Relations of derives the qualques function to base

#I. MO signal(x) = 
$$\frac{1}{2}$$
 (tanh ( $\frac{1}{2}$ -x) = 1)

signal(x) =  $\frac{1}{2}$  (tanh ( $\frac{1}{2}$ ) + 1) ( $\frac{1}{2}$ ) 2 signal(x) - 1 = tanh ( $\frac{1}{2}$ )

2 signal(x) =  $\frac{2-1-\exp(x)}{1+\exp(x)}$  =  $\frac{1-\exp(x)}{1+\exp(x)}$  =  $\frac{\exp(\frac{x}{2})}{1+\exp(\frac{x}{2})}$  ( $\frac{1-\exp(x)}{1+\exp(x)}$ )

=  $\frac{2\exp(\frac{x}{2})-\exp(\frac{x}{2})\exp(x)}{1+\exp(x)}$  =  $\frac{\exp(\frac{x}{2})-\exp(\frac{x}{2})}{1+\exp(x)}$  =  $\frac{1}{2}$  tanh ( $\frac{\pi}{2}$ )

=  $\exp(\frac{x}{2})+\exp(\frac{x}{2})$  +  $\exp(\frac{x}{2})$  +  $\exp(\frac{x}{2})$  =  $\frac{1}{2}$  tanh ( $\frac{\pi}{2}$ )

#13. MQ la signal(x) =  $\ln(\frac{1}{1+\exp(-x)})$  =  $\ln(1)-\ln(1+\exp(-x))$  =  $\frac{1}{2}$  -  $\frac{1}{2}$ 

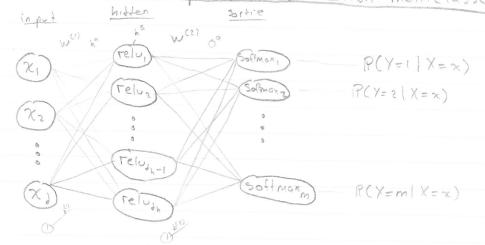
#8. 
$$\|x\|_2^2 = \sum_i x_i^2$$
 Ecrire  $\frac{\partial \|x\|_2^2}{\partial x}$ 

$$\frac{\partial \Sigma_{1} \chi_{1}^{2}}{\partial \chi} = \left(\frac{\partial (\chi_{1}^{2} + \chi_{2}^{2} + \chi_{1}^{2} + \ldots + \chi_{n}^{2})}{\partial \chi_{1}}\right) = \left(\frac{\partial \chi_{1}}{\partial \chi_{2}}\right) = \left(\frac{\partial \chi_{1}}{\partial \chi_{2}$$

$$\frac{\partial \Sigma_{1}|x_{1}|}{\partial x} = \frac{\partial (|x_{1}|+|x_{2}|+...+|x_{n}|)}{\partial x_{1}} = \frac{\operatorname{sign}(x_{1})}{\operatorname{sign}(x_{2})} = \operatorname{sign}(x)$$

$$\frac{\partial (|x_{1}|+|x_{2}|+...+|x_{n}|)}{\partial x_{n}} = \frac{\operatorname{sign}(x_{1})}{\operatorname{sign}(x_{n})} = \frac{\operatorname{sign}(x_{1})}{\operatorname{sign}(x_{n})}$$

2. Calcul du gradient pour l'optimisation des paramètres dun réseau de neurones pour la classification multiclasse



#1. W(1) daxd, entre conche input et hidden, b ER?

Donner la formule de coicul du vecteur de préactivations des neurones de la hidden layer h<sup>a</sup> given x as input. first in a metrix form (h°=...). Then how to compute one element (h°=...)

Ecrime le vecter des sortier des neurones de la hidden layer his en fonction de his

#2 Donner dimension de W (2) et 6(2) (entre hidden and output).

