

Une approche tensorielle pour entraîner des réseaux de neurones flexibles*

Yassine ZNIYED¹, Konstantin USEVICH², Sebastian MIRON², David BRIE²

¹Université de Toulon, Aix Marseille Université, CNRS, LIS, Toulon, France

²Université de Lorraine, CNRS, CRAN, F-54000 Nancy, France

yassine.zniyed@univ-tln.fr, {konstantin.usevich, sebastian.miron, david.brie}@univ-lorraine.fr

Résumé – Dans ce travail, nous nous concentrons sur la compression des réseaux de neurones (RNs), en considérant une approche basée sur les tenseurs. En particulier, nous nous intéressons à l’estimation de fonctions d’activation flexibles à l’aide de solutions basées sur des tenseurs. La méthode proposée est basée sur une factorisation couplée matrice-tenseur (FCMT), fusionnant les informations d’ordre 0 et d’ordre 1 des RNs. L’algorithme d’apprentissage proposé est basé sur une approche contrainte des moindres carrés alternés. Les fonctions d’activation (FAs) sont exprimées sous la forme d’une somme pondérée d’éléments de base prédéfinis. Notre méthode permet de compresser fortement des grandes couches des RNs, avec une très légère dégradation de la précision de la classification. Cela donne également un aperçu du fonctionnement des RNs, qui est un problème du plus haut intérêt dans ce domaine.

Abstract – In this work, we focus on the compression of Neural Networks (NNs), considering a tensor-based approach. In particular, we are interested in the estimation of flexible activation functions using tensor-based solutions. The proposed method is based on a coupled matrix-tensor factorization (CMTF), fusing the zeroth and first order information of the NN. The proposed learning algorithm is based on a constrained alternating least squares (ALS) approach. The activation functions (AFs) are expressed as a weighted sum of pre-defined basis elements. Our method allows to strongly compress large NN layers, with a very slight degradation of the classification accuracy. This also provides insights into the functioning of NNs, which is a problem of utmost interest in NN domain.

1 Introduction

Les réseaux de neurones (RNs) sont des outils puissants, qui ont la capacité d’apprendre et de modéliser des fonctions non linéaires complexes, tout en montrant de bonnes performances de prédiction. Avant d’entraîner un réseau de neurones, il faut décider (*i*) combien de neurones et de couches cachées utiliser ? et (*ii*) quelles fonctions d’activation (FAs) utiliser ? Ces choix contrôlent la capacité d’apprentissage du RN à partir des données. Les architectures des RNs modernes de pointe sont trop volumineuses et comportent de nombreux paramètres. Les exécuter dans des systèmes avec une capacité de calcul limitée est une tâche difficile. Pour cette raison, plusieurs travaux dans la littérature se sont intéressés à la compression de ces modèles.

Les approches tensorielles les plus populaires [1] pour la compression visent principalement à compresser les couches des réseaux convolutifs, qui peuvent être considérés comme des tenseurs. Par une décomposition polyadique canonique (CPD) [2], ils remplacent les convolutions multidimensionnelles par des convolutions unidimensionnelles. Une autre direction de recherche se concentre sur la relation entre les décompositions de tenseurs et les réseaux de neurones avec des unités de produit (au lieu d’unités de sommation) [3]; ce type de représentations est cependant peu utilisé en pratique.

Dans ce travail, nous considérons une approche totalement différente. Tout en conservant la structure traditionnelle du réseau de neurones (poids linéaires + fonctions d’activation non linéaires), nous visons à ajouter de la flexibilité [4] aux FAs, par opposition aux FAs fixes utilisées classiquement. En particulier, les fonctions d’activation peuvent être différentes dans différents nœuds. Une telle architecture est particulièrement intéressante grâce à la théorie de l’identifiabilité (unicité) disponible dans le cas polynomial [5]. Les propriétés d’identifiabilité peuvent donner un aperçu du fonctionnement de ces RNs et aider à renforcer la stabilité de la représentation.

