



Reconnaissance des formes pour l'analyse et l'interprétation d'images

Rapport TP 3-4 : Introduction aux réseaux de neurones

Étudiant :

DJEGHRI Amine

MAMOU Idles

Numéro :

3801757

3803676

Novembre 2020

Partie 1 – Formalisation mathématique

1.1 Jeu de données

1. L'ensemble d'apprentissage sert à entraîner le modèle
L'ensemble de validation sert à choisir les meilleurs hyper-paramètres de notre modèle
L'ensemble de test sert à évaluer notre modèle et voir ses performances si on l'applique sur un jeu de donnée inconnu
2. L'influence du nombre N d'exemples :
Plus N le nombre d'exemple est grand, plus le modèle sera mieux entraîné et par conséquent le modèle pourra mieux de généraliser.

1.2 Architecture du réseau (phase forward)

3. Les fonctions d'activation permettent de transformer nos fonctions linéaires en non linéaire et ainsi rendre notre modèle complexe et pouvoir traiter des problèmes complexes (par exemple des données non linéairement séparable) qui ne peuvent être bien résolu avec des fonctions linéaires.
4. $n_x = 2$ (la taille n'est pas choisie car elle dépend de la dimension des entrées)
 $n_h = 4$ (C'est le nombre de neurones de la couche cachée, c'est un hyper-paramètre,)
 $n_y = 2$ (la taille n'est pas choisie, elle dépend du nombre de classes)
5. y est la valeur réelle, le label réel
 \hat{y} est la valeur prédite par notre modèle
6. On utilise une fonction SoftMax en sortie, pour avoir une distribution de probabilités des différentes classes (la fonction SoftMax est continue et dérivable), la classe dont la probabilité est la plus élevée est choisie comme la classe la plus probable
7. Les équations mathématiques permettant d'effectuer la passe forward :

$$\begin{aligned}\tilde{h} &= W_h * X + B_h \\ h &= \text{tanh}(\tilde{h}) \\ \tilde{y} &= W_y * h + B_y \\ \hat{y} &= \text{SoftMax}(\tilde{y})\end{aligned}$$

1.3 Fonction de coût

8. Les \hat{y} doivent être proche de y pour faire diminuer la loss
9. Les deux fonctions sont convexes
Le cout MSE pour la regression car elle mesure l'écart de distance
L'entropie croisée(cross-entropy) pour la classification car elle a pour but de mesurer la différence de distribution de probabilité

1.4 Méthode d'apprentissage

10. Les avantages et inconvénients des diverses variantes de descente de gradient entre les versions :

- Classique : avantage : la descente est stable
Inconvénient : met trop de temps à converger
- Stochastique sur mini-batch : Avantage : temps de calcul efficace
Inconvénient : Convergence moins stable (un peu de variance)
- Stochastique online : Avantage : Rapide dans le temps de calcul
Inconvénient : pas stable, trop de variance

11. Si le learning rate est trop petit, le modèle va prendre beaucoup de temps à converger

Si le learning est grand, le modèle risque de ne pas converger à cause des grands sauts du gradient d'une direction à une autre.

12. L'algorithme de Backprop est moins coûteux que l'approche naïve, car le gradient calculé par rapport à la sortie d'une couche peut être réutilisé pour calculer les gradients par rapport à l'entrée et aux paramètres de cette même couche et donc on aura une complexité égale au nombre de couches du réseau, par contre pour l'approche naïve aura une complexité bien plus élevée.

13. Il faut que les fonctions d'activation et les couches soient dérivables

14.

$$\begin{aligned}\hat{y} &= \frac{e^{\tilde{y}_i}}{\sum_j e^{\tilde{y}_j}} \\ \ell(y, \hat{y}) &= -\sum_i y_i \log \hat{y}_i \\ &= -\sum_i y_i \log \left(\frac{e^{\tilde{y}_i}}{\sum_j e^{\tilde{y}_j}} \right) \\ &= -\sum_i y_i (\log e^{\tilde{y}_i} - \log \sum_j e^{\tilde{y}_j}) \\ &= -\sum_i y_i \tilde{y}_i + \log \left(\sum_j e^{\tilde{y}_j} \right)\end{aligned}$$

15.

$$\begin{aligned}\frac{\partial \ell}{\partial \tilde{y}_i} &= -y_i + \frac{e^{\tilde{y}_i}}{\sum_j e^{\tilde{y}_j}} \\ &= \hat{y}_i - y_i\end{aligned}$$

16.

$$\begin{aligned}\frac{\partial \ell}{\partial w_{y,i,j}} &= \sum_k \frac{\partial \ell}{\partial \tilde{y}_k} \cdot \frac{\partial \tilde{y}_k}{\partial w_{y,i,j}} = \sum_k (-y_k + \hat{y}_k) h_k \\ * \quad \frac{\partial \ell}{\partial b_y} &= \sum_k \frac{\partial \ell}{\partial \tilde{y}_k} \cdot \frac{\partial \tilde{y}_k}{\partial b_y} \\ &= \sum_k (\hat{y}_k - y_k)\end{aligned}$$

17.

$$\begin{aligned}17) \quad 1) \quad \frac{\partial \ell}{\partial \tilde{h}_i} &= \frac{\partial \ell}{\partial h_i} \cdot \frac{\partial h_i}{\partial \tilde{h}_i} \\ &= (1 - h_i^2) w_{y,i} \left(-y_i + \frac{e^{\tilde{y}_i}}{\sum_j e^{\tilde{y}_j}} \right) \\ 2) \quad \frac{\partial \ell}{\partial w_{h,i}} &= \frac{\partial \ell}{\partial \tilde{h}_i} \cdot \frac{\partial \tilde{h}_i}{\partial w_{h,i}} \\ &= w_{y,i} \left(-y_i + \frac{e^{\tilde{y}_i}}{\sum_j e^{\tilde{y}_j}} \right) \underbrace{(1 - \tanh(\tilde{h}_i))}_{1 - h_i^2} x_i \\ 3) \quad \frac{\partial \ell}{\partial b_h} &= \frac{\partial \ell}{\partial \tilde{h}_i} \cdot \frac{\partial \tilde{h}_i}{\partial b_h} \\ &= w_{y,i} \left(-y_i + \frac{e^{\tilde{y}_i}}{\sum_j e^{\tilde{y}_j}} \right) \underbrace{(1 - \tanh(\tilde{h}_i))}_{1 - h_i^2}\end{aligned}$$

