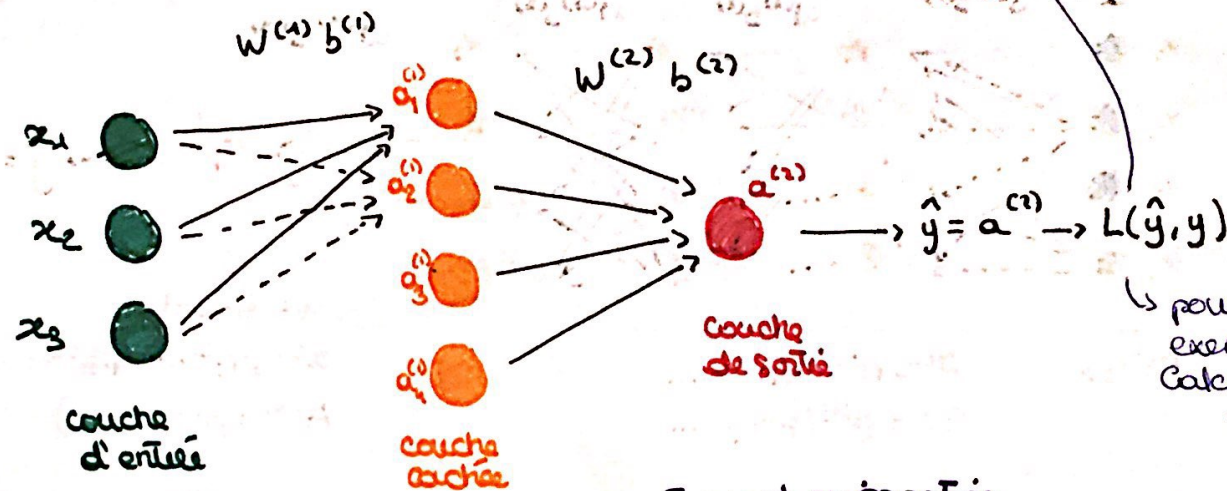


Réseau de neurones (2 couches, 1 couche cachée)



Classification binaire

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

pour tous les
exemples on
calculer le coût
global.

forward propagation

$$A^{(0)} = X$$

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

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

$$z^{(2)} = A^{(1)} W^{(2)} + b^{(2)}$$

$$A^{(2)} = g^{(2)}(z^{(2)})$$

Pour tous les
exemples.

Backward propagation

$$dz^{(1)} = dA^{(1)} \odot g^{(1)'}(z^{(1)})$$

$$dW^{(1)} = \frac{1}{m} A^{(0)T} dz^{(1)}$$

$$db^{(1)} = \frac{1}{m} \sum_{i=1}^m dz^{(1)}$$

$$dz^{(2)} = A^{(2)} - y$$

$$dW^{(2)} = \frac{1}{m} A^{(1)T} dz^{(2)}$$

$$db^{(2)} = \frac{1}{m} \sum_{i=1}^m dz^{(2)}$$

$$dA^{(1)} = dz^{(2)} W^{(2)T}$$

update

$$W^{(1)} = W^{(1)} - \alpha dW^{(1)}$$

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

$$W^{(2)} = W^{(2)} - \alpha dW^{(2)}$$

$$b^{(2)} = b^{(2)} - \alpha db^{(2)}$$

Réseaux de neurones profonds (plus de 2 couches cachées)

Fonctions d'activation $a = g(z)$ Sortie.

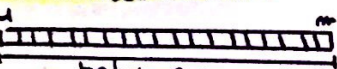
- Sigmoide (binaire)
- Softmax (multiclasse)

Fonctions d'activation couches.

=> non linéaire

- Sigmoide
- Tangente hyperbolique } Vanishing gradient.
- ReLU + variantes

Différents types d'entraînement.

