

# Correction matricielle de l'indétermination d'échelle pour l'optimisation alternée

Massyl MOUDOUD Céline MEILLIER Vincent MAZET

ICube, Université de Strasbourg, CNRS (UMR 7357), 300 Bd Sébastien Brant, 67400 Illkirch-Graffenstaden, France

**Résumé** – L'apprentissage de dictionnaire souffre du problème d'indétermination d'échelle (multiplication d'un terme par une matrice diagonale et de l'autre terme par l'inverse de la matrice). La méthode AMORS répond à ce problème en intégrant un scalaire qui compense l'indétermination. Nous proposons une extension de cette méthode en intégrant une matrice diagonale qui permet de compenser chaque atome du dictionnaire. La méthode proposée est plus rapide et plus stable que AMORS pour l'apprentissage de dictionnaire dans le contexte de la connectivité fonctionnelle en IRMf.

**Abstract** – Dictionary learning suffers from the problem of scale indetermination (multiplication of one term by a diagonal matrix and the other term by the inverse of the matrix). The AMORS method addresses this problem by integrating a scalar that compensates for the indetermination. We propose an extension to this method by integrating a diagonal matrix allowing to compensate each atom of the dictionary. The proposed method is faster and more stable compared to AMORS for the dictionary learning in the context of functional connectivity in fMRI.

L'apprentissage de dictionnaire est un problème inverse qui apparaît dans de nombreuses applications [2, 7]. Il existe une indétermination d'échelle dans la résolution de ce problème. En effet, le problème ne peut être résolu qu'à un facteur d'échelle près. La majorité des méthodes d'apprentissage de dictionnaire [2, 7] ignorent ce facteur d'échelle et se contentent de normaliser la solution *a posteriori* [1, 8].

Thé et. al [6] ont montré que les termes de régularisation du problème sont affectés par le facteur d'échelle. Ils ont présenté la méthode AMORS qui permet de prendre en compte le facteur d'échelle, de trouver sa valeur optimale et de l'utiliser pour ajuster les hyperparamètres de régularisation durant la résolution du problème, ce qui améliore la convergence.

Dans cet article, nous proposons d'étendre cette méthode en faisant en sorte que chaque atome du dictionnaire ait son propre facteur d'échelle. Dans la première partie, nous présentons la preuve de la nécessité de prendre en compte un facteur d'échelle par atome dans les problèmes d'apprentissage de dictionnaire. Ensuite, nous montrons comment tirer avantage de ces facteurs d'échelle pour améliorer la convergence des algorithmes d'optimisation. Enfin, nous présentons une application de notre méthode et un comparatif avec AMORS.

## 1 Méthode

Nous reprenons ci-après la méthode AMORS [6] en l'adaptant au cas d'un facteur d'échelle par atome du dictionnaire.

### 1.1 Le modèle d'apprentissage de dictionnaire

On considère le modèle d'apprentissage de dictionnaire

$$\mathbf{C} \approx \mathbf{D}\mathbf{A} \quad (1)$$

avec  $\mathbf{C} \in \mathbb{R}^{E \times T}$  une matrice d'observations. L'objectif est d'exprimer cette matrice comme un produit des deux ma-

trices  $\mathbf{D} \in \mathbb{D} \subseteq \mathbb{R}^{E \times P}$  et  $\mathbf{A} \in \mathbb{A} \subseteq \mathbb{R}^{P \times T}$ . On peut voir  $\mathbf{D}$  comme un dictionnaire d'atomes et  $\mathbf{A}$  comme leurs amplitudes. Les ensembles  $\mathbb{D}$  et  $\mathbb{A}$  définissent l'espace des solutions qui peuvent être restreintes aux valeurs positives par exemple. Le modèle (1) n'étant pas exact, on cherche à trouver une paire de matrices  $(\hat{\mathbf{D}}, \hat{\mathbf{A}})$  minimisant l'erreur de reconstruction de  $\mathbf{C}$  :

$$\hat{\mathbf{D}}, \hat{\mathbf{A}} = \arg \min_{\mathbf{D}, \mathbf{A}} \|\mathbf{C} - \mathbf{D}\mathbf{A}\|_F^2. \quad (2)$$