Donner la formule de vecteur d'activation des neurones de jartier o en fonction de h', puis 0 à  $W^{(2)}$  est  $m \times dh$  et  $b^{(2)} \in \mathbb{R}^m$   $0^a = b^{(2)} + W^{(2)}h^s$ 

$$O^{\alpha} = \begin{pmatrix} b_{1} \\ b_{2} \\ b_{3} \end{pmatrix} + \begin{bmatrix} \omega_{11} \\ \omega_{22} \\ \omega_{32} \\ \omega_{33} \end{pmatrix} = \begin{pmatrix} c_{2} \\ b_{3} \\ c_{23} \\ \omega_{34} \end{pmatrix} = \begin{pmatrix} c_{2} \\ b_{34} \\ c_{33} \\ c_{34} \\ c_{35} \\ c_{$$

#3. La sortie des neurones de sortie est 0°= soft max (0°). Préciser 0x.

Demontrer que les ox sont positifs et somment à 1. Pourquoi est-ce important?

$$O_{k}^{s} = softmax(o^{q})_{k} = \frac{exp(o_{k}^{q})}{\sum_{i=1}^{m} exp(o_{i}^{q})}$$

Les Ox sont positifs per définitions de exp(x) (i.e. YxER, exp(x) > 0).

$$\sum_{k=1}^{m} o_{k}^{s} = \sum_{k=1}^{m} \frac{\exp(o_{k}^{a})}{\sum_{i=1}^{m} \exp(o_{i}^{a})} = \sum_{k=1}^{m} \exp(o_{i}^{a})$$

Il est important que les où soient positif et qu'ils somment à 1, as cela permet d'interpréter où comme P(Y=k|X=x) (c'est-à-dire qu'on interprète où comme étant la probabilité que l'entrée x soit de la classe k).

$$L(x,y) = -\log\left(\frac{\exp(O_y^a)}{\sum_{i=1}^{m} \exp(O_i^a)}\right) = \log\left(\sum_{i=1}^{m} \exp(O_i^a)\right) - \log\left(\exp(O_y^a)\right)$$

$$= \log\left(\sum_{i=1}^{m} \exp(O_i^a)\right) - O_y^a$$

Formuler  $\widehat{R}$ . Intiques precisement lensemble  $\widehat{\Theta}$  des paramètres du réseau. Indiquer à combien de paramètres scalaires no cela correspond. Formuler le problème d'aptimisation qui correspond à l'entrainement du réseau pour trouver une voleir paramètres.  $\widehat{R} = \widehat{h} \ \widehat{\Sigma}[\widehat{E}] \ L(\widehat{X}^{(i)}, \widehat{g}^{(i)}) = \widehat{h} \ \widehat{\Sigma}[\widehat{E}] \ (\log(\widehat{\Sigma}_{j=1}^{m} \exp(\widehat{O}_{j}^{\alpha}(\widehat{X}^{(i)}))) - \widehat{O}_{g}^{(i)}(\widehat{X}^{(i)}))$   $\widehat{\Theta} = \{W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}\}$ 

no = dhxd +dh+ mxdh+ m

Le problème d'optimisation qui correspond à l'entremement du réseau permettant de trouver une voleur optimale des paramètres est argmin R(A, D+min).

#6. Exprimer avec un bref pseudo-code la descente de gradient pour ce problème initialize 0

for N itérations:  $\theta \leftarrow \theta - \eta \cdot \left(\frac{1}{n} \cdot \sum_{i=1}^{n} \frac{\partial}{\partial \theta} \left(\log \left(\sum_{j=1}^{m} \exp(O_j^2(x^{(j)})\right) - O_j^{(j)}(x^{(j)})\right)\right)$ 

H7. MQ  $\frac{\partial L}{\partial Q} = O^{S} - Onehot_{m}(y)$ Pour  $K \neq y$ ,  $\frac{\partial L(X_{1}y)}{\partial Q_{K}^{2}} = \frac{\partial log(\Sigma_{1}^{m}, exp(Q_{1}^{Q})) - Q_{2}^{Q}}{\partial Q_{K}^{2}} = \frac{\partial log(\Sigma_{1}^{m}, exp(Q_{1}^{Q}))}{\partial Q_{K}^{2}}$ 