Contrairement aux méthodes existantes pour les FAs flexibles (FAFs) qui utilisent des techniques d’entraînement conventionnelles [4], nous utilisons un cadre original développé dans la communauté d’identification de systèmes non linéaires. Les travaux de [6] ont montré qu’une architecture avec une couche flexible cachée peut être identifiée comme une CPD d’un tenseur jacobien. Cependant, il n’est pas directement applicable dans la configuration d’apprentissage; en particulier, il n’y a pas de moyen simple d’estimer les fonctions d’activation. Dans ce travail, nous proposons une nouvelle méthode de compression de réseaux de neurones pré-entraînés basée sur la factorisation couplée matrice-tenseur (FCMT). Cette technique fusionne les informations d’ordre 0 et 1 du RN, où les informations de premier ordre sont contenues dans un tenseur jacobien,

*Travail soutenu par l’ANR LeaFleT (ANR-19-CE23-0021).

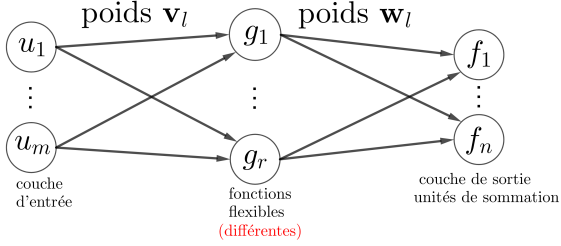


FIGURE 1 – Représentation graphique du bloc de base à une couche du RN flexible.

suivant une CPD contrainte.

Le nouveau RN flexible a des FAs différentes et flexibles dans chaque neurone. Les FAFs sont exprimées sous forme de somme pondérée de fonctions de base prédéfinies. Le but de ce modèle est d'augmenter l'expressivité du modèle pour pouvoir apprendre des fonctions avec un nombre de paramètres inférieur par rapport aux modèles classiques. Ce genre de résultat est susceptible d'accélérer l'exécution des réseaux de neurones et permettra leur utilisation sur des systèmes à faibles ressources de calcul.

2 Modèle de FAFs

Modèle flexible Le modèle proposé dans ce travail est un RN flexible, où chaque neurone a une FA différente et flexible. Les FAs sont exprimées sous la forme d'une combinaison linéaire de fonctions de base prédéfinies. Une illustration graphique d'un bloc de base à une couche de ce modèle est donnée dans la Figure 1. La fonction multivariée qui représente la relation entrée-sortie est notée \mathbf{f} . Cette fonction $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n$ est donnée par :

$$\mathbf{f}(\mathbf{u}) = [f_1(\mathbf{u}) \quad f_n(\mathbf{u})]^T, \quad \text{with } \mathbf{u} = [u_1 \quad u_m]^T. \quad (1)$$

On dit que \mathbf{f} admet un modèle RN flexible s'il peut s'exprimer sous la forme :

$$\mathbf{f}(\mathbf{u}) = \mathbf{W}\mathbf{g}(\mathbf{V}^T \mathbf{u}), \quad (2)$$

où $\mathbf{V} \in \mathbb{R}^{m \times r}$, $\mathbf{W} \in \mathbb{R}^{n \times r}$ sont des matrices de transformation, \mathbf{v}_ℓ et \mathbf{w}_ℓ , pour $1 \leq \ell \leq r$, sont respectivement les colonnes de \mathbf{V} et \mathbf{W} , et $\mathbf{g} : \mathbb{R}^r \rightarrow \mathbb{R}^r$ est exprimée comme

$$\mathbf{g}(t_1, \dots, t_r) = [g_1(t_1) \quad g_r(t_r)]^T, \quad (3)$$

avec $g_\cdot : \mathbb{R} \rightarrow \mathbb{R}$ est une fonction non linéaire univariée.

Fonctions de base Les fonctions non linéaires univariées g_\cdot ($1 \leq \ell \leq r$) dans (2) sont une combinaison linéaire de fonctions de base prédéfinies $\{\phi_1, \dots, \phi_d\}$:

$$g_\cdot(t) = c_0 \cdot 1 + c_1 \cdot \phi_1(t) + \dots + c_d \cdot \phi_d(t). \quad (4)$$

Les fonctions ϕ_k , pour $1 \leq k \leq d$, dans (4) étant *a priori* choisies, les paramètres c_k , pour $0 \leq k \leq d$ et $1 \leq \ell \leq r$, sont les paramètres apprenables.

