

Ordonnés: $\frac{\partial E^{(p)}}{\partial \tau_k} = \left(\delta_k^{sortie(p)} \right) = \begin{pmatrix} 2(\hat{y}_0^{(p)} - y_0^{(p)}) g'_s(\tau_0) \\ \vdots \\ 2(\hat{y}_k^{(p)} - y_k^{(p)}) g'_s(\tau_k) \\ \vdots \\ 2(\hat{y}_s^{(p)} - y_s^{(p)}) g'_s(\tau_s) \end{pmatrix} = \begin{pmatrix} \delta_0^{sortie(p)} \\ \vdots \\ \delta_k^{sortie(p)} \\ \vdots \\ \delta_s^{sortie(p)} \end{pmatrix} = \text{vecteur à } N_s \text{ composantes } (k=0, \dots, s)_{N_s-1}$

Alors $\frac{\partial E^{(p)}}{\partial \sigma_j} = \sum_{k=0}^s \frac{\partial E^{(p)}}{\partial \tau_k} \frac{\partial \tau_k}{\partial \sigma_j}$ mais $\tau_k = \sum_j U_j z_{jk} = \sum_j g(\sigma_j) \cdot z_{jk}$

$= \sum_{k=0}^s \frac{\partial E^{(p)}}{\partial \tau_k} \cdot g'(\sigma_j) z_{jk} \Leftarrow \text{d'où } \frac{\partial \tau_k}{\partial \sigma_j} = g'(\sigma_j) z_{jk}$

$= \sum_{k=0}^s [g'(\sigma_j) z_{jk}] \times \frac{\partial E^{(p)}}{\partial \tau_k} = \text{vecteur à } N_c \text{ composantes } (j=0, \dots, N_c-1)$

$= \delta_j^{caché(p)}$

$\vec{\delta}^{caché(p)} = \begin{pmatrix} \delta_0^{caché(p)} \\ \vdots \\ \delta_j^{caché(p)} \\ \vdots \\ \delta_{N_c-1}^{caché(p)} \end{pmatrix} = \begin{pmatrix} g'(\sigma_0) z_{00} & \dots & g'(\sigma_0) z_{0k} & \dots & g'(\sigma_0) z_{0s} \\ g'(\sigma_j) z_{j0} & \dots & g'(\sigma_j) z_{jk} & \dots & g'(\sigma_j) z_{js} \\ \vdots & & \vdots & & \vdots \\ g'(\sigma_{N_c-1}) z_{N_c-1/0} & \dots & g'(\sigma_{N_c-1}) z_{N_c-1/k} & \dots & g'(\sigma_{N_c-1}) z_{N_c-1/s} \end{pmatrix} \begin{pmatrix} \delta_0^{sortie(p)} \\ \vdots \\ \delta_k^{sortie(p)} \\ \vdots \\ \delta_s^{sortie(p)} \end{pmatrix}$

$\delta_j^{caché(p)}$: nb neurones cachées

$\delta_s^{sortie(p)}$: nb neurones de sortie

rq: si $g = \text{logistique}$, $g'(\sigma_j) = g(\sigma_j) \cdot [1 - g(\sigma_j)]$

$g_s = \text{logistique}$, $g'_s(\tau_k) = g_s(\hat{y}_k) \cdot [1 - g_s(\hat{y}_k)]$

Algorithme:

- 1) calculer toutes les prédictions: (phase forward) 1 échantillon = $(\vec{x}^{(p)}, \vec{y}^{(p)})$
 $p = 0, \dots, P-1$ (pour tous les échantillons du jeu de données), faire: a) $X_i^{(p)}$ = données de la couche d'entrée, $i = 1, \dots, N_{\text{entrée}} - 1$
 et $X_0^{(p)} = 1, 0$ (neurone de biais)

b) calculer $\sigma_j^{(p)} = \sum_{i=0}^{N_e-1} w_{ij} X_i^{(p)}$ et $U_j^{(p)} = g(\sigma_j^{(p)})$
 $j = 1, \dots, N_c - 1$
 $\sigma_0^{(p)} = 1, 0$ (neurone de biais)

c) calculer $\tau_k^{(p)} = \sum_{j=0}^{N_c-1} z_{jk} U_j^{(p)}$ et $\hat{y}_k^{(p)} = g_s(\tau_k^{(p)})$
 $k = 0, \dots, N_{\text{sortie}} - 1$
 $\sigma_0^{(p)} = 1, 0$ (neurone de biais)

d) Conserver les valeurs de: $X_i^{(p)}, U_j^{(p)}, \hat{y}_k^{(p)}$
 si logistique

