

# Optimisation bi-niveaux pour l'apprentissage d'hyperparamètres

Jordan Frécon

Equipe Apprentissage - LITIS - INSA de Rouen



1. Introduction : définition et exemples
2. Rappels sur la sélection de modèles
3. Formalisme bi-niveaux pour l'apprentissage d'hyperparamètres
4. Algorithmes basés sur le gradient
  - 4.1 Principe
  - 4.2 Cas différentiable
  - 4.3 Cas non différentiable
5. Application à la découverte de structure
5. Conclusion et perspectives

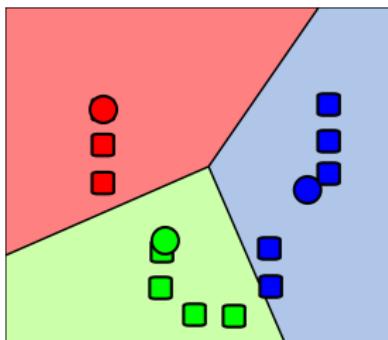
# Introduction

# Qu'est-ce qu'un hyperparamètre ?

## Hyperparamètre

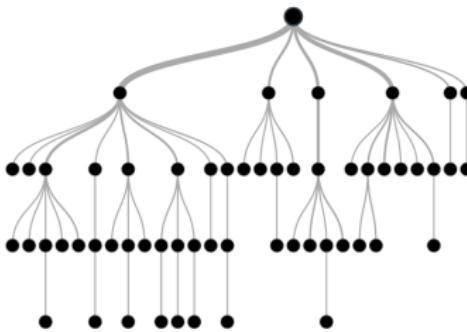
Un hyperparamètre  $\lambda$  est un paramètre défini avant le début du processus d'apprentissage. Ces (hyper)-paramètres sont réglables et peuvent affecter directement les performances d'un modèle et la façon dont il s'apprend.

## Exemple 1 : Partitionnement en k-moyennes



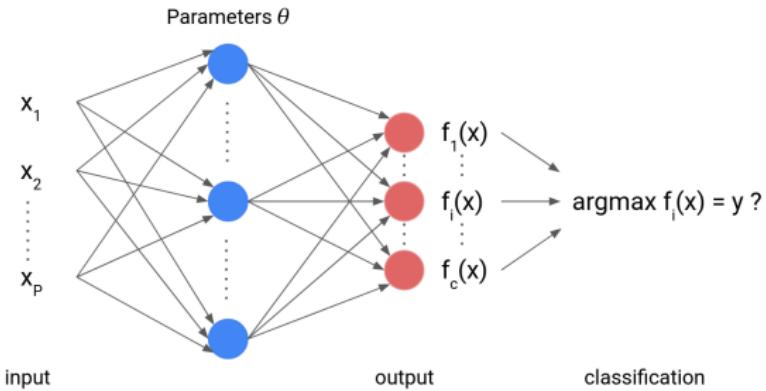
Hyperparamètre  $\lambda = \text{nombre de clusters } k$

## Exemple 2 : Arbre de décision



Hyperparamètre  $\lambda$  = nombre de branches

## Exemple 3 : Réseaux de neurones



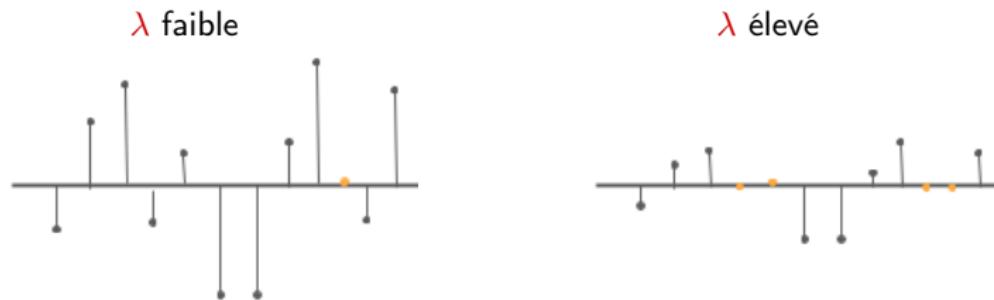
Hyperparamètre  $\lambda$  = nombre d'epochs, taux d'apprentissage, ...

## Exemple 4 : Lasso

Soient des données  $\{x_i, y_i\}_{i=1}^n$ , on souhaite trouver un **modèle de régression linéaire parcimonieux** entre  $x$  et  $y$ .

En d'autres termes, on cherche  $w$  parcimonieux de telle sorte que  $y_i \approx x_i w$ .

$$\min_w \left\{ \mathcal{L}(w; \lambda) \triangleq \frac{1}{2} \|y - Xw\|^2 + \lambda \|w\|_1 \right\}$$



Hyperparamètre  $\lambda$  = paramètre de régularisation

## Exemple 5 : Détection de changements

$$\min_w \left\{ \mathcal{L}(w; \lambda) \triangleq \frac{1}{2} \|y - w\|^2 + \lambda \text{TV}(w) \right\}$$

Hyperparamètre  $\lambda$  = paramètre de régularisation

## Entiers

- nombre de clusters dans un algorithme de partitionnement
- nombre de branches dans un arbre de décision
- nombre d'epochs

## Réels

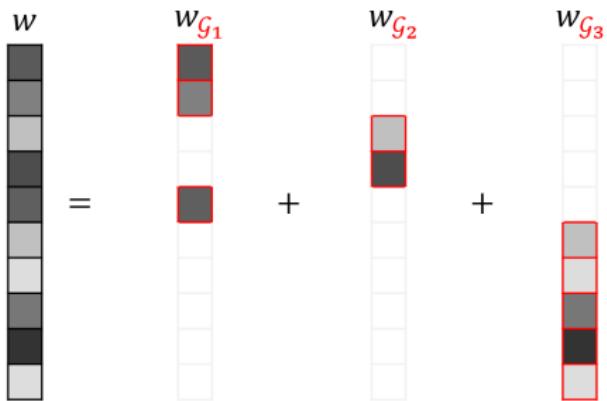
- taux d'apprentissage
- moments
- paramètres de régularisation

Particularité mon domaine de recherche : **Structures** (groupes, hiérarchies, graphes)

## Exemple 6 : Group Lasso

La solution du *Group Lasso* [Yuan and Lin (2006)] dépend des groupes  $\{\mathcal{G}_1, \dots, \mathcal{G}_L\}$

$$\hat{w}(\mathcal{G}_1, \dots, \mathcal{G}_L) = \underset{w}{\operatorname{argmin}} \frac{1}{2} \|y - Xw\|^2 + \gamma \sum_{l=1}^L \|w_{\mathcal{G}_l}\|_2$$



Hyperparamètre  $\lambda = \text{Groupes } \{\mathcal{G}_1, \dots, \mathcal{G}_L\}$

## Rappels sur la sélection de modèles

Pour chaque hyperparamètre  $\lambda$ , nous avons un modèle  $\hat{w}(\lambda)$  différent

Exemples :

$\hat{w}(\lambda)$  = partitionnement en  $\lambda$  moyennes ( $k$ -means clustering)

$\hat{w}(\lambda)$  = réseau de neurones appris avec un taux d'apprentissage valant  $\lambda$

Choisir  $\lambda \Leftrightarrow$  Sélection de modèle  $\{\hat{w}(\lambda)\}_\lambda$

Questions :

1. Dans quel espace choisir  $\lambda$  ?
2. Comment évaluer la qualité du modèle  $\hat{w}(\lambda)$  ?

Pour chaque hyperparamètre  $\lambda$ , nous avons un modèle  $\hat{w}(\lambda)$  différent

**Exemples :**

$\hat{w}(\lambda)$  = partitionnement en  $\lambda$  moyennes ( $k$ -means clustering)

$\hat{w}(\lambda)$  = réseau de neurones appris avec un taux d'apprentissage valant  $\lambda$

Choisir  $\lambda \Leftrightarrow$  Sélection de modèle  $\{\hat{w}(\lambda)\}_\lambda$

**Questions :**

1. Dans quel espace choisir  $\lambda$  ?
2. Comment évaluer la qualité du modèle  $\hat{w}(\lambda)$  ?

Une stratégie populaire consiste à choisir  $\lambda$  sur une grille  $\Lambda$  prédéfinie

La qualité du modèle  $\hat{w}(\lambda)$  est évaluée selon un critère  $\mathcal{E}$  :

- critère d'information d'Akaike (AIC)
- critère d'information bayésien (BIC)
- norme induisant la parcimonie
- :
- erreur de validation

## Principe de la recherche sur grille

Pour chaque  $\lambda \in \Lambda = \{\lambda_{\min}, \dots, \lambda_{\max}\}$

  | Calculer  $\hat{w}(\lambda)$

  | Evaluer  $\mathcal{E}(\hat{w}(\lambda))$

Choisir  $\hat{\lambda}$  qui minimise  $\{\mathcal{E}(\hat{w}(\lambda))\}_{\lambda \in \Lambda}$



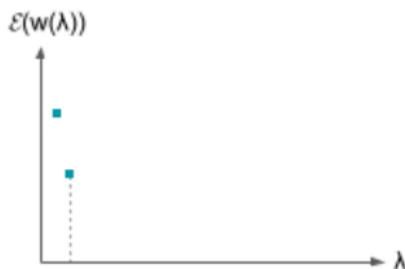
## Principe de la recherche sur grille

Pour chaque  $\lambda \in \Lambda = \{\lambda_{\min}, \dots, \lambda_{\max}\}$

  | Calculer  $\hat{w}(\lambda)$

  | Evaluer  $\mathcal{E}(\hat{w}(\lambda))$

  | Choisir  $\hat{\lambda}$  qui minimise  $\{\mathcal{E}(\hat{w}(\lambda))\}_{\lambda \in \Lambda}$



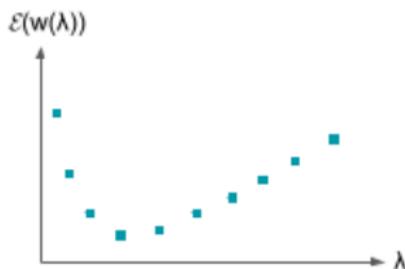
## Principe de la recherche sur grille

Pour chaque  $\lambda \in \Lambda = \{\lambda_{\min}, \dots, \lambda_{\max}\}$

  | Calculer  $\hat{w}(\lambda)$

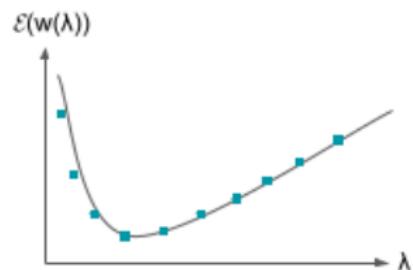
  | Evaluer  $\mathcal{E}(\hat{w}(\lambda))$

  | Choisir  $\hat{\lambda}$  qui minimise  $\{\mathcal{E}(\hat{w}(\lambda))\}_{\lambda \in \Lambda}$

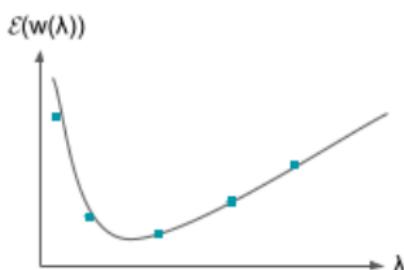


## Limitations de la recherche sur grille

Grille fine

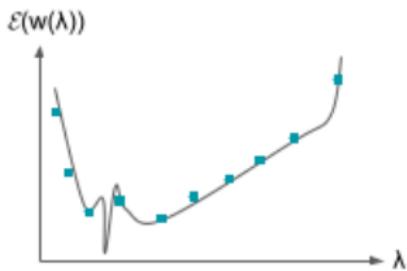


Grille grossière



La grille doit être choisie selon les propriétés de  $\lambda \mapsto \mathcal{E}(\hat{w}(\lambda))$

## Limitations de la recherche sur grille



La fonction  $\lambda \mapsto E(\hat{w}(\lambda))$  peut être non convexe

# Formalisme bi-niveaux pour l'apprentissage d'hyperparamètres

### Problème bi-niveaux

$$\min_{\lambda \in \Lambda} \mathcal{E}(\hat{w}(\lambda), \lambda) \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

Le critère  $\mathcal{E}$  peut également dépendre de  $\lambda$  et l'espace de recherche  $\Lambda$  est un ensemble continu

#### Questions :

1. Est-ce que le problème est bien posé ?
2. Est-ce qu'il possède des solutions ?
3. Si oui, comment trouver ces solutions ?