3 Décomposition d'une fonction multivariée basée sur un tenseur

Représenter une fonction non-linéaire \mathbf{f} par une représentation plus simple, telle que (2) a été abordée dans plusieurs travaux, puisqu'elle permet d'approximer la fonction avec un nombre réduit de paramètres et peut donner un aperçu de la compréhension des mécanismes derrière la fonction non linéaire. Ces méthodes supposent que la fonction \mathbf{f} est connue à travers une représentation paramétrique polynomiale, et ont été développées dans le contexte du problème de découplage pour les polynômes réels multivariés. L'une de ces méthodes est la méthode jacobienne [6] qui vise à décomposer $\mathbf{f}(\mathbf{u})$ selon (2) à partir de son information de premier ordre. Cette approche repose sur l'observation que la matrice jacobienne de \mathbf{f} , évaluée en un point donné \mathbf{u} s'exprime par

$$\begin{aligned} \mathbf{J}_{\mathbf{f}}(\mathbf{u}) &:= \begin{bmatrix} \frac{\partial f_1}{\partial u_1}(\mathbf{u}) & \dots & \frac{\partial f_1}{\partial u_m}(\mathbf{u}) \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial u_1}(\mathbf{u}) & \dots & \frac{\partial f_n}{\partial u_m}(\mathbf{u}) \end{bmatrix} \quad (5) \\ &= \mathbf{W} \text{diag} \{g_1^{\ell}(\mathbf{v}_1^T \mathbf{u}), \dots, g_r^{\ell}(\mathbf{v}_r^T \mathbf{u})\} \mathbf{V}^T. \quad (6) \end{aligned}$$

En empilant les matrices jacobienes évaluées en N points d'échantillonnage différents $\mathbf{u}^{(j)} \in \mathbb{R}^m$, pour $1 \leq j \leq N$, tel que $\mathbf{J}_{\dots j} = \mathbf{J}_{\mathbf{f}}(\mathbf{u}^{(j)})$, nous construisons un tenseur jacobien d'ordre 3 \mathcal{J} de taille $n \times m \times N$ suivant une CPD exprimée comme

$$\mathcal{J} = \underset{j=1}{\times} \mathbf{w}_\cdot \cdot \mathbf{v}_\cdot \cdot \mathbf{h}_\cdot, \quad (7)$$

où $\underset{j=1}{\times}$ dénote le produit tensoriel. De manière équivalente, on note $\mathcal{J} = [\mathbf{j}\mathbf{W}, \mathbf{V}, \mathbf{H}]$. Comme montré dans [6, 7], les colonnes de $\mathbf{H} \in \mathbb{R}^{n \times r}$ ont la structure suivante

$$\mathbf{h}_\cdot = [g_1^{\ell}(\mathbf{v}_1^T \mathbf{u}^{(1)}), \dots, g_r^{\ell}(\mathbf{v}_r^T \mathbf{u}^{(N)})]^T, \quad (8)$$

ce qui montre que le facteur matriciel \mathbf{H} dépend de \mathbf{V} . Il convient de noter que l'approche originale de [6] est assez sensible à la non-unicité de la CPD et au bruit [8, 9], puisqu'une solution ALS relaxée est utilisée pour estimer les facteurs.

Suivant [8, 9], cela peut être attribué au fait que l'algorithme ne tient pas compte de la structure particulière du tenseur jacobien (8), *i.e.* la relation entre les colonnes de \mathbf{H} et \mathbf{V} . Cela devient problématique, en particulier lorsque l'unicité de la CPD est perdue, auquel cas l'ALS relâché ne parviendra pas à récupérer des FAs significatives. Cela a motivé le développement d'algorithmes structurés pour le découplage de polynômes multivariés [8, 9]. De plus, il a été montré dans [7] que la prise en compte de la structure (8) conduit à des conditions d'unicité moins restrictives par rapport au cas non structuré, élargissant ainsi le domaine d'application de la méthode.