- 2) Calculer les erreurs pour $p = 0, \dots, P-1$ (i.e. chaque échantillon).

a) en couche de sortie: $\delta_k^{(p)\text{sortie}} = 2g'_s(\tau_k) (\hat{y}_k^{(p)} - y_k^{(p)}) = 2g_s(\hat{y}_k^{(p)}) [1 - g_s(\hat{y}_k^{(p)})] \times [\hat{y}_k^{(p)} - y_k^{(p)}]$
 $k = 0, \dots, N_{\text{sortie}} - 1$
 pour chaque $p = 0, \dots, P-1$

b) en couche cachée: $\delta_j^{(p)\text{cachée}} = g'(\sigma_j) \sum_{k=0}^{N_s-1} z_{jk} \delta_k^{(p)\text{sortie}}$
 $j = 1, \dots, N_{\text{caché}} - 1$
 (car le neurone caché $j=0$ est celui de biais donc pas d'erreur.)
 (prédiction) (vraie valeur voir (a))
 si logistique

- 3) calculer les gradients et les corrections:

a) pour $p = 0, \dots, P-1$, calculer: $\frac{\partial E^{(p)}}{\partial z_{jk}} = \delta_k^{(p)\text{sortie}} \cdot U_j$
 et $\frac{\partial E^{(p)}}{\partial w_{ij}} = X_i \times \delta_j^{(p)\text{caché}}$ (voir page d)

b) d'où les gradients totaux moyennés: $\frac{\partial E}{\partial z_{jk}} = \frac{1}{P} \sum_{p=0}^{P-1} \frac{\partial E^{(p)}}{\partial z_{jk}} = \frac{1}{P} \sum_{p=0}^{P-1} \delta_k^{(p)\text{sortie}} \cdot U_j$
 $\frac{\partial E}{\partial w_{ij}} = \frac{1}{P} \sum_{p=0}^{P-1} \frac{\partial E^{(p)}}{\partial w_{ij}} = \frac{1}{P} \sum_{p=0}^{P-1} \delta_j^{(p)\text{caché}} \cdot X_i$

c) d'où les corrections: $\Delta z_{jk} = -\alpha \frac{\partial E}{\partial z_{jk}}$ et $\Delta w_{ij} = -\alpha \frac{\partial E}{\partial w_{ij}}$, à ajouter à w_{ij} et z_{jk} .

On a:
$$\frac{\partial E^{(p)}}{\partial w_{ij}} = 2 \sum_k (\hat{y}_k^{(p)} - y_k^{(p)}) \frac{\partial y_k^{(p)}}{\partial w_{ij}}$$

$$\frac{\partial y_k}{\partial w_{ij}} = \frac{\partial g_s(\tau_k)}{\partial w_{ij}} = g'_s(\tau_k) \frac{\partial \tau_k}{\partial w_{ij}}$$

$$\frac{\partial \tau_k}{\partial w_{ij}} = \frac{\partial [\sum_j z_{jk} u_j]}{\partial w_{ij}}$$

$$= \sum_j z_{jk} \frac{\partial u_j}{\partial w_{ij}} = \sum_j z_{jk} \frac{\partial g(\sigma_j)}{\partial w_{ij}}$$