## 1. Etude du problème

### Problème bi-niveaux

$$\min_{\lambda \in \Lambda} \mathcal{E}(\hat{w}(\lambda), \lambda) \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

### Partitionnement en k-moyennes.

Soient des données  $\{x_i\}_{i=1}^n$ . On souhaite trouver  $k = \lambda$  clusters, notés  $w_1, \dots, w_\lambda$ , via la minimisation de

$$\mathcal{L}(w; \lambda) = \sum_{i=1}^{\lambda} \sum_{x \in w_i} \|x - \mu_i\|^2$$

où est la moyenne des points dans le cluster  $w_i$  est noté  $\mu_i$ .

Problème NP difficile !

## 1. Etude du problème

### Problème bi-niveaux

$$\min_{\lambda \in \Lambda} \mathcal{E}(\hat{w}(\lambda), \lambda) \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

#### Probleme Lasso.

Soient des données  $\{x_i, y_i\}_{i=1}^n$ , on souhaite trouver un modèle de régression linéaire parcimonieux entre  $x$  et  $y$ . En d'autres termes, on cherche  $w$  parcimonieux de telle sorte que  $y_i \approx x_i w$ . Pour cela, on considère la minimisation de

$$\mathcal{L}(w; \lambda) = \frac{1}{2} \|y - xw\|^2 + \lambda \|w\|_1$$

L'objectif  $\mathcal{L}(\cdot; \lambda)$  est convexe mais non différentiable

## 1. Etude du problème

### Problème bi-niveaux

$$\min_{\lambda \in \Lambda} \mathcal{E}(\hat{w}(\lambda), \lambda) \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

On souhaite :

- Qu'il existe une seule solution  $\hat{w}(\lambda)$
- La solution  $\hat{w}(\lambda)$  soit *facile* à trouver

$\Rightarrow \mathcal{L}(\cdot; \lambda)$  strictement convexe

## 2. Existence de solutions

### Problème bi-niveaux

$$\min_{\lambda \in \Lambda} \mathcal{E}(\hat{w}(\lambda), \lambda) \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

#### Théorème informel

Si  $\Lambda$  est compacte,  $\mathcal{L}(\cdot; \lambda)$  strictement convexe et  $\mathcal{E}$  est continu, alors  $\lambda \mapsto \hat{w}(\lambda)$  est continu et le problème bi-niveaux admet des solutions.

« J. Frecon, S. Salzo, and M. Pontil. [Bilevel learning of the group lasso structure](#). In Advances in Neural Information Processing Systems 31 (NeurIPS), pages 8301–8311. 2018 »

### 3. Méthodes de résolution

#### Problème bi-niveaux

$$\min_{\lambda \in \Lambda} \mathcal{E}(\hat{w}(\lambda), \lambda) \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

- Recherche sur grille
- Recherche aléatoire
- Algorithmes évolutionnaires
- :
- Algorithmes basés sur le gradient —> sujet de la suite de cet exposé !

# Optimisation bi-niveaux basée sur le gradient

## Problème bi-niveaux

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}(\lambda) \triangleq \mathcal{E}(\hat{w}(\lambda), \lambda) \right\} \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

On peut voir le problème supérieur comme visant à minimiser un objectif  $\mathcal{U}$

Choix initial  $\lambda_0 \in \Lambda$

Pour  $q = 0, \dots,$

$$| \quad \lambda_{q+1} = \operatorname{Proj}_{\Lambda} (\lambda_q - \eta \nabla \mathcal{U}(\lambda_q))$$

## Problèmes potentiels

- Cela suppose que l'on connaisse  $\hat{w}(\lambda)$  de manière exacte !
- $\mathcal{U}$  n'est pas forcément différentiable
- Comment choisir le pas  $\eta$ ? Est-ce que  $\mathcal{U}$  est localement lisse ?

## Problème bi-niveaux

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}(\lambda) \triangleq \mathcal{E}(\hat{w}(\lambda), \lambda) \right\} \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

On peut voir le problème supérieur comme visant à minimiser un objectif  $\mathcal{U}$

Choix initial  $\lambda_0 \in \Lambda$

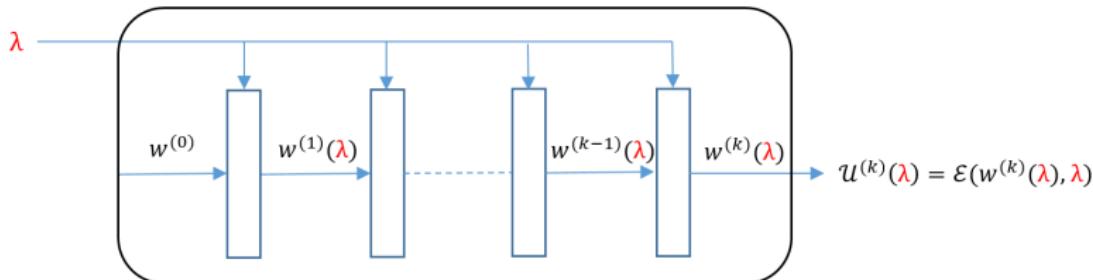
Pour  $q = 0, \dots,$

$$|\lambda_{q+1} = \operatorname{Proj}_{\Lambda}(\lambda_q - \eta \nabla \mathcal{U}(\lambda_q))$$

## Problèmes potentiels

- Cela suppose que l'on connaisse  $\hat{w}(\lambda)$  de manière exacte !
- $\mathcal{U}$  n'est pas forcément différentiable
- Comment choisir le pas  $\eta$ ? Est-ce que  $\mathcal{U}$  est localement lisse ?

## Problème bi-niveaux approché



$$w^{(k)}(\lambda) \xrightarrow[k \rightarrow \infty]{} \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w, \lambda)$$

## Problème bi-niveaux approché

$$\min_{\lambda \in \Lambda} \mathcal{E}(w^{(k)}(\lambda), \lambda) \quad \text{s.c.} \quad w^{(k)}(\lambda) \xrightarrow[k \rightarrow \infty]{} \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

# Qualité de l'approximation

## Problème exact

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}(\lambda) \triangleq \mathcal{E}(\hat{w}(\lambda), \lambda) \right\}$$

s.c.  $\hat{w}(\lambda) = \operatorname{argmin}_{w \in \mathcal{W}} \mathcal{L}(w; \lambda)$

## Problème approché

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}^{(k)}(\lambda) \triangleq \mathcal{E}(w^{(k)}(\lambda), \lambda) \right\}$$

for  $i = 0, 1, \dots, k - 1$   
s.c.  $w^{(i+1)}(\lambda) = \mathcal{A}(w^{(i)}(\lambda), \lambda)$   
 $w^{(k)}(\lambda) \rightarrow \hat{w}(\lambda)$

# Qualité de l'approximation

## Problème exact

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}(\lambda) \triangleq \mathcal{E}(\hat{\mathbf{w}}(\lambda), \lambda) \right\}$$

s.c.  $\hat{\mathbf{w}}(\lambda) = \operatorname{argmin}_{w \in \mathcal{W}} \mathcal{L}(w; \lambda)$

## Problème approché

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}^{(k)}(\lambda) \triangleq \mathcal{E}(\mathbf{w}^{(k)}(\lambda), \lambda) \right\}$$

for  $i = 0, 1, \dots, k - 1$   
s.c.  $\begin{aligned} \mathbf{w}^{(i+1)}(\lambda) &= \mathcal{A}(\mathbf{w}^{(i)}(\lambda), \lambda) \\ \mathbf{w}^{(k)}(\lambda) &\rightarrow \hat{\mathbf{w}}(\lambda) \end{aligned}$

## Théorème informel

Si  $\mathbf{w}^{(k)}(\lambda) \xrightarrow{k \rightarrow \infty} \hat{\mathbf{w}}(\lambda)$  uniformément, alors les minima et la valeur infimum de  $\mathcal{U}^{(k)}$  convergent vers ceux de  $\mathcal{U}$

« J. Frecon, S. Salzo, and M. Pontil. [Bilevel learning of the group lasso structure](#). In Advances in Neural Information Processing Systems 31 (NeurIPS), pages 8301–8311. 2018 »

# Qualité de l'approximation

## Problème exact

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}(\lambda) \triangleq \mathcal{E}(\hat{\mathbf{w}}(\lambda), \lambda) \right\}$$

s.c.  $\hat{\mathbf{w}}(\lambda) = \operatorname{argmin}_{w \in \mathcal{W}} \mathcal{L}(w; \lambda)$

## Problème approché

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}^{(k)}(\lambda) \triangleq \mathcal{E}(\mathbf{w}^{(k)}(\lambda), \lambda) \right\}$$

for  $i = 0, 1, \dots, k - 1$   
s.c.  $\begin{aligned} \mathbf{w}^{(i+1)}(\lambda) &= \mathcal{A}(\mathbf{w}^{(i)}(\lambda), \lambda) \\ \mathbf{w}^{(k)}(\lambda) &\rightarrow \hat{\mathbf{w}}(\lambda) \end{aligned}$

Les propriétés de  $\mathcal{U}$  dépendent de l'objectif  $\mathcal{L}$

Les propriétés de  $\mathcal{U}^{(k)}$  dépendent de l'algorithme  $\mathcal{A}$

Certaines propriétés de  $\mathcal{A}$  dépendent de  $\mathcal{L}$

Dans la suite, nous allons considérer les problèmes bi-niveaux exacts et approchés selon les cas où ils sont différentiables ou non

# Optimisation bi-niveaux basée sur le gradient

- Cas différentiable -

## Problème bi-niveaux exact et différentiable

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}(\lambda) \triangleq \mathcal{E}(\hat{w}(\lambda), \lambda) \right\} \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

## Différentiation implicite

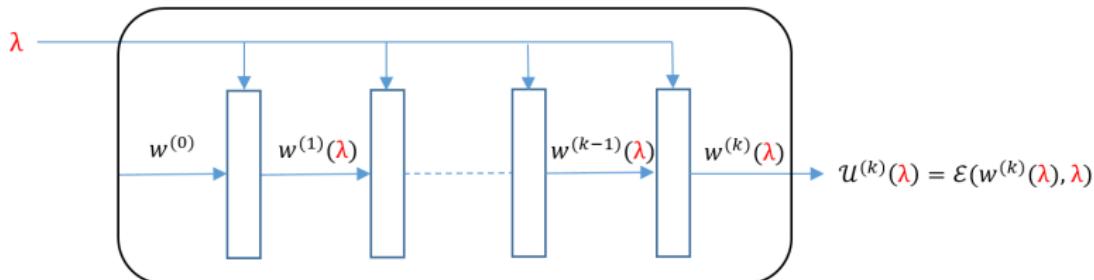
« Almeida, L. [A learning rule for asynchronous perceptrons with feedback in a combinatorial environment](#). In Proceedings, 1st First International Conference on Neural Networks 1987 »

« Pineda, F. J. [Generalization of back-propagation to recurrent neural networks](#) . In Physical Review Letters 1987 »

« Domke, J. [Generic methods for optimization-based modeling](#). In AISTATS 2012 »

« Pedregosa, F. [Hyperparameter optimization with approximate gradient](#). In ICML 2016 »

## Problème bi-niveaux approché et différentiable

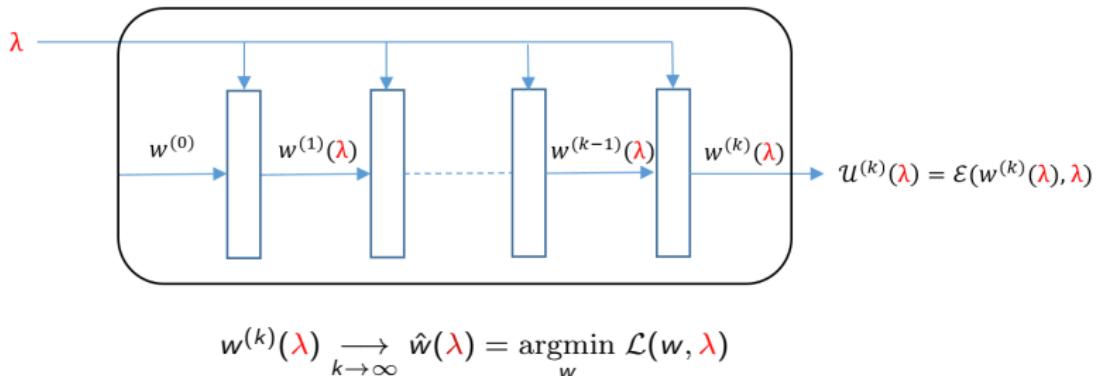


$$w^{(k)}(\lambda) \xrightarrow[k \rightarrow \infty]{} \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w, \lambda)$$

$\mathcal{L}$  différentiable  $\Rightarrow$  Algorithme généralement différentiable  
 $\Rightarrow$  Rétropropagation du gradient

$$\nabla \mathcal{U}^{(k)}(\lambda) = ? \longrightarrow \text{Dérivation en chaîne de } w^{(i+1)}(\lambda) = \mathcal{A}(w^{(i)}(\lambda), (\lambda))$$

