

---

# Conditional Meta-Learning of Linear Representations

---

**Giulia Denevi\***

Leonardo Labs (Italy)  
giulia.denevi.ext@leonardo.com

**Massimiliano Pontil**

Istituto Italiano di Tecnologia (Italy) & University College of London (UK)  
massimiliano.pontil@iit.it

**Carlo Ciliberto**

University College of London (UK) & Istituto Italiano di Tecnologia (Italy)  
c.ciliberto@ucl.ac.uk

## Abstract

Standard meta-learning for representation learning aims to find a common representation to be shared across multiple tasks. The effectiveness of these methods is often limited when the nuances of the tasks’ distribution cannot be captured by a single representation. In this work we overcome this issue by inferring a conditioning function, mapping the tasks’ side information (such as the tasks’ training dataset itself) into a representation tailored to the task at hand. We study environments in which our conditional strategy outperforms standard meta-learning, such as those in which tasks can be organized in separate clusters according to the representation they share. We then propose a meta-algorithm capable of leveraging this advantage in practice. In the unconditional setting, our method yields a new estimator enjoying faster learning rates and requiring less hyper-parameters to tune than current state-of-the-art methods. Our results are supported by preliminary experiments.

## 1 Introduction

Learning a shared representation among a class of machine learning problems is a well-established approach used both in multi-task learning [3, 20, 11] and meta-learning [18, 15, 5, 17, 35, 24, 30, 9, 7]. The idea behind this methodology is to consider two nested problem: at the within-task level an empirical risk minimization is performed on each task, using inputs transformed by the current representation, on the outer-task (meta-) level, such a representation is updated taking into account the errors of the within-task algorithm on previous tasks.

Such a technique was shown to be advantageous in contrast to solving each task independently when the tasks share a low dimensional representation, see e.g. [27, 25, 15, 24, 35, 5, 22, 9]. However, in real world applications we often deal with heterogeneous classes of learning tasks, which may overall be only loosely related. Consequently, the tasks’ commonalities may not be captured well by a single representation shared among all the tasks. This is for instance the case in which the tasks can be organized in different groups (clusters), where only tasks belonging to the same cluster share the same low-dimensional representation.

In order to overcome this issue, previous authors developed non-convex methods (or convex relaxations) attempting at clustering the tasks, see e.g. [4, 26, 2, 20, 28, 40, 38, 31]. In this work, we follow

---

\*Work done while the first author was with Istituto Italiano di Tecnologia (Italy).

the recent literature on heterogeneous meta-learning [37, 36, 32, 21, 10, 39, 14, 7] and propose a so-called *conditional meta-learning* approach for meta-learning a representation. Our algorithm learns a conditioning function mapping available tasks’ side information into a *linear* representation that is tuned to that task at hand. Our approach borrows from [14], where the authors proposed a conditional meta-learning approach for fine tuning and biased regularization. In those cases however, the tasks’ target vectors are assumed to be all close to a common bias vector rather than sharing the same low-dimensional linear representation, as instead explored in this work. As we explain in the following, working with the representation setting requires a significant contribution with respect to the bias one in order to give a new formulation of the problem, to consider a different meta-objective and a different interpretation of the results. In addition, the representation setting is known to be a more relevant and effective framework in many scenarios in comparison to the bias one (see e.g. [28]).

In this work, we propose for the first time an online conditional method for linear representation learning with strong theoretical guarantees. In particular, we show that the method is advantageous over standard (unconditional) representation learning methods used in meta-learning when the environment of observed tasks is heterogeneous.

**Contributions and organization.** The contributions of this work are the following. First, in Sec. 2, we design a conditional meta-learning approach to infer a linear representation that is tuned to the task at hand. Second, in Sec. 3, we formally characterize circumstances under which our conditional framework brings advantage with respect to the standard unconditional approach. In particular, we argue that this is the case when the tasks are organized in different clusters according to the support pattern or linear representation their target vectors’ share. Third, in Sec. 4, we design a convex meta-algorithm providing a comparable gain as the number of the tasks it observes increases. In the unconditional setting, the proposed method is able to recover faster rates and it requires to tune one less hyper-parameter with respect to the state-of-the-art unconditional methods. Finally, in Sec. 5, we present numerical experiments supporting our theoretical claims. We conclude our work in Sec. 6 and we postpone the missing proofs to the supplementary material.

## 2 Conditional representation learning

In this section we introduce our conditional meta-learning setting for representation learning. Then, we proceed to identify the differences with respect to (with respect to) the standard unconditional counterpart. We begin our overview by first introducing the class of inner learning algorithms we use in this work.

**Within-task algorithms.** We consider the standard linear supervised learning setting over  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$  with  $\mathcal{X} \subseteq \mathbb{R}^d$  and  $\mathcal{Y} \subseteq \mathbb{R}$  input and output spaces, respectively. We denote by  $\mathcal{P}(\mathcal{Z})$  the set of probability distributions (tasks) over  $\mathcal{Z}$ . For any task  $\mu \in \mathcal{P}(\mathcal{Z})$  and a given loss function  $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ , we aim at finding a weight vector  $w_\mu \in \mathbb{R}^d$  minimizing the *expected risk*

$$\min_{w \in \mathbb{R}^d} \mathcal{R}_\mu(w) \quad \mathcal{R}_\mu(w) = \mathbb{E}_{(x,y) \sim \mu} \ell(\langle x, w \rangle, y), \quad (1)$$

where,  $\langle \cdot, \cdot \rangle$  represents the Euclidean product in  $\mathbb{R}^d$ . In practice,  $\mu$  is only partially observed through a dataset  $Z = (x_i, y_i)_{i=1}^n \sim \mu^n$ , namely, a collection of  $n$  identically independently distributed (i.i.d.) points sampled from  $\mu$ . Thus, the goal becomes to use a learning algorithm in order to estimate a candidate weight vector with a small expected risk converging to the ideal  $\mathcal{R}_\mu(w_\mu)$  as the sample size  $n$  grows. Specifically, in this work we will consider as candidate estimators, the family of regularized empirical risk minimizers for linear feature learning [3]. Formally, denoting by  $\mathcal{D} = \bigcup_{n \in \mathbb{N}} \mathcal{Z}^n$  the space of all datasets on  $\mathcal{Z}$ , for a given  $\theta \in \Theta$  in  $\Theta = \mathbb{S}_+^d$  the set of positive definite  $d \times d$  matrices, we will consider the following learning algorithms  $A(\theta, \cdot) : \mathcal{D} \rightarrow \mathbb{R}^d$ :

$$A(\theta, Z) = \operatorname{argmin}_{w \in \operatorname{Ran}(\theta) \subseteq \mathbb{R}^d} \mathcal{R}_{Z,\theta}(w), \quad \mathcal{R}_{Z,\theta}(w) = \frac{1}{n} \sum_{i=1}^n \ell(\langle x_i, w \rangle, y_i) + \frac{1}{2} \langle w, \theta^\dagger w \rangle \quad (2)$$

where  $\operatorname{Ran}(\theta)$  denotes the range of  $\theta$ . Here  $\theta^\dagger$  denotes the pseudoinverse of  $\theta$ . Throughout this work we will denote by  $\mathcal{R}_Z(\cdot) = 1/n \sum_{i=1}^n \ell(\langle x_i, \cdot \rangle, y_i)$  the empirical risk associated to  $Z$ . Here,  $\theta$  plays the role of a linear feature representation that is learned during the meta-learning process (see [3]).

**Remark 1** (Within-task regularization parameter). Differently to previous work, see e.g. [15], we do not impose any constraints on the trace of  $\theta$  (e.g.  $\text{Tr}(\theta) \leq 1$ ). This allows us to absorb the regularization parameter  $\lambda$  typically used to control  $\lambda \langle w, \theta^\dagger w \rangle$  in  $\theta$ . As we will discuss later, this choice reduces the number of hyper-parameter to tune and it allows to enjoy faster learning rates.

