z_l}$ et $\frac{\partial L}{\partial a_l}$, pour $l \in \llbracket 1, L-1 \rrbracket$

$\frac{\partial L}{\partial \hat{y}} = \frac{1}{m} \left(\frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}} \right)$ (division coordonnée à coordonnée, np.divide)

$\frac{\partial L}{\partial a_l} = \frac{1}{m} ((1-y) \times \hat{y} - y \times (1-\hat{y}))$ (* : multiplication coordonnée à coordonnée, np.multiply)

$\forall l \in \llbracket 1, L \rrbracket : \frac{\partial L}{\partial z_l} = \frac{\partial L}{\partial a_l} * 1_{\mathbb{R}_+^*}(z_l)$

$\forall l \in \llbracket 1, L-1 \rrbracket : \frac{\partial L}{\partial a_l} = {}^T W_{l+1} \frac{\partial L}{\partial z_{l+1}}$ (produit matriciel np.dot)

$\frac{\partial L}{\partial W_l} = \frac{\partial L}{\partial z_l} {}^T a_{l-1}$

d'où $\forall l \in \llbracket 1, L \rrbracket$

$\frac{\partial L}{\partial b_l} = \text{np.sum} \left(\frac{\partial L}{\partial z_l}, \text{axis}=1 \right)$

Démonstration :

• $\frac{\partial L}{\partial \hat{y}} = \left(\frac{\partial L}{\partial \hat{y}_{li}} \right)_{\substack{1 \leq l \leq m \\ 1 \leq i \leq 10}}$ \hat{y} matrice de taille $m \times 10$ de même pour y
 \hat{y}_{li} ses coordonnées

$$\left(\frac{\partial L}{\partial \hat{y}_{li}} \right) = - \frac{1}{m} \sum_{k=1}^m \sum_{j=1}^{10} y_{kj} \underbrace{\frac{\partial \log \hat{y}_{kj}}{\partial \hat{y}_{li}}}_{\frac{y_{ki}}{\hat{y}_{li}} \delta_{kl}} + \sum_{j=1}^{10} (1 - y_{kj}) \underbrace{\frac{\partial \log(1 - \hat{y}_{kj})}{\partial \hat{y}_{li}}}_{-\frac{1 - y_{kj}}{1 - \hat{y}_{li}} \delta_{kl}} = \frac{1}{m} \left(\frac{1 - y_{li}}{1 - \hat{y}_{li}} - \frac{y_{li}}{\hat{y}_{li}} \right)$$

d'où $\boxed{\frac{\partial L}{\partial \hat{y}} = \frac{1}{m} \left(\frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} \right)}$ (np. divide)

• $\frac{\partial L}{\partial a_L} = \left(\frac{\partial L}{\partial a_{Li}} \right)_{1 \leq i \leq 10}$

$$\frac{\partial L}{\partial a_{Li}} = \sum_{k=1}^{10} \frac{\partial L}{\partial \hat{y}_k} \times \underbrace{\frac{\partial \hat{y}_k}{\partial a_{Li}}}_{\downarrow} = \frac{\partial L}{\partial \hat{y}_i} * \hat{y}_i * (1 - \hat{y}_i)$$

$$\begin{aligned} \frac{\partial}{\partial a_{Li}} (\text{softmax } a_{Li}) &= \text{softmax } a_{Li} * (1 - \text{softmax } a_{Li}) \\ &= \hat{y}_i * (1 - \hat{y}_i) \delta_{ki} \end{aligned}$$

d'où $\boxed{\frac{\partial L}{\partial a_L} = \frac{1}{m} ((1 - y) * \hat{y} - y * (1 - \hat{y}))}$ (np. multiply)