## Problème bi-niveaux approché et différentiable



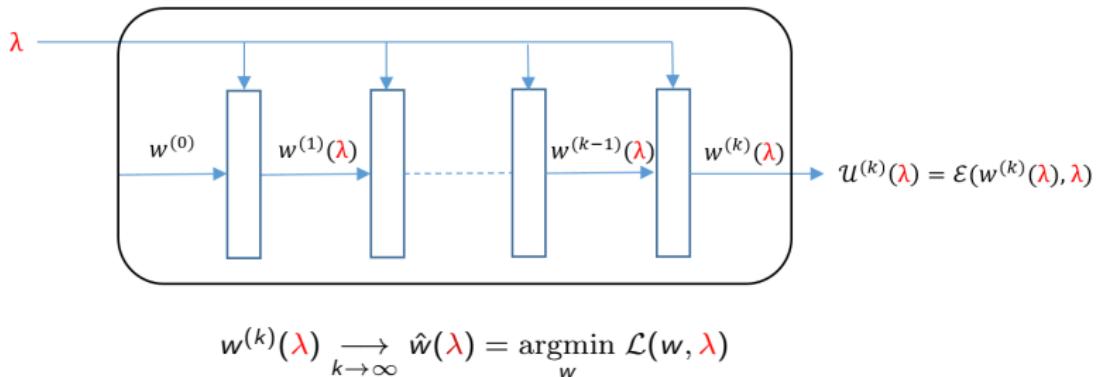
### Différentiation déroulée

« Griewank, A. W. [Evaluating Derivatives : Principles and Techniques of Algorithmic Differentiation](#). In Cambridge press 2008 »

« Maclaurin, D., Duvenaud, D., and Adams, R. P. [Gradient-based hyperparameter optimization through reversible learning](#). In ICML 2015 »

« Franceschi, L., Donini, M., Frasconi, P., and Pontil, M. [Forward and reverse gradient-based hyper-parameter optimization](#). In ICML 2017 »

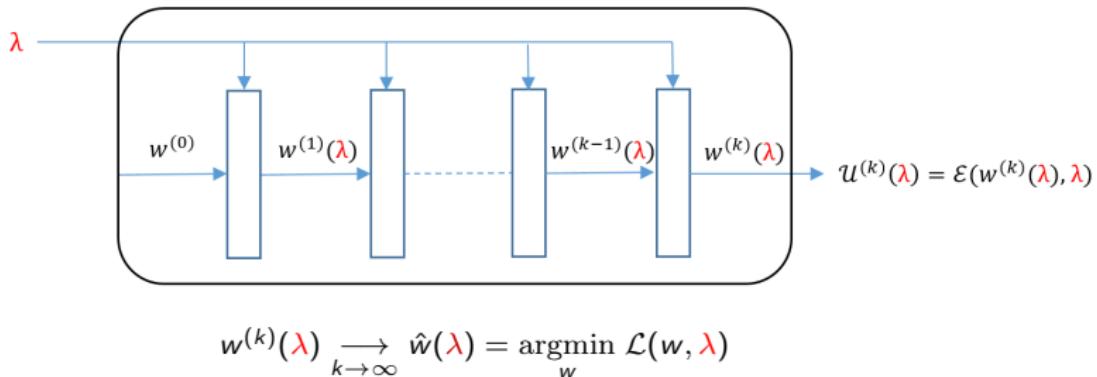
## Problème bi-niveaux approché et différentiable



Différentiation déroulée tronquée

« Shaban, A., Cheng, C.-A., Hatch, N., and Boots, B. [Truncated back-propagation for bilevel optimization](#). In AISTATS 2019 »

## Problème bi-niveaux approché et différentiable



Approximation à la dernière itération

« Lorraine, J., Vicol, P., and Duvenaud, D. [Optimizing millions of hyperparameters by implicit differentiation](#). In [AISTATS 2020](#) »

# Optimisation bi-niveaux basée sur le gradient

- Cas non différentiable -

## Problème bi-niveaux exact et non différentiable

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}(\lambda) \triangleq \mathcal{E}(\hat{w}(\lambda), \lambda) \right\} \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w; \lambda)$$

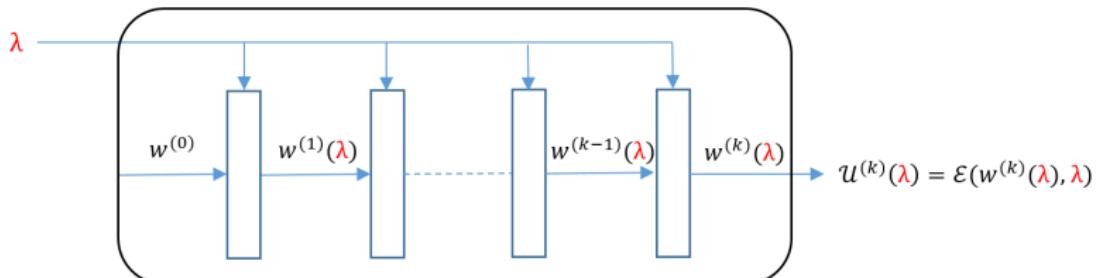
Analyse variationnelle de l'application  $\lambda \mapsto \hat{w}(\lambda)$

Dérivation des conditions d'optimalité

« B. Mordukhovich. [Variational analysis and applications](#). In Springer Monographs in Mathematics, 2018. »

« J.C. De los Reyes, D. Villacís. [Optimality Conditions for Bilevel Imaging Learning Problems with Total Variation Regularization](#). arXiv preprint. 2021. »

## Problème bi-niveaux approché et non-différentiable

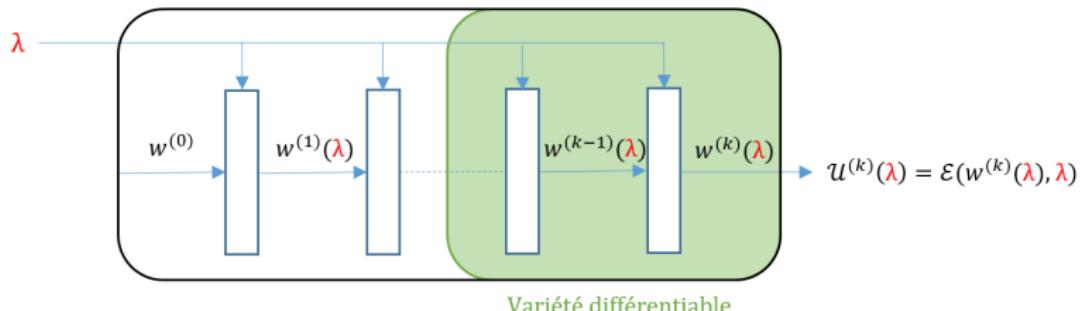


$$w^{(k)}(\lambda) \xrightarrow[k \rightarrow \infty]{} \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w, \lambda)$$

$\mathcal{L}$  non-différentiable  $\Rightarrow$  Algorithme généralement non-différentiable  
 $\Rightarrow$  Rétropropagation du gradient instable

« S. M. Kakade, and J. D. Lee [Provably correct automatic subdifferentiation for qualified programs](#). In NIPS 2018 »

## Problème bi-niveaux approché et non-différentiable



$$w^{(k)}(\lambda) \xrightarrow{k \rightarrow \infty} \hat{w}(\lambda) = \operatorname{argmin}_w \mathcal{L}(w, \lambda)$$

Differentiation dans les variétés différentiables

« Q. Bertrand, Q. Klopfenstein, M. Blonde, S. Vaiter, A. Gramfort, and J. Salmon [Implicit differentiation of Lasso-type models for hyperparameter optimization](#). In ICML 2020 »

## Aparté : exemple de variété différentiable

$$\hat{w}(\lambda) = \operatorname{argmin}_w \underbrace{\frac{1}{2} \|y - Xw\|^2}_{f(w)} + \underbrace{\lambda \|w\|_1}_{g(w)}$$

Résolution par minimisation successive de majorants quadratiques

Choix initial  $w^{(0)}$

Pour  $k = 0, \dots$

$$w^{(k+1)} = \operatorname{argmin}_w f(w^{(k)}) + \nabla f(w^{(k)})^\top (w - w^{(k)}) + \frac{1}{2\tau} \|w - w^{(k)}\|^2 + g(w)$$

## Aparté : exemple de variété différentiable

$$\hat{w}(\lambda) = \operatorname{argmin}_w \underbrace{\frac{1}{2} \|y - Xw\|^2}_{f(w)} + \underbrace{\lambda \|w\|_1}_{g(w)}$$

Se réduit à un algorithme forward-backward

Choix initial  $w^{(0)}$

Pour  $k = 0, \dots$

$$w^{(k+1)} = \operatorname{prox}_{\tau g}(w^{(k)} - \tau \nabla f(w^{(k)}))$$

## Aparté : exemple de variété différentiable

$$\hat{w}(\lambda) = \operatorname{argmin}_w \underbrace{\frac{1}{2} \|y - Xw\|^2}_{f(w)} + \underbrace{\lambda \|w\|_1}_{g(w)}$$

Le tout est en fait l'algorithme de descente de gradient seuillé

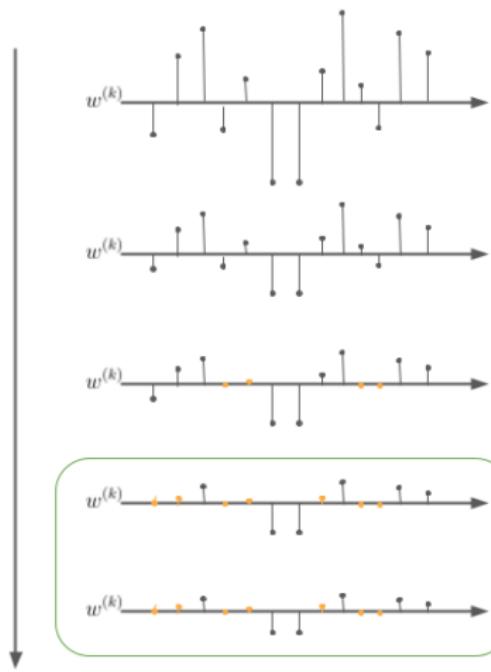
Choix initial  $w^{(0)}$

Pour  $k = 0, \dots$

$$w^{(k+1)} = \text{Soft}_{\tau\lambda}(w^{(k)} - \tau \nabla f(w^{(k)}))$$

## Aparté : exemple de variété différentiable

En appliquant cet algorithme, on observe ...

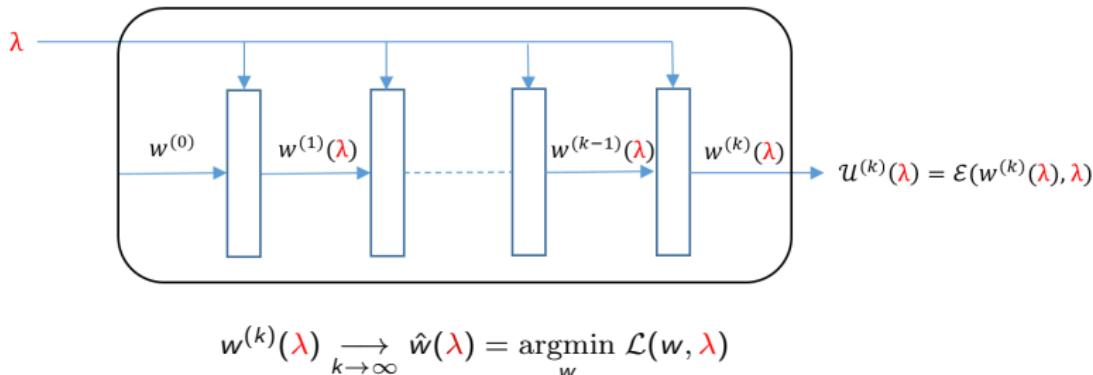


A partir de cette itération  
le support reste inchangé

Variété différentiable

itérations  $k$

## Problème bi-niveaux approché et non-différentiable



### Algorithme différentiables pour résoudre des problèmes non-différentiables

« J. Frecon, S. Salzo, and M. Pontil. [Bilevel learning of the group lasso structure](#). In Advances in Neural Information Processing Systems 31 (NeurIPS), pages 8301–8311. 2018 »

« J. Frecon, S. Salzo, and M. Pontil. [Unveiling groups of related tasks in multi-task learning](#). In Proceedings of the International Conference on Pattern Recognition (ICPR). 2020 »