$$= \sum_j z_{jk} g'(\sigma_j) \frac{\partial \sigma_j}{\partial w_{ij}}$$

$$\frac{\partial (\sum_i w_{ij} x_i)}{\partial w_{ij}}$$

$$= z_{jk} g'(\sigma_j) x_i$$

$$\Rightarrow \frac{\partial E^{(p)}}{\partial w_{ij}} = 2 \sum_k (\hat{y}_k^{(p)} - y_k^{(p)}) g'_s(\tau_k) z_{jk} g'(\sigma_j) x_i$$

$$= 2 x_i \sum_k g'(\sigma_j) z_{jk} \times (\hat{y}_k^{(p)} - y_k^{(p)}) g'_s(\tau_k)$$

$$= 2 x_i \cdot \delta_j^{\text{caché}(p)}$$

Couche sortie

↓ n° couche -

$$\delta_d^m = (\hat{y}_d - y_d) g'_0(a_d^m)$$

$\hat{y} = g_0(a_1^m)$

indice du cas dans l'échantillon (Wellennummer)

Couche cachée k.

$$\delta_j^k = o_j^k (1 - o_j^k) \sum_{l=1}^{k+1} w_{jl}^{k+1} \delta_l^{k+1}$$

et 1 couche

$$\delta_j^{m-1} = o_j^{m-1} (1 - o_j^{m-1}) \sum_{l=1}^m w_{jl}^m \delta_l^m$$

nb nodes
en couche (k+1)

$$\delta_j^k = \frac{\partial E}{\partial a_j^k} = \sum_{l=1}^{k+1} \frac{\partial E}{\partial a_l^{k+1}} \cdot \frac{\partial a_l^{k+1}}{\partial a_j^k}$$

$$= \sum_{l=1}^{k+1} \delta_l^{k+1} \frac{\partial a_l^{k+1}}{\partial a_j^k}$$

(l ≠ 0 car
node biais
n'intervient pas)

mais $a_l^{k+1} = \sum_{j=1}^k w_{jl}^{k+1} g(a_j^k) \Rightarrow \frac{\partial a_l^{k+1}}{\partial a_j^k} = w_{jl}^{k+1} g'(a_j^k)$

w_{ij}^k = poids du node j de la couche k, des nœuds i

o_i^k = output du node i de la couche k $\Rightarrow o_i^k = g(a_i^k)$

$a_i^k = \sum \text{pondérées } w_{ji}^k g(a_j^{k-1})$

1) forward: $\forall d=1, \dots, N (\hat{y}_d, y_d) \rightarrow \hat{y}_d, a_j^k, o_j^k = g(a_j^k)$

\hat{y}_d : indice de couche
 a_j^k : indice de node
 o_j^k : indice de couche
 $j=0, \dots, m$: input | out
 m : hidden

2) backward $\forall d=1, \dots, N$:

faire calculer $\frac{\partial \mathcal{E}_d}{\partial w_{ij}^k}$ pour chaque poids w_{ij}^k

δ_j^k : couche $= m, m-1, \dots, 0$
 o_j^k : dans la couche
 i : node de couche k
 j : node de couche $k-1$

pour cela: a) calcul erreur couche sortie: $\delta_j^m = g'_j(a_j^m) (\hat{y}_d - y_d)$