• $\frac{\partial L}{\partial z_L} = \left(\frac{\partial L}{\partial z_{Li}} \right)_{1 \leq i \leq \text{taille couche}}$ l fixé, $l \in \llbracket 1, L \rrbracket$

$$\frac{\partial L}{\partial z_{Li}} = \sum_k \frac{\partial L}{\partial a_{Lk}} \underbrace{\frac{\partial a_{Lk}}{\partial z_{Li}}}_{\downarrow} = 1_{\mathbb{R}_*^+}(z_{Li}) * \frac{\partial L}{\partial a_{Li}}$$

$$\frac{\partial}{\partial z_{Li}} \text{relu}(z_{Lk}) = 1_{\mathbb{R}_*^+}(z_{Li}) \delta_{ik}$$

d'où $\boxed{\frac{\partial L}{\partial z_L} = \frac{\partial L}{\partial a_L} * 1_{\mathbb{R}_*^+}(z_L)}$ (np multiply) $\forall l \in \llbracket 1, L \rrbracket$

• $\frac{\partial L}{\partial a_l} = \left(\frac{\partial L}{\partial a_{li}} \right)_{i \leq \text{taille couche } l}$ $l \text{ fixé, } l \in [1, L-1]$

$\frac{\partial L}{\partial a_{li}} = \sum_k \frac{\partial L}{\partial z_{li+1,k}} \cdot \boxed{\frac{\partial z_{li+1,k}}{\partial a_{li}}} = \sum_k \left(\frac{\partial L}{\partial z_{li+1,k}} \right) w_{ki}^{l+1} = \sum_k ({}^t w_{ki}^{l+1}) \left(\frac{\partial L}{\partial z_{li+1,k}} \right) = L_i \left({}^t w_{ki}^{l+1} \frac{\partial L}{\partial z_{li+1,k}} \right)$

\downarrow
 w_{ki}^{l+1} car $(z_{li+1})_k = \sum_{t=1}^{l+1} w_{kt}^{l+1} (a_l)_t + b_{l+1,k}$

d'où $\boxed{\frac{\partial L}{\partial a_l} = {}^t w_{l+1} \frac{\partial L}{\partial z_{l+1}}} \text{ (mp.dot)} \quad \forall l \in [1, L-1]$

Avec ces 4 résultats, on peut obtenir le $\frac{\partial L}{\partial w_l}$ et $\frac{\partial L}{\partial b_l}$

• $\frac{\partial L}{\partial w_l} = \left(\frac{\partial L}{\partial w_{ij}^l} \right)_{i,j}$, w_l matrice dont la taille dépend de la taille des couches $l-1$ et l $\forall l \in [1, L]$

$\frac{\partial L}{\partial w_{ij}^l} = \sum_k \frac{\partial L}{\partial z_k^{l+1}} \cdot \boxed{\frac{\partial z_k^{l+1}}{\partial w_{ij}^l}} = \frac{\partial L}{\partial z_i^{l+1}} \cdot (a_{l-1})_j$

\downarrow
 $\frac{\partial}{\partial w_{ij}^l} \left(\sum_u w_{ku} a_u^{l-1} \right) = \delta_{ik} a_j^{l-1}$

donc $\boxed{\frac{\partial L}{\partial w_l} = \frac{\partial L}{\partial z_l} \cdot {}^t a_{l-1}} \text{ (produit matriciel mp.dot)} \quad \forall l \in [1, L]$
 $a_0 = x.$

• $\frac{\partial L}{\partial b_l} = \left(\frac{\partial L}{\partial b_i^l} \right)_i$ b_l vecteur de taille de la couche l

$\frac{\partial L}{\partial b_i^l} = \sum_k \frac{\partial L}{\partial z_k^{l+1}} \cdot \boxed{\frac{\partial z_k^{l+1}}{\partial b_i^l}} = \frac{\partial L}{\partial z_i^{l+1}} \rightarrow \text{taille } [1 \times m] : \text{ il faut sommer les lignes}$

\downarrow
 $= \delta_{ki}$

$\rightarrow \sum_{j=1}^m \frac{\partial L}{\partial z_j^{l+1}}$ car le biais s'ajoute à w_a par broadcasting (pté de la librairie numpy)

d'où $\frac{\partial L}{\partial b_l} = \sum_{k=1}^m \left(\frac{\partial L}{\partial z_k} \right) \text{ (vecteur de taille couche } l)$ encore noté $\boxed{\frac{\partial L}{\partial b_l} = \text{np.sum} \left(\frac{\partial L}{\partial z_l}, \text{axis} = 1 \right)}$