« J. Frecon, S. Salzo, and M. Pontil. [Bilevel optimization of groupwise penalties](#). En préparation. 2022 »

$$\hat{w} = \operatorname{argmin}_w f(w) + g(w)$$

Résolution par minimisation successive de majorants *bien choisis*

Choix initial  $w^{(0)}$

Pour  $k = 0, \dots$

$$| \quad w^{(k+1)} = \operatorname{argmin}_w f(w^{(k)}) + \nabla f(w^{(k)})^\top (w - w^{(k)}) + \frac{1}{\tau} D(w, w^{(k)}) + g(w)$$

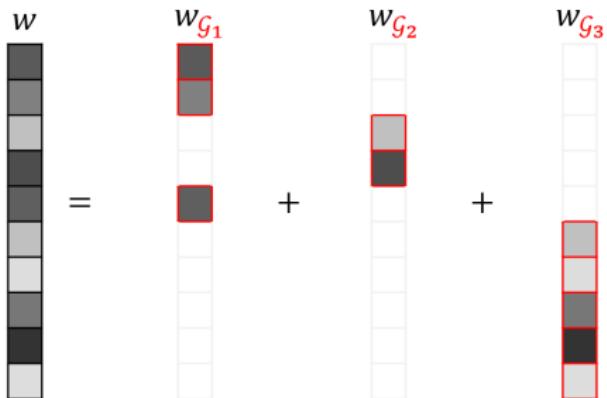
Le choix de la distance  $D(w, w^{(k)})$  peut permettre d'obtenir un algorithme différentiable

## Application à la découverte de structure

## Application à la découverte de structure

La solution du *Group Lasso* [Yuan and Lin (2006)] dépend des groupes  $\{\mathcal{G}_1, \dots, \mathcal{G}_L\}$

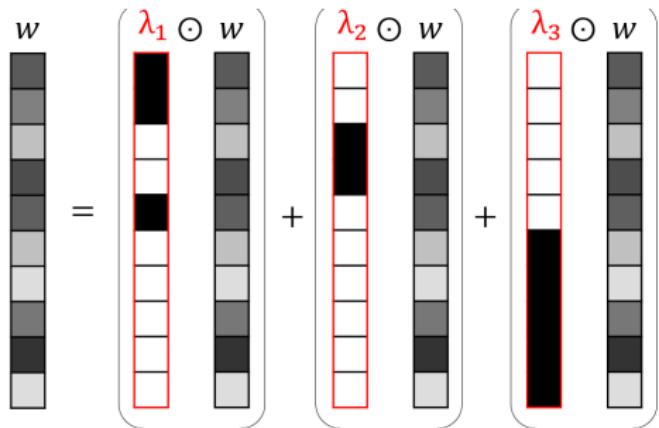
$$\hat{w}(\mathcal{G}_1, \dots, \mathcal{G}_L) = \operatorname{argmin}_{w \in \mathcal{W}} \ell_{\text{trn}}(w) + \gamma \sum_{l=1}^L \|w_{\mathcal{G}_l}\|_2$$

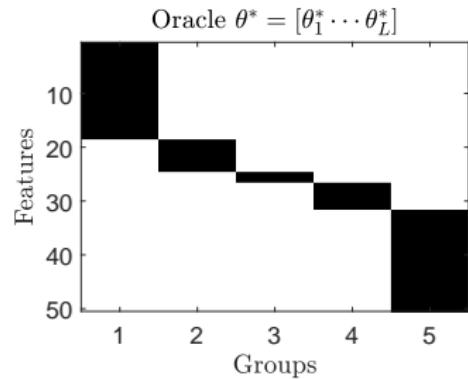


## Application à la découverte de structure

Formalisme bi-niveaux pour apprendre l'indicatrice  $\lambda = [\lambda_1 \cdots \lambda_L]$  des groupes

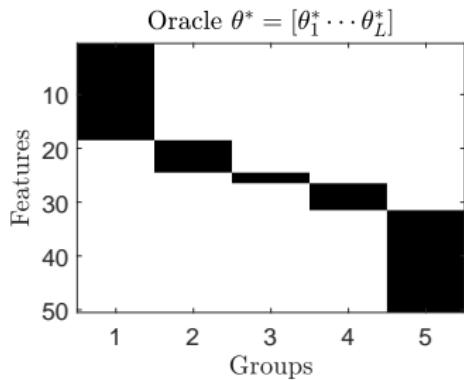
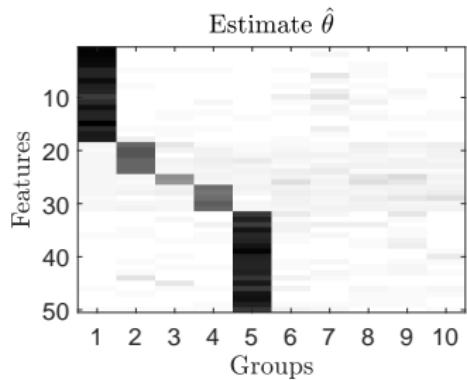
$$\min_{\lambda=[\lambda_1 \cdots \lambda_L] \in \Lambda} \ell_{\text{val}}(\hat{w}(\lambda)) \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_{w \in \mathcal{W}} \ell_{\text{trn}}(w) + \gamma \sum_{l=1}^L \|\lambda_l \odot w\|_2$$





On retrouve les bons groupes ! (modulo une permutation)

## Quand le nombre de groupes est inconnu



Identifie correctement les groupes inutiles

## Régularisation structurée : $\lambda$ = groupes de variables

« J. Frecon, S. Salzo, and M. Pontil. [Bilevel optimization of groupwise penalties](#). En préparation.  
2022 »

$$\min_{\lambda=[\lambda_1 \cdots \lambda_L] \in \Lambda} \ell_{\text{val}}(\hat{w}(\lambda)) \quad \text{s.c.} \quad \hat{w}(\lambda) = \operatorname{argmin}_{w \in \mathcal{W}} \ell_{\text{trn}}(w) + \underbrace{\sum_{l=1}^L \Omega_l(w, \lambda_l)}_{\text{apprendre la structure de la régularisation}}$$

### Théorème informel

Les itérées convergent vers un point stationnaire de l'erreur de validation

# Contributions en apprentissage bi-niveaux

## Régularisation structurée : $\lambda$ = groupes de variables

« J. Frecon, S. Salzo, and M. Pontil. [Bilevel optimization of groupwise penalties](#). En préparation. 2021 »

## Apprentissage multi-tâches : $\lambda$ = groupes de tâches

« J. Frecon, S. Salzo, and M. Pontil. [Unveiling groups of related tasks in multi-task learning](#). In Proceedings of the International Conference on Pattern Recognition (ICPR). 2020 »

Tâche 1



Tâche 2



Tâche 3



Tâche 4



Tâche 5



# Contributions en apprentissage bi-niveaux

## Régularisation structurée : $\lambda$ = groupes de variables

« J. Frecon, S. Salzo, and M. Pontil. [Bilevel optimization of groupwise penalties](#). En préparation. 2021 »

## Apprentissage multi-tâches : $\lambda$ = groupes de tâches

« J. Frecon, S. Salzo, and M. Pontil. [Unveiling groups of related tasks in multi-task learning](#). In Proceedings of the International Conference on Pattern Recognition (ICPR). 2020 »

Tâche 1



Tâche 2



Tâche 3



Tâche 4



Tâche 5



# Contributions en apprentissage bi-niveaux

## Régularisation structurée : $\lambda$ = groupes de variables

« J. Frecon, S. Salzo, and M. Pontil. [Bilevel optimization of groupwise penalties](#). En préparation. 2021 »

## Apprentissage multi-tâches : $\lambda$ = groupes de tâches

« J. Frecon, S. Salzo, and M. Pontil. [Unveiling groups of related tasks in multi-task learning](#). In Proceedings of the International Conference on Pattern Recognition (ICPR). 2020 »



### Théorème informel

L'opérateur proximal de la norme spectrale est différentiable pour une classe de distances de Bregman

## Régularisation structurée : $\lambda$ = groupes de variables

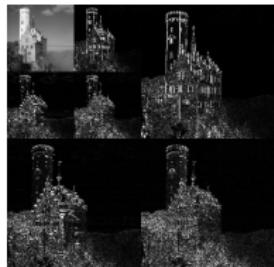
« J. Frecon, S. Salzo, and M. Pontil. [Bilevel optimization of groupwise penalties](#). En préparation. 2021 »

## Apprentissage multi-tâches : $\lambda$ = groupes de tâches

« J. Frecon, S. Salzo, and M. Pontil. [Unveiling groups of related tasks in multi-task learning](#). In Proceedings of the International Conference on Pattern Recognition (ICPR). 2020 »

## Apprentissage d'ondelette : $\lambda$ = filtre d'ondelette

« J. Frecon, R. Grazzi, S. Salzo, and M. Pontil. [Smooth optimization of wavelet bases](#). Submitted to Conférence sur l'Apprentissage Automatique (CAp). 2021 »



double objectif :

Réduire le nombre de coefficients d'ondelettes

Diminuer l'erreur de reconstruction

## Conclusion et perspectives

# 1. Etude des problèmes bi-niveaux non différentiables

## Problème exact

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}(\lambda) \triangleq \mathcal{E}(\hat{\mathbf{w}}(\lambda), \lambda) \right\}$$

s.c.  $\hat{\mathbf{w}}(\lambda) = \operatorname{argmin}_{w \in \mathcal{W}} \mathcal{L}(w; \lambda)$

## Problème approché

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}^{(k)}(\lambda) \triangleq \mathcal{E}(\mathbf{w}^{(k)}(\lambda), \lambda) \right\}$$

for  $i = 0, 1, \dots, k - 1$   
s.c.  $\begin{cases} \mathbf{w}^{(i+1)}(\lambda) = \mathcal{A}(\mathbf{w}^{(i)}(\lambda), \lambda) \\ \mathbf{w}^{(k)} \rightarrow \hat{\mathbf{w}}(\lambda) \end{cases}$

$$\lambda^{(n+1)} = \operatorname{Proj}_{\Lambda}(\lambda^{(n)} - \mu \nabla \mathcal{U}^{(k)}(\lambda^{(n)}))$$

- Convergence de  $\{\mathcal{U}^{(k)}\}_k$  vers  $\mathcal{U}$ ? ( $\inf$  ✓,  $\operatorname{argmin}$  ✓, points stationnaires, ...)
- Convergence de  $\{\nabla \mathcal{U}^{(k)}(\lambda)\}_k$  vers un sous-gradient de  $\partial \mathcal{U}(\lambda)$ ?
- Convergence de  $\{\lambda^{(n)}\}_n$  vers un point critique de  $\mathcal{U}$ ?
- Optimisation stochastique
- Calcul efficace d'hypergradient  $\nabla \mathcal{U}^{(k)}(\lambda^{(n)})$

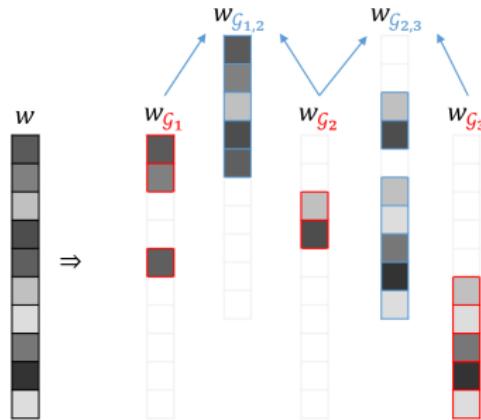
## 2. Découverte de structure dans les données

Extension du formalisme bi-niveaux

⇒ apprendre des structures de graphe  $\{\mathcal{G}_I, \mathcal{G}_{I,I'}\}$  sur les paramètres à estimer

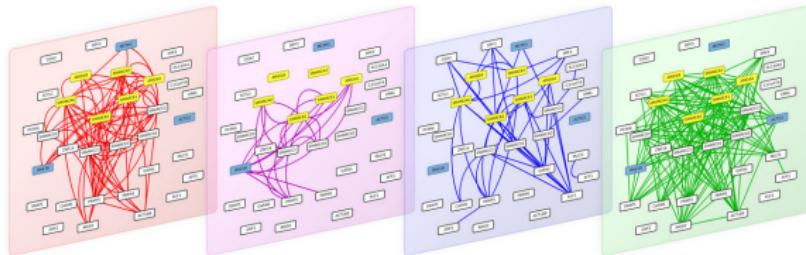
$$\min_{\{\mathcal{G}_I, \mathcal{G}_{I,I'}\}} \mathcal{E}(\hat{w}(\{\mathcal{G}_I, \mathcal{G}_{I,I'}\}))$$

$$\text{s.c. } \hat{w}(\{\mathcal{G}_I, \mathcal{G}_{I,I'}\}) = \operatorname{argmin}_{w \in \mathcal{W}} \ell_{\text{trn}}(w) + \sum_{I=1}^L \gamma_I \Omega_I(w_{\mathcal{G}_I}) + \sum_{I \neq I'} \gamma_{I,I'} \Omega_{I,I'}(w_{\mathcal{G}_{I,I'}})$$



