

# Bilevel Learning of the Group Lasso Structure

Jordan Frecon<sup>1</sup>, Saverio Salzo<sup>1</sup>, Massimiliano Pontil<sup>1,2</sup>

<sup>1</sup> CSML - Istituto Italiano di Tecnologia

<sup>2</sup> Dept of Computer Science - University College London

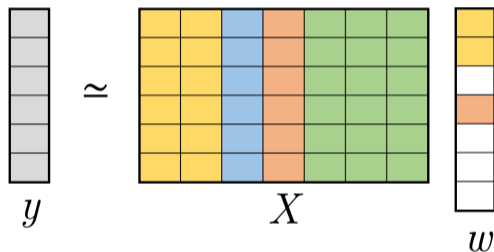


Thirty-second Conference on Neural Information Processing Systems, Montreal, Canada

# Linear Regression and Group Sparsity

**Problem:** Predict  $y \in \mathbb{R}^N$  from  $X \in \mathbb{R}^{N \times P}$

**Linear Regression:** Find  $w \in \mathbb{R}^P$  such that



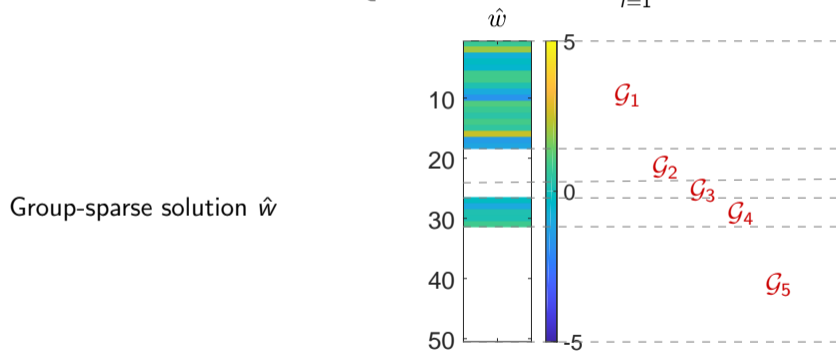
In many applications, **few groups are relevant** to predict  $y \Rightarrow$  **Group Sparse**  $w$

- Predict psychiatric disorder from activities in regions of the brain
- Predict protein functions from their molecular composition

# Group Lasso

Given  $\lambda > 0$  and a group-structure  $\{\mathcal{G}_1, \dots, \mathcal{G}_L\}$ , find

$$\hat{w} \in \operatorname{argmin}_{w \in \mathbb{R}^P} \frac{1}{2} \|y - Xw\|^2 + \lambda \sum_{l=1}^L \|w_{\mathcal{G}_l}\|_2,$$

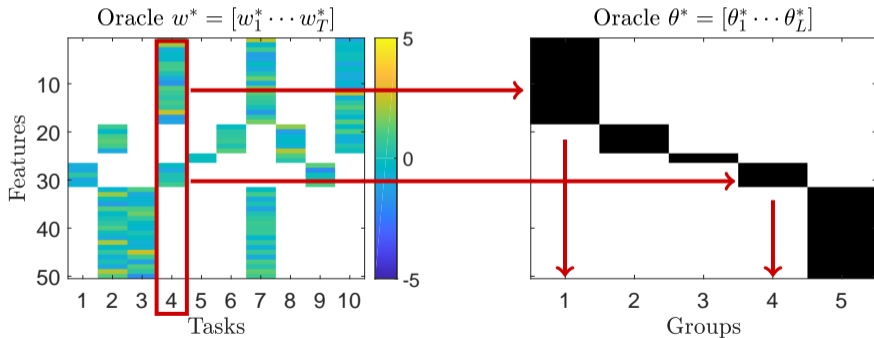


**Limitation:** The group-structure  $\{\mathcal{G}_1, \dots, \mathcal{G}_L\}$  may be unknown

# Setting

**Setting:**  $T$  Group Lasso problems with shared group-structure

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



**Goal:** Estimation of the optimal group-structure  $\theta^*$

# A Bilevel Programming Approach

## Upper-level Problem:

$$\underset{[\theta_1 \dots \theta_L] \in \Theta}{\text{minimize}} \mathcal{U}(\theta) := \sum_{t=1}^T \mathcal{E}_t(\hat{w}_t(\theta)) \quad (\text{e.g., validation error})$$