**Remark 2** (Online variant of Eq. (2)). Paying additional negligible logarithmic factors, our analysis and results extend also to the setting in which the minimizer in Eq. (2) is replaced by a pre-conditioned variant of online gradient descent on  $\mathcal{R}_{Z,\theta}$  with starting point  $w_0 = 0$  and appropriate step size:

$$A(\theta, Z) = \frac{1}{n} \sum_{i=1}^n w_i, \quad w_{i+1} = w_i - \frac{\theta p_i}{i}, \quad p_i = s_i x_i + \theta^\dagger w_i, \quad s_i \in \partial \ell(\cdot, y_i)(\langle x_i, w_i \rangle). \quad (3)$$

**Unconditional Meta-Learning.** The standard unconditional meta-learning setting assumes there exist a meta-distribution  $\rho \in \mathcal{P}(\mathcal{M})$  – also called *environment* in [6] – over a family  $\mathcal{M} \subseteq \mathcal{P}(\mathcal{Z})$  of distributions (tasks)  $\mu$  and it aims at selecting an inner algorithm in the family above that is well suited to solve tasks  $\mu$  sampled from  $\rho$ . This target can be reformulated as finding a linear representation  $\theta_\rho \in \Theta$  such that the corresponding algorithm  $A(\theta_\rho, \cdot)$  minimizes the *transfer risk*

$$\min_{\theta \in \Theta} \mathcal{E}_\rho(\theta), \quad \mathcal{E}_\rho(\theta) = \mathbb{E}_{\mu \sim \rho} \mathbb{E}_{Z \sim \mu^n} \mathcal{R}_\mu(A(\theta, Z)). \quad (4)$$

In practice, this stochastic problem is usually tackled by iteratively sampling a task  $\mu \sim \rho$  and a corresponding dataset  $Z \sim \mu^n$ , and, then, performing a step of stochastic gradient descent on an empirical approximation of Eq. (4) computed from  $Z$ . This approach has proven effective for instance when the tasks of the environment share a simple common linear representation, see e.g. [18, 5, 22, 15, 16, 12, 17, 9]. However, when a single linear representation is not sufficient for the entire environment of tasks (e.g. multi-clusters), this homogeneous approach is expected to fail. In order to overcome this limitation, some recent works have adopted the following conditional approach to the problem, see e.g. [37, 36, 32, 21, 10, 39, 14].

**Conditional Meta-learning.** Analogously to [14], we assume that any task  $\mu \sim \rho$  is provided of additional side information  $s \in \mathcal{S}$ . In such a case, we consider the environment  $\rho$  as a distribution  $\rho \in \mathcal{P}(\mathcal{M}, \mathcal{S})$  over the set  $\mathcal{M}$  of tasks and the set  $\mathcal{S}$  of possible side information. Moreover, as usual, we assume  $\rho$  to decompose in  $\rho(\cdot|s)\rho_S(\cdot)$  and  $\rho(\cdot|\mu)\rho_{\mathcal{M}}(\cdot)$  the conditional and marginal distributions with respect to  $\mathcal{S}$  and  $\mathcal{M}$ . For instance, we observe that the side information  $s$  could contain descriptive features of the associated task, for example attributes in collaborative filtering [1], or additional information about the users in recommendation systems [19]). Moreover  $s$  could be formed by a portion of the dataset sampled from  $\mu$  (see [37, 14]). Conditional meta-learning leverages this additional side information in order to adapt (or condition) the linear representation  $\theta \in \Theta$  on the associated task at hand, by learning a linear-representation-valued function  $\tau$  solving the problem

$$\min_{\tau \in \mathcal{T}} \mathcal{E}_\rho(\tau), \quad \mathcal{E}_\rho(\tau) = \mathbb{E}_{(\mu,s) \sim \rho} \mathbb{E}_{Z \sim \mu^n} \mathcal{R}_\mu(A(\tau(s), Z)) \quad (5)$$

over the space  $\mathcal{T}$  of measurable functions  $\tau : \mathcal{S} \rightarrow \Theta$ . Notice that we retrieve the unconditional meta-learning problem in Eq. (4) if we restrict Eq. (5) to the set of functions  $\mathcal{T}^{\text{const}} = \{\tau \mid \tau(\cdot) \equiv \theta, \theta \in \Theta\}$ , mapping all the side information into the same constant linear representation.

In the next section, we will investigate the theoretical advantages of adopting such a conditional perspective and, then, we will introduce a convex meta-algorithm to tackle Eq. (5).

### 3 The advantage of conditional representation learning

In order to characterize the behavior of the optimal solution of Eq. (5) and to investigate the potential advantage of conditional meta-learning, we analyze the generalization properties of a given conditioning function  $\tau$ . Formally, we compare the error  $\mathcal{E}_\rho(\tau)$  with respect to the optimal minimum risk

$$\mathcal{E}_\rho^* = \mathbb{E}_{\mu \sim \rho} \mathcal{R}_\mu(w_\mu) \quad w_\mu = \underset{w \in \mathbb{R}^d}{\text{argmin}} \mathcal{R}_\mu(w). \quad (6)$$

In order to do this, we first need to introduce the following standard assumptions used also in previous literature. Throughout this work we will denote by  $\cdot^\top$  the standard transposition operation.

**Assumption 1.** Let  $\ell$  be a convex and  $L$ -Lipschitz loss function in the first argument. Additionally, there exist  $R > 0$  such that  $\|x\| \leq R$  for any  $x \in \mathcal{X}$ .

**Theorem 1** (Excess risk with generic conditioning function  $\tau$ ). *Let Asm. 1 hold. For any  $s \sim \rho_S$ , introduce the conditional covariance matrices  $W(s) = \mathbb{E}_{\mu \sim \rho(\cdot|s)} w_\mu w_\mu^\top$  and  $C(s) = \mathbb{E}_{\mu \sim \rho(\cdot|s)} \mathbb{E}_{x \sim \eta_\mu} x x^\top$ , where,  $\eta_\mu$  denotes the inputs' marginal distribution of the task  $\mu$ . Let  $\tau \in \mathcal{T}$  such that  $\text{Ran}(W(s)) \subseteq \text{Ran}(\tau(s))$  for any  $s \sim \rho_S$  and let  $A(\tau(s), \cdot)$  be the associated inner algorithm from Eq. (2). Then,*

$$\mathcal{E}_\rho(\tau) - \mathcal{E}_\rho^* \leq \frac{\mathbb{E}_{s \sim \rho_S} \text{Tr}(\tau(s)^\dagger W(s))}{2} + \frac{2L^2 \mathbb{E}_{s \sim \rho_S} \text{Tr}(\tau(s) C(s))}{n}. \quad (7)$$

*Proof.* For any  $(\mu, s) \sim \rho$ , consider the decomposition  $\mathcal{E}_\rho(\tau) - \mathcal{E}_\rho^* = \mathbb{E}_{(\mu, s) \sim \rho} [\mathbf{B}_{\mu, s} + \mathbf{C}_{\mu, s}]$ , with

$$\begin{aligned} \mathbf{B}_{\mu, s} &= \mathbb{E}_{Z \sim \mu^n} [\mathcal{R}_\mu(A(\tau(s), Z)) - \mathcal{R}_Z(A(\tau(s), Z))] \\ \mathbf{C}_{\mu, s} &= \mathbb{E}_{Z \sim \mu^n} [\mathcal{R}_Z(A(\tau(s), Z)) - \mathcal{R}_\mu(w_\mu)]. \end{aligned}$$

$\mathbf{B}_{\mu, s}$  is the generalization error of the inner algorithm  $A(\tau(s), \cdot)$  on the task  $\mu$ . Hence, applying stability arguments (see Prop. 6 in App. A), we can write  $\mathbf{B}_{\mu, s} \leq 2L^2 \text{Tr}(\tau(s) \mathbb{E}_{x \sim \eta_\mu} x x^\top) n^{-1}$ . Regarding the term  $\mathbf{C}_{\mu, s}$ , for any conditioning function  $\tau$  such that  $w_\mu \in \text{Ran}(\tau(s))$ , we can write  $\mathbf{C}_{\mu, s} \leq \mathbb{E}_{Z \sim \mu^n} [\mathcal{R}_{Z, \tau(s)}(w_\mu) - \mathcal{R}_\mu(w_\mu)] = 2^{-1} \text{Tr}(\tau(s)^\dagger w_\mu w_\mu^\top)$ , where, the inequality exploits the definition of the algorithm in Eq. (2) as minimum of the regularized empirical risk. The desired statement follows by combining the two bounds above and rewriting  $\mathbb{E}_{(\mu, s) \sim \rho} = \mathbb{E}_{s \sim \rho_S} \mathbb{E}_{\mu \sim \rho(\cdot|s)}$ .  $\square$

Thm. 1 suggests that the conditioning function  $\tau_\rho$  minimizing the right hand side of Eq. (7) is a good candidate to solve the meta-learning problem. The following result explores this question by showing that such a minimizer admits a closed form solution. The proof is reported in App. B. In the following, we will denote by  $\|\cdot\|_F$  and  $\|\cdot\|_*$  the Frobenius and trace norm of a matrix, respectively.