 $= \frac{1}{\sum_{j=1}^{m} \exp(o_{j}^{\alpha})} \cdot \frac{\partial \sum_{j=1}^{m} \exp(o_{k}^{\alpha})}{\partial o_{k}^{\alpha}} = \frac{\exp(o_{k}^{\alpha})}{\sum_{j=1}^{m} \exp(o_{j}^{\alpha})} = o_{k}^{\alpha}$ 

 $\frac{\partial L(x,y)}{\partial o_{3}^{2}} = \frac{\partial \left(\log \left(\sum_{j=1}^{m} \exp(o_{j}^{2})\right) - o_{3}^{2}\right)}{\log \left(\sum_{j=1}^{m} \exp(o_{j}^{2})\right)} = \frac{\exp(o_{3}^{2})}{\sum_{j=1}^{m} \exp(o_{j}^{2})}$ 

Alors,  $\frac{\partial L(x,y)}{\partial o^{\alpha}} = \begin{pmatrix} \frac{\partial L}{\partial o^{\alpha}} \\ \frac{\partial L}{\partial o^{\alpha}} \end{pmatrix} = \begin{pmatrix} \frac{o^{S}}{o^{S}} \\ \frac{o^{S}}{o^{S}} \end{pmatrix} - \begin{pmatrix} \frac{o}{o} \\ \frac{o}{o} \end{pmatrix} = o^{S} - one hotm(y)$ 

#8. Donner l'expression correspondente en numpy

gradoa = 0.5 - np. eye cm) [y-1]

car y & El,..., m3 et le vecteur one hotm à des index de 0 à m-1

#II. Colcoler 
$$\frac{\partial L}{\partial h^{S}}$$

$$\frac{\partial L}{\partial h^{S}} = \sum_{k=1}^{m} \frac{\partial L}{\partial O_{k}^{\alpha}} \cdot \frac{\partial O_{k}^{\alpha}}{\partial h^{S}} = \sum_{k=1}^{m} \frac{\partial L}{\partial O_{k}^{\alpha}} \cdot \frac{\partial (b^{(2)} + \sum_{i=1}^{d_{1}} W_{ki}^{(2)} h^{S})}{\partial O_{k}^{\alpha}} \cdot \frac{\partial O_{k}^{\alpha}}{\partial h^{S}} = \sum_{k=1}^{m} \frac{\partial L}{\partial O_{k}^{\alpha}} \cdot \frac{\partial O_{k}^{\alpha}}{\partial h^{S}} \cdot \frac{\partial O_{k}^{\alpha}}{\partial O_{k}^{\alpha}} \cdot \frac{\partial O_{k}^{\alpha}}{\partial h^{S}} + \dots + O_{m}^{S} W_{mj}^{(2)}$$

$$= \sum_{k=1}^{m} \frac{\partial L}{\partial O_{k}^{\alpha}} \cdot \frac{\partial V_{kj}^{(2)}}{\partial O_{k}^{\alpha}} - W_{kj}^{(2)} + \dots + O_{m}^{S} W_{mj}^{(2)} + \dots + O_{m}^{S} W_{mj}^{(2)}$$

$$= \sum_{k=1}^{m} \frac{\partial L}{\partial O_{k}^{\alpha}} \cdot \frac{\partial V_{kj}^{(2)}}{\partial O_{k}^{\alpha}} - W_{kj}^{(2)} - W_{kj}^{(2)} + \dots + O_{m}^{S} W_{mj}^{(2)} + \dots + O_{m}^{S} W_{mj}^{(2)}$$

$$= \sum_{k=1}^{m} \frac{\partial L}{\partial O_{k}^{\alpha}} \cdot \frac{\partial V_{kj}^{(2)}}{\partial O_{k}^{\alpha}} - W_{kj}^{(2)} - W_{kj}^{(2)} + \dots + O_{m}^{S} W_{mj}^{(2)} + \dots + O_{m}^{S} W_$$

#12. Exprimer sous forme matricielle, Preiciser les dimensions, en numpy grad-hs=?