L'erreur de reconstruction est invariante à la multiplication de  $\mathbf{A}$  par une matrice diagonale  $\mathbf{\Gamma}$  et  $\mathbf{D}$  par son inverse :

$$\begin{aligned} \|\mathbf{C} - \mathbf{D}\mathbf{A}\|_F^2 &= \|\mathbf{C} - (\mathbf{D}\mathbf{\Gamma}^{-1})(\mathbf{\Gamma}\mathbf{A})\|_F^2 \\ \text{avec } \mathbf{\Gamma} \in \mathbb{R}^{P \times P} \text{ et } \mathbf{\Gamma}_{ij} &= \begin{cases} \gamma_i \in \mathbb{R}_+^*, & \text{si } i = j \\ 0, & \text{sinon} \end{cases} \end{aligned} \quad (3)$$

Cela signifie que si chaque atome du dictionnaire (colonne  $i$  de  $\mathbf{D}$ ) est divisé par un facteur  $\gamma_i$ , alors le problème (2) reste inchangé tant que l'amplitude correspondante dans la ligne  $i$  de  $\mathbf{A}$  est multipliée par ce même facteur  $\gamma_i$ .

Le problème (2) est mal posé au sens d'Hadarnard. Pour le résoudre, il est courant d'ajouter des termes de régularisation sur  $\mathbf{D}$  et  $\mathbf{A}$  issus de connaissances *a priori*. Le problème (2) est donc remplacé par  $(\boldsymbol{\lambda}, \boldsymbol{\mu} \in \mathbb{R}_+^{*P})$  :

$$\hat{\mathbf{D}}, \hat{\mathbf{A}} = \arg \min_{\mathbf{D}, \mathbf{A}} \mathcal{L}(\mathbf{D}, \mathbf{A}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \quad (4)$$

$$\mathcal{L}(\mathbf{D}, \mathbf{A}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \|\mathbf{C} - \mathbf{D}\mathbf{A}\|_F^2 + \boldsymbol{\lambda}^T \mathcal{J}(\mathbf{D}) + \boldsymbol{\mu}^T \mathcal{K}(\mathbf{A})$$

$$\boldsymbol{\lambda}^T \mathcal{J}(\mathbf{D}) + \boldsymbol{\mu}^T \mathcal{K}(\mathbf{A}) = \sum_{i=1}^P (\lambda_i j(\mathbf{D}_i) + \mu_i k(\mathbf{A}^i))$$

où  $\mathbf{D}_i$  est la colonne  $i$  de  $\mathbf{D}$  et  $\mathbf{A}^i$  la ligne  $i$  de  $\mathbf{A}$ . Contrairement à la formulation habituelle où les hyperparamètres de régularisation sont des scalaires, nous choisissons ici le cas plus général où les paramètres  $\boldsymbol{\lambda}$  et  $\boldsymbol{\mu}$  sont des vecteurs. L'intérêt est d'affecter une régularisation adaptée à chaque atome

Cette recherche a été financée par l'Agence Nationale de la Recherche au titre du projet DynaSTI ANR-22-CE45-0008.

du dictionnaire en réglant plus précisément les connaissances a priori. Par exemple, si un atome est supposé devoir être d'amplitude plus faible que les autres, alors son paramètre de régularisation propre sera augmenté. Il reste toutefois possible d'avoir le même poids de régularisation sur l'ensemble des inconnues en choisissant des vecteurs  $\lambda$  et  $\mu$  constants.

On impose que les termes de régularisations  $\mathcal{J}$  et  $\mathcal{K}$  soient séparables en fonction des colonnes de  $\mathbf{D}$  et des lignes de  $\mathbf{A}$  ce qui est le cas de la majorité des fonctions utilisées (normes  $\ell_1, \ell_2$ , variation totale anisotrope...):

$$\mathcal{J}(\mathbf{D}) = \left[ j(\mathbf{D}_i) \right]_{i \in [1, P]} \quad \mathcal{K}(\mathbf{A}) = \left[ k(\mathbf{A}^i) \right]_{i \in [1, P]} \quad (5)$$