**Proposition 2** (Best conditioning function in hindsight). *The conditioning function minimizer and the minimum of the bound presented in Thm. 1 over the set  $\{\tau \in \mathcal{T} \mid \text{Ran}(W(s)) \subseteq \text{Ran}(\tau(s)), \rho_S\text{-almost surely}\}$ , are respectively*

$$\begin{aligned} \tau_\rho(s) &= (2L)^{-1} n^{1/2} C(s)^\dagger / 2 (C(s)^{1/2} W(s) C(s)^{1/2})^{1/2} C(s)^\dagger / 2 \\ \mathcal{E}_\rho(\tau_\rho) - \mathcal{E}_\rho^* &\leq 2L \mathbb{E}_{s \sim \rho_S} \|W(s)^{1/2} C(s)^{1/2}\|_* n^{-1/2}. \end{aligned} \quad (8)$$

We observe that, in comparison to [14], the numerator term in the bound above describes a different kind of tasks' similarity assumption: the conditional variance term  $\mathbb{E}_{(\mu, s) \sim \rho} \|w_\mu - \mathbb{E}_{\mu \sim \rho(\cdot|s)} w_\mu\|^2$  present in [14] is now substituted by the trace norm term  $\mathbb{E}_{s \sim \rho_S} \|W(s)^{1/2} C(s)^{1/2}\|_*$  above. Additionally, the bound above allows us to quantify the benefits of adopting the conditional feature learning strategy.

**Conditional vs. unconditional Meta-Learning.** Applying Prop. 2 to  $\mathcal{T}^{\text{const}}$ , we obtain the optimal (constant) meta-parameter and the corresponding excess risk bound for unconditional meta-learning

$$\begin{aligned} \tau &\equiv \theta_\rho = (2L)^{-1} n^{1/2} C_\rho^\dagger / 2 (C_\rho^{1/2} W_\rho C_\rho^{1/2})^{1/2} C_\rho^\dagger / 2 \\ \mathcal{E}_\rho(\theta_\rho) - \mathcal{E}_\rho^* &\leq 2L \|W_\rho^{1/2} C_\rho^{1/2}\|_* n^{-1/2} \end{aligned} \quad (9)$$

with unconditional covariance matrices  $W_\rho = \mathbb{E}_{\mu \sim \rho} w_\mu w_\mu^\top$  and  $C_\rho = \mathbb{E}_{\mu \sim \rho} \mathbb{E}_{x \sim \eta_\mu} x x^\top$ . We observe that in the previous literature [13, 15] the authors restricted the unconditional problem over the smaller class of linear representation  $\hat{\Theta} = \{\theta \in \mathbb{S}_+^d : \text{Ran}(W_\rho) \subseteq \text{Ran}(\theta), \text{Tr}(\theta) \leq 1\}$  and they considered as the best unconditional representation, the matrix minimizing only a part of the previous bound, namely,

$$\hat{\theta}_\rho = \underset{\theta \in \hat{\Theta}}{\text{argmin}} \text{Tr}(\theta^\dagger W_\rho) = W_\rho^{1/2} (\text{Tr}(W_\rho^{1/2}))^{-1}. \quad (10)$$

On the other hand, the unconditional oracle we introduce above in Eq. (9) allows us to recover a tighter bound which is able to recover the best performance between independent task learning (ITL) and the oracle considered in previous literature [15]. Indeed, by exploiting the duality between the trace norm  $\|\cdot\|_*$  and the operator norm  $\|\cdot\|_\infty$  of a matrix, we can upper bound the right-side-term in Eq. (9) by the quantity

$$2L \min \left\{ \|W_\rho^{1/2}\|_* \|C_\rho^{1/2}\|_\infty, \|W_\rho^{1/2}\|_F \|C_\rho^{1/2}\|_F \right\} n^{-1/2},$$

namely, the minimum between the bound for independent task learning and the bound for unconditional oracle obtained by previous authors. Notice that the unconditional quantity in Eq. (9) is always bigger than the conditional quantity in Eq. (8), since Eq. (9) coincides with the minimum over a smaller class of function. In order to quantify the gap between these two quantities – namely, the advantage in using the conditional approach with respect to the unconditional one – we have to compare the term  $\|W_\rho^{1/2} C_\rho^{1/2}\|_*$  with the term  $\mathbb{E}_{s \sim \rho_S} \|C(s)^{1/2} W(s)^{1/2}\|_*$ .

We report below a setting that can be considered illustrative for many real-world scenarios in which such a gap in performance is significant.

**Example 1 (Clusters).** Let  $S = \mathbb{R}^q$  be the side information space, for some integer  $q > 0$ . Let  $\rho$  be such that the side information marginal distribution  $\rho_S$  is given by a uniform mixture of  $m$  uniform distributions. More precisely, let  $\rho_S = \frac{1}{m} \sum_{i=1}^m \rho_S^{(i)}$ , with  $\rho_S^{(i)} = \mathcal{U}(\mathcal{B}(a_i, 1/2))$  the uniform distribution on the ball of radius  $1/2$  centered at  $a_i \in S$ , characterizing the cluster  $i$ . For a given side information  $s$ , a task  $\mu \sim \rho(\cdot|s)$  is sampled such that: 1) its inputs' marginal  $\eta_\mu$  is a distribution with constant covariance matrix  $C(s) = \mathbb{E}_{\mu \sim \rho(\cdot|s)} \mathbb{E}_{x \sim \eta_\mu} x x^\top = C$ , for some  $C \in \mathbb{S}_+^d$ , 2)  $w_\mu$  is sampled from a distribution with conditional covariance matrix  $W(s) = \mathbb{E}_{\mu \sim \rho(\cdot|s)} w_\mu w_\mu^\top$ , with  $W(s)$  such that  $(C^{1/2} W(s) C^{1/2})(C^{1/2} W(p) C^{1/2}) = 0$  if  $s \neq p$ . Then,  $\mathbb{E}_{s \sim \rho_S} \|C(s)^{1/2} W(s)^{1/2}\|_* = \frac{1}{\sqrt{m}} \|W_\rho^{1/2} C_\rho^{1/2}\|_*$ .

The inequality above tells us that, in the setting of Ex. 1, the conditional approach gains a  $\sqrt{m}$  factor in comparison to the unconditional approach. Therefore, the larger the number of clusters is, the more pronounced the advantage of conditional approach with respect to the unconditional one will be. The  $\sqrt{m}$  gain factor follows from the fact that the weight vectors  $w_\mu$  sampled from the different clusters share disjoint supports (they share orthogonal representations). This allows us to rewrite the overall clusters weight vectors' covariance as the average of the intra clusters weight vectors' covariances. The  $\sqrt{m}$  term comes from this rewriting and the quadratic behavior of the covariance matrix. We refer to App. C for more details and the deduction. We also observe that a particular case of the setting above could be that one in which  $q = 1$  and the side information are *noisy* observations of the index of the cluster the tasks belong to. In our experiments, in Sec. 5, we consider a more interesting and realistic variant of the setting above, in which we will use as task's side information a training dataset sampled from that task. In the next section, we introduce a convex meta-algorithm mimicking this advantage also in practice.