[ = 0 × W Kdh - W y dh

où west drxm, os est mxl et W(2) Ey, :] Test drxl

#13. Calcules 
$$\frac{\partial L}{\partial h_{i}^{\alpha}}$$

$$\frac{\partial L}{\partial h_{i}^{\alpha}} = \frac{\partial L}{\partial h_{i}^{\beta}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i} - W_{y,i}^{(2)}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha} - W_{y,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} = \left(\sum_{k=1}^{m} O_{k}^{\beta} W_{k,i}^{\alpha}\right) \cdot \frac{\partial \left(\operatorname{rect}(h_{i}^{\alpha})\right)}{\partial h_{i}^{\alpha}} =$$

#14. Exprimer sous forme motricielle, Preciser les informations, donner l'équivolent rumpy
$$\frac{\partial L}{\partial h^{a}} = \left(\frac{\partial L}{\partial h^{s}}\right) \odot \begin{pmatrix} \mathbb{I}_{Eha} > 0.3 & (h^{a}) \\ \mathbb{I}_{Eha} > 0.3 & (h^{a}) \end{pmatrix}$$

où de de vecteur contenant les fonctions indicatrices sont du x l

rector\_indicator = np.array ([life > 0 else O for e in hs])
grad-ha = np. multiply (grad\_hs, vector\_indicator)

#15.	Cajculer 2 W(1) et 2 b(1)
	$\frac{\partial L}{\partial W_{3,k}^{(1)}} = \frac{\partial L}{\partial h_{j}^{\alpha}} = \frac{\partial L}{\partial W_{j,k}^{(2)}} = \frac{\partial L}{\partial W_{j,k}^{(2$
	$= \left(\sum_{k=1}^{\infty} O_{k}^{3} W_{kj}^{(2)} - W_{kj}^{(2)}\right) 1_{\xi h_{j}^{2} > 03} (h_{j}^{3}) \chi_{p} = \frac{\partial L}{\partial h_{j}^{2}} \chi_{p}^{2} = \frac{\partial L}{\partial h_{j}^{2}} \chi_{p}^{2}$
	$\frac{\partial L}{\partial b_{j}} = \frac{\partial L}{\partial h_{j}} = \frac{\partial L}{\partial h_{j}} = \frac{\partial L}{\partial h_{j}} = \frac{\partial L}{\partial h_{j}} = \frac{\partial L}{\partial h_{j}}$
	$\frac{\partial h^{\alpha}}{\partial \Gamma} \cdot x^{1} \frac{\partial h^{\alpha}}{\partial \Gamma} \cdot x^{2} = 0 \cdot 0 \cdot \frac{\partial h^{\alpha}}{\partial \Gamma} \cdot x^{3}$ $\frac{\partial h^{\alpha}}{\partial \Gamma} \cdot x^{1} \frac{\partial h^{\alpha}}{\partial \Gamma} \cdot x^{2} = 0 \cdot 0 \cdot \frac{\partial h^{\alpha}}{\partial \Gamma} \cdot x^{3}$ $\frac{\partial h^{\alpha}}{\partial \Gamma} \cdot x^{1} \frac{\partial h^{\alpha}}{\partial \Gamma} \cdot x^{2} = 0 \cdot 0 \cdot \frac{\partial h^{\alpha}}{\partial \Gamma} \cdot x^{3}$
	$\frac{\partial P_{(1)}}{\partial \Gamma} = \frac{\partial \Gamma}{\partial \Gamma}$
#16.	Exprimer sous forme matricielle, définir les dimensions, donner légaliselent en numpy
	$\frac{\partial P_{CD}}{\partial \Gamma} = \frac{\partial \Gamma}{\partial $
	$\frac{\partial L}{\partial W^{(1)}} = \frac{\partial L}{\partial h^{\alpha}} \cdot x^{T}  \text{est}  \partial_{h} x \partial  \text{car}  \frac{\partial L}{\partial h^{\alpha}} = \frac{\partial L}{\partial h^{\alpha}} \cdot x^{T} \cdot est  1 \times \partial u$
	grad-bl = grad-ha
	grad_W1 = np.outer (grad_ha, x)

