

Réseaux de neurones simples et Deep learning

Soit une image de (28×28) pixels. On représente cette image comme un élément de $[0, 1]^{28 \times 28}$, soit un 784 uplet de valeurs dans $[0, 1]$.

L'image représente au plus un élément de l'ensemble des caractères hexadécimaux.

$$\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F\} = \mathcal{C}$$

on dit que ces 16 éléments sont les classes du problème.

On suppose qu'il existe une application $f: [0, 1]^{28 \times 28} \rightarrow \mathcal{C}$

un réseau neuronal est en fait une approximation de f .

• Fonctionnement d'un réseau de neurone standard: Forward pass

Un réseau de neurone est une composition de fonctions linéaires et non linéaires. Un neurone standard reçoit un set d'échantillons S chacun contenant F caractéristiques.

un ensemble de neurones constituent une couche. un neurone recevant une matrice $(S \times F)$ nous renseigne sur les couches avant: F neurones.

$$S = \begin{bmatrix} x_{s1} & \dots & x_{sF} \\ \vdots & & \vdots \\ x_{m1} & \dots & x_{mF} \end{bmatrix}$$

$m1 \quad \quad mF$

Pour chaque neurone, on associe une matrice de poids W , et un vecteur de bias B . Il est nécessaire que W soit de dimensions $(F \times 1)$ et B soit de dimension $(S \times 1)$.

le neurone effectue les opérations matricielles suivantes:

$XW + B = Z$. une fois Z calculé, on applique une fonction non linéaire dite d'activation à Z , élément par élément. Z est de dimension $(S \times 1)$. Par exemple: tanh, sigmoid $(\frac{1}{1+e^{-x}})$, ReLU ... etc. Il est nécessaire (ou plutôt fortement conseillé) que la fonction d'activation soit monotone et injective (on autorise toutefois 0 en valeur multiple). S'il y'a m neurones dans un niveau, les neurones de la couche suivantes reçoivent des matrices de dimension $(s \times m)$

$$S \begin{bmatrix} x_{s1} & x_{s2} & \dots & x_{sF} \\ \vdots & \vdots & & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mF} \end{bmatrix} \times F \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_F \end{bmatrix} + S \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_S \end{bmatrix} = S \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_S \end{bmatrix} \Rightarrow \text{Activation} \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_S \end{bmatrix}$$
$$\Rightarrow \begin{bmatrix} \text{Activation}(z_1) \\ \vdots \\ \text{Activation}(z_S) \end{bmatrix} \quad (S \times 1)$$

Soit un neurone de la couche suivante. et soit m neurones dans la couche qui le précède. le neurone reçoit une matrice de dimension $(S \times m)$: c'est la concaténation des vecteurs colonnes renvoyés par chaque neurone de la couche antérieure.

Le forward pass se termine par une couche contenant 16 neurones; Soit la cardinalité de nos classes. Cette couche renvoie une matrice de dimension $(S \times 16)$

$S \begin{bmatrix} C_{00} & C_{01} & C_{02} & \dots & C_{015} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ C_{S0} & C_{S1} & C_{S2} & \dots & C_{S15} \end{bmatrix}$ chaque C_{ij} est une probabilité que l'échantillon i soit de classe $j \in C$.
On suppose que pour chaque échantillon, la vraie classe soit connue, on peut alors regrouper ces classes en une matrice Y de même dimension

$P = \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 1 & \dots & 0 & 0 \end{bmatrix}$ si P est la matrice des prédictions, et Y la matrice d'observation, on calcule le taux d'erreur, une valeur numérique donnée par une fonction coût opérant ligne à ligne. on obtient un vecteur d'erreurs: chaque erreur correspond à un échantillon
($S \times 16$)

Fonction coût 1: MSE: sur L_i : $\sum_{j=0}^{15} (Y_{ij} - P_{ij})^2$ (borne sup incluse)

Fonction coût 2: CE : sur L_i : $\sum_{j=0}^{15} -Y_{ij} \times \ln(P_{ij}) - (1 - Y_{ij}) \times \ln(1 - P_{ij})$

la fonction CE (cross entropy loss) est plus efficace sur des lignes ayant été altérées par la fonction softmax: $\text{softmax}(L_i) = \frac{e^{L_{ij}}}{\sum_l e^{L_{il}}}$

une fois le vecteur d'erreur calculé, de dimension $(S \times 1)$, le forward pass est terminé. (CES = CE(softmax))

• backward pass

Rappel mathématique: le gradient d'une fonction $g(x_1, x_2, \dots, x_m)$ en x

est un vecteur des dérivées partielles de g en x : $\frac{\partial g(x)}{\partial x} = \left[\frac{\partial g(x)}{\partial x_1}, \frac{\partial g(x)}{\partial x_2}, \dots, \frac{\partial g(x)}{\partial x_m} \right]$