## 4 Conditional representation Meta-Learning algorithm

To tackle conditional meta-learning in practice we consider a parametrization where the conditioning functions that are modeled with respect to a given feature map  $\Phi : S \rightarrow \mathbb{R}^k$  (with  $k \in \mathbb{N}$ ) on the side information space. In other words, we consider  $\tau : S \rightarrow \mathbb{S}_+^d$ ,

$$\tau(\cdot) = (M\Phi(\cdot))^\top M\Phi(\cdot) + C, \quad (11)$$

for some tensor  $M \in \mathbb{R}^{p \times d \times k}$  ( $p \in \mathbb{N}$ ) and matrix  $C \in \mathbb{S}_+^d$ . By construction, the above parametrization guarantees us to learn functions taking values in the set of positive semi-definite matrices. However, directly addressing the meta-learning problem poses two issues: first, dealing with tensorial structures might become computationally challenging in practice and second, such parametrization is quadratic in  $M$  and would lead to a non-convex optimization functional in practice. To tackle this issue, the following results shows that we can rewrite the conditioning function in the form of Eq. (11) by using a matrix in  $\mathbb{S}_+^{dk}$ . This will allow us to implement our method working with matrices in  $\mathbb{S}_+^{dk}$ , instead of tensors in  $\mathbb{R}^{p \times d \times k}$ . Throughout this work, we will denote by  $\otimes$  the Kronecker product.

**Proposition 3** (Matricial re-formulation of  $\tau_M(s)$ ). *Let  $\tau$  be as in Eq. (11). Then,*

$$\tau(s) = (I_d \otimes \Phi(s)^\top) H_M (I_d \otimes \Phi(s)) + C, \quad (12)$$

where  $I_d$  is the identity in  $\mathbb{R}^{d \times d}$  and  $H_M$  is the matrix in  $\mathbb{R}^{dk \times dk}$  defined by the entries

$$(H_M)_{(i-1)k+h, (j-1)k+z} = \langle M(:, i, h), M(:, j, z) \rangle, \quad i, j = 1, \dots, d, \quad h, z = 1, \dots, k.$$

The arguments above motivate us to consider the following set of conditioning functions:

$$\mathcal{T}_\Phi = \left\{ \tau(\cdot) = (I_d \otimes \Phi(\cdot)^\top) H (I_d \otimes \Phi(\cdot)) + C \mid \text{such that } H \in \mathbb{S}_+^{dk}, C \in \mathbb{S}_+^d \right\}. \quad (13)$$

To highlight the dependency of a function  $\tau \in \mathcal{T}_\Phi$  with respect to its parameter  $H$  and  $C$ , we will denote  $\tau = \tau_{H,C}$ . Evidently,  $\mathcal{T}_\Phi$  contains the space of all unconditional estimators  $\mathcal{T}^{\text{const}}$ . We consider  $\mathcal{T}_\Phi$  equipped with the canonical norm  $\|\tau_{H,C}\|^2 = \|(H, C)\|_F^2 = \|H\|_F^2 + \|C\|_F^2$ , where, recall,  $\|\cdot\|_F$  denotes the Frobenius norm. The following two standard assumptions will allow us to design and analyse our method.

**Assumption 2.** *The optimal function  $\tau_\rho$  belongs to  $\mathcal{T}_\Phi$ , namely there exist  $H_\rho \in \mathbb{S}_+^{dk}$  and  $C_\rho \in \mathbb{S}_+^d$ , such that  $\tau_\rho(\cdot) = \tau_{H_\rho, C_\rho}(\cdot) = (I_d \otimes \Phi(\cdot)^\top) H_\rho (I_d \otimes \Phi(\cdot)) + C_\rho$ .*

**Assumption 3.** *There exists  $K > 0$  such that  $\|\Phi(s)\| \leq K$  for any  $s \in \mathcal{S}$ .*

Asm. 2, known as well-specified setting assumption (see e.g. [33]), is a standard assumption in learning theory and it allows us to restrict the conditional meta-learning problem in Eq. (5) to  $\mathcal{T}_\Phi$ , rather than to the entire space  $\mathcal{T}$  of measurable functions. Asm. 3 ensures that the meta-objective is Lipschitz (see below).

**The convex surrogate problem.** We start from observing that, exploiting the generalization properties of the within-task algorithm (see Prop. 6 in App. A), we can write the following

$$\mathcal{E}_\rho(\tau) \leq \mathbb{E}_{(\mu,s) \sim \rho} \mathbb{E}_{Z \sim \mu^n} F_Z(\tau(s)), \quad F_Z(\theta) = \mathcal{R}_{Z,\theta}(A(\theta, Z)) + \frac{2L^2}{n} \text{Tr}\left(\theta \frac{X^\top X}{n}\right)$$

where  $X \in \mathbb{R}^{n \times d}$  is the matrix with the inputs vectors  $(x_i)_{i=1}^n$  as rows. The inequality above suggests us to introduce the surrogate problem

$$\min_{\tau \in \mathcal{T}} \hat{\mathcal{E}}_\rho(\tau), \quad \hat{\mathcal{E}}_\rho(\tau) = \mathbb{E}_{(\mu,s) \sim \rho} \mathbb{E}_{Z \sim \mu^n} F_Z(\tau(s)). \quad (14)$$

We stress that the surrogate problem we take here is different from the one considered in previous work [12, 15, 14, 9], where the authors considered as meta-objective only a part of the function above, namely,  $\mathbb{E}_{(\mu,s) \sim \rho} \mathbb{E}_{Z \sim \mu^n} [\mathcal{R}_{Z,\tau(s)}(A(\tau(s), Z))]$ . As we will see in the following, such a choice is more appropriate for the problem at hand, since, differently from the meta-objective used in previous literature, it will allow us to develop a conditional meta-learning method that is theoretically grounded also for linear representation learning.

Exploiting Asm. 2, the surrogate problem in Eq. (14) can be restricted to the class of linear functions  $\mathcal{T}_\Phi$  in Eq. (13) and it can be rewritten more explicitly as

$$\min_{H \in \mathbb{S}_+^{dk}, C \in \mathbb{S}_+^d} \mathbb{E}_{(\mu,s) \sim \rho} \mathbb{E}_{Z \sim \mu^n} \mathcal{L}(H, C, s, Z), \quad \mathcal{L}(H, C, s, Z) = F_Z(\tau_{H,C}(s)). \quad (15)$$

In the following proposition we outline some useful properties of the meta-loss  $\mathcal{L}(\cdot, \cdot, s, Z)$  introduced above (such as convexity) supporting its choice as surrogate meta-loss.

**Proposition 4** (Properties of the surrogate meta-loss  $\mathcal{L}$ ). *For any  $Z \in \mathcal{D}$  and  $s \in \mathcal{S}$ , the function  $\mathcal{L}(\cdot, \cdot, s, Z)$  is convex and one of its subgradients is given, for any  $H \in \mathbb{S}_+^{dk}$  and  $C \in \mathbb{S}_+^d$ , by*

$$\nabla \mathcal{L}(H, \cdot, s, Z)(C) = \hat{\nabla}, \quad \nabla \mathcal{L}(\cdot, C, s, Z)(H) = (I_d \otimes \Phi(s)) \hat{\nabla} (I_d \otimes \Phi(s)^\top),$$

where

$$\hat{\nabla} = -\frac{\lambda}{2} \tau_{H,C}(s)^\dagger w_{\tau_{H,C}(s)}^\top w_{\tau_{H,C}(s)}^\top \tau_{H,C}(s)^\dagger + \frac{2L^2 X^\top X}{n^2}.$$

Moreover, by Asm. 1 and Asm. 3,

$$\|\nabla \mathcal{L}(\cdot, \cdot, s, Z)(H, C)\|_F \leq (1 + K^2)(LR)^2(2^{-1} + 2n^{-1}).$$

The proof of Prop. 4 is reported in App. D.2. It follows from combining results from [15] with the composition of the linear parametrization of the functions  $\tau_{H,C} \in \mathcal{T}_\Phi$ .