La fonction  $\mathcal{L}$  est minimisée par une optimisation alternée des sous-problèmes sur  $\mathbf{D}$  et  $\mathbf{A}$ :

$$\hat{\mathbf{D}} = \arg \min_{\mathbf{D}} \|\mathbf{C} - \mathbf{D}\mathbf{A}\|_F^2 + \lambda^T \mathcal{J}(\mathbf{D}), \quad (6a)$$

$$\hat{\mathbf{A}} = \arg \min_{\mathbf{A}} \|\mathbf{C} - \mathbf{D}\mathbf{A}\|_F^2 + \mu^T \mathcal{K}(\mathbf{A}). \quad (6b)$$

Les sous-problèmes sont convexes tant que  $\mathcal{J}$  et  $\mathcal{K}$  sont convexes et peuvent donc être optimisés avec des méthodes simples. Notons que des contraintes de positivité ou de support peuvent être ajoutées sans compromettre la méthode proposée.

## 1.2 Exploiter le facteur d'échelle

Afin de pouvoir utiliser l'indétermination de facteur d'échelle à notre avantage, on impose l'homogénéité de degré  $r > 0$  et  $q > 0$  des régularisations [6]:

$$\begin{aligned} j(\gamma_i \mathbf{D}_i) &= \gamma_i^{-q} j(\mathbf{D}_i) \implies \mathcal{J}(\mathbf{D}\mathbf{\Gamma}^{-1}) = \mathbf{\Gamma}^{-q} \mathcal{J}(\mathbf{D}) \\ k(\gamma_i \mathbf{A}_i) &= \gamma_i^r k(\mathbf{A}_i) \implies \mathcal{K}(\mathbf{\Gamma}\mathbf{A}) = \mathbf{\Gamma}^r \mathcal{K}(\mathbf{A}) \end{aligned} \quad (7)$$

$\forall \gamma_i \in \mathbb{R}_+^* \quad \text{et} \quad \forall \mathbf{D} \in \mathbb{D}, \mathbf{A} \in \mathbb{A}$

D'après (4) et (7), on constate que la fonction de coût  $\mathcal{L}$  avec les inconnues  $\mathbf{D}$  et  $\mathbf{A}$  multipliées par le facteur d'échelle est équivalente à la remise à l'échelle des hyperparamètres de régularisations:

$$\mathcal{L}(\mathbf{D}\mathbf{\Gamma}^{-1}, \mathbf{\Gamma}\mathbf{A}, \lambda, \mu) = \mathcal{L}(\mathbf{D}, \mathbf{A}, \mathbf{\Gamma}^{-q}\lambda, \mathbf{\Gamma}^r\mu) \quad (8)$$

En d'autres termes, le facteur d'échelle change la valeur des hyperparamètres de régularisations. Cela signifie que, si au cours des itérations de l'optimisation alternée, le facteur d'échelle évolue, alors le problème résolu évolue également. La matrice qui augmente en valeur sera plus pénalisée et celle qui diminue le sera moins. Cet effet va donc changer la forme des matrices estimées étant donné qu'elles sont la solution à une autre formulation du problème et non celle souhaitée. Cet effet a pour conséquence de déplacer les minima locaux.

Pour résoudre ce problème, on cherche la matrice diagonale  $\mathbf{\Gamma}$  qui minimise la fonction de coût  $\mathcal{L}$ . On définit:

$$\begin{aligned} \mathcal{E}(\mathbf{\Gamma}) &= \mathcal{L}(\mathbf{D}\mathbf{\Gamma}^{-1}, \mathbf{\Gamma}\mathbf{A}, \lambda, \mu) \\ &= \|\mathbf{C} - \mathbf{D}\mathbf{A}\|_F^2 + \sum_{i=1}^P (\lambda_i \gamma_i^{-q} j(\mathbf{D}_i) + \mu_i \gamma_i^r k(\mathbf{A}^i)) \end{aligned} \quad (9)$$

Les termes  $\epsilon(\gamma_i) = \lambda_i \gamma_i^{-q} j(\mathbf{D}_i) + \mu_i \gamma_i^r k(\mathbf{A}^i)$  sont chacun dépendant d'un élément  $\gamma_i$  de  $\mathbf{\Gamma}$  et indépendants les uns des autres. Il est alors possible de minimiser par rapport à  $\gamma_i$  la fonction globale  $\mathcal{E}(\mathbf{\Gamma})$  en minimisant chaque terme  $\epsilon(\gamma_i)$ .