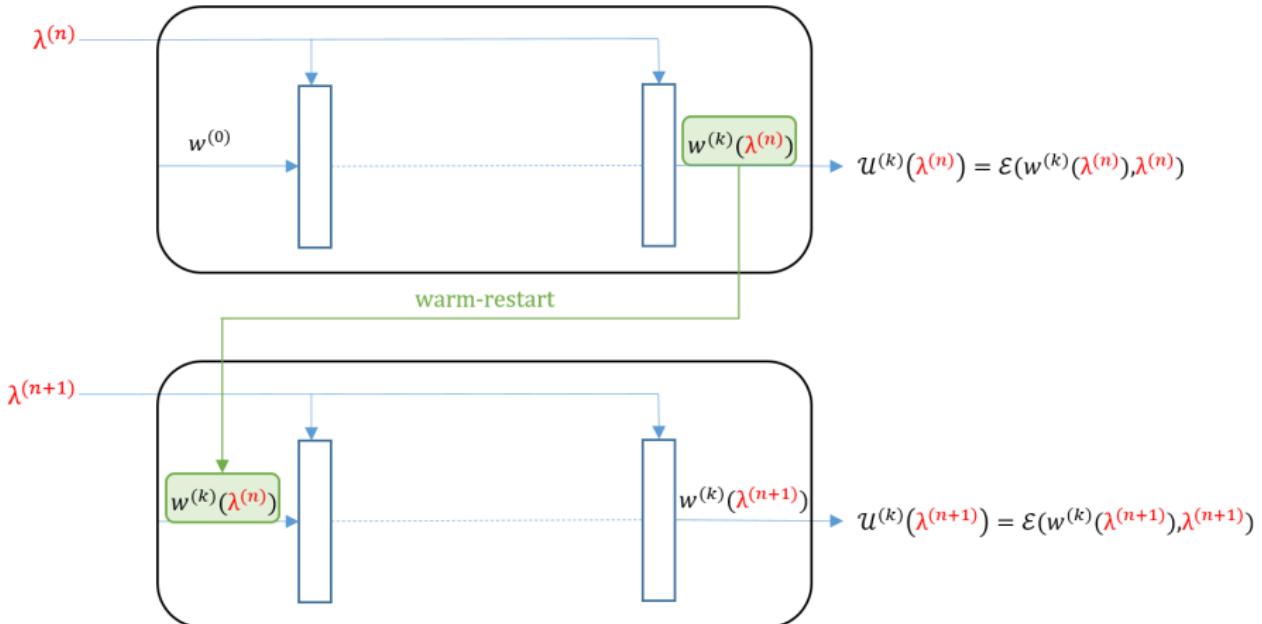
### 3. Application

**Application :** identification de communautés dans les réseaux biologiques



Communautés de protéines associées au syndrome de Coffin-Siris [Didier et al., 2015]

### 3. Impact du « warm-restart »



## 4. Optimisation bi-niveaux : impact du jeu de données

$$\min_{\lambda \in \Lambda} \mathbb{E}_{x \sim val} [\mathcal{E}(\hat{w}(\lambda), \lambda; \textcolor{red}{x})] \quad \text{s.c.} \quad \hat{w}(\lambda) \in \operatorname{argmin}_{w \in \mathcal{W}} \mathbb{E}_{x \sim trn} [\mathcal{L}(w, \lambda; \textcolor{red}{x})]$$

Garanties numériques pour l'optimisation stochastique

Extension à la validation croisée (k-fold, Monte-Carlo)

Développement de bornes statistiques

Merci pour votre attention

## Problème exact

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}(\lambda) \triangleq \mathcal{E}(\hat{w}(\lambda), \lambda) \right\}$$

s.c.  $\hat{w}(\lambda) = \operatorname{argmin}_{w \in \mathcal{W}} \mathcal{L}(w; \lambda)$

$\hat{w}(\lambda)$  sans expression explicite

## Problème approché

$$\min_{\lambda \in \Lambda} \left\{ \mathcal{U}^{(k)}(\lambda) \triangleq \mathcal{E}(w^{(k)}(\lambda), \lambda) \right\}$$

$w^{(0)}(\lambda)$  choisi arbitrairement  
 for  $i = 0, 1, \dots, k-1$   
 s.c.  $\begin{cases} w^{(i+1)}(\lambda) = \mathcal{A}(w^{(i)}(\lambda), \lambda) \\ w^{(k)}(\lambda) \rightarrow \hat{w}(\lambda) \end{cases}$

$\mathcal{U}^{(k)}$  lisse si  $\mathcal{A}$  lisse

$$\lambda^{(n+1)} = \operatorname{Proj}_{\Lambda}(\lambda^{(n)} - \mu \nabla \mathcal{U}^{(k)}(\lambda^{(n)}))$$

- $\inf_{\lambda \in \Lambda} \mathcal{U}^{(k)}(\lambda) \xrightarrow{k \rightarrow +\infty} \inf_{\lambda \in \Lambda} \mathcal{U}(\lambda) ?$
- $\operatorname{argmin}_{\lambda \in \Lambda} \mathcal{U}^{(k)}(\lambda) \xrightarrow{k \rightarrow +\infty} \operatorname{argmin}_{\lambda \in \Lambda} \mathcal{U}(\lambda) ?$
- $\lim_{k \rightarrow \infty} \nabla \mathcal{U}^{(k)}(\lambda) \in \partial \mathcal{U}(\lambda) ?$

- Impact de "warm-restart" sur  $w^{(0)}$  ?
- Calcul efficace de  $\nabla \mathcal{U}^{(k)}$  ?
- $\lim_{n \rightarrow \infty} \lambda^{(n)} \in \partial \mathcal{U}^{-1}(0) ?$

$$\min_{\lambda \in \Lambda} \mathbb{E}_{x \sim val} [\mathcal{E}(\hat{w}(\lambda), \lambda; \textcolor{red}{x})] \quad \text{s.c.} \quad \hat{w}(\lambda) \in \operatorname{argmin}_{w \in \mathcal{W}} \mathbb{E}_{x \sim trn} [\mathcal{L}(w, \lambda; \textcolor{red}{x})]$$

Développement de bornes statistiques

Garanties numériques pour l'optimisation stochastique

Extension à la validation croisée (k-fold, Monte-Carlo)

**Setting :**  $T$  linear regression tasks arranged in  $L$  groups of related tasks  $\{\mathcal{G}_1, \dots, \mathcal{G}_L\}$

$$\begin{array}{lllll} \text{find } w_1 & \text{find } w_2 & \text{find } w_3 & \text{find } w_4 & \text{find } w_5 \\ y_1 \approx X_1 w_1 & y_2 \approx X_2 w_2 & y_3 \approx X_3 w_3 & y_4 \approx X_4 w_4 & y_5 \approx X_5 w_5 \end{array}$$

$$W_{\mathcal{G}_1} = [w_1 \ w_2] \\ \text{low-rank}$$

$$W_{\mathcal{G}_2} = [w_3 \ w_4 \ w_5] \\ \text{low-rank}$$

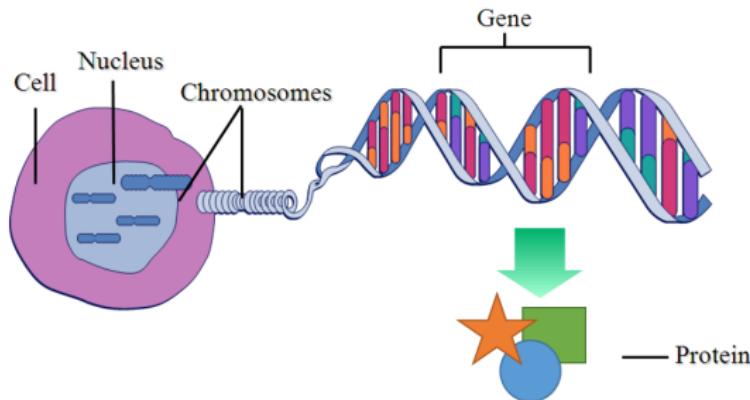
$$\hat{W} \in \underset{W=[w_1 \cdots w_T]}{\operatorname{argmin}} \sum_{t=1}^T \frac{1}{2} \|y_t - X_t w_t\|^2 + \lambda \sum_{l=1}^L \|W_{\mathcal{G}_l}\|_{\text{tr}}$$

**Issue :** In practice we don't know how tasks are related  $\rightarrow$  need to estimate  $\{\mathcal{G}_1, \dots, \mathcal{G}_L\}$

# Motivation : genes expression analysis

**Goal :** Predict the function of proteins from regulatory patterns

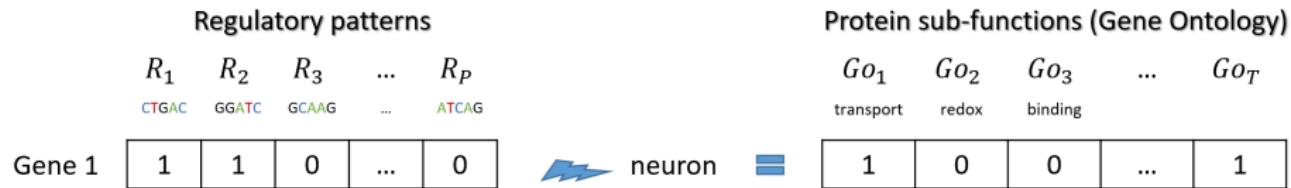
*Collaboration with Giorgio Valentini (Universita degli Studi di Milano)*



Gene = long sequence with regulatory patterns (dictate gene expression)

Proteins perform various functions (transport, redox, binding ...)

## Motivation : genes expression analysis



Sequences  $R_1$  and  $R_2$  are present in Gene 1

Gene 1 produces neurons

Neurons perform transport of electrons

## Motivation : genes expression analysis

Regulatory patterns						Protein sub-functions (Gene Ontology)					
	$R_1$	$R_2$	$R_3$	...	$R_P$		$Go_1$	$Go_2$	$Go_3$	...	$Go_T$
	CTGAC	GGATC	GCAAG	...	ATCAG						
Gene 1	1	1	0	...	0		transport	redox	binding		
Gene 2	1	1	1	...	0		neuron	=		1	
							hormone	=		0	

Sequences  $R_1$ ,  $R_2$  and  $R_3$  are present in Gene 2

Gene 2 produces hormones

Hormones perform transport of particles and reduction–oxidation

# Motivation : genes expression analysis

Regulatory patterns

$R_1 \quad R_2 \quad R_3 \quad \dots \quad R_P$

CTGAC    GGATC    GCAAG    ...    ATCAG

	$R_1$	$R_2$	$R_3$	...	$R_P$
Gene 1	1	1	0	...	0
Gene 2	1	1	1	...	0
Gene 3	0	0	0	...	1

Protein sub-functions (Gene Ontology)

$Go_1 \quad Go_2 \quad Go_3 \quad \dots \quad Go_T$

transport    redox    binding

- neuron
- hormone
- antibody

	$Go_1$	$Go_2$	$Go_3$	...	$Go_T$
	1	0	0	...	1
	1	1	0	...	0
	0	0	1	...	0

Sequence  $R_P$  is present in Gene 3

Gene 3 produces antibodies

Antibodies perform binding of particles

# Motivation : genes expression analysis

Regulatory patterns

$R_1 \quad R_2 \quad R_3 \quad \dots \quad R_P$

CTGAC    GGATC    GCAAG    ...    ATCAG

	$R_1$	$R_2$	$R_3$	...	$R_P$
Gene 1	1	1	0	...	0
Gene 2	1	1	1	...	0
Gene 3	0	0	0	...	1
⋮	⋮	⋮	⋮	⋮	⋮
Gene N	1	1	0	...	1

$$X \in \{0,1\}^{N \times P}$$

Protein sub-functions (Gene Ontology)

$Go_1 \quad Go_2 \quad Go_3 \quad \dots \quad Go_T$

transport    redox    binding

- neuron
- hormone
- antibody

	$Go_1$	$Go_2$	$Go_3$	...	$Go_T$
Gene 1	1	0	0	...	1
Gene 2	1	1	0	...	0
Gene 3	0	0	1	...	0
⋮	⋮	⋮	⋮	⋮	⋮
Gene N	0	0	1	...	1

$$y = [y_1 \cdots y_T] \in \{0,1\}^{N \times T}$$

# Motivation : genes expression analysis

Regulatory patterns

$R_1 \quad R_2 \quad R_3 \quad \dots \quad R_P$

CTGAC    GGATC    GCAAG    ...    ATCAG

	$R_1$	$R_2$	$R_3$	...	$R_P$
Gene 1	1	1	0	...	0
Gene 2	1	1	1	...	0
Gene 3	0	0	0	...	1
⋮	⋮	⋮	⋮	⋮	⋮
Gene N	1	1	0	...	1

$$X \in \{0,1\}^{N \times P}$$

Protein sub-functions (Gene Ontology)