δ_j^m : couche sortie
 j : 1 node en sortie

b) calcul erreur couche cachée $k = m-1, m-2, \dots$

$\delta_j^k = g'_j(a_j^k) \sum_{l=k+1}^m w_{jl}^{k+1} \delta_l^{k+1}$

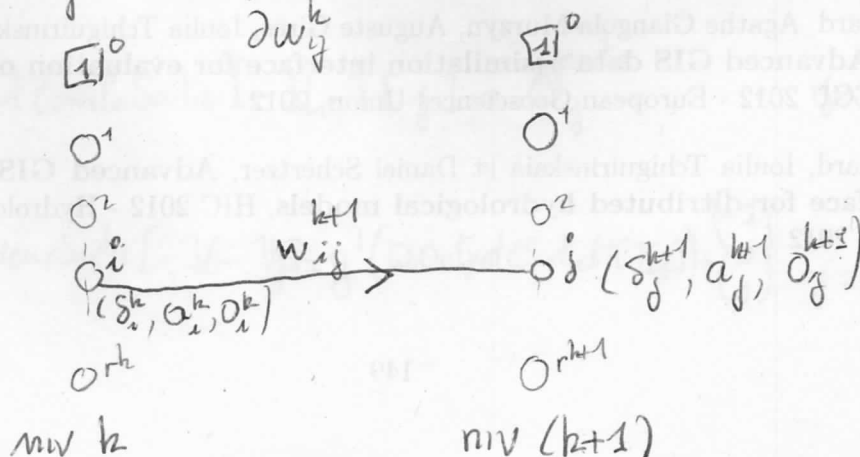
δ_j^k : couche k
 j : node j
 l : car on ne prend pas node biais ($l=0$)
 w_{jl}^{k+1} : nb nodes couche $k+1$
 δ_l^{k+1} : couche $k+1$

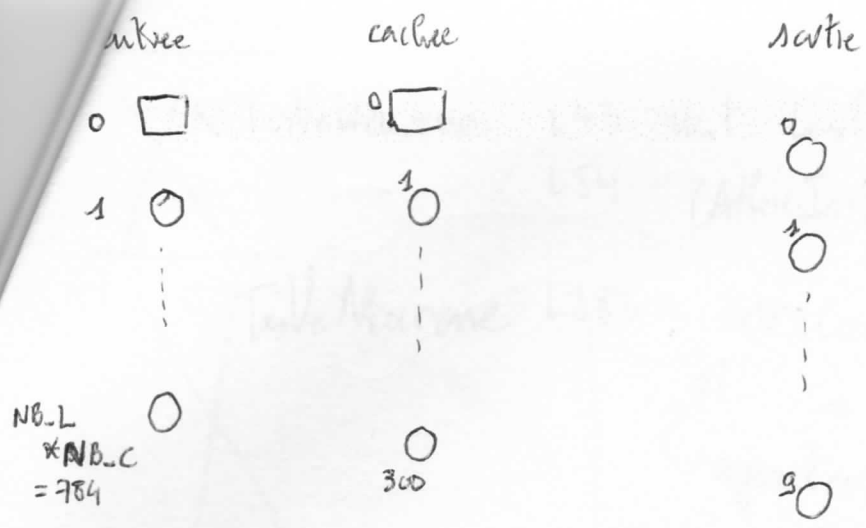
c) $\frac{\partial \mathcal{E}_d}{\partial w_{ij}^k} = \delta_j^k o_i^{k-1}$

3) Combinaisons gradients individuels

$$\frac{\partial \mathcal{E}}{\partial w_{ij}^k} = \frac{1}{N} \sum_{d=1}^N \frac{\partial}{\partial w_{ij}^k} \left(\frac{1}{2} (\hat{y}_d - y_d)^2 \right) = \frac{1}{N} \sum_{d=1}^N \frac{\partial \mathcal{E}_d}{\partial w_{ij}^k}$$