On distingue trois cas dans le calcul de la dérivée de  $\epsilon(\gamma_i)$ :

1. Si  $j(\mathbf{D}_i) \neq 0$  et  $k(\mathbf{A}^i) \neq 0$  alors:

$$\frac{\partial \mathcal{E}}{\partial \gamma_i} = \frac{d\epsilon(\gamma_i)}{d\gamma_i} = -q\lambda_i \gamma_i^{-q-1} j(\mathbf{D}_i) + r\mu_i \gamma_i^{r-1} k(\mathbf{A}^i) \quad (10)$$

Dans ce cas  $\epsilon(\gamma_i)$  est strictement convexe et son minimum est:

$$\hat{\gamma}_i = \left( \frac{q\lambda_i j(\mathbf{D}_i)}{r\mu_i k(\mathbf{A}^i)} \right)^{\frac{1}{r+q}} \quad (11)$$

2. Si  $j(\mathbf{D}_i) = 0$  et  $k(\mathbf{A}^i) = 0$ , alors  $\gamma_i$  peut prendre une valeur arbitraire.
3. Si  $j(\mathbf{D}_i) = 0$  ou  $k(\mathbf{A}^i) = 0$ , mais pas les deux, alors c'est un cas dégénéré qui n'a pas de solution. Il peut apparaître pendant une estimation alternée de  $\mathbf{D}$  et  $\mathbf{A}$ . Un atome du dictionnaire peut avoir une activation nulle mais quand même exister dans  $\mathbf{D}$  ou vice-versa. Nous pouvons négliger ce cas et assigner une valeur arbitraire à  $\gamma_i$  car dès l'estimation suivante d'une des deux matrices, on retombe dans le cas 2.

L'algorithme 1 s'inspire de la méthode AMORS [6] pour incorporer la connaissance du facteur d'échelle optimal dans la procédure d'optimisation alternée. Après chaque estimation de l'une des deux matrices  $\mathbf{D}$  et  $\mathbf{A}$ , le facteur d'échelle optimal est calculé pour chaque atome du dictionnaire avec (9). Lors de l'estimation suivante, le sous-problème (6a) ou (6b) est résolu en pondérant l'autre matrice (qui reste constante) par le facteur d'échelle optimal. Dans AMORS, c'est l'hyperparamètre de régularisation qui est pondéré. En théorie, selon (8), c'est équivalent. Cependant, nous avons constaté des problèmes numériques avec cette approche du fait de l'apparition de valeurs extrêmes à la limite de la précision machine.

---

**Algorithme 1 :** Minimisation alternée avec matrice de correction d'échelle optimale.

---

**Entrées :**  $\mathbf{D}^{[0]}, \lambda, \mu$

```

1  $\mathbf{\Gamma} \leftarrow \mathbf{I}_P$ ;
2  $k \leftarrow 0$ ;
3 tant que pas convergé faire
4    $\mathbf{A} \leftarrow \arg \min_{\mathbf{A}} \|\mathbf{C} - \mathbf{D}\mathbf{\Gamma}^{-1}\mathbf{A}\|_F^2 + \mu^T \mathcal{K}(\mathbf{A})$ ;
5    $\tilde{\mathbf{\Gamma}} \leftarrow \arg \min_{\mathbf{\Gamma}} \mathcal{L}(\mathbf{D}\mathbf{\Gamma}^{-1}, \mathbf{\Gamma}\mathbf{A}, \lambda, \mu)$ ;
6   si  $k = 0$  et  $\max_i (\log(\tilde{\Gamma}_{ii}) - \log(\Gamma_{ii})) > \epsilon$  alors
7      $\mathbf{\Gamma} \leftarrow \tilde{\mathbf{\Gamma}}$ ;
7     Retourner à la ligne : 4
8   fin
9    $\mathbf{D} \leftarrow \arg \min_{\mathbf{D}} \|\mathbf{C} - \mathbf{D}\tilde{\mathbf{\Gamma}}\mathbf{A}\|_F^2 + \lambda^T \mathcal{J}(\mathbf{D})$ ;
10   $\mathbf{\Gamma} \leftarrow \arg \min_{\mathbf{\Gamma}} \mathcal{L}(\mathbf{D}\mathbf{\Gamma}^{-1}, \mathbf{\Gamma}\mathbf{A}, \lambda, \mu)$ ;
11   $k \leftarrow k + 1$ ;
12 fin
13 retourner  $(\mathbf{D}, \mathbf{A})$ 
```