$Go_1 \quad Go_2 \quad Go_3 \quad \dots \quad Go_T$

transport    redox    binding

- neuron
- hormone
- antibody

	$Go_1$	$Go_2$	$Go_3$	...	$Go_T$
Gene 1	1	0	0	...	1
Gene 2	1	1	0	...	0
Gene 3	0	0	1	...	0
⋮	⋮	⋮	⋮	⋮	⋮
Gene N	0	0	1	...	1

$$y = [y_1 \cdots y_T] \in \{0,1\}^{N \times T}$$

**Goal 1 :** Predict each  $y_t$  from  $X$

# Motivation : genes expression analysis

Regulatory patterns						Protein sub-functions (Gene Ontology)					
	$R_1$	$R_2$	$R_3$	$\dots$	$R_P$		$Go_1$	$Go_2$	$Go_3$	$\dots$	$Go_T$
	CTGAC	GGATC	GCAAG	$\dots$	ATCAG		transport	redox	binding	$\dots$	$\dots$
Gene 1				$\dots$						$\dots$	
Gene 2				$\dots$						$\dots$	
Gene 3				$\dots$						$\dots$	
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$					$\vdots$	$\vdots$
Gene N				$\dots$						$\dots$	

$X \in \mathbb{R}^{N \times P}$

$y = [y_1 \dots y_T] \in \mathbb{R}^{N \times T}$

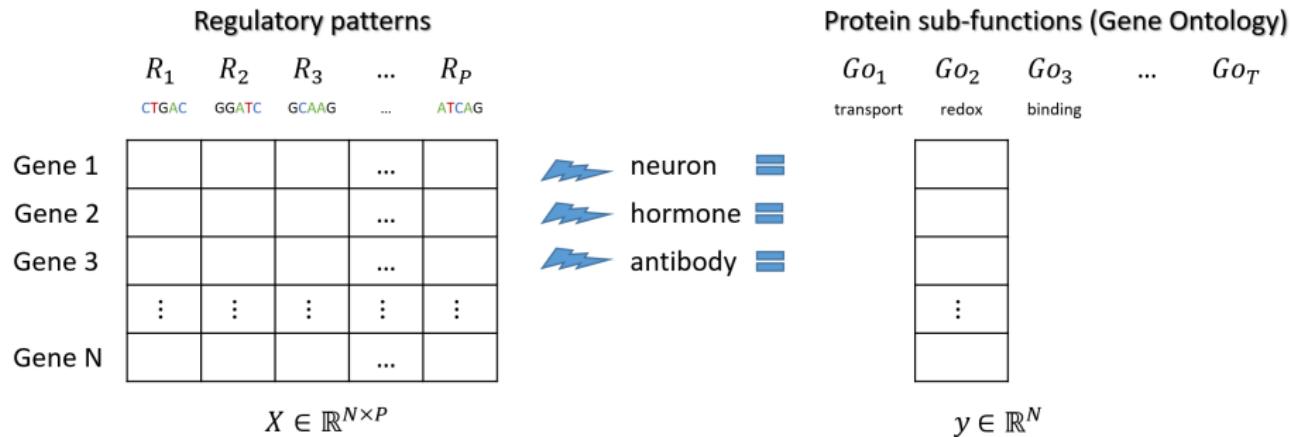
Legend:

- neuron
- hormone
- antibody

**Goal 1 :** Predict each  $y_t$  from  $X$

→ to generalize :  $X$  and  $y$  are not only made of 0's and 1's

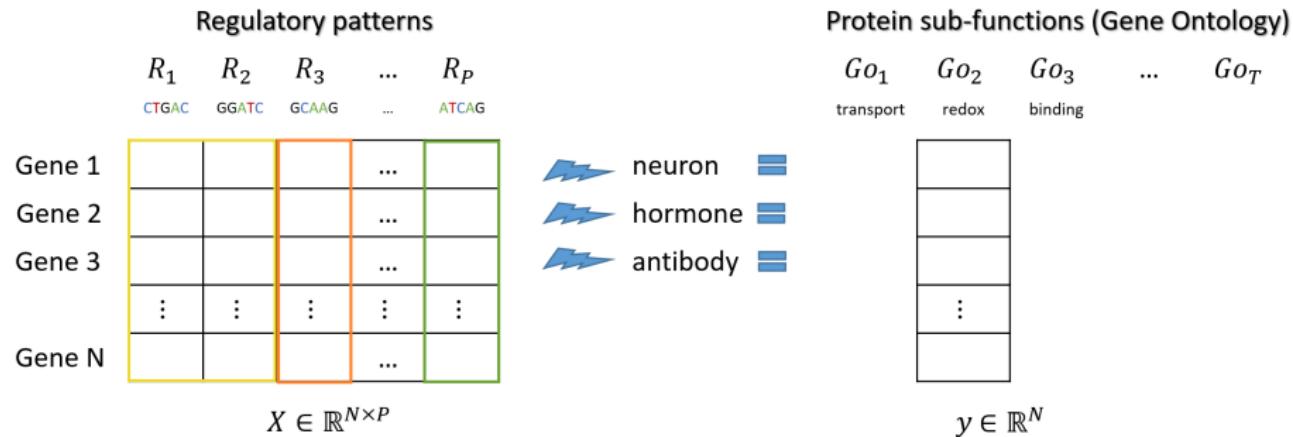
## Motivation : genes expression analysis



**Goal 1 :** Predict  $y$  from  $X$

- to generalize :  $X$  and  $y$  are not only made of 0's and 1's
- to simplify : we first assume that  $T = 1$  and omit the index  $t$

## Motivation : genes expression analysis



**Goal 1 :** Predict  $y$  from  $X$

- to generalize :  $X$  and  $y$  are not only made of 0's and 1's
- to simplify : we first assume that  $T = 1$  and omit the index  $t$

**Goal 2 :** Discover if there exist some groups in  $X$

ex :  $R_1$  and  $R_2$  are both equally relevant to predict  $y$

## Assumptions

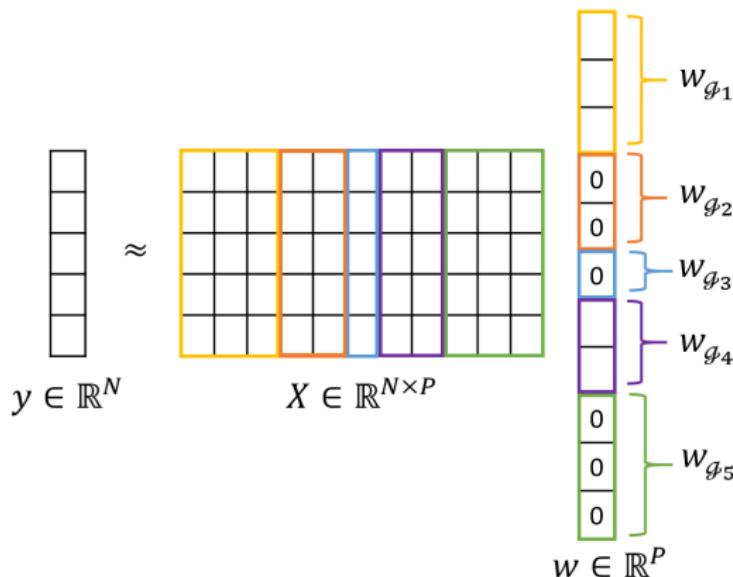
### Model the observations

⇒ linear model + Gaussian distribution : there exists  $w$  such that  $y \sim \mathcal{N}(Xw, \sigma)$

### Model the group structure

few groups of features in  $X$  are relevant to predict  $y$

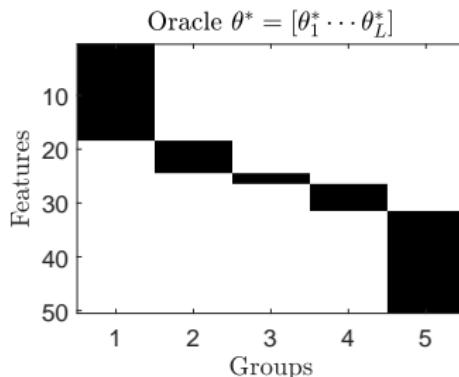
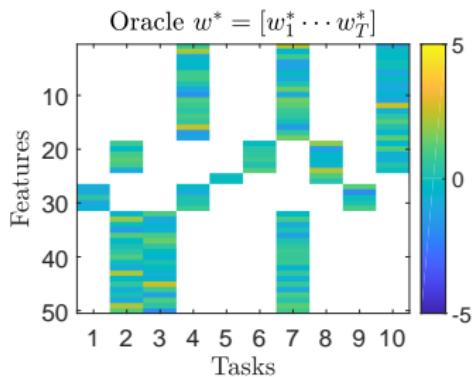
⇒ group sparsity : some groups of variables in  $w$  are zero while others are non-zero



## Group Lasso

**Configuration :**  $T$  problèmes Group Lasso partageant la même structure de groupe

$$(\forall t \in \{1, \dots, T\}) \quad \hat{w}_t(\theta) \in \operatorname{argmin}_{w_t \in \mathbb{R}^P} \frac{1}{2} \|y_t - X_t w_t\|^2 + \lambda \sum_{l=1}^L \|w_t \odot \theta_l\|_2,$$



**But :** Estimer la structure de groupe optimale  $\theta^*$

# Algorithme lisse

## Problème primal

$$\underset{w \in \mathbb{R}^P}{\text{minimize}} \left\{ \ell(w, \theta) := \underbrace{\frac{1}{2} \|y - Xw\|_2^2 + \frac{\epsilon}{2} \|w\|_2^2}_{f(w) \text{ différentiable}} + \underbrace{\lambda \sum_{l=1}^L \|\theta_l \odot w\|_2}_{g(A_\theta w) \text{ non différentiable}} \right\}$$

où  $A_\theta : w \in \mathbb{R}^P \mapsto (\theta_1 \odot w, \dots, \theta_L \odot w) \in \mathbb{R}^{P \times L}$

## Algorithme forward-backward

$$\left\{ \begin{array}{l} \text{for } q = 0, 1, \dots, Q-1 \\ \quad \left[ \begin{array}{l} w^{(q+1)}(\theta) = \text{prox}_{\beta g \circ A_\theta}(w^{(q)}(\theta) - \beta \nabla f(w^{(q)}(\theta))) \end{array} \right] \end{array} \right.$$

**Difficulté :**  $\text{prox}_{g \circ A_\theta}(v) = \underset{w \in \mathbb{R}^P}{\text{argmin}} g(A_\theta w) + \frac{1}{2} \|w - v\|_2^2$  sans expression explicite

## Problème primal

$$\underset{w \in \mathbb{R}^P}{\text{minimize}} \left\{ \ell(w, \theta) := \underbrace{\frac{1}{2} \|y - Xw\|_2^2 + \frac{\epsilon}{2} \|w\|_2^2}_{f(w) \text{ différentiable}} + \underbrace{\lambda \sum_{l=1}^L \|\theta_l \odot w\|_2}_{g(A_\theta w) \text{ non différentiable}} \right\}$$

où  $A_\theta : w \in \mathbb{R}^P \mapsto (\theta_1 \odot w, \dots, \theta_L \odot w) \in \mathbb{R}^{P \times L}$

## Algorithme forward-backward

$$\left\{ \begin{array}{l} \text{for } q = 0, 1, \dots, Q-1 \\ \quad \left[ \begin{array}{l} w^{(q+1)}(\theta) = \text{prox}_{\beta g \circ A_\theta}(w^{(q)}(\theta) - \beta \nabla f(w^{(q)}(\theta))) \end{array} \right] \end{array} \right.$$

**Difficulté :**  $\text{prox}_{g \circ A_\theta}(v) = \underset{w \in \mathbb{R}^P}{\text{argmin}} g(A_\theta w) + \frac{1}{2} \|w - v\|_2^2$  sans expression explicite

## Problème primal

$$\underset{w \in \mathbb{R}^P}{\text{minimize}} \left\{ \ell(w, \theta) := \underbrace{\frac{1}{2} \|y - Xw\|_2^2 + \frac{\epsilon}{2} \|w\|_2^2}_{f(w)} + \lambda \underbrace{\sum_{l=1}^L \|\theta_l \odot w\|_2}_{g(\mathbf{A}_\theta w)} \right\}$$

$\text{prox}_{g \circ A_\theta}$  n'a pas de forme explicite  $\Rightarrow$  considère le problème dual pour déplacer  $A_\theta$

## Problème dual

$$\underset{u \in \mathbb{R}^{P \times L}}{\text{minimize}} \left\{ \tilde{\ell}(u, \theta) := f^*(-\mathbf{A}_\theta^* u) + g^*(u) \right\}$$

où

$A_\theta^*$  est l'adjoint de  $A_\theta$

$f^*$  et  $g^*$  sont les conjuguées de Fenchel de  $f$  et  $g$ , respectivement.