4 Problème d'apprentissage

Pour surmonter les inconvénients susmentionnés des méthodes existantes, nous proposons dans cette section un algorithme

pour la CPD structurée du tenseur jacobien. Il consiste à formuler la CPD structurée comme une factorisation couplée matrice-tenseur contrainte, utilisant conjointement les informations d'ordre 0 et 1.

Information d'ordre 0 Soit F une matrice de taille $n \times N$ dont les lignes sont les évaluations des n sorties du sous-réseau d'origine à N points d'échantillonnage différents $\mathbf{u}^{(j)} \in \mathbb{R}^m$:

$$F = \begin{bmatrix} f_1(\mathbf{u}^{(1)}) & \dots & f_1(\mathbf{u}^{(N)}) \\ \vdots & \ddots & \vdots \\ f_n(\mathbf{u}^{(1)}) & \dots & f_n(\mathbf{u}^{(N)}) \end{bmatrix}. \quad (9)$$

En se basant sur (9) et (2), nous pouvons voir que F est sujet à une factorisation donnée par

$$F = \begin{bmatrix} f_1(\mathbf{u}^{(1)}) & \dots & f_1(\mathbf{u}^{(N)}) \\ \vdots & \ddots & \vdots \\ f_n(\mathbf{u}^{(1)}) & \dots & f_n(\mathbf{u}^{(N)}) \end{bmatrix} = W \begin{bmatrix} g_1(\mathbf{v}_1^T \mathbf{u}^{(1)}) & \dots & g_1(\mathbf{v}_r^T \mathbf{u}^{(N)}) \\ \vdots & \ddots & \vdots \\ g_r(\mathbf{v}_1^T \mathbf{u}^{(1)}) & \dots & g_r(\mathbf{v}_r^T \mathbf{u}^{(N)}) \end{bmatrix}, \quad (10)$$

$$= W \begin{bmatrix} g_1(\mathbf{v}_1^T \mathbf{u}^{(1)}) & \dots & g_1(\mathbf{v}_r^T \mathbf{u}^{(N)}) \\ \vdots & \ddots & \vdots \\ g_r(\mathbf{v}_1^T \mathbf{u}^{(1)}) & \dots & g_r(\mathbf{v}_r^T \mathbf{u}^{(N)}) \end{bmatrix} \underbrace{\begin{bmatrix} \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_r \end{bmatrix}}_{\mathbf{Z}^T}, \quad (11)$$

où \mathbf{Z} est de taille $N \times r$. On peut voir à partir de (11) que la matrice F partage le même facteur W avec le tenseur \mathcal{J} . Sur la base de cette observation, nous construisons un problème FCMT contraint, impliquant la matrice F et le tenseur \mathcal{J} .

Algorithme d'apprentissage L'idée de l'approche basée sur la FCMT contrainte est de formuler l'apprentissage des paramètres, W , V et c_k , pour $0 \leq k \leq d$ et $1 \leq \ell \leq r$, comme suit.

$$\min_{W, V, H, Z} \mathcal{J} [W, V, H] + \lambda \|F - WZ^T\|^2 \quad (12)$$

$$\text{t.q. } h_{j,\ell} = c_{1,\ell} \cdot \phi_1^\ell(\mathbf{v}_1^T \mathbf{u}^{(j)}) + \dots + c_{d,\ell} \cdot \phi_d^\ell(\mathbf{v}_d^T \mathbf{u}^{(j)}),$$

$$z_{j,\ell} = c_{0,\ell} + c_{1,\ell} \cdot \phi_1(\mathbf{v}_1^T \mathbf{u}^{(j)}) + \dots + c_{d,\ell} \cdot \phi_d(\mathbf{v}_d^T \mathbf{u}^{(j)}).$$

