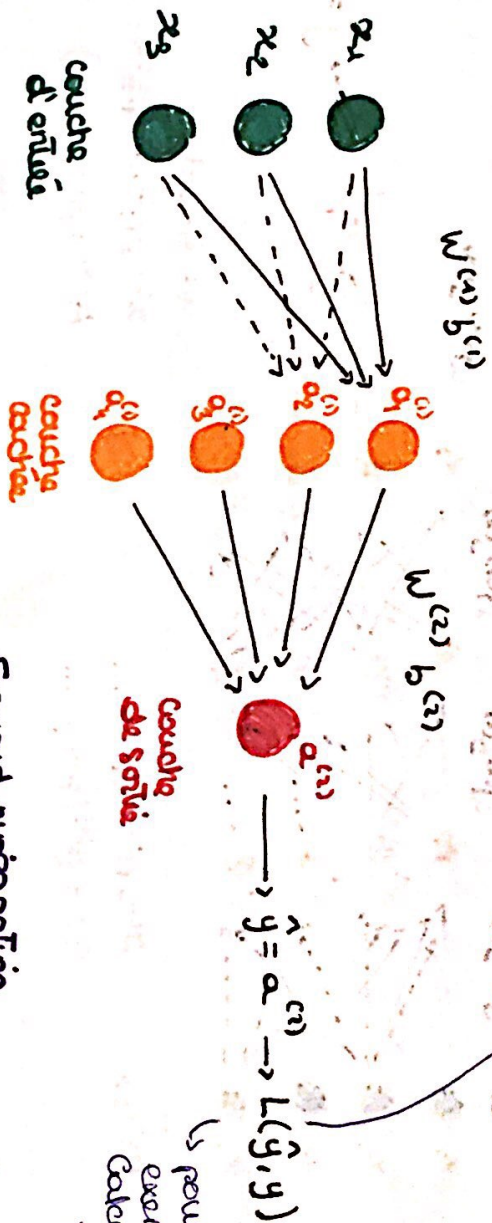


# Réseaux de neurones (2 couches, 1 couche cachée)



**Classification binaire**  
 $y \sim \mathcal{B}(p)$  on a  $\begin{cases} P(y=1)=p \\ P(y=0)=1-p \end{cases}$   
 on cherche  $\theta$  qui maximise le  
 log likelihood max  $\log R(y|x)$   
 $\Leftrightarrow$  minimiser la BCE  
 $\Leftrightarrow$  minimiser le coût global.

pour tous les  
exemples on  
calculer le coût  
global.

## Forward propagation

$$A^{(0)} = X \quad \begin{matrix} z^{(1)} = A^{(0)} w^{(1)} + b^{(1)} & z^{(2)} = A^{(1)} w^{(2)} + b^{(2)} & \hat{y} = A^{(2)} \\ A^{(1)} = g^{(1)}(z^{(1)}) & A^{(2)} = g^{(2)}(z^{(2)}) & \end{matrix} \quad \text{Pour tous les exemples.}$$

## Backward propagation

$$\begin{aligned} dz^{(2)} &= f'(z^{(2)}) - y \\ dw^{(2)} &= \frac{1}{n} A^{(2)T} dz^{(2)} \\ db^{(2)} &= \frac{1}{n} \sum_{i=1}^n dz^{(2)} \\ \left[ dA^{(1)} \right] &= dz^{(2)} w^{(2)T} \\ dw^{(1)} &= \frac{1}{n} A^{(1)T} dz^{(1)} \\ db^{(1)} &= \frac{1}{n} \sum_{i=1}^n dz^{(1)} \end{aligned}$$

**Update**

$$\begin{aligned} w^{(1)} &= w^{(1)} - \alpha dw^{(1)} & w^{(2)} &= w^{(2)} - \alpha dw^{(2)} \\ b^{(1)} &= b^{(1)} - \alpha db^{(1)} & b^{(2)} &= b^{(2)} - \alpha db^{(2)} \end{aligned}$$