where  $\hat{w}(\theta) = [\hat{w}_1(\theta) \dots \hat{w}_T(\theta)]$  solves

## Lower-level Problem: ( $T$ Group Lasso problems)

$$\underset{w \in \mathbb{R}^{P \times T}}{\text{minimize}} \mathcal{L}(w, \theta) := \sum_{t=1}^T \left( \frac{1}{2} \|y_t - X_t w_t\|^2 + \lambda \sum_{l=1}^L \|\theta_l \odot w_t\|_2 \right)$$

## Difficulties:

- $\hat{w}(\theta)$  not available in closed form
- $\theta \mapsto \hat{w}(\theta)$  is nonsmooth  $[\Rightarrow \mathcal{U}$  is nonsmooth]

# Approximate Bilevel Problem

## Upper-level Problem:

$$\underset{[\theta_1 \dots \theta_L] \in \Theta}{\text{minimize}} \mathcal{U}_K(\theta) := \sum_{t=1}^T \mathcal{E}_t(w_t^{(K)}(\theta))$$

$$\text{where } w_t^{(K)}(\theta) \rightarrow \hat{w}_t(\theta)$$

## Dual Algorithm:

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

for  $k = 0, 1, \dots, K - 1$

$$\lfloor u^{(k+1)}(\theta) = \mathcal{A}(u^{(k)}(\theta), \theta)$$

*dual update*

$$\lfloor w_1^{(K)}(\theta) \cdots w_T^{(K)}(\theta) \rfloor = \mathcal{B}(u^{(K)}(\theta), \theta)$$

*primal dual relationship*

## Goals:

- Find  $\mathcal{A}$  and  $\mathcal{B}$  smooth [ $\Rightarrow w^{(K)}$  is smooth  $\Rightarrow \mathcal{U}_K$  is smooth]
- Prove that the approximate bilevel scheme converges

- **Bilevel Framework for Estimating the Group Lasso Structure**
- **Design of a Dual Forward-Backward Algorithm with Bregman Distances** such that
  - 1  $\mathcal{A}$  and  $\mathcal{B}$  are smooth  $\Rightarrow \mathcal{U}_K$  is smooth
  - 2  $\begin{cases} \min \mathcal{U}_K \rightarrow \min \mathcal{U} \\ \operatorname{argmin} \mathcal{U}_K \rightarrow \operatorname{argmin} \mathcal{U} \end{cases}$

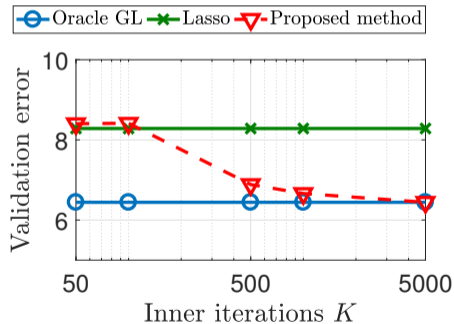
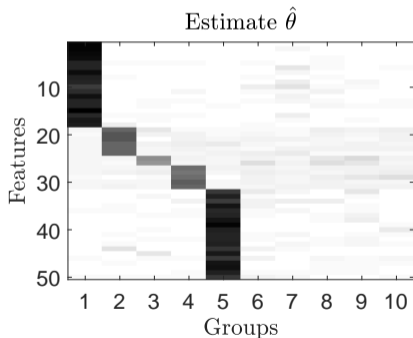
Implementation of proxSAGA algorithm: nonconvex stochastic variant of

$$\theta^{(q+1)} = \mathcal{P}_{\Theta}(\theta^{(q)} - \gamma \nabla \mathcal{U}_K(\theta^{(q)}))$$

# Numerical Experiment

**Setting:**  $T = 500$  tasks,  $N = 25$  noisy observations,  $P = 50$  features.

Estimate and group the features into, at most,  $L = 10$  groups.





# Thank You

Our poster AB #92 will be presented in Room 210 & 230 at 5pm