alors $\Delta w_{ij}^k = -\alpha \frac{\partial \mathcal{E}}{\partial w_{ij}^k}$

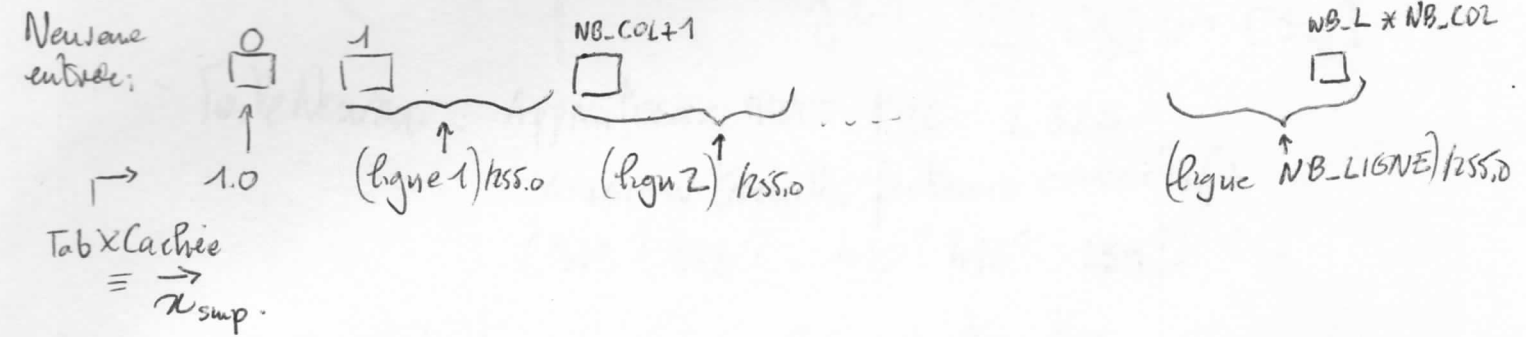




NB_NEURONES
 - C. ENTREE = $28 \times 28 + 1 = 785$
 NB_L = 784
 NB_C. CACHEE = 301
 NB_SORTIE = 10

Pour l'échantillon smp ($= 0, \dots, 60000 - 1$)

Entree: $(TabDonnee MNIST [\frac{smp}{TBLOC}][\frac{smp}{TBLOC}]. Image [ligne][col]. / 255) \cdot 1 \rightarrow 784$
 TabDonnee $[\text{---}] [\text{---}]$



Cachee: $lePNC. TabOutputCoucheCachee [-] [-] [0] = 1.0$
 $[-] [-] [j] = \sigma_j^{cachee}$

Sortie: $lePNC. TabOutputCoucheSortie [-] [-] [j] = \sigma_j^{sortie}$

$j = 1, \dots, lePNC_0$
 couche cachee
 - 1
 $= 300$
 $j = 0, \dots, lePNC_1$
 - 1
 $= 9$

En sortie: $lePNC \rightarrow TabNeurSortie [-] [-] [j] = g' (TabOutputSortie [-] [-] [j]) \cdot \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}$ $j = 0, \dots, 9$

X - Test Neureme. L43 init les Poids $\begin{bmatrix} 1 \\ 2 \end{bmatrix} [i][j] = 0,5$ et
L54 (Alloc Init les Poids) non 0.00005

X - Test Neureme L283 TabXCachee $[j * NB_colonnes + h]$
= (double) $\sim Image [j][h]$
↑ ajouter/neureme de biais/couche entrée
- TabXCachee a 1 case de plus

L520/L333 $\begin{bmatrix} 0 \\ 1 \end{bmatrix} T \sim [0] = 1.0$
 $\begin{bmatrix} 0 \\ 1 \end{bmatrix} T \sim [1 + j * NB_colonnes + h]$
= $\sim Image [j][h] / 255$
normalisation à $[0, 1]$ ↑

X - Test Neureme L288 Le PNL Tab Output Couche Cachee
et 255 ↑ ajouter un neurone de biais = 1.0
en cas [0] et tout de calcul
 j de 1 à 500 $[j] \mapsto [1+j]$

X - Test Neureme. Apprentissage MNIST PNL L324-
↳ revoir tous les facteurs correctifs
(379? 386? 410? 416? 489?)

X - MNIST-Test.h ← voir les constantes

- ① neurone de biais
- ② normaliser 0-255 → 0.0 / 1.0
- ③ facteurs correctifs ↗ gradient
↳ retro propag.

Traitement avec la fonction de coût logarithmique et la fonction d'activation Softmax en sortie.

la fonction softmax est adaptée comme fonction d'activation de la couche de sortie car elle permet de normaliser les productions T_0, \dots, T_9 des neurones de sortie afin d'obtenir des probabilités $\hat{y}_0, \hat{y}_1, \dots, \hat{y}_9$ dont la somme vaut 1.