---

Il reste la question de l'initialisation du facteur d'échelle. Nous reprenons la méthode proposée dans [6]. Tout d'abord, le dictionnaire  $\mathbf{D}$  est initialisé à une valeur définie par l'utilisateur et la matrice des facteurs d'échelle  $\mathbf{\Gamma}$  est initialisée à l'identité. Les activations  $\mathbf{A}$  sont ensuite estimées grâce à

l'équation (6b), puis la matrice des facteurs d'échelle est mise à jour. Cette opération est répétée tant que le facteur d'échelle change significativement, c'est-à-dire tant que :

$$\max_i \left( \log(\tilde{\Gamma}_{ii}) - \log(\Gamma_{ii}) \right) < \epsilon. \quad (12)$$

Concernant le cas où un atome est à zéro et qu'on tombe dans le cas 2 ou 3 du calcul du facteur d'échelle optimal, on choisit de garder la valeur de  $\gamma_i$  précédente qui a fait que cet atome est estimé à zéro. Ce choix garantit la stabilité de la méthode, surtout pendant l'initialisation du facteur d'échelle.

Il faut bien noter que l'algorithme 1 ne résout pas le facteur d'échelle : il y aura toujours une indétermination de la solution à un facteur près par atome du dictionnaire. Notre méthode propose d'annuler l'effet du facteur d'échelle sur la résolution du problème. On garantit que les deux matrices sont régularisées au niveau souhaité, peu importe l'évolution du facteur d'échelle au cours des itérations. Grâce à la prise en compte du facteur d'échelle par atome du dictionnaire, nous allons encore plus loin et garantissons que tous les atomes sont pénalisés aux niveaux souhaités, même si les atomes évoluent vers différentes échelles de valeurs.

## 2 Application

La méthode proposée est testée pour la résolution du problème de dictionnaire lié à l'extraction d'unités de connectivité fonctionnel (UCF) dans l'étude de la dynamique de la connectivité fonctionnelle en IRMf [4, 5]. L'objectif est de décomposer les corrélations dans  $\mathbf{C}$  en réseaux cérébraux définis dans  $\mathbf{D}$  dont l'activité temporelle est estimée dans  $\mathbf{A}$ . Chaque ligne de  $\mathbf{C}$  représente l'évolution dans le temps de la corrélation entre deux régions du cerveau. Dans nos travaux sur les UCF [4], les sous-problèmes (6a) et (6b) s'écrivent :

$$\min_{\mathbf{D} \in \mathbb{R}_+^P} \|\mathbf{C} - \mathbf{D}\mathbf{A}\|_F^2 + \sum_{i=1}^P \left( \lambda_i \|\mathbf{D}_i\|_2^2 \right) + \mathcal{I}_{\tilde{\mathbf{D}}}(\mathbf{D}) \quad (13a)$$

$$\min_{\mathbf{A} \in \mathbb{R}_+^P} \|\mathbf{C} - \mathbf{D}\mathbf{A}\|_F^2 + \sum_{i=1}^P \left( \mu_i \|\mathbf{A}^i\|_1 + \nu_i TV(\mathbf{A}^i) \right) \quad (13b)$$

Dans le sous-problème (13a), la pénalisation  $\ell_2$  définit le terme  $\mathcal{J}$ . Elle évite l'apparition de valeurs aberrantes dans le dictionnaire. Il y a aussi un ensemble de contraintes d'égalité avec le terme  $\mathcal{I}_{\tilde{\mathbf{D}}}(\mathbf{D})$ . Elles sont définies par la matrice binaire  $\tilde{\mathbf{D}}$  qui impose la structure des réseaux (atomes du dictionnaire), c.-à-d.  $\mathbf{D}_{ij} = 0$  si  $\tilde{\mathbf{D}}_{ij} = 0$ . Elle est fournie par les biologistes qui définissent les réseaux d'intérêts pour l'étude.