**The conditional Meta-Learning estimator.** The meta-learning strategy we propose consists in applying Stochastic Gradient Descent (SGD) on the surrogate problem in Eq. (15). Such a meta-algorithm is implemented in Alg. 1: we assume to observe a sequence of i.i.d. pairs  $(Z_t, s_t)_{t=1}^T$  of training datasets and side information, and at each iteration we update the conditional parameters  $(H_t, C_t)$  by performing a step of constant size  $\gamma > 0$  in the direction of  $-\nabla \mathcal{L}(\cdot, \cdot, s_t, Z_t)(H_t, C_t)$  and a projection step on  $\mathbb{S}_+^{dk} \times \mathbb{S}_+^d$ . Finally, we output the conditioning function  $\tau_{\bar{H}, \bar{C}}$  parametrized by  $(\bar{H}, \bar{C})$ , the average across all the iterates  $(H_t, C_t)_{t=1}^T$ . The theorem below analyzes the generalization properties of such a conditioning function.



---

**Algorithm 1** Meta-algorithm, SGD on Eq. (15)

---

**Input**  $\gamma > 0$  meta-step size,  $H_0 \in \mathbb{S}_+^{dk}$ ,  $C_0 \in \mathbb{S}_+^d$   
**Initialization**  $H_1 = H_0 \in \mathbb{S}_+^{dk}$ ,  $C = C_0 \in \mathbb{S}_+^d$   
**For**  $t = 1$  to  $T$   
  Receive  $(\mu_t, s_t) \sim \rho$  and  $Z_t \sim \mu_t^n$   
  Let  $\theta_t = (I_d \otimes \Phi(s_t))H_t(I_d \otimes \Phi(s_t)^\top) + C_t$  and compute  $w_{\theta_t} = A(\theta_t, Z_t)$  by Eq. (2)  
  Compute  $\nabla \mathcal{L}(\cdot, \cdot, s_t, Z_t)(H_t, C_t)$  as in Prop. 4 with  $w_{\theta_t}$   
  Update  $(H_{t+1}, C_{t+1}) = \text{proj}_\Theta((H_t, C_t) - \gamma \nabla \mathcal{L}(\cdot, \cdot, s_t, Z_t)(H_t, C_t))$   
**Return**  $\bar{H} = \frac{1}{T} \sum_{t=1}^T H_t$ ,  $\bar{C} = \frac{1}{T} \sum_{t=1}^T C_t$ 

---

**Theorem 5** (Excess risk bound for the conditioning function returned by Alg. 1). *Let Asm. 1 and Asm. 3 hold. For any  $s \sim \rho_S$ , recall the conditional covariance matrices  $W(s)$  and  $C(s)$  introduced in Thm. 1. Let  $\tau_{H,C}$  be a fixed function in  $\mathcal{T}_\Phi$  such that  $\text{Ran}(W(s)) \subseteq \text{Ran}(\tau_{H,C}(s))$  for any  $s \sim \rho_S$ . Let  $\bar{H}$  and  $\bar{C}$  be the outputs of Alg. 1 applied to a sequence  $(Z_t, s_t)_{t=1}^T$  of i.i.d. pairs sampled from  $\rho$  with an appropriate meta-step size  $\gamma$ . Then, in expectation with respect to the sampling of  $(Z_t, s_t)_{t=1}^T$ ,*

$$\mathbb{E} \mathcal{E}_\rho(\tau_{\bar{H}, \bar{C}}) - \mathcal{E}_\rho^* \leq \frac{\mathbb{E}_{s \sim \rho_S} \text{Tr}(\tau_{H,C}(s)^\dagger W(s))}{2} + \frac{2L^2 \mathbb{E}_{s \sim \rho_S} \text{Tr}(\tau_{H,C}(s) C(s))}{n} \\ + \left( \frac{1}{2} + \frac{2}{n} \right) \frac{(1 + K^2)(LR)^2 \|(H - H_0, C - C_0)\|_F}{\sqrt{T}}.$$

*Proof (Sketch).* The detailed proof is reported in App. D.4. Exploiting the fact that, for any  $\tau \in \mathcal{T}$ ,  $\mathcal{E}_\rho(\tau) \leq \hat{\mathcal{E}}_\rho(\tau)$  and adding  $\pm \hat{\mathcal{E}}_\rho(\tau_{H,C})$ , we can write the following

$$\mathbb{E}_{\mathbf{Z}} \mathcal{E}_\rho(\tau_{\bar{H}, \bar{C}}) - \mathcal{E}_\rho^* \leq A(\tau_{H,C}) + B(\tau_{H,C}) \quad (16) \\ A(\tau_{H,C}) = \mathbb{E}_{\mathbf{Z}} \hat{\mathcal{E}}_\rho(\tau_{\bar{H}, \bar{C}}) - \hat{\mathcal{E}}_\rho(\tau_{H,C}) \quad B(\tau_{H,C}) = \hat{\mathcal{E}}_\rho(\tau_{H,C}) - \mathcal{E}_\rho^*.$$

The term  $A(\tau_{H,C})$  can be controlled according to the convergence properties of the meta-algorithm in Alg. 1 as described in Prop. 12. Regarding the term  $B(\tau_{H,C})$ , exploiting the definition of the within-task algorithm in Eq. (2) as minimum, for any  $\tau \in \mathcal{T}$  such that  $\text{Ran}(\mathbb{E}_{\mu \sim \rho(\cdot|s)} w_\mu w_\mu^\top) \subseteq \text{Ran}(\tau(s))$  for any  $s \sim \rho_S$ , we can rewrite

$$B(\tau) \leq \frac{\mathbb{E}_{(\mu,s) \sim \rho} \text{Tr}(\tau(s)^\dagger w_\mu w_\mu^\top)}{2} + \frac{2L^2 \mathbb{E}_{(\mu,s) \sim \rho} \text{Tr}(\tau(s) \mathbb{E}_{x \sim \eta_\mu} x x^\top)}{n}.$$

The desired statement then derives from combining the two parts above and optimizing with respect to  $\gamma$ .  $\square$

**Remark 3** (Online variant of Eq. (2)). *Similarly to the bias regularization framework in [14], for the online inner family in Rem. 2, we approximate the meta-subgradient in Prop. 4 by replacing the batch minimizer  $A(\tau_{H,C}(s), Z)$  in Eq. (2) with the last iterate of the online algorithm in Eq. (3).*

**Proposed vs. optimal conditioning function.** Specializing the bound in Thm. 5 to the best conditioning function  $\tau_\rho$  in Prop. 2, thanks to Asm. 2, we get, up to constants, the following bound for our estimator,

$$\mathbb{E}_{s \sim \rho_S} \|W(s)^{1/2} C(s)^{1/2}\|_* n^{-1/2} + \|(H_\rho - H_0, C_\rho - C_0)\|_F T^{-1/2}.$$

From such a bound, we can state that our proposed meta-algorithm achieves comparable performance to the best conditioning function  $\tau_\rho$  in hindsight, when the number of observed tasks is sufficiently large. Moreover, recalling the unconditional oracle  $\hat{\theta}_\rho$  in Eq. (10) used in previous literature, regarding the second term vanishing with  $T$ , we observe that our conditional meta-learning approach incurs a cost of  $\|(H_\rho - H_0, C_\rho - C_0)\|_F T^{-1/2}$  as opposed to the cost of  $\|\hat{\theta}_\rho - \theta_0\| T^{-1/4}$  associated to state-of-the-art unconditional meta-learning approaches (see [15, 5, 22, 9]). Thus, our conditional approach presents a faster convergence rate with respect to  $T$  than such unconditional methods, but a complexity term that is expected to be larger due to the larger complexity of the class of functions we

