

# Résultats concernant l'arithmétique d'intervalle appliquée aux réseaux de neurones

Capitaine de corvette Guillaume Berthelot

Aout 2024

On considère l'arithmétique d'intervalle appliquée aux réseaux de neurones.

**Théorème 1.** *Soit un réseau de neurones quelconque. Soit  $A_\epsilon$  un tenseur de symboles de bruits pour l'arithmétique d'intervalle. Soit  $f$  la fonction telle que l'image du tenseur à travers le réseau est donnée par  $f(a_{\epsilon_i})$  pour chacun des éléments de  $A_\epsilon$ .*

*S'il existe une base de dimension  $n$  sur laquelle la projection de  $A_\epsilon$  est injective et s'il existe une matrice  $W$  représentant  $f$  sur cette base, c'est-à-dire pour tout  $i$  appartenant à  $\{1, \dots, n\}$ ,  $f(a_{\epsilon_i}) = W a_{\epsilon_i} e_i$ ,*

*Alors,*

$$Z^+ = |W| |A_\epsilon|$$

*est un vecteur bornant pour l'arithmétique d'intervalle. Si  $C$  est le vecteur centre, alors les bornes sont données par  $C \pm Z^+$ .*

*Proof.* Soit  $p$  le nombre d'éléments de  $A_\epsilon$ , on a

$$Z^+ = \sum_{i=1}^p |f(a_{\epsilon_i})| = \sum_{i=1}^p |W a_{\epsilon_i} e_i| = |W| \sum_{i=1}^p |a_{\epsilon_i} e_i|$$

d'où le résultat. □

**Lemme 1.** *Il existe une base canonique dans laquelle les opérations dites linéaires des réseaux de neurones sont représentables par une matrice. Il s'agit de la base générée par l'opération d'aplatissement.*

**Lemme 2.** *Si  $p$  est le tenseur des coefficients d'approximation d'une fonction d'activation, alors l'image de  $p$  dans la base canonique est l'aplatissement de  $p$ .*

**Corollaire 1.** *Si  $L_1, R_1, L_2, R_2, \dots, L_n, R_n$  est une succession de couches linéaires et d'activations, et si  $A_\epsilon$  est un tenseur de symboles de bruits pour l'arithmétique d'intervalle se projetant injectivement sur la base canonique de  $L_1$ ,*

*Si  $p_1, p_2, \dots, p_n$  sont les tenseurs d'approximation projetés dans leurs bases canoniques respectives, alors l'image  $A_s$  de  $A_\epsilon$  pour l'arithmétique d'intervalle est donnée par*

$$Z_s = |W_r| |A_\epsilon| = |((p_n W_n) \otimes (p_{n-1} W_{n-1}) \otimes \dots \otimes (p_1 W_1))| |A_\epsilon|$$

Autrement dit, il est possible de réduire à un produit matriciel le calcul de l'arithmétique d'intervalle appliquée à l'ensemble du réseau.

*Proof.* Soit  $L_n, R_n$  les deux dernières couches linéaire et d'activation,

Soit  $A_{n-1}$  le tenseur de symboles d'entrée,  $A_s$  l'image du tenseur d'entrée par l'approximation linéaire du doublet  $L_n, R_n$ . Alors

$$A_s = f_n(p_n A_{n-1}) = p_n W_n A_{n-1}$$

puis récursivement

$$\begin{aligned} &= p_n W_n p_{n-1} W_{n-1} A_{n-2} = \\ &\dots = ((p_n W_n) \otimes (p_{n-1} W_{n-1}) \otimes \dots \otimes (p_1 W_1)) A_\epsilon \end{aligned}$$

Or  $A_\epsilon$  se projecte injectivement sur son espace canonique.

D'où le résultat. □

**Théorème 2.** Soit  $d$  le tenseur de bruit généré par l'approximation d'une couche d'activation. Alors  $d$  se projecte injectivement sur la base canonique.

*Proof.* L'opération d'approximation est définie. □

En appliquant les résultats précédents, il est possible de réduire l'évaluation de l'arithmétique d'intervalle à un produit matriciel.

**Théorème 3.** Soit  $L_1, R_1, L_2, R_2, \dots, L_n, R_n$  une succession de couches linéaires et d'activations, soit  $j$  une couche intermédiaire  $L_j$ . Alors les bornes pour l'arithmétique d'intervalle pour la couche d'activation  $R_j$  sont données par

$$C_j \pm \sum_{l=1}^j |W_{r_l}| |A_{d_l}|$$

où

$$W_{r_i} = (W_j) \otimes (p_{j-1} W_{j-1}) \otimes \dots \otimes (p_i W_i) \quad \forall i \in \{1, \dots, j\}$$

*Proof.* Ce résultat est immédiat. □

**Corollaire 2.** Soit un réseau constitué de couches linéaires et de fonctions d'activations. L'algorithme 1 est un algorithme d'approximation affine pour ce réseau en temps polynomial

---

**Algorithm 1** Affine Approximation Algorithm for the Network

---

```
1: Step 1:
2: for each linear layer do
3:   Calculate  $W_l$ 
4:   Create an empty list  $L_W$ 
5: end for
6: Step 2:
7: Choose an input  $x$ 
8: for each dimension of  $x$  do
9:   Establish a noise level  $\delta$ 
10: end for
11: Create the vector  $A_\delta$ 
12: Create a list  $A$  with the first element  $A_\delta$ 
13: Initialize a unit approximation vector  $p$ 
14: for each layer of the network, in increasing order do
15:   if linear layer then
16:     Calculate the center
17:     for each element of the list  $L_W$  do
18:       Multiply it by  $p \times W_l$  on left side
19:     end for
20:     Add  $W_l$  to the list  $L_W$ 
21:     Create a copy  $|L_W|$ 
22:     for each pair of elements  $(|W_L|, |A|)$  do
23:       Stack the sum of the products: result  $Z^+$ 
24:     end for
25:     Overwrite  $p$  with a unit vector of the output dimension
26:   else if activation layer then
27:     Calculate the bounds and store the result
28:     Define  $p, q, d$ 
29:     Shift the center
30:     Add  $A_d$  to the list  $A$ 
31:   end if
32: end for
```

---