# Différents types d'entraînement

batch GD

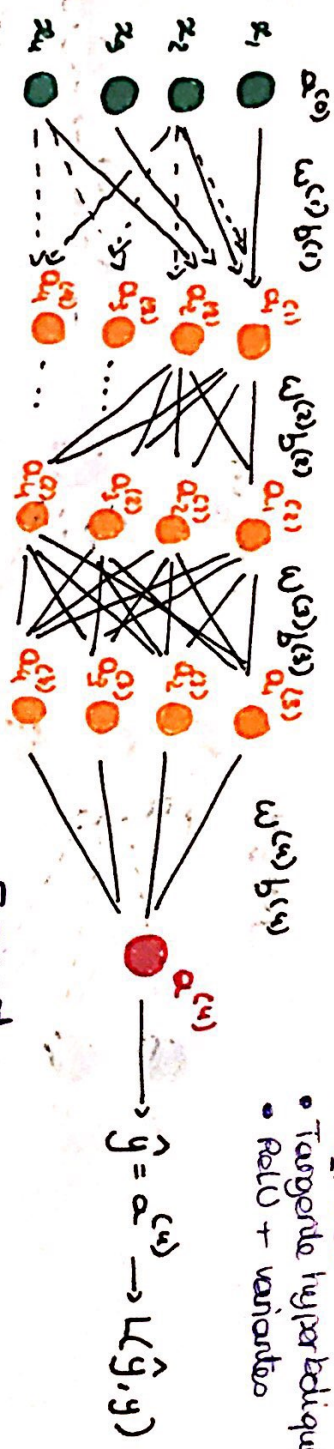
mini batch GD

Stochastic GD

Différents types de descentes de gradient

- mini batch GD
- Momentum
- NAG
- Adagrad
- AdaDelta
- Adam
- learning rate

## Réseaux de neurones profonds (plus de 2 couches cachées)



Forward

$$z^{(1)} = A^{(0)} W^{(0)} + b^{(1)}$$

$$A^{(1)} = g^{(1)}(z^{(1)})$$

$$z^{(4)} = A^{(3)} W^{(3)} + b^{(4)}$$

$$A^{(4)} = g^{(4)}(z^{(4)})$$

Pour chaque couche, on met en cache certains éléments pour la backpropagation.

Backward

$$dw^{(1)}$$

$$db^{(1)}$$

$$dz^{(3)} = dA^{(3)} b g^{(3)'}(z^{(3)})$$

$$dA^{(3)} = \dots$$

$$dz^{(4)} = dA^{(3)} W^{(3)T}$$

$$dA^{(4)} = \dots$$

la dérivée de la fonction de coût comprise entre 0 et 1/2 donc plus le réseau est profond, plus db de du est petit or disparaît.

- Fonctions d'activation  $a = g(z)$  Sortie
- Sigmoïde (binaire)
  - Softmax (multiclasse)
- Fonctions d'activation cachées
- non linéaire
- Sigmoïde
  - Tangente hyperbolique
  - ReLU + variantes
- Vanishing gradient

### Initialisation des poids

He init = random(size E, size E-1)  $\times \sqrt{\frac{2}{size E-1}}$

Xavier =  $\times \sqrt{\frac{1}{size E-1}}$

Yengio =  $\times \sqrt{\frac{2}{size E-1 \times size E}}$

### Update

$$w^{(1)} = w^{(1)} - \alpha dw^{(1)}$$

$$b^{(1)} = b^{(1)} - \alpha db^{(1)}$$

Normalisation = on peut normaliser les entrées par la min et max pour éviter des valeurs similaires de poids, éviter l'explosion des activations et l'instabilité numérique.

### Dropout

on cache aléatoirement certains neurones ce qui évite des coconnations.

Regulation des réseaux réduire l'overfitting.