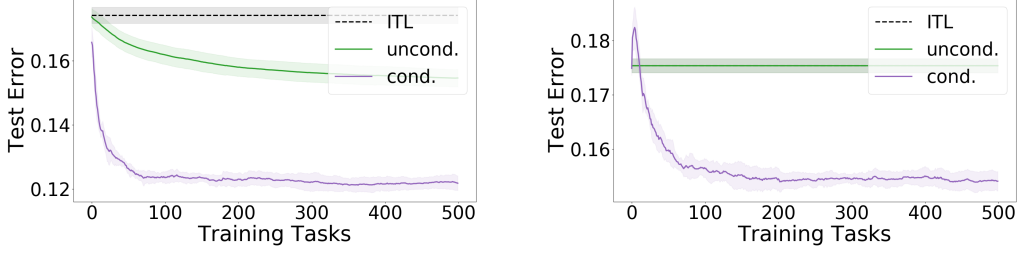


Figure 1: Mean test error (5 generations) on synthetic data. 2 (Left) and 6 (Right) clusters.

are working with. Such a faster rate with respect to  $T$  is essentially due to our formulation of the problem on the entire set of positive-semidefinite matrices (with no trace constraints). This in fact allows us to incorporate the within-task regularization parameter  $\lambda$  directly in the linear representation and to gain a  $\sqrt{T}$  order that was lost in previous literature when tuning with respect to the parameter  $\lambda$ . At the same time, this allows us to develop also a method requiring to tune just one hyper-parameter, while previous unconditional approaches requires to tune two hyper-parameters.

**Comparison to unconditional Meta-Learning.** Specializing Thm. 5 to the best unconditional estimator  $\tau_{H,C} \equiv \theta_\rho$  we introduced in Eq. (9), the bound for our estimator becomes, up to constants,

$$\|W_\rho^{1/2} C_\rho^{1/2}\|_* n^{-1/2} + \|\theta_\rho - C_0\| T^{-1/2}.$$

From the bound above, we can conclude that the conditional approach provides, at least, the same guarantees as its unconditional counterpart. Moreover, we stress again that the bound above presents a faster rate with respect to  $T$  in comparison to the state-of-the-art unconditional methods.

## 5 Experiments

We now present preliminary experiments in which we compare the proposed conditional meta-learning approach in Alg. 1 (cond.) with the unconditional counterpart (uncond.) and solving the tasks independently (ITL, namely, running the inner algorithm separately across the tasks with the constant linear representation  $\theta = I_d \in \mathbb{S}_+^d$ ). We considered regression problems and we evaluated the errors by  $\ell$  the absolute loss. We implemented the online variant of the within-task algorithm introduced in Eq. (3). The hyper-parameter  $\gamma$  was chosen by (meta-)cross validation on separate  $T_{tr}$ ,  $T_{va}$  and  $T_{te}$  respectively meta-train, -validation and -test sets. Each task is provided with a training dataset  $Z_{tr}$  of  $n_{tr}$  points and a test dataset  $Z_{te}$  of  $n_{te}$  points used to evaluate the performance of the within-tasks algorithm. In App. E we report the details of this process in our experiments.

**Synthetic clusters.** We considered two variants of the setting described in Ex. 1 with side information corresponding to the training datasets  $Z_{tr}$  associated to each task. In both settings, we sampled  $T_{tot} = 900$  tasks from a uniform mixture of  $m$  clusters. For each task  $\mu$ , we generated the target vector  $w_\mu \in \mathbb{R}^d$  with  $d = 20$  as  $w_\mu = P(j_\mu) \tilde{w}_\mu$ , where,  $j_\mu \in \{1, \dots, m\}$  denotes the cluster from which the task  $\mu$  was sampled and with the components of  $\tilde{w}_\mu \in \mathbb{R}^{d/(10)}$  sampled from the Gaussian distribution  $\mathcal{G}(0, 1)$  and then  $\tilde{w}_\mu$  normalized to have unit norm, with  $P(j_\mu) \in \mathbb{R}^{d \times d/(10)}$  a matrix with orthonormal columns. We then generated the corresponding dataset  $(x_i, y_i)_{i=1}^{n_{tot}}$  with  $n_{tot} = 80$  according to the linear equation  $y = \langle x, w_\mu \rangle + \epsilon$ , with  $x$  sampled uniformly on the unit sphere in  $\mathbb{R}^d$  and  $\epsilon$  sampled from a Gaussian distribution,  $\epsilon \sim \mathcal{G}(0, 0.1)$ . In this setting, the operator norm of the inputs' covariance matrix is small (equal to  $1/d$ ) and the weight vectors' covariance matrix of each single cluster is low-rank (its rank is  $d/(10) = 2$ ). We implemented our conditional method using the feature map  $\Phi : \mathcal{D} \rightarrow \mathbb{R}^{2d}$  defined by  $\Phi(Z) = \frac{1}{n_{tr}} \sum_{i=1}^{n_{tr}} \phi(z_i)$ , with  $\phi(z_i) = \text{vec}(x_i(y_i, 1)^\top)$ , where, for any matrix  $A = [a_1, a_2] \in \mathbb{R}^{d \times 2}$  with columns  $a_1, a_2 \in \mathbb{R}^d$ ,  $\text{vec}(A) = (a_1, a_2)^\top \in \mathbb{R}^{2d}$ . In Fig. 1, we report the results we got on an environment of tasks generated as above with  $m = 2$  (Left) and  $m = 6$  (Right) clusters, respectively. As we can see, when the clusters are two, the unconditional approach outperforms ITL (as predicted from previous literature), but the unconditional method is in turn outperformed by our conditional counterpart. When the number of clusters raises to six, the performance of unconditional meta-learning degrades to the same performance of ITL,



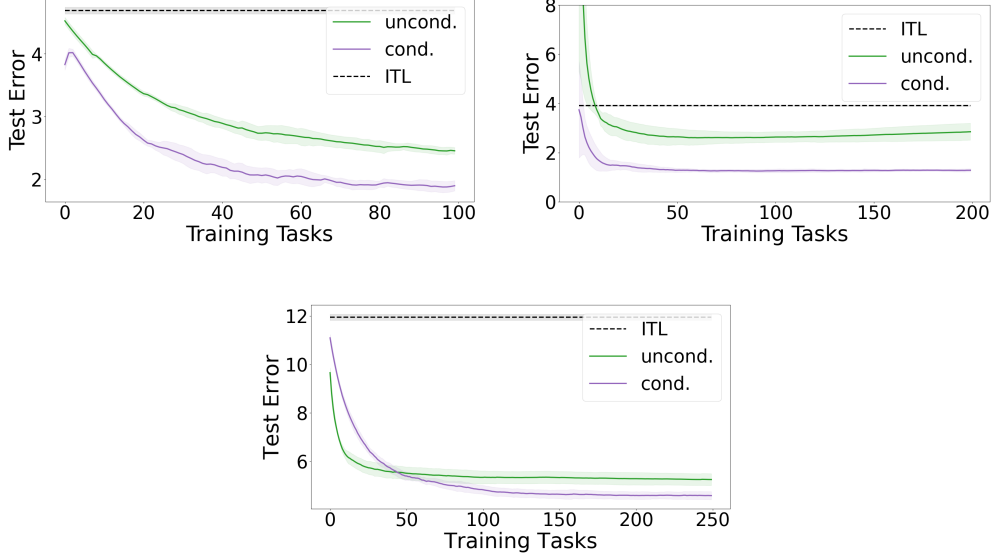


Figure 2: Mean test error (5 splits) on Lenk (Top-Left), Movielens-100k (Top-Right), Jester-1 (Bottom) dataset.

while conditional meta-learning outperforms both methods. Summarizing, the more the heterogeneity of the environment (number of clusters) is significant, the more the conditional approach brings advantage with respect to the unconditional one. This is in line with Ex. 1.

