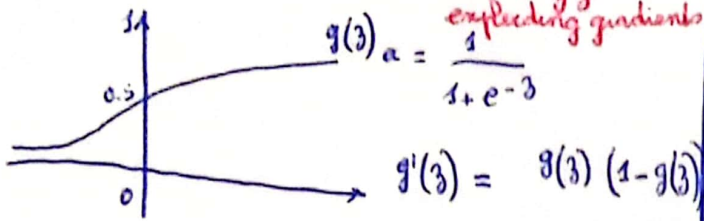


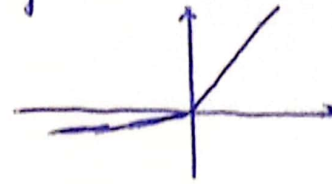
activation functions : fonctions de seuillage qui peuvent se décomposer en 3 parties [partie non active : au dessus du seuil
phase de transition : aux alentours du seuil
partie active : au dessous du seuil]

① sigmoid



les fonctions
vanishing gradients
exploding gradients

② ReLU partie active : au dessous du seuil

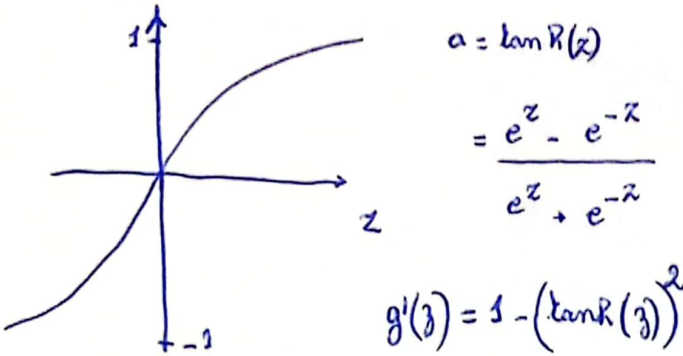


not very used in practice.

$$g(z) = \begin{cases} \max(0, z) & \end{cases}$$

$$g'(z) = \begin{cases} 0, & \text{if } z < 0 \\ 1, & \text{if } z > 0 \end{cases}$$

③



④ $ELU = \begin{cases} \alpha(e^x - 1) & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$

almost works better than sigmoid because it centers the data \Rightarrow makes learning easier

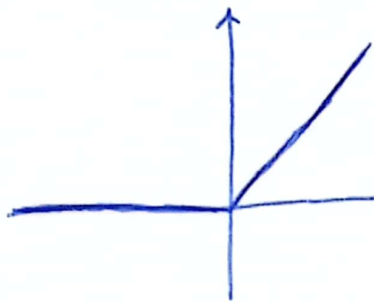
disadvantages of these 2 :

if $z \rightarrow +\infty$ ou $z \rightarrow -\infty$:

the slope $\rightarrow 0 \Rightarrow$ gradient descent is very slow

③ Rectified linear unit : ReLU

$$g(z) = \max(0, z)$$



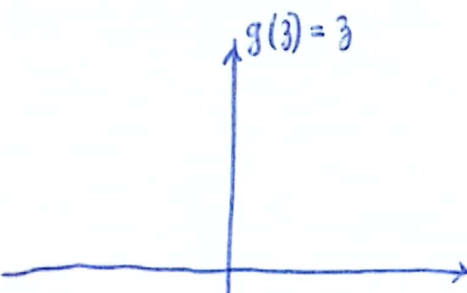
$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z > 0 \\ \text{non defined} & \text{in } 0 \end{cases}$$

in programming need

define g'(z) when

$$z = 0 = \begin{cases} 0 \\ 1 \end{cases}$$

④ Linear activation function



Fonctions d'activation usuelles

| sigmoid | Image | centré en 0 | saturation | évanescence | |
|---|----------------------|-------------|----------------|-----------------------------------|--------|
| sigmoid + cross entropy (log loss) for binary classification | $[0, 1]$ | Non | valeurs + et - | oui | Border |
| Tanh | $[-1, 1]$ | Oui | valeurs + et - | oui | Border |
| ReLU + MSE for regression (+ values) | $[0, +\infty[$ | Non | valeurs - | oui (moins que Border et tanh) | |
| Leaky ReLU | $[-\infty, +\infty[$ | non | non | non | facile |
| Linear + MSE for regression | | | | | |

softmax + categorical cross entropy loss for multiclass

choosing activation functions for output layer

binary classification : $y \in \{0, 1\}$

\Rightarrow sigmoid : used only in this case
(output of binary classification)

regression : $y = + / -$

\Rightarrow linear activation function

regression : y ne prend que des valeurs positives

\Rightarrow ReLU

for hidden layer :

most common choice : ReLU

faster computation \Rightarrow faster learning

(gradient doesn't need to be zero if we have a lot of flat proportions because $\frac{\partial}{\partial w} J(w, b) = 0$ when $g(z)$ is flat

! don't use linear activation in hidden layers

all linear + output linear

\Rightarrow linear regression

hidden layers linear + output sigmoid

\Rightarrow logistic regression without hidden layers

sigmoid

La dérivée de la fonction sigmoïde est à valeurs dans $[0, \frac{1}{4}]$

\Rightarrow diminue l'amplitude des gradients
propagés à travers les couches du réseau
de neurones

\Rightarrow pb d'extinction des gradients

ReLU

gradient égal à $\begin{cases} 1 & \text{dans la zone activée} \end{cases}$

\Rightarrow améliore la convergence

Mais

pas dérivable en 0

dérivée nulle sur \mathbb{R}^-

à valeurs dans \mathbb{R}^+

\Rightarrow optimisation plus compliquée si
beaucoup de couches activées

leaky ReLU

moins d'occurrences où le gradient
est nul

\Rightarrow meilleure convergence de la descente
de gradient