## Problème primal

$$\underset{w \in \mathbb{R}^P}{\text{minimize}} \left\{ \ell(w, \theta) := \underbrace{\frac{1}{2} \|y - Xw\|_2^2 + \frac{\epsilon}{2} \|w\|_2^2}_{f(w)} + \lambda \underbrace{\sum_{l=1}^L \|\theta_l \odot w\|_2}_{g(\mathbf{A}_\theta w)} \right\}$$

prox <sub>$g \circ A_\theta$</sub>  n'a pas de forme explicite  $\Rightarrow$  considère le problème dual pour déplacer  $A_\theta$

## Problème dual

$$\underset{u \in \mathbb{R}^{P \times L}}{\text{minimize}} \left\{ \tilde{\ell}(u, \theta) := f^*(-\mathbf{A}_\theta^* u) + g^*(u) \right\}$$

où

$A_\theta^*$  est l'adjoint de  $A_\theta$

$f^*$  et  $g^*$  sont les conjuguées de Fenchel de  $f$  et  $g$ , respectivement.

## Algorithme lisse

### Problème dual

$$\underset{u \in \mathbb{R}^{P \times L}}{\text{minimize}} \left\{ \tilde{\ell}(u, \theta) := f^*(-A_\theta^* u) + g^*(u) \right\}$$

où  $f^*$  est différentiable et  $g^* = \iota_{\mathcal{B}(\lambda \mathbb{1}_P)^L}$  est non différentiable

### Algorithme dual forward-backward

$$\begin{cases} \text{for } q = 0, 1, \dots, Q-1 \\ \quad \left| \begin{array}{l} u^{(q+1)}(\theta) = \text{prox}_{\beta g^*}(u^{(q)}(\theta) + \beta A_\theta \nabla f^*(-A_\theta^* u^{(q)}(\theta))) \\ w^{(Q)}(\theta) = \nabla f^*(-A_\theta u^{(Q)}(\theta)). \end{array} \right. \end{cases}$$

Cependant  $\text{prox}_{\beta g^*} = \text{Proj}_{\mathcal{B}(\lambda)^L}$  est non différentiable

### Algorithme dual forward-backward avec distances de Bregman

$$\begin{cases} \text{for } q = 0, 1, \dots, Q-1 \\ \quad \left| \begin{array}{l} u^{(q+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi(\nabla \Phi(u^{(q)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^* u^{(q)}(\theta))) \\ w^{(Q)}(\theta) = \nabla f^*(-A_\theta u^{(Q)}(\theta)). \end{array} \right. \end{cases}$$

où l'opérateur proximal de Bregman avec la fonction de Legendre  $\Phi$  :

$$\text{prox}_h^\Phi(v) = \underset{u \in \mathbb{R}^{P \times L}}{\text{argmin}} h(u) + \Phi(u) - \langle u, v \rangle$$

## Problème dual

$$\underset{u \in \mathbb{R}^{P \times L}}{\text{minimize}} \left\{ \tilde{\ell}(u, \theta) := f^*(-A_\theta^* u) + g^*(u) \right\}$$

où  $f^*$  est différentiable et  $g^* = \iota_{\mathcal{B}(\lambda \mathbb{1}_P)^L}$  est non différentiable

## Algorithme dual forward-backward avec distances de Bregman

$$\begin{cases} \text{for } q = 0, 1, \dots, Q-1 \\ \quad \left| \begin{array}{l} u^{(q+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi(\nabla \Phi(u^{(q)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^* u^{(q)}(\theta))) \\ w^{(Q)}(\theta) = \nabla f^*(-A_\theta u^{(Q)}(\theta)). \end{array} \right. \end{cases}$$

où l'opérateur proximal de Bregman avec la fonction de Legendre  $\Phi$  :

$$\text{prox}_h^\Phi(v) = \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} h(u) + \Phi(u) - \langle u, v \rangle$$

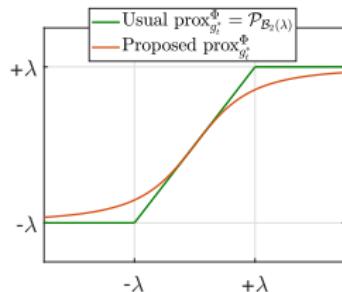
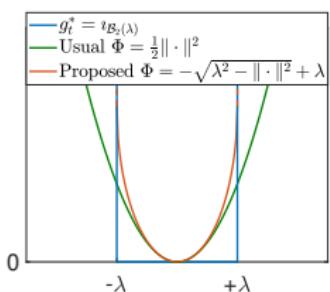
## Choix de $\Phi$ pour lisser les mises à jour

### Algorithme dual forward-backward avec distances de Bregman

$$\left\{ \begin{array}{l} \text{for } q = 0, 1, \dots, Q-1 \\ \quad \left| \begin{array}{l} u^{(q+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi(\nabla \Phi(u^{(q)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^* u^{(q)}(\theta))) \\ w^{(Q)}(\theta) = \nabla f^*(-A_\theta u^{(Q)}(\theta)). \end{array} \right. \end{array} \right.$$

où

$$\text{prox}_{\beta g^*}^\Phi(v) = \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \beta g^*(u) + \Phi(u) - \langle u, v \rangle$$



## Choix de $\Phi$ pour lisser les mises à jour

### Algorithme dual forward-backward avec distances de Bregman

$$\left\{ \begin{array}{l} \text{for } q = 0, 1, \dots, Q-1 \\ \quad \left| \begin{array}{l} u^{(q+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi(\nabla \Phi(u^{(q)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^* u^{(q)}(\theta))) \\ w^{(Q)}(\theta) = \nabla f^*(-A_\theta u^{(Q)}(\theta)). \end{array} \right. \end{array} \right.$$

où

$$\text{prox}_{\beta g^*}^\Phi(v) = \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \sum_{l=1}^L \iota_{\mathcal{B}(\lambda)}(u_l) + \Phi(u) - \langle u, v \rangle$$

## Choix de $\Phi$ pour lisser les mises à jour

### Algorithme dual forward-backward avec distances de Bregman

$$\left\{ \begin{array}{l} \text{for } q = 0, 1, \dots, Q-1 \\ \quad \left| \begin{array}{l} u^{(q+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi(\nabla \Phi(u^{(q)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^* u^{(q)}(\theta))) \\ w^{(Q)}(\theta) = \nabla f^*(-A_\theta u^{(Q)}(\theta)). \end{array} \right. \end{array} \right.$$

où

$$\text{prox}_{\beta g^*}^\Phi(v) = \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \sum_{l=1}^L \left( \iota_{\mathcal{B}(\lambda)}(u_l) - \sqrt{\lambda^2 - \|u_l\|^2} - \langle u_l, v_l \rangle \right)$$

$$\text{pour } \Phi(u) = \sum_{l=1}^L -\sqrt{\lambda^2 - \|u_l\|^2}$$

## Choix de $\Phi$ pour lisser les mises à jour

### Algorithme dual forward-backward avec distances de Bregman

```
{ for q = 0, 1, ..., Q - 1
    |   u(q+1)(θ) = proxΦβg*(∇Φ(u(q)(θ)) + βAθ∇f*(-Aθ*u(q)(θ)))
    |   w(Q)(θ) = ∇f*(-Aθu(Q)(θ)).
```

où

$$\text{prox}_{\beta g^*}^\Phi(v) = \left( \frac{\lambda v_l}{\sqrt{1 + \|v\|_2^2}} \right)_{l=1, \dots, L}$$

for  $\Phi(u) = \sum_{l=1}^L -\sqrt{\lambda^2 - \|u_l\|^2}$  ⇒ mise à jour différentiable

[

plain,c]

# Algorithmic Solution

## Optimization problem

$$\underset{w \in \mathbb{R}^P}{\text{minimize}} \quad \underbrace{\frac{1}{2} \|y - Xw\|_2^2}_{f(w) \text{ differentiable}} + \underbrace{\lambda \sum_{l=1}^L \|\theta_l \odot w\|_2}_{g(A_\theta w) \text{ non differentiable}}$$

where  $A_\theta : w \in \mathbb{R}^P \mapsto (\theta_1 \odot w, \dots, \theta_L \odot w) \in \mathbb{R}^{P \times L}$

Forward-backward algorithm [Combettes and Wajs (2005)]

$$\begin{cases} \text{for } i = 0, \dots, k-1 \\ \quad \left[ \begin{array}{l} w^{(i+1)}(\theta) = \text{prox}_{\beta g \circ A_\theta}(w^{(i)}(\theta) - \beta \nabla f(w^{(i)}(\theta))) \end{array} \right] \end{cases}$$

*proximity operator (prox) ?*  $\longrightarrow$  see next slide

## Optimization problem

$$\underset{w \in \mathbb{R}^P}{\text{minimize}} \quad \underbrace{\frac{1}{2} \|y - Xw\|_2^2}_{f(w) \text{ differentiable}} + \underbrace{\lambda \sum_{l=1}^L \|\theta_l \odot w\|_2}_{g(A_\theta w) \text{ non differentiable}}$$

where  $A_\theta : w \in \mathbb{R}^P \mapsto (\theta_1 \odot w, \dots, \theta_L \odot w) \in \mathbb{R}^{P \times L}$

**Forward-backward algorithm** [Combettes and Wajs (2005)]

$$\left\{ \begin{array}{l} \text{for } i = 0, \dots, k-1 \\ \quad \left[ \begin{array}{l} w^{(i+1)}(\theta) = \text{prox}_{\beta g \circ A_\theta}(w^{(i)}(\theta) - \beta \nabla f(w^{(i)}(\theta))) \end{array} \right] \end{array} \right.$$

*proximity operator (prox) ? —> see next slide*

## Proximity operator

In the 1960s, [Moreau (1962)] proposed an extension of the notion of projection operator to any convex function  $h$ , leading to the so-called proximity operator

$$\begin{aligned}\text{Proj}_{\mathcal{C}}(v) &= \underset{w \in \mathcal{C}}{\operatorname{argmin}} \frac{1}{2} \|w - v\|_2^2 \\ &= \underset{w \in \mathbb{R}^P}{\operatorname{argmin}} \iota_{\mathcal{C}}(w) + \frac{1}{2} \|w - v\|_2^2 \quad \text{where} \quad \iota_{\mathcal{C}}(w) = \begin{cases} 0 & \text{if } w \in \mathcal{C} \\ +\infty & \text{otherwise} \end{cases}\end{aligned}$$

$$\text{prox}_h(v) = \underset{w \in \mathbb{R}^P}{\operatorname{argmin}} h(w) + \frac{1}{2} \|w - v\|_2^2$$

# Group Lasso solver $\mathcal{A}$

## Optimization problem

$$\underset{w \in \mathbb{R}^P}{\text{minimize}} \quad \underbrace{\frac{1}{2} \|y - Xw\|_2^2}_{f(w) \text{ differentiable}} + \underbrace{\lambda \sum_{l=1}^L \|\theta_l \odot w\|_2}_{g(A_\theta w) \text{ non differentiable}}$$

where  $A_\theta : w \in \mathbb{R}^P \mapsto (\theta_1 \odot w, \dots, \theta_L \odot w) \in \mathbb{R}^{P \times L}$

## Forward-backward algorithm [Combettes and Wajs (2005)]

$$\left\{ \begin{array}{l} \text{for } i = 0, \dots, k-1 \\ \quad \boxed{w^{(i+1)}(\theta) = \text{prox}_{\beta g \circ A_\theta}(w^{(i)}(\theta) - \beta \nabla f(w^{(i)}(\theta)))} \end{array} \right.$$

*generalization of projected gradient descent  
projection  $\rightarrow$  proximity operator*

$$\text{prox}_{\beta g \circ A_\theta}(v) = \underset{w \in \mathbb{R}^P}{\operatorname{argmin}} \beta g(A_\theta w) + \frac{1}{2} \|w - v\|_2^2 \quad \triangleleft \text{ without closed form}$$

## Group Lasso solver $\mathcal{A}$

### Optimization problem

$$\underset{w \in \mathbb{R}^P}{\text{minimize}} \quad \underbrace{\frac{1}{2} \|y - Xw\|_2^2}_{f(w) \text{ differentiable}} + \underbrace{\lambda \sum_{l=1}^L \|\theta_l \odot w\|_2}_{g(A_\theta w) \text{ non differentiable}}$$

where  $A_\theta : w \in \mathbb{R}^P \mapsto (\theta_1 \odot w, \dots, \theta_L \odot w) \in \mathbb{R}^{P \times L}$

