

# Unveiling Groups of Related Tasks in Multi-Task Learning

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



25<sup>th</sup> International Conference on Pattern Recognition, Milan, Italy

# Multi-task learning

**Setting:**  $T$  linear regression tasks

$$\begin{array}{cccccccc} \text{find } w_1 & \text{find } w_2 & \text{find } w_3 & \text{find } w_4 & \text{find } w_5 & \text{find } w_6 & \cdots & \text{find } w_T \\ y_1 \approx X_1 w_1 & y_2 \approx X_2 w_2 & y_3 \approx X_3 w_3 & y_4 \approx X_4 w_4 & y_5 \approx X_5 w_5 & y_6 \approx X_6 w_6 & \cdots & y_T \approx X_T w_T \end{array}$$

$$W = [w_1 \cdots w_T]$$

low-rank

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

# Multi-task learning

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

$$\begin{array}{ccccccccccc} \text{find } w_1 & \text{find } w_2 & \text{find } w_3 & \text{find } w_4 & \text{find } w_5 & \text{find } w_6 & \dots & \text{find } w_T \\ y_1 \approx X_1 w_1 & y_2 \approx X_2 w_2 & y_3 \approx X_3 w_3 & y_4 \approx X_4 w_4 & y_5 \approx X_5 w_5 & y_6 \approx X_6 w_6 & \dots & y_T \approx X_T w_T \end{array}$$

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

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

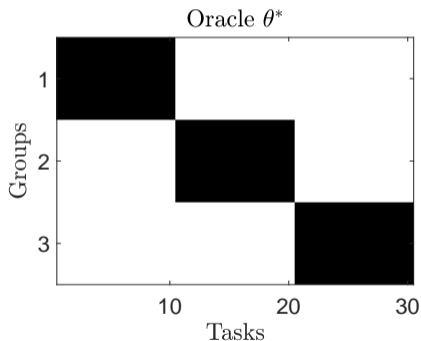
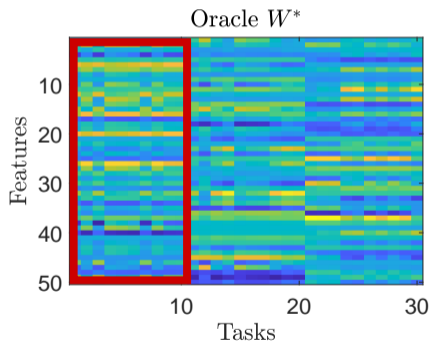
$$W_{\mathcal{G}_3} = [w_6 \dots w_T] \\ \text{low-rank}$$

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

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

# Parametrization of related tasks

$$\hat{W}(\theta_1 \cdots \theta_L) \in \underset{W=[w_1 \cdots w_T]}{\operatorname{argmin}} \sum_{t=1}^T \frac{1}{2} \|y_t - X_t w_t\|^2 + \lambda \sum_{l=1}^L \|\theta_l \odot W\|_{\operatorname{tr}},$$



**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:

$$\underset{W=[w_1 \dots w_T]}{\text{minimize}} \mathcal{L}(W, \theta) := \sum_{t=1}^T \frac{1}{2} \|y_t - X_t w_t\|^2 + \lambda \sum_{l=1}^L \|\theta_l \odot W\|_{\text{tr}}$$

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

$$\lfloor U^{(k+1)}(\theta) = \mathcal{A}(U^{(k)}(\theta), \theta) \quad \text{dual update}$$

$$w^{(K)}(\theta) = \mathcal{B}(U^{(K)}(\theta), \theta) \quad \text{primal dual relationship}$$

## 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 finding groups of related tasks
- 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 a projected gradient descent algorithm (and some variants)

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

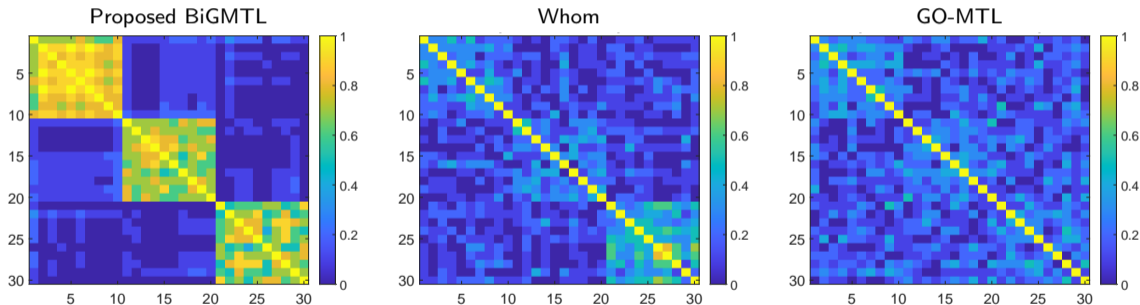
**Technicality:**  $\mathcal{U}_K$  involves generalized matrix functions  $\rightarrow$  computing  $\nabla \mathcal{U}_K$  requires attention

# Numerical Experiment

**Setting:**  $T = 30$  tasks arranged of 3 groups made of 10 tasks each.

$N = 10$  noisy observations and  $P = 20$  features per task.

→ Estimate and group the features into, at most,  $L = 6$  groups.



**Figure 1:** Mean group covariance matrix  $\theta^\top \theta$  on the synthetic experiment. Only the proposed method manages to clearly estimate the three groups of tasks.

# Thank You

A Matlab toolbox will be available at <https://github.com/jordanFrecon>