**Real datasets.** We tested the performance of the methods also on the regression problem on the computer survey data from [23] (see also [28]).  $T_{\text{tot}} = 180$  people (tasks) rated the likelihood of purchasing one of  $n_{\text{tot}} = 20$  computers. The input represents  $d = 13$  computers’ characteristics and the label is a rate in  $\{0, \dots, 10\}$ . In this case, we used as side information the training datapoints  $Z = (z_i)_{i=1}^{n_{\text{tr}}}$  and the feature map  $\Phi : \mathcal{D} \rightarrow \mathbb{R}^{d+1}$  defined by  $\Phi(Z) = w_Z$ , with  $w_Z$  the solution of Tikhonov regularization with the squared loss, namely, the vector satisfying  $(\hat{X}^\top \hat{X} + I_{d+1})w_Z = \hat{X}^\top y$ , where,  $\hat{X} \in \mathbb{R}^{(d+1) \times n}$  is the matrix obtained by adding to the matrix  $X \in \mathbb{R}^{n \times d}$  one column of ones at the end. Fig. 2 (Top-Left) shows that also in this case, the unconditional approach outperforms ITL, but the performance of its conditional counterpart is much better.

We also tested the performance of the methods on the Movielens-100k and Jester-1 real-world datasets, containing ratings of users (tasks) to movies and jokes (points), respectively. Recommendation system settings with  $d$  items can be interpreted within the meta-learning setting by considering each data point  $(x, y)$  to have input  $x \in \mathbb{R}^d$  to be the one-hot encoding of the current item to be rated (e.g. a movie or a joke) and  $y \in \mathbb{R}$  the corresponding score, see e.g. [12] for more details. We restricted the original dataset to the  $n_{\text{tot}} = 20$  most voted movies/jokes (as a consequence, by formulation,  $d = 20$ ). We guaranteed each user voted at least 5 movies/jokes, which led to a total of  $T_{\text{tot}} = 400/450$  tasks (i.e. users). In both cases, we used as side information the training datapoints  $Z = (z_i)_{i=1}^{n_{\text{tr}}}$ . For the Movielens-100k dataset we used the same feature map described for the synthetic clusters experiments in Fig. 1. For the Jester-1 dataset, let  $M$  and  $m$  denote the maximum and minimum rating value that can be assigned to a joke. We adopted the feature map  $\Phi : \mathcal{D} \rightarrow \mathbb{R}^{2d+1}$  such that, for any dataset  $Z = (x_i, y_i)_{i=1}^n$ , we have  $\Phi(Z) = (\text{vec}(\tilde{\Phi}(Z)); 1)$ , where  $\text{vec}$  denotes the vectorization operator (i.e. mapping a matrix in the vector concatenating all its columns) and  $\tilde{\Phi} : Z \rightarrow \mathbb{R}^{d \times 2}$  is such that  $\tilde{\Phi}(Z) = (\cos(\alpha(Z)), \sin(\alpha(Z))) \odot (\sum_{i=1}^n x_i)$  with  $\alpha(Z) = \sum_{i=1}^n x_i \left( \frac{\pi}{4} \frac{M - y_i}{M - m} \right)$  and  $\odot$  denoting the Hadamard (entry-wise) product broad-casted across both columns. The rationale behind this feature map is to represent as similar vectors those users with similar scores for the same movies. In particular, each item-score pair observed in training is represented as a unitary vector in  $\mathbb{R}_{++}^2$ , with the angle depending on the score attributed to that item (the vector corresponds to zero if that movie was not observed at the training time). As it can be noticed in Fig. 2 (Top Right and Bottom), the

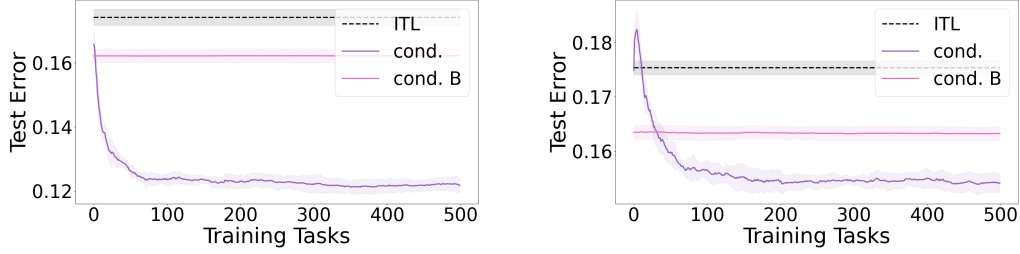


Figure 3: Mean test error (5 generations) of our method and [14] on 2 (Left) and 6 (Right) clusters.

proposed approach performs significantly better than ITL and its unconditional counterpart also on these two benchmarks. This suggests that groups of users might rely each on similar features (but different from those of other groups) to rate an item in the dataset (respectively a movie or a joke).

**Comparison with [14].** We conclude the experimental section by comparing the performance of our method with the conditional meta-learning approach for biased regularization proposed in [14]. In that case, the tasks’ target vectors are assumed to be all close to a common bias vector rather than sharing the same low-dimensional linear representation, as instead in our method. The representation setting is known to be a more relevant and effective framework in many scenarios in comparison to the bias one (see e.g. [28]). This is confirmed in Fig. 3 where our conditional representation learning method (‘cond.’) significantly outperforms the one in [14] (‘cond. B’) in the synthetic settings used in Fig. 1, since the tasks’ similarity leveraged by [14] is not appropriate in these cases.

## 6 Conclusion

We proposed a conditional meta-learning approach aiming at learning a function mapping task’s side information into a linear representation that is well suited for the task at hand. We theoretically and experimentally showed that the proposed conditional approach is advantageous with respect to the standard unconditional counterpart when the observed tasks share heterogeneous linear representations. As a consequence of our analysis we also developed a new unconditional meta-learning variant requiring tuning less hyper-parameters and relying on faster rates with respect to state-of-the-art unconditional approaches. We identify two future directions addressing the limitations of our method. A first question left opened is how to design a suitable feature map  $\Phi$  when the tasks’ training data is used as side information. Following [32, 37], we adopted a mean embedding representation. However, given the importance played by such feature map in Thm. 5, it will be worth investigating better alternatives. Secondly, it will be valuable to investigate how to predict non-linear conditioning functions (similarly to e.g. [7, 16, 32]) and use less expensive algorithms to update the positive matrices, such as the Frank-Wolfe algorithm used in [9] for unconditional settings. In this last case, according to our analysis, applying standard convergence rates for Frank-Wolfe algorithm, we expect the meta-learning algorithm based on Frank-Wolfe iteration to incur a slower rate of order  $T^{-1/4}$  (instead of  $T^{-1/2}$  for the SGD method proposed here) in Thm. 5, paying the computational benefits in terms of statistical performance.

## Acknowledgments and Disclosure of Funding

This work was supported in part by SAP SE and EPSRC Grant N. EP/P009069/1. C.C. acknowledges the support of the Royal Society (grant SPREM RGS\R1\201149) and Amazon.com Inc. (Amazon Research Award – ARA). G.D. acknowledges Leonardo SpA for funding her participation to the conference.