On a:

$$\hat{y}_k = \frac{e^{T_k}}{\sum_{i=0}^9 e^{T_i}} = \hat{y}_k(T_0, \dots, T_9) \quad (\text{vrai au plus } \hat{y}_k = g_s(T_k))$$

la fonction softmax a ceci de particulier qu'elle prend en donnée d'entrée non seulement la production T_k du neurone de sortie k , mais aussi celle de tous les autres neurones (T_0, \dots, T_9). Elle diffère en cela notablement des autres fonctions d'activation (par exemple la fonction sigmoïde), ^{lorsqu'elles sont appliquées à un neurone k} qui ne prennent en donnée d'entrée que la seule valeur T_k produite par le neurone de sortie k .

la fonction de coût logarithmique est mathématiquement construite pour pénaliser les erreurs de prédictions, en s'appuyant sur le comportement de la fonction Log. (voir les pages pour justifications avec les mains).

On a:

$$C = - \sum_{p=0}^{P-1} \sum_{k=0}^9 y_k^{(p)} \log(\hat{y}_k^{(p)}) = \sum_{p=0}^{P-1} C^{(p)}$$

\uparrow vraie valeur \uparrow prédiction pour l'échantillon p d'être le chiffre k .
 $y_k^{(p)} = 1$ si $k = \text{vrai chiffre}$
 $= 0$ si $k \neq \text{vrai chiffre}$

avec $C^{(p)} = \text{coût logarithmique de l'échantillon } p = - \sum_{k=0}^9 y_k^{(p)} \log(\hat{y}_k^{(p)})$

On note que si le chiffre vrai est k (alors $y_k^{(p)} = 1$ et $y_{j \neq k}^{(p)} = 0$) et si la prédiction est bonne (alors $\hat{y}_k^{(p)} \rightarrow 1$), le seul terme non nul ($-y_k^{(p)} \log(\hat{y}_k^{(p)})$) $\rightarrow 0$ car $\log x \rightarrow 0$ quand $x \rightarrow 1$.

En revanche, si la prédiction est mauvaise (alors $\hat{y}_k^{(p)} \rightarrow 0$ et $\hat{y}_j^{(p)} = 1$ pour une valeur de $j \neq k$), alors le seul terme non nul dans $C^{(p)}$, qui est $(-y_j^{(p)} \log \hat{y}_j^{(p)}) \rightarrow +\infty$ car $\log x \rightarrow -\infty$ quand $x \rightarrow 0$.

dérivées partielles de la fonction de coût (pour un échantillon).

$$\frac{\partial C^{(p)}}{\partial x} = \sum_{k=0}^9 - \left(y_k^{(p)} / \hat{y}_k^{(p)} \right) \times \frac{\partial \hat{y}_k^{(p)}}{\partial x} \quad \text{(on mettra l'indice (p) dans la suite des calculs)}$$

$$\hat{y}_k = \text{softmax}(\tau_0, \tau_k, \tau_9) = e^{\tau_k} / \left(\sum_{e=0}^9 e^{\tau_e} \right)$$

On calcule les dérivées partielles de \hat{y}_k par rapport à τ_q , $q=0, \dots, 9$.

• si $q=k$, τ_k apparaît 2 fois dans softmax (au numérateur et au dénominateur).

$$\frac{\partial \hat{y}_k}{\partial \tau_k} = \frac{(e^{\tau_k} \times (\sum_e e^{\tau_e}) - e^{\tau_k} \times e^{\tau_k})}{(\sum_e e^{\tau_e})^2}$$

$$= \frac{(e^{\tau_k} \times (e^{\tau_k} + \sum_{e \neq k} e^{\tau_e}) - e^{\tau_k} \times e^{\tau_k})}{(\sum_e e^{\tau_e})^2}$$

$$= \frac{e^{\tau_k} \times (\sum_{e \neq k} e^{\tau_e})}{(\sum_e e^{\tau_e})^2}$$

$$= \frac{e^{\tau_k}}{(\sum_e e^{\tau_e})} \times \frac{(\sum_e e^{\tau_e} - e^{\tau_k})}{\sum_e e^{\tau_e}} = \hat{y}_k \times (1 - \hat{y}_k)$$