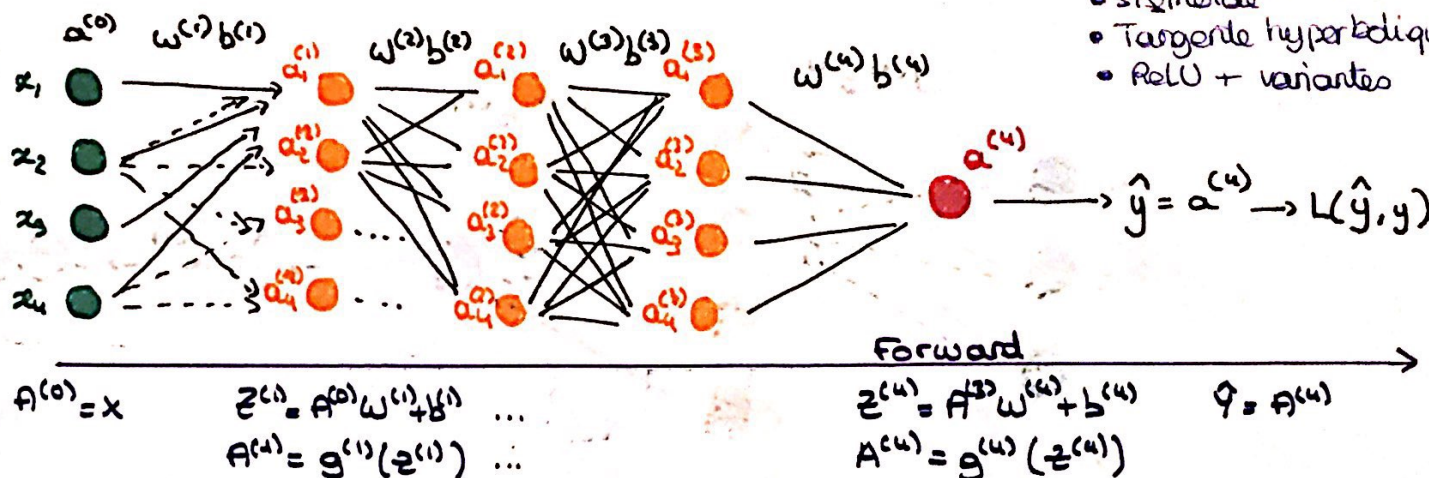


mini batch GD

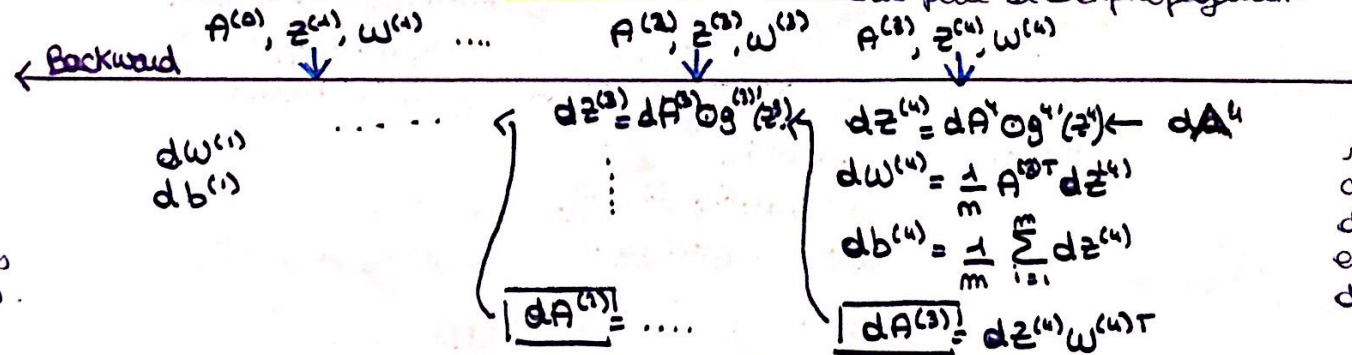
Stochastic GD

Différents types de descentes de grad.

- mini batch GD
- Momentum
- NAG
- AdaGrad
- AdaDelta
- Adam
- learning rate.



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



la backpropagation est un algorithme efficace qui permet de calculer les dérivées en chaîne des composées en stockant les résultats intermédiaires.

la dérivée de la sigmoïde est comprise entre 0 et 1/4 donc plus le réseau est profond, plus db ou dW est petit et disparaît.

Update $W^{(i)} = W^{(i)} - \alpha dW^{(i)}$
 $b^{(i)} = b^{(i)} - \alpha db^{(i)}$

Dropout = on retire aléatoirement certains neurones ce qui retire des connexions.
=> joue le rôle de régularisation

Régularisation des réseaux réduit l'overfitting.

Initialisation des poids.

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

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

Bengio = $\times \sqrt{\frac{2}{\text{size l-1} \times \text{size l}}}$

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