• la règle de la chaîne: $\frac{d(g(h(u)))}{du} = \frac{d(g(h(u)))}{dh(u)} \times \frac{dh(u)}{du}$

On a vu dit qu'un réseau neuronal est une approximation de f , la fonction de $[0, 1]^{28 \times 28}$ dans C . le taux d'erreur calculé par CES est une mesure de cette approximation. De plus cette approximation n'est autre qu'une composition de fonctions, avec chaque neurone étant: Activation($XW+B$)

où X est un set de vecteurs lignes contenant les résultats des neurones de la couche antérieure. On discute à présent l'algorithme de backpropagation, qui permet d'ajuster tous les W et B de chaque neurone de réseau neuronal.

Notons déjà que la fonction composée CES = CE(softmax) est une application de $\{0, 1\}^{16} \times \{0, 1\}^{16}$ dans \mathbb{R} opérant sur les lignes de Y et P resp.

on note $X = (x_1, x_2, \dots, x_{16}) \in [0, 1]^{16}$ les éléments de la ligne i de P
 $Z = Y$.

On calcule les dérivées partielles de CES par rapport à chaque $x_i \in X$

on trouve: $\frac{\partial CE(Z, \text{softmax}(X))}{\partial x_1} = \text{softmax}(x_1) - Z_1 = \frac{e^{x_1}}{e^{x_1} + \dots + e^{x_6}} - Z_1$

donc, pour la matrice P , le gradient de la ligne i est simplement $\text{softmax}(P_i) - Y_i$, et de manière plus générale, pour de multiples échantillons

$$\begin{bmatrix} \left[\frac{e^{c_{11}}}{\sum e^c} - P_{11}, \frac{e^{c_{12}}}{\sum e^c} - P_{12}, \dots, \frac{e^{c_{16}}}{\sum e^c} - P_{16} \right] \\ \vdots \\ \left[\frac{e^{c_{51}}}{\sum e^c} - P_{51}, \frac{e^{c_{52}}}{\sum e^c} - P_{52}, \dots, \frac{e^{c_{56}}}{\sum e^c} - P_{56} \right] \end{bmatrix}$$

on nomme cette matrice la matrice des gradients des échantillons

Pour chaque neurone de la dernière couche, on ajuste les poids W et B

Par abstraction: le neurone N reçoit le vecteur colonne:

$$\begin{bmatrix} \frac{e^{c_{11}}}{\sum e^c} - P_{11} \\ \vdots \\ \frac{e^{c_{51}}}{\sum e^c} - P_{51} \end{bmatrix} \text{ qu'on renomme } G_l \text{ (gradient local)} \quad \begin{bmatrix} G_{l1} \\ \vdots \\ G_{ls} \end{bmatrix} = G_l$$

on effectue la somme des G_{li} , ce qui donne un gradient global du neurone en question. on nomme cette quantité

$\frac{\partial E}{\partial A(Z)}$; graphe de computation:

l'algorithme de back propagation se déroule ainsi:

$$\left. \begin{array}{l} \frac{\partial E}{\partial A(Z)} \text{ donnée} \\ \frac{\partial A(Z)}{\partial Z} = A'(Z) \\ \frac{\partial Z}{\partial B} = 1 \\ \frac{\partial Z}{\partial xw} = 1 \\ \frac{\partial xw}{\partial x} = w \\ \frac{\partial xw}{\partial w} = x \end{array} \right\} \begin{array}{l} \frac{\partial E}{\partial w} = \frac{\partial E}{\partial A(Z)} \times \frac{\partial A(Z)}{\partial Z} \times \frac{\partial Z}{\partial xw} \times \frac{\partial xw}{\partial w} =: g_w \\ \text{ajustement } w = w - \alpha g_w \quad \alpha = \text{learning rate} \\ \frac{\partial E}{\partial B} = \frac{\partial E}{\partial A(Z)} \times \frac{\partial A(Z)}{\partial Z} \times \frac{\partial Z}{\partial B} =: g_b \\ B = B - \alpha g_b \quad (\text{ajustement}) \\ \text{à renvoyer à la couche antérieure} \\ \frac{\partial E}{\partial x} = \frac{\partial E}{\partial A(Z)} \times \frac{\partial A(Z)}{\partial Z} \times \frac{\partial Z}{\partial xw} \times \frac{\partial xw}{\partial x} \end{array}$$

on peut, étant donné une couche entière, exprimer ces opérations en termes d'opérations bien plus élémentaires, Transposition des vecteurs w et x .