



## INITIALIZATION

Vincent Barra  
LIMOS, UMR 6158 CNRS, Université Clermont Auvergne

## What we know

- ▶ Initialization should break symmetry
- ▶ the relative scale of weights is fundamental.
- ▶ and so is initialization.

## First bad idea

$$\forall i \ w_i = 0$$

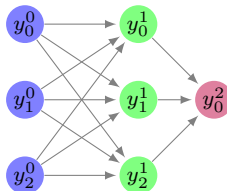
1 hidden layer MLP. ( $\forall l \in \{1, 2\}\})(\forall i \in \{0, 2\})$

$$y_i^l = \sigma \left( \sum_{j=0}^2 w_{ij}^l y_j^{l-1} \right)$$

Using backpropagation on error  $E$

$$\frac{\partial E}{\partial w_{0j}^1} = \frac{\partial E}{\partial w_{1j}^1} = \frac{\partial E}{\partial w_{2j}^1}$$

Hidden neurons learn the same parameters  
 $\Rightarrow$  redundancy.



## Second bad idea

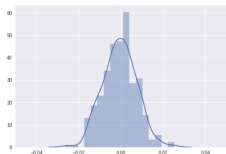
$\forall i \ w_i = \delta$  : Same result!

## Second bad idea

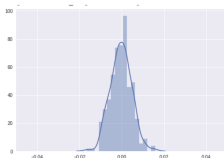
$\forall i \ w_i = \delta$  : Same result!

## A better idea ?

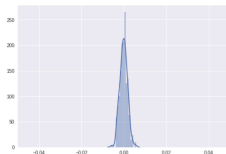
$$\forall i \ w_i \rightsquigarrow 10^{-m} \mathcal{N}(0, 1) \ m > 0$$



After 10 epochs  
 $\Rightarrow$  vanishing gradient (See dedicated slides).



After 20 epochs



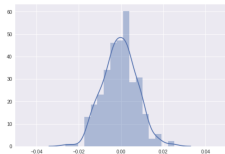
After 50 epochs

## Second bad idea

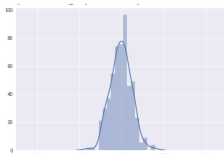
$\forall i \ w_i = \delta$  : Same result!

## A better idea ?

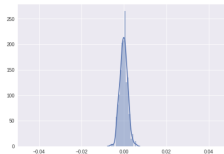
$$\forall i \ w_i \rightsquigarrow 10^{-m} \mathcal{N}(0, 1) \ m > 0$$



After 10 epochs  
 $\Rightarrow$  vanishing gradient (See dedicated slides).



After 20 epochs



After 50 epochs

## A second better idea ?

$$\forall i \ w_i \rightsquigarrow \text{high random values}$$

The scalar products of the post synaptic potentials saturate the activation functions  $\sigma$ . The derivatives tend towards 0.

$\Rightarrow$  vanishing gradient (See dedicated slides).

## Xavier/Glorot initialization

MLP  $\mathbf{w} \in \mathbb{R}^d$  i.i.d  $\rightsquigarrow \mathcal{N}(0, 1)$  with input  $\mathbf{x} \in \mathbb{R}^d$ ,  $x_i$  i.i.d  $\rightsquigarrow \mathcal{N}(0, 1)$ .  
Post synaptic potential of a neuron in the first hidden layer:  $\mathbf{w}^T \mathbf{x}$ .

$$\begin{aligned} \text{Var}(\mathbf{w}^T \mathbf{x}) &= \sum_{i=1}^d \text{Var}(w_i x_i) \\ &= \sum_{i=1}^d (\mathbb{E}(w_i)^2 \text{Var}(x_i) + \mathbb{E}(x_i)^2 \text{Var}(w_i) + \text{Var}(w_i) \text{Var}(x_i)) \end{aligned}$$

Since  $\mathbb{E}(w_i) = \mathbb{E}(x_i) = 0$

$$\text{Var}(\mathbf{w}^T \mathbf{x}) = \sum_{i=1}^d \text{Var}(w_i) \text{Var}(x_i)$$

$w_i, x_i$  i.i.d  $\Rightarrow \text{Var}(\mathbf{w}^T \mathbf{x}) = d \text{Var}(w_i) \text{Var}(x_i)$

More generally

$$\text{Var}(y^l) = (d \text{Var}(w_i))^l \text{Var}(x_i)$$

$\Rightarrow$  Each neuron can vary about  $d$  times the variation of its input.

## Xavier/Glorot initialization

- ▶ if  $dVar(w_i) > 1$  the gradient increases when going deeper in the network
  - ▶ if  $dVar(w_i) < 1$  vanishing gradient as  $l$  increases
- $\Rightarrow$  impose  $dVar(w_i) = 1$ , thus  $Var(w_i) = 1/d$  and

$$w_{ij}^l \rightsquigarrow \frac{1}{\sqrt{m^{l-1}}} \mathcal{N}(0, 1)$$

where  $m^{l-1}$  : number of neurons of layer  $l - 1$ .

If  $\sigma = ReLU$  one can multiply by  $\frac{\sqrt{2}}{\sqrt{m^{l-1}}}$  to take into account the negative part that does not participate to the variance.