 $\frac{\partial L}{\partial x_{\ell}} = \sum_{j=1}^{d_h} \frac{\partial L}{\partial h_j^{\alpha}} = \sum_{j=1}^{d_h} \frac{\partial L}{\partial L} \frac{\partial (b_j^{\alpha})^{\alpha} + \sum_{j=1}^{d_h} W_{j,l}^{\alpha} x_i)}{\partial h_j^{\alpha}}$ = Zj= ( Em 0 x W kj - Wy; 2) 1 Eng >03 (hg) W; e Alors,  $dL = \left(\sum_{j=1}^{d_h} \left(\sum_{k=1}^{m} o_k W_{kj}^{(2)} - W_{yj}^{(2)}\right) 1_{\xi h_j^2 > 0_3} (h_j^2) W_{j1}^{(1)}\right)$   $\left(\sum_{j=1}^{d_h} \left(\sum_{k=1}^{m} o_k W_{kj}^{(2)} - W_{yj}^{(2)}\right) 1_{\xi h_j^2 > 0_3} (h_j^2) W_{j3}^{(1)}\right)$  $= \left(\sum_{j=1}^{d_{A}} \frac{\partial L}{\partial h_{j}^{\alpha}} \cdot W_{j}^{(1)}\right)$   $= \left(\sum_{j=1}^{d_{A}} \frac{\partial L}{\partial h_{j}^{\alpha}} \cdot W_{j}^{(1)}\right)$ #18. Comment minimises \( \tilde{R} = \tilde{R} + \lambda\_{ii} \left( \Sinj \wight) \wight) + \lambda\_{i2} \left( \Sinj \wight) \varphi \rangle \lambda\_{ij} \left( \Wight)^2 \right) + \lambda\_{i2} \left( \Sinj \wight) \varphi \rangle \lambda\_{ij} \left( \Wight) \varphi^2 \right) \right) « de R change le gradient par sapport aux différents paramètres? Il ny a pas de différence pour b'ét b(2) con de(0) = de(0) = 0  $\frac{\partial \mathcal{L}(\theta)}{\partial \mathcal{M}^{(1)}} = \lambda_{11} \begin{bmatrix} sign(W_{11}^{(1)}) & sign(W_{12}^{(1)}) & ... & sign(W_{12}^{(1)}) \\ sign(W_{21}^{(1)}) & sign(W_{21}^{(1)}) & ... & sign(W_{21}^{(1)}) \end{bmatrix} + \lambda_{12} \begin{bmatrix} 2W_{11}^{(1)} & 2W_{12}^{(1)} & ... & 2W_{12}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ Sign(W_{21}^{(1)}) & Sign(W_{21}^{(1)}) & ... & Sign(W_{21}^{(1)}) \end{bmatrix} + \lambda_{12} \begin{bmatrix} 2W_{11}^{(1)} & ... & 2W_{12}^{(1)} \\ \vdots & \ddots & \vdots \\ 2W_{21}^{(1)} & ... & 2W_{21}^{(1)} \end{bmatrix}$ = ln sign (W.(1)) + 2 ln W(1) Alors, 2R = 2R + his sign(W(1)) + 2 hiz W(1)

et de la même façon, dR = dR + lz sign (W(2)) + 2 lz W(2)