• si $q \neq k$, τ_q apparaît une seule fois dans softmax (au dénominateur)

$$\frac{\partial \hat{y}_k}{\partial \tau_q} = \frac{e^{\tau_k} \times (-e^{\tau_q})}{(\sum_e e^{\tau_e})^2} = \frac{e^{\tau_k}}{(\sum_e e^{\tau_e})} \times \frac{(-e^{\tau_q})}{(\sum_e e^{\tau_e})} = -\hat{y}_k \times \hat{y}_q$$

Ceci permet d'obtenir les dérivées partielles suivantes:

$$\frac{\partial C^{(p)}}{\partial \tau_q} = \sum_{k=0}^9 - \left(\frac{y_k}{\hat{y}_k} \right) \frac{\partial \hat{y}_k}{\partial \tau_q} = \sum_{k \neq q} - \frac{y_k}{\hat{y}_k} \times (-\hat{y}_k \times \hat{y}_q) - \frac{y_q}{\hat{y}_q} \times \hat{y}_q (1 - \hat{y}_q)$$

$$= \hat{y}_q \underbrace{\sum_{k \neq q} y_k - y_q (1 - \hat{y}_q)}_{=1} = \hat{y}_q \underbrace{\sum_{k=0}^9 y_k - y_q}_{=1} = \hat{y}_q - y_q, \quad q=0, \dots, 9.$$

Donc: $\frac{\partial C^{(p)}}{\partial \tau_k} = \hat{y}_k - y_k$, vecteur à 10 composantes

$$\boxed{\frac{\partial C^{(p)}}{\partial z_{jk}}} = \frac{\partial C^{(p)}}{\partial \tau_k} \times \frac{\partial \tau_k}{\partial z_{jk}}$$

$$= (\hat{y}_k - y_k) \times U_j$$

$$\text{avec: } \tau_k = \sum_j U_j \cdot z_{jk}$$

$$\boxed{\frac{\partial C^{(p)}}{\partial U_j}} = \sum_k \frac{\partial C^{(p)}}{\partial \tau_k} \times \frac{\partial \tau_k}{\partial U_j} = \sum_k (\hat{y}_k - y_k) \times z_{jk}$$

Posons: $\frac{\partial C^{(p)}}{\partial \tau_k} = \left(\delta_k^{sortie(p)} \right) = \begin{pmatrix} \hat{y}_0^{(p)} - y_0^{(p)} \\ \vdots \\ \hat{y}_g^{(p)} - y_g^{(p)} \end{pmatrix} = \begin{pmatrix} \delta_0^{sortie(p)} \\ \vdots \\ \delta_g^{sortie(p)} \end{pmatrix} = \text{vecteur à } N_s \text{ composés}$
 $(k=0, \dots, g)$
 $N_s - 1$.

Alors: $\frac{\partial C^{(p)}}{\partial \sigma_j} = \sum_{k=0}^g \frac{\partial C^{(p)}}{\partial \tau_k} \times \frac{\partial \tau_k}{\partial \sigma_j}$
 $= \sum_{k=0}^g (\hat{y}_k^{(p)} - y_k^{(p)}) \times g'(\sigma_j) \times z_{jk}$
 $= \sum_{k=0}^g [g'(\sigma_j) \cdot z_{jk}] \times \delta_k^{sortie(p)}$
 $= \delta_j^{cache(p)}$.

mais $\tau_k = \sum_j U_j z_{jk} = \sum_j g(\sigma_j) z_{jk}$
 fonction d'activation de la couche ~~cache~~ cachée, qui prend le seul σ_j pour calculer U_j .
 donc: $\frac{\partial \tau_k}{\partial \sigma_j} = g'(\sigma_j) \times z_{jk}$.