### Forward-backward algorithm [Combettes and Wajs (2005)]

$$\left\{ \begin{array}{l} \text{for } i = 0, \dots, k-1 \\ \quad \left[ \begin{array}{l} w^{(i+1)}(\theta) = \text{prox}_{\beta g \circ A_\theta}(w^{(i)}(\theta) - \beta \nabla f(w^{(i)}(\theta))) \end{array} \right] \end{array} \right.$$

*generalization of projected gradient descent  
projection  $\rightarrow$  proximity operator*

$$\text{prox}_{\beta g \circ A_\theta}(v) = \underset{w \in \mathbb{R}^P}{\operatorname{argmin}} \beta g(A_\theta w) + \frac{1}{2} \|w - v\|_2^2 \quad \triangleq \text{without closed form}$$

## Duality in convex optimization

The ideas of duality and transforms are ubiquitous in mathematics

- Harmonics analysis → Fourier transform
- Convex analysis → Fenchel conjugate :  $h^*(x) = \sup_w \langle w, x \rangle - h(w)$

[Rockafellar (1970)]

**Example :**  $h: w \mapsto \|w\|_2 \implies h^*: x \mapsto \iota_{B(1)}(x) = \begin{cases} 0 & \text{if } \|x\|_2 \leq 1 \\ +\infty & \text{otherwise} \end{cases}$

## Duality in convex optimization

The ideas of duality and transforms are ubiquitous in mathematics

- Harmonics analysis → Fourier transform
- Convex analysis → Fenchel conjugate :  $h^*(x) = \sup_w \langle w, x \rangle - h(w)$

[Rockafellar (1970)]

**Example :**  $h: w \mapsto \lambda \|w\|_2 \quad \Rightarrow \quad h^*: x \mapsto i_{B(\lambda)}(x) = \begin{cases} 0 & \text{if } \|x\|_2 \leq \lambda \\ +\infty & \text{otherwise} \end{cases}$

# Duality in convex optimization

The ideas of duality and transforms are ubiquitous in mathematics

- Harmonics analysis → Fourier transform
- Convex analysis → Fenchel conjugate :  $h^*(x) = \sup_w \langle w, x \rangle - h(w)$

[Rockafellar (1970)]

## Primal problem

↔

## Dual problem

$$\underset{w \in \mathbb{R}^P}{\text{minimize}} f(w) + g(\mathbf{A}_\theta w)$$

$$\mathbf{A}_\theta : \mathbb{R}^P \rightarrow \mathbb{R}^{P \times L}$$

$$g(v_1 \dots v_L) = \sum_{l=1}^L \underbrace{\lambda \|v_l\|_2}_{\text{norm}}$$

$$\underset{u \in \mathbb{R}^{P \times L}}{\text{minimize}} f^*(-\mathbf{A}_\theta^\top u) + g^*(u)$$

$$\mathbf{A}_\theta^\top : \mathbb{R}^{P \times L} \rightarrow \mathbb{R}^P$$

$$g^*(u_1 \dots u_L) = \sum_{l=1}^L \underbrace{\iota_{\mathcal{B}(\lambda)}(u_l)}_{\substack{\text{indicator} \\ \text{dual norm ball}}}$$

## Link

$$w = \nabla f^*(-\mathbf{A}_\theta^\top u)$$

$\text{prox}_{g \circ \mathbf{A}_\theta}$  without closed form  $\Rightarrow$  solve dual problem to move  $\mathbf{A}_\theta$  in smooth part

# Group Lasso solver $\mathcal{A}$ : dual approach

## Dual problem

$$\underset{u \in \mathbb{R}^{P \times L}}{\text{minimize}} \quad \underbrace{f^*(-A_\theta^\top u)}_{\text{differentiable}} + \underbrace{g^*(u)}_{\text{non differentiable}}$$

## Dual forward-backward algorithm

$$\left\{ \begin{array}{l} \text{for } i = 0, \dots, k-1 \\ \quad \left| \begin{array}{l} u^{(i+1)}(\theta) = \text{prox}_{\beta g^*}(u^{(i)}(\theta) + \beta A_\theta \nabla f^*(-A_\theta^\top u^{(i)}(\theta))) \\ w^{(k)}(\theta) = \nabla f^*(-A_\theta^* u^{(k)}(\theta)) \quad (\text{link}) \end{array} \right. \end{array} \right.$$

## Group Lasso solver $\mathcal{A}$ : dual approach

### Dual problem

$$\underset{u \in \mathbb{R}^{P \times L}}{\text{minimize}} \quad \underbrace{f^*(-A_\theta^\top u)}_{\text{differentiable}} + \underbrace{g^*(u)}_{\text{non differentiable}}$$

### Dual forward-backward algorithm

$$\begin{cases} \text{for } i = 0, \dots, k-1 \\ \quad \left| \begin{array}{l} u^{(i+1)}(\theta) = \text{prox}_{\beta g^*}(u^{(i)}(\theta) + \beta A_\theta \nabla f^*(-A_\theta^\top u^{(i)}(\theta))) \\ w^{(k)}(\theta) = \nabla f^*(-A_\theta^\top u^{(k)}(\theta)) \quad (\text{link}) \end{array} \right. \end{cases}$$

where the proximal operator reads :

$$\begin{aligned} \text{prox}_{\beta g^*}(v) &= \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \beta g^*(u) + \frac{1}{2} \|u - v\|^2 \\ &= \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \beta \sum_{l=1}^L \iota_{\mathcal{B}(\lambda)}(u_l) + \frac{1}{2} \|u - v\|^2 \\ &= \text{Proj}_{\mathcal{B}(\lambda)^L}(v) \quad \triangleq \text{not differentiable} \end{aligned}$$

## Reminder : why differentiability is important

$\theta^{(0)}$  chosen arbitrarily

for  $n = 0, 1, \dots$

$w^{(0)}(\theta^{(n)})$  chosen arbitrarily  
for  $i = 0, \dots, k - 1$  (inner algorithm = dual forward-backward)  
 $w^{(i+1)}(\theta^{(n)}) = \mathcal{A}(w^{(i)}(\theta^{(n)}))$   
 $\theta^{(n+1)} = \text{Proj}_{\Theta}(\theta^{(n)} - \gamma \nabla \mathcal{U}^{(k)}(\theta^{(n)}))$  where  $\mathcal{U}^{(k)}(\theta^{(n)}) = \mathcal{E}(w^{(k)}(\theta^{(n)}))$

We want a **differentiable dual forward-backward algorithm**  
because it inside a bilevel algorithm !

### Dual forward-backward algorithm

$$\left\{ \begin{array}{l} \text{for } i = 0, \dots, k-1 \\ \quad \left| \begin{array}{l} u^{(i+1)}(\theta) = \text{prox}_{\beta g^*} \left( u^{(i)}(\theta) + \beta A_\theta \nabla f^*(-A_\theta^\top u^{(i)}(\theta)) \right) \\ w^{(k)}(\theta) = \nabla f^*(-A_\theta u^{(k)}(\theta)). \end{array} \right. \end{array} \right.$$

where

$$\begin{aligned} \text{prox}_{\beta g^*}(v) &= \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \beta g^*(u) + \frac{1}{2} \|u - v\|^2 \\ &= \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \beta g^*(u) + \frac{1}{2} \|u\|^2 - \langle u, v \rangle + \text{cst} \end{aligned}$$

### Dual forward-backward algorithm

$$\left\{ \begin{array}{l} \text{for } i = 0, \dots, k-1 \\ \quad \left| \begin{array}{l} u^{(i+1)}(\theta) = \text{prox}_{\beta g^*} \left( u^{(i)}(\theta) + \beta A_\theta \nabla f^*(-A_\theta^\top u^{(i)}(\theta)) \right) \\ w^{(k)}(\theta) = \nabla f^*(-A_\theta u^{(k)}(\theta)). \end{array} \right. \end{array} \right.$$

where

$$\text{prox}_{\beta g^*}(v) = \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \beta g^*(u) + \frac{1}{2} \|u\|^2 - \langle u, v \rangle + \text{cst}$$

## Group Lasso solver $\mathcal{A}$ : dual approach

Dual forward-backward algorithm **with Bregman distances** [Bauschke et al. (2016)]

$$\left\{ \begin{array}{l} \text{for } i = 0, \dots, k-1 \\ \quad \left| \begin{array}{l} u^{(i+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi(\nabla \Phi(u^{(i)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^\top u^{(i)}(\theta))) \\ w^{(k)}(\theta) = \nabla f^*(-A_\theta u^{(k)}(\theta)). \end{array} \right. \end{array} \right.$$

where the Bregman proximal operator associated to  $\Phi$  :

$$\text{prox}_{\beta g^*}^\Phi(v) = \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \beta g^*(u) + \Phi(u) - \langle u, v \rangle$$

## Choice of $\Phi$ to smooth the updates

### Dual forward-backward algorithm with Bregman distances

$$\left\{ \begin{array}{l} \text{for } i = 0, \dots, k-1 \\ \quad \left| \begin{array}{l} u^{(i+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi (\nabla \Phi(u^{(i)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^\top u^{(i)}(\theta))) \\ w^{(k)}(\theta) = \nabla f^*(-A_\theta u^{(k)}(\theta)). \end{array} \right. \end{array} \right.$$

where

$$\text{prox}_{\beta g^*}^\Phi(v) = \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \beta g^*(u) + \Phi(u) - \langle u, v \rangle$$

## Choice of $\Phi$ to smooth the updates

### Dual forward-backward algorithm with Bregman distances

$$\left\{ \begin{array}{l} \text{for } i = 0, \dots, k-1 \\ \quad \left| \begin{array}{l} u^{(i+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi(\nabla \Phi(u^{(i)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^\top u^{(i)}(\theta))) \\ w^{(k)}(\theta) = \nabla f^*(-A_\theta u^{(k)}(\theta)). \end{array} \right. \end{array} \right.$$

where

$$\text{prox}_{\beta g^*}^\Phi(v) = \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \sum_{l=1}^L \varphi_{B(\lambda)}(u_l) + \Phi(u) - \langle u, v \rangle$$

## Choice of $\Phi$ to smooth the updates

### Dual forward-backward algorithm with Bregman distances

$$\begin{cases} \text{for } i = 0, \dots, k-1 \\ \quad \left| \begin{array}{l} u^{(i+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi (\nabla \Phi(u^{(i)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^\top u^{(i)}(\theta))) \\ w^{(k)}(\theta) = \nabla f^*(-A_\theta u^{(k)}(\theta)). \end{array} \right. \end{cases}$$

where

$$\text{prox}_{\beta g^*}^\Phi(v) = \underset{u \in \mathbb{R}^P \times L}{\operatorname{argmin}} \sum_{l=1}^L (\varphi_{\mathcal{B}(\lambda)}(u_l) + \phi(u_l) - \langle u_l, v_l \rangle)$$

for  $\Phi(u) = \sum_{l=1}^L \phi(u_l)$

## Choice of $\Phi$ to smooth the updates

### Dual forward-backward algorithm with Bregman distances

$$\left\{ \begin{array}{l} \text{for } i = 0, \dots, k-1 \\ \quad \left| \begin{array}{l} u^{(i+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi(\nabla \Phi(u^{(i)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^\top u^{(i)}(\theta))) \\ w^{(k)}(\theta) = \nabla f^*(-A_\theta u^{(k)}(\theta)). \end{array} \right. \end{array} \right.$$

where

$$\text{prox}_{\beta g^*}^\Phi(v) = \underset{u \in \mathbb{R}^{P \times L}}{\operatorname{argmin}} \sum_{l=1}^L \left( \iota_{\mathcal{B}(\lambda)}(u_l) - \sqrt{\lambda^2 - \|u_l\|^2} - \langle u_l, v_l \rangle \right)$$

for  $\phi(u_l) = -\sqrt{\lambda^2 - \|u_l\|^2} \Rightarrow \text{dom } \phi = \mathcal{B}(\lambda)$   
 $\Rightarrow \iota_{\mathcal{B}(\lambda)}(u_l)$  always equal to 0 !

⚠️ trick for a differentiable algorithm

## Choice of $\Phi$ to smooth the updates

### Dual forward-backward algorithm with Bregman distances

$$\begin{cases} \text{for } i = 0, \dots, k-1 \\ \quad \left[ u^{(i+1)}(\theta) = \text{prox}_{\beta g^*}^\Phi (\nabla \Phi(u^{(i)}(\theta)) + \beta A_\theta \nabla f^*(-A_\theta^\top u^{(i)}(\theta))) \right. \\ \quad \left. w^{(k)}(\theta) = \nabla f^*(-A_\theta u^{(k)}(\theta)). \right. \end{cases}$$

where

$$\text{prox}_{\beta g^*}^\Phi(v) = \left( \frac{\lambda v_l}{\sqrt{1 + \|v\|_2^2}} \right)_{l=1, \dots, L}$$