Concernant (13b),  $\mathcal{K}$  est composé de deux régularisations : une contrainte de parcimonie de type  $\ell_1$  et une régularisation de variation totale [3] sur les lignes de  $\mathbf{A}$ . Cette pénalisation combinée sur  $\mathbf{A}$  n'empêche pas l'utilisation de notre méthode car la contrainte  $\ell_1$  et la variation totale sont toutes les deux homogènes d'ordre  $r = 1$ . On peut donc les regrouper dans une seule pénalisation en contrôlant l'influence relative des deux termes et l'influence globale avec deux hyperparamètres. Notons qu'il n'aurait pas été possible d'appliquer deux pénalisations homogènes d'ordres différents sur la même matrice, variation totale et  $\ell_2$  par exemple, car dans ce cas la solution explicite du facteur d'échelle n'est plus valide.

Les sous-problèmes (13a) et (13b) sont convexes, ils peuvent être résolus avec n'importe quelle méthode d'opti-

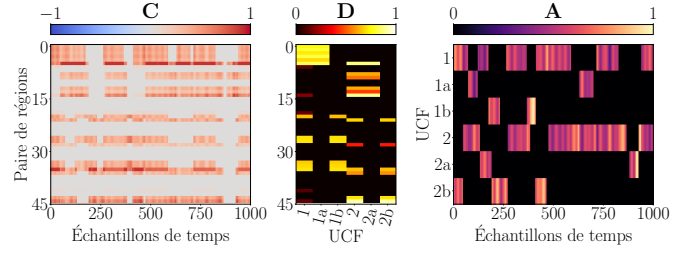


FIGURE 1 : Jeu de données du modèle direct  $\mathbf{C} = \mathbf{D}\mathbf{A}$  avec  $(E, T) = (45, 1000)$  et  $P = 6$  atomes.

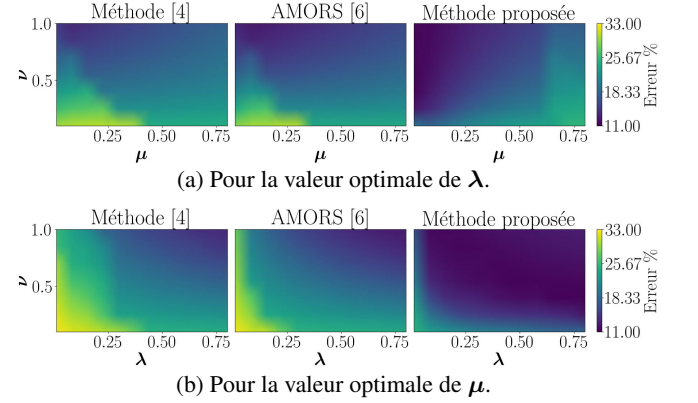


FIGURE 2 : Erreur de reconstruction de  $\mathbf{C}$  en fonction des couples  $(\nu, \mu)$  (a) et  $(\nu, \lambda)$  (b).

misation sous contraintes. Dans notre précédent travail [4], l'algorithme ADMM (*alternating direction method of multipliers*) a été utilisé. Dans ce travail, nous préférons utiliser l'algorithme FISTA (*fast iterative shrinkage thresholding algorithm*) qui donne la même solution que ADMM mais converge plus rapidement et avec une plus faible complexité algorithmique.

## 3 Résultats

Dans cette section, nous comparons la méthode proposée à l'algorithme AMORS [6] et à notre précédente méthode [4] (optimisation alternée classique sans correction du facteur d'échelle, mais dans laquelle l'algorithme ADMM est remplacé par FISTA). Nous utilisons le jeu de données synthétiques figure 1, généré selon la procédure présentée dans [4]. Ces données sont créées selon le modèle direct (1) et fournissent donc la vérité terrain pour  $\mathbf{C}$ ,  $\mathbf{D}$  et  $\mathbf{A}$ .

La structure de  $\mathbf{D}$  étant connue dans ce problème grâce à  $\tilde{\mathbf{D}}$ , les algorithmes sont initialisés avec

$$\mathbf{D}_i^{[0]} = \sum_{t=1}^T \mathbf{C}_t \circ \tilde{\mathbf{D}}_i, \quad (14)$$

c.-à-d. la moyenne temporelle de la matrice d'observations bruitée masquée par  $\tilde{\mathbf{D}}$  ( $\circ$  est le produit élément par élément). Pour assurer la bonne convergence, l'optimisation alternée est calculée avec 500 itérations et chaque sous-problème effectuée 500 itérations de FISTA.