$\vec{\delta}^{cache(p)} = \begin{pmatrix} \delta_0^{cache(p)} \\ \vdots \\ \delta_j^{cache(p)} \\ \vdots \\ \delta_{N_c-1}^{cache(p)} \end{pmatrix} = \begin{pmatrix} g'(\sigma_0) z_{00} \dots g'(\sigma_0) z_{0g} \\ g'(\sigma_1) z_{10} \dots g'(\sigma_1) z_{1g} \\ \vdots \\ g'(\sigma_{N_c-1}) z_{N_c-1,0} \dots g'(\sigma_{N_c-1}) z_{N_c-1,g} \end{pmatrix} \begin{pmatrix} \delta_0^{sortie(p)} \\ \vdots \\ \delta_k^{sortie(p)} \\ \vdots \\ \delta_g^{sortie(p)} \end{pmatrix}$
 $N_c = \text{nb neurones cachés}$
 $g = \text{nb neurones de sortie}$

ex: si $g = \text{sigmoïde}$, $g'(\sigma_j) = g(U_j) \times [1 - g(U_j)]$.

ex: $\frac{\partial C^{(p)}}{\partial w_{ij}} = \frac{\partial C^{(p)}}{\partial U_j} \times \frac{\partial U_j}{\partial w_{ij}} = \sum_k (\hat{y}_k - y_k) z_{jk} \times g'(\sigma_j) \times \frac{\partial \sigma_j}{\partial w_{ij}}$
 $= X_i \times \sum_k [g'(\sigma_j) z_{jk} \times (\hat{y}_k - y_k)]$
 $= X_i \times \delta_j^{cache(p)}$

Méthode:

1) calculer les prédictions:

a) \hookrightarrow si $X_i^{(p)}$ = donnée du neurone d'entrée i pour l'échantillon p .
(rq: $X_0^{(p)} = 1$ car neurone de biais)

b) \hookrightarrow $\sigma_j^{(p)} = \sum_i w_{ij} X_i^{(p)}$ et $U_d^{(p)} = g(\sigma_j^{(p)})$
 \hookrightarrow activation couche cachée

c) \hookrightarrow $\tau_k^{(p)} = \sum_j z_{jk} U_j^{(p)}$ et $\hat{y}_k^{(p)} = g_s(\tau_0^{(p)}, \dots, \tau_9^{(p)})$
 \hookrightarrow softmax

d) \hookrightarrow conserver les valeurs de: $X_i^{(p)}$,

2) calculer toutes les erreurs:

a) couche sortie: $\delta_k^{soutie(p)} = \frac{\partial C^{(p)}}{\partial \tau_k} = \underbrace{\hat{y}_k^{(p)} - y_k^{(p)}}_{\substack{\text{probabilité} \\ \text{produite par softmax}}} \cdot \underbrace{1}_{\substack{\text{si } k = \text{vrai chiffre} \\ 0 \text{ si } k \neq \text{vrai chiffre}}}$

b) couche cachée: $\delta_j^{caché(p)} = g'(\sigma_j) \times \sum_k z_{jk} \delta_k^{soutie(p)}$
 \downarrow
 $= g(U_j) \times [1 - g(U_j)]$ si $g = \text{sigmoïde}$

3) calculer les gradients et les corrections

a) pour $p=0, \dots, P-1$ calculer: $\frac{\partial C^{(p)}}{\partial z_{jk}} = \delta_k^{soutie(p)} \times U_j$

et $\frac{\partial C^{(p)}}{\partial w_{ij}} = X_i \times \delta_j^{caché(p)}$

b) d'où les gradients cumulés, totaux:

$$\frac{\partial C}{\partial z_{jk}} = \sum_{p=0}^{P-1} \frac{\partial C^{(p)}}{\partial z_{jk}} = \sum_{p=0}^{P-1} \delta_k^{soutie(p)} \times U_j$$

$$\frac{\partial C}{\partial w_{ij}} = \sum_{p=0}^{P-1} \frac{\partial C^{(p)}}{\partial w_{ij}} = \sum_{p=0}^{P-1} \delta_j^{caché(p)} \times X_i$$

c) d'où les corrections:

$$\Delta z_{jk} = -\alpha \frac{\partial C}{\partial z_{jk}} \quad \text{et} \quad \Delta w_{ij} = -\alpha \frac{\partial C}{\partial w_{ij}}$$

à ajouter à w_{ij} et z_{jk} .

Finalement, changer \hat{y} en C et la fonction d'activation de sortie en Softmax change peu de choses