Ces contraintes s'expriment de manière équivalente comme

$$\mathbf{h}_\ell = \mathbf{X}_\ell \cdot \mathbf{c}_\ell, \quad \mathbf{z}_\ell = \mathbf{Y}_\ell \cdot \mathbf{c}_\ell, \quad \text{for } 1 \leq \ell \leq r,$$

où les matrices \mathbf{X}_ℓ et \mathbf{Y}_ℓ , de taille $N \times (d+1)$, sont données par

$$\mathbf{X}_\ell = \begin{bmatrix} \phi_1^\ell(\mathbf{v}_1^T \mathbf{u}^{(1)}) & \dots & \phi_d^\ell(\mathbf{v}_d^T \mathbf{u}^{(1)}) \\ \vdots & \ddots & \vdots \\ \phi_1^\ell(\mathbf{v}_1^T \mathbf{u}^{(N)}) & \dots & \phi_d^\ell(\mathbf{v}_d^T \mathbf{u}^{(N)}) \end{bmatrix}, \quad (13)$$

$$\mathbf{Y}_\ell = \begin{bmatrix} 1 & \phi_1(\mathbf{v}_1^T \mathbf{u}^{(1)}) & \dots & \phi_d(\mathbf{v}_d^T \mathbf{u}^{(1)}) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(\mathbf{v}_1^T \mathbf{u}^{(N)}) & \dots & \phi_d(\mathbf{v}_d^T \mathbf{u}^{(N)}) \end{bmatrix}, \quad (14)$$

et le ℓ -ième vecteur de coefficient \mathbf{c}_ℓ est défini comme $\mathbf{c}_\ell = [c_{0,\ell}, c_{1,\ell}, \dots, c_{d,\ell}]^T$.

Pour résoudre le problème (12), nous proposons un algorithme des moindres carrés alternés à facteurs H et Z contraints, dans lequel nous mettons à jour les coefficients des FAs g_1, \dots, g_r selon (4). La procédure complète est donnée dans Algorithme 1, où \odot dénote le produit de Khatri-Rao.

Algorithme 1 Algorithme d'apprentissage basé sur la FCMT

Entrée : Tenseur \mathcal{J} de dimension $n \times m \times N$, matrice F de taille $n \times N$, les fonctions $f, \phi_1, \dots, \phi_d, g$, paramètre d'ajustement λ , rang r .

Sortie : Facteurs W , V , H et Z , et les coefficients \mathbf{c}_ℓ .

- 1: Initialiser V , H , Z
- 2: **répéter**
- 3: Mettre à jour W avec $\min_W \text{unfold}_1 \mathcal{J} \cdot W \cdot H + \lambda \|F - WZ^T\|^2$
- 4: Mettre à jour V avec $\min_V \text{unfold}_2 \mathcal{J} \cdot V \cdot H \cdot W$
- 5: Mettre à jour H avec $\min_H \text{unfold}_3 \mathcal{J} \cdot H \cdot V \cdot W$
- 6: Mettre à jour Z avec $\min_Z \|F - WZ^T\|^2$
- 7: **pour** $\ell = 1 \dots r$ **faire**
- 8: Construire \mathbf{X}_ℓ selon (13).
- 9: Construire \mathbf{Y}_ℓ selon (14).
- 10: **fin pour**
- 11: Calculer \mathbf{c}_ℓ comme $\min_{\mathbf{c}_\ell} \text{vec}(\mathbf{X}_\ell) \cdot \mathbf{c}_\ell + \lambda \text{vec}(\mathbf{Y}_\ell) \cdot \mathbf{c}_\ell$
- 12: **pour** $\ell = 1 \dots r$ **faire**
- 13: Calculer \mathbf{h}_ℓ t.q. : $\mathbf{h}_\ell = \mathbf{X}_\ell \cdot \mathbf{c}_\ell$
- 14: Calculer \mathbf{z}_ℓ t.q. : $\mathbf{z}_\ell = \mathbf{Y}_\ell \cdot \mathbf{c}_\ell$
- 15: **fin pour**
- 16: **jusqu'à** un test de convergence est satisfait ou le nombre maximal d'itérations est atteint
- 17: **retourne** W , V , \mathbf{c}_ℓ .

5 Compression des réseaux de neurones pré-entraînés

Pour illustrer l'efficacité de l'apprentissage de réseaux de neurones flexibles basé sur des tenseurs, la compression de réseaux de neurones pré-formés est considérée dans cette section. L'un des principaux objectifs de la compression est de réduire la complexité d'un réseau de neurones afin de permettre son implémentation sur des systèmes embarqués avec des capaci-