L'effet des hyperparamètres de régularisations est étudié en exécutant les algorithmes pour des triplets d'hyperparamètres pris dans  $\Theta = (\lambda, \mu, \nu) \in [0, 0.8] \times [0, 0.8] \times [0, 1]$  pour un total de 1008 combinaisons. L'erreur de reconstruction de  $\mathbf{C}$ ,

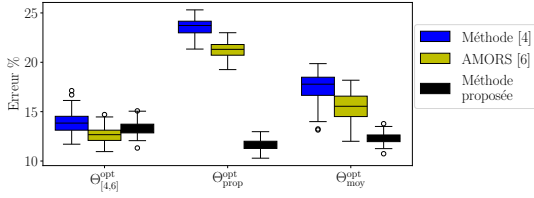


FIGURE 3 : Distribution de l’erreur de reconstruction de  $\mathbf{C}$  normaliser en pourcentage pour les méthodes [4], AMORS [6] et la méthode proposée, avec trois jeux d’hyperparamètres.

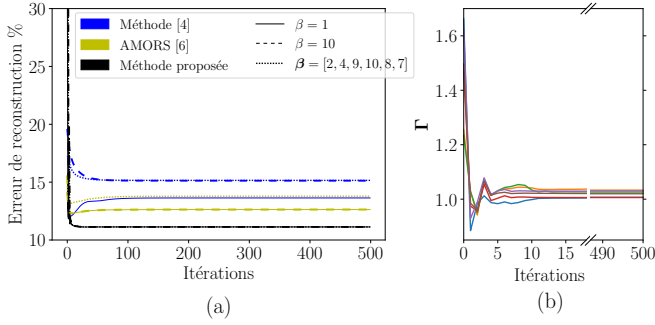


FIGURE 4 : (a) Courbes de convergences de l’erreur de reconstruction de la matrice  $\mathbf{C}$  estimée par les trois méthodes : classique [4], AMORS ainsi que la méthode proposée, (b) évolution du facteur d’échelle optimal par atome du dictionnaire estimé par la méthode proposée.

normalisée en pourcentage, pour les trois algorithmes en fonction des couples  $(\nu, \mu)$  et  $(\nu, \lambda)$  est affichée dans les figures 2a et 2b pour la meilleure valeur du paramètre restant.

Nous constatons que AMORS a des résultats très similaires à l’optimisation alternée classique. La méthode proposée minimise l’erreur pour une grande plage d’hyperparamètres. Elle est en moyenne meilleure que les deux autres méthodes.

La stabilité d’estimation des méthodes est étudiée en générant 100 matrices  $\mathbf{A}$  différentes pour un dictionnaire  $\mathbf{D}$  fixé. Pour chaque matrice  $\mathbf{C}$  résultante, les trois méthodes sont utilisées pour estimer la décomposition. Trois jeux d’hyperparamètres sont considérés : le triplet  $\Theta_{[4,6]}^{\text{opt}}$  minimisant l’erreur de  $\mathbf{C}$  pour l’optimisation alternée et AMORS (qui se sont avérées avoir le même minimiseur sur ce jeu de données), celui minimisant l’erreur de la méthode proposée  $\Theta_{\text{prop}}^{\text{opt}}$  ainsi que la moyenne de ces minimiseurs  $\Theta_{\text{moy}}^{\text{opt}}$ .

La distribution de l’erreur pour chaque méthode et jeu de paramètres est affichée figure 3. Dans la majorité des cas, notre méthode a une erreur inférieure et sa variance est plus faible, à l’exception du cas  $\Theta_{[4,6]}^{\text{opt}}$  où AMORS a de meilleurs résultats. De plus, notre méthode est plus stable à une variation des valeurs des hyperparamètres. Pour l’estimation de  $\mathbf{D}$  et  $\mathbf{A}$ , l’erreur est significativement plus faible, y compris pour  $\Theta_{[4,6]}^{\text{opt}}$ .

L’influence de l’échelle de la matrice d’initialisation  $\mathbf{D}^{[0]}$  sur la convergence est évaluée. Pour ce faire, l’initialisation de (14) est multipliée par un facteur  $\beta = 10$  (scalaire) et  $\beta = [2, 4, 9, 10, 8, 7]^T$  où l’échelle de chaque atome est modifiée. Ces deux cas sont comparés au cas  $\beta = 1$  (14). Pour ces trois conditions, la valeur des hyperparamètres correspond au  $\Theta^{\text{opt}}$  de chaque méthode déterminé dans le cas  $\beta = 1$  (cf. figure 2) Les courbes de convergence de l’erreur de reconstruction de  $\mathbf{C}$  en pourcentage sont affichées figure 4 (a). Notre méthode

garde la même vitesse de convergence et n’est pas affectée par l’initialisation contrairement à la méthode [4] et à AMORS qui convergent toutes les deux vers une moins bonne solution en fonction de l’initialisation.

Enfin, la figure 4 (b) montre l’évolution du facteur d’échelle optimal estimé par notre méthode pour chaque atome du dictionnaire. On voit clairement que chaque atome converge vers un facteur différent, confortant la nécessité d’optimiser le facteur d’échelle par atome du dictionnaire et non pour toute la matrice comme dans AMORS.

## 4 Conclusion

Nous avons proposé une extension de la méthode AMORS pour prendre en compte l’indétermination de facteur d’échelle par atome dans les problèmes d’apprentissage de dictionnaire. Cela permet de réduire le risque de tomber dans des cas dégénérés où un atome converge vers de très grandes valeurs et écrase les autres du fait de la régularisation.

La méthode proposée a été comparée à AMORS et à l’optimisation alternée classique sur des données synthétiques issues de l’étude de la connectivité fonctionnelle en IRMf. Les résultats ont montré que notre méthode améliore la convergence et obtient de meilleurs résultats. Enfin, notre méthode n’est pas affectée par l’échelle de l’initialisation. Un cas pratique de cette dernière propriété est le démélange hyperspectral [3]. Si les différents spectres initiaux des *end-members* sont estimés avec des instruments différents alors le résultat n’est pas affecté par le réglage de ces instruments.

Une perspective à ce travail est d’évaluer notre méthode sur des opérateurs bilinéaires autres que le produit matriciel, notamment le produit de convolution pour la résolution du problème de déconvolution.

Le code pour reproduire les résultats présentés ainsi que la génération des données des simulations est disponible à l’adresse : [https://github.com/massylmoudoud/DynaSTI\\_Gretsi2025](https://github.com/massylmoudoud/DynaSTI_Gretsi2025).

## Références

- [1] D. BENACHIR, Y. DEVILLE, S. HOSSEINI, M. S. KAROUI et A. HAMEURLAIN. “Hyperspectral image unmixing by non-negative matrix factorization initialized with modified independent component analysis”. In : *WHISPERS*. 2013.
- [2] P. COMON et C. JUTTEN, éd. *Handbook of blind source separation. Independent component analysis and applications*. Academic Press, 2010.
- [3] M.-D. IORDACHE, J. M. BIOUSCAS-DIAS et A. PLAZA. “Total Variation Spatial Regularization for Sparse Hyperspectral Unmixing”. In : *IEEE Trans. Geosci. Remote Sens.* 50.11 (2012), p. 4484-4502.
- [4] M. MOUDOUD, C. MEILLIER, M. SOURTY et V. MAZET. “Spatio-Temporal Model for Dynamic Functional Connectivity in Resting State fMRI Analysis”. In : *EUSIPCO*. 2024.
- [5] V. PORTMANN, C. MEILLIER et V. MAZET. “Analyse de la dynamique spatio-temporelle de la connectivité fonctionnelle cérébrale : données synthétiques et modélisation”. In : *GRETSI*, 2023.
- [6] S. THÉ, É. THIÉBAUT, L. DENIS et F. SOULEZ. “Exploiting the scaling indetermination of bi-linear models in inverse problems”. In : *EUSIPCO*. 2020.
- [7] I. TOSIC et P. FROSSARD. “Dictionary Learning”. In : *IEEE Signal Process Mag.* (2011).
- [8] H. YANG et C. SEOIGHE. “Impact of the Choice of Normalization Method on Molecular Cancer Class Discovery Using Nonnegative Matrix Factorization”. In : *PLOS ONE* (2016).