## References

- [1] Jacob Abernethy, Francis Bach, Theodoros Evgeniou, and Jean-Philippe Vert. A new approach to collaborative filtering: Operator estimation with spectral regularization. *Journal of Machine Learning Research*, 10(Mar):803–826, 2009.
- [2] Andreas Argyriou, Stéphan Cléménçon, and Rucong Zhang. Learning the graph of relations among multiple tasks. 2013.
- [3] Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. Convex multi-task feature learning. *Machine Learning*, 73(3):243–272, 2008.
- [4] Andreas Argyriou, Andreas Maurer, and Massimiliano Pontil. An algorithm for transfer learning in a heterogeneous environment. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 71–85. Springer, 2008.
- [5] Maria-Florina Balcan, Mikhail Khodak, and Ameet Talwalkar. Provable guarantees for gradient-based meta-learning. In *International Conference on Machine Learning*, pages 424–433, 2019.
- [6] Jonathan Baxter. A model of inductive bias learning. *J. Artif. Intell. Res.*, 12(149–198):3, 2000.
- [7] Luca Bertinetto, Joao F Henriques, Philip HS Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. *arXiv preprint arXiv:1805.08136*, 2018.
- [8] Olivier Bousquet and André Elisseeff. Stability and generalization. *Journal of machine learning research*, 2(Mar):499–526, 2002.
- [9] Brian Bullins, Elad Hazan, Adam Kalai, and Roi Livni. Generalize across tasks: Efficient algorithms for linear representation learning. In *Algorithmic Learning Theory*, pages 235–246, 2019.
- [10] T Tony Cai, Tengyuan Liang, and Alexander Rakhlin. Weighted message passing and minimum energy flow for heterogeneous stochastic block models with side information. *Journal of Machine Learning Research*, 21(11):1–34, 2020.
- [11] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- [12] Giulia Denevi, Carlo Ciliberto, Riccardo Grazi, and Massimiliano Pontil. Learning-to-learn stochastic gradient descent with biased regularization. In *International Conference on Machine Learning*, pages 1566–1575, 2019.
- [13] Giulia Denevi, Carlo Ciliberto, Dimitris Stamos, and Massimiliano Pontil. Incremental learning-to-learn with statistical guarantees. In *Proc. 34th Conference on Uncertainty in Artificial Intelligence (UAI)*, 2018.
- [14] Giulia Denevi, Massimiliano Pontil, and Carlo Ciliberto. The advantage of conditional meta-learning for biased regularization and fine tuning. *Advances in Neural Information Processing Systems*, 33, 2020.
- [15] Giulia Denevi, Dimitris Stamos, Carlo Ciliberto, and Massimiliano Pontil. Online-within-online meta-learning. In *Advances in Neural Information Processing Systems*, pages 13089–13099, 2019.
- [16] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1126–1135. PMLR, 2017.
- [17] Chelsea Finn and Sergey Levine. Meta-learning and universality: Deep representations and gradient descent can approximate any learning algorithm. In *International Conference on Learning Representations*, 2018.
- [18] Chelsea Finn, Aravind Rajeswaran, Sham Kakade, and Sergey Levine. Online meta-learning. In *International Conference on Machine Learning*, pages 1920–1930, 2019.

- [19] F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- [20] Laurent Jacob, Jean-philippe Vert, and Francis R Bach. Clustered multi-task learning: A convex formulation. In *Advances in neural information processing systems*, pages 745–752, 2009.
- [21] Ghassen Jerfel, Erin Grant, Tom Griffiths, and Katherine A Heller. Reconciling meta-learning and continual learning with online mixtures of tasks. In *Advances in Neural Information Processing Systems*, pages 9119–9130, 2019.
- [22] Mikhail Khodak, Maria-Florina F Balcan, and Ameet S Talwalkar. Adaptive gradient-based meta-learning methods. In *Advances in Neural Information Processing Systems*, pages 5915–5926, 2019.
- [23] Peter J Lenk, Wayne S DeSarbo, Paul E Green, and Martin R Young. Hierarchical bayes conjoint analysis: Recovery of partworth heterogeneity from reduced experimental designs. *Marketing Science*, 15(2):173–191, 1996.
- [24] Andreas Maurer. Transfer bounds for linear feature learning. *Machine Learning*, 75(3):327–350, 2009.
- [25] Andreas Maurer, Massi Pontil, and Bernardino Romera-Paredes. Sparse coding for multitask and transfer learning. In *International Conference on Machine Learning*, 2013.
- [26] Andreas Maurer and Massimiliano Pontil. Transfer learning in a heterogeneous environment. In *2012 3rd International Workshop on Cognitive Information Processing (CIP)*, pages 1–6. IEEE, 2012.
- [27] Andreas Maurer, Massimiliano Pontil, and Bernardino Romera-Paredes. The benefit of multitask representation learning. *The Journal of Machine Learning Research*, 17(1):2853–2884, 2016.
- [28] Andrew M McDonald, Massimiliano Pontil, and Dimitris Stamos. New perspectives on k-support and cluster norms. *Journal of Machine Learning Research*, 17(155):1–38, 2016.
- [29] Charles A Micchelli, Jean M Morales, and Massimiliano Pontil. Regularizers for structured sparsity. *Advances in Computational Mathematics*, 38(3):455–489, 2013.
- [30] Anastasia Pentina and Christoph Lampert. A PAC-Bayesian bound for lifelong learning. In *International Conference on Machine Learning*, pages 991–999, 2014.
- [31] Bernardino Romera-Paredes, Hane Aung, Nadia Bianchi-Berthouze, and Massimiliano Pontil. Multilinear multitask learning. In *International Conference on Machine Learning*, pages 1444–1452. PMLR, 2013.
- [32] Andrei A Rusu, Dushyant Rao, Jakub Sygnowski, Oriol Vinyals, Razvan Pascanu, Simon Osindero, and Raia Hadsell. Meta-learning with latent embedding optimization. In *International Conference on Learning Representations*, 2018.
- [33] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.
- [34] Ingo Steinwart and Andreas Christmann. *Support vector machines*. Springer Science & Business Media, 2008.
- [35] Nilesh Tripuraneni, Chi Jin, and Michael I Jordan. Provable meta-learning of linear representations. *arXiv preprint arXiv:2002.11684*, 2020.
- [36] Risto Vuorio, Shao-Hua Sun, Hexiang Hu, and Joseph J Lim. Multimodal model-agnostic meta-learning via task-aware modulation. In *Advances in Neural Information Processing Systems*, pages 1–12, 2019.
- [37] Ruohan Wang, Yiannis Demiris, and Carlo Ciliberto. A structured prediction approach for conditional meta-learning. *Advances in Neural Information Processing Systems*, 2020.

- [38] Kishan Wimalawarne, Masashi Sugiyama, and Ryota Tomioka. Multitask learning meets tensor factorization: task imputation via convex optimization. In *NIPS*, pages 2825–2833. Citeseer, 2014.
- [39] Huaxiu Yao, Ying Wei, Junzhou Huang, and Zhenhui Li. Hierarchically structured meta-learning. *arXiv preprint arXiv:1905.05301*, 2019.
- [40] Jiayu Zhou, Jianhui Chen, and Jieping Ye. Clustered multi-task learning via alternating structure optimization. *Advances in neural information processing systems*, 2011:702, 2011.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#) See the abstract and Lines 44 – 54.
  - (b) Did you describe the limitations of your work? [\[Yes\]](#) See Sec. 6.
  - (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#) Our work does not imply any negative societal impact.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) See Asm. 1, Asm. 2, Asm. 3.
  - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) See Sec. 3, Sec. 4 and Supplementary material.
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) See Sec. 5, App. E and Supplementary material.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) See Sec. 5, App. E and Supplementary material.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#) See Sec. 5, App. E and Supplementary material.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See App. E and Supplementary material.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#) In Fig. 3 we compared with the method in [14] and we cited that work.
  - (b) Did you mention the license of the assets? [\[Yes\]](#) See Fig. 3.
  - (c) Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#) We attached our code in the Supplementary material.
  - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[N/A\]](#) We used benchmark data and we cited the original source where one can find all the details.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#) The data we are using does not contain personally identifiable information or offensive content.
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#) We did not use crowdsourcing or conducted research with human subjects.
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#) We did not use crowdsourcing or conducted research with human subjects.

- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] We did not use crowdsourcing or conducted research with human subjects.



## Appendix

The supplementary material is organized as follows. In App. A we give the bound on the generalization error of the algorithm in Eq. (2) that we used in various proofs. In App. B we report the proof to get the closed form of the best conditioning function  $\tau_\rho$  outlined in Prop. 2. In App. C we report the proof of the statement in Ex. 1. In App. D, we report the proofs of the statements we used in Sec. 4 in order to prove the expected excess risk bound in Thm. 5 for Alg. 1. Finally, in App. E we report the experimental details we missed in the main body.

### A Generalization bound of the within-task algorithm

We now study the generalization error of the within-task algorithm in Eq. (2), i.e. the discrepancy between the (true) risk and the empirical risk of the corresponding estimator. This is done in the following result where we exploit stability arguments, more precisely the so-called hypothesis stability, see [8, Def. 3].

**Proposition 6** (Generalization error of the within-task algorithm in Eq. (2)). *Let Asm. 1 hold. For a distribution  $\mu \sim \rho$ , fix a dataset  $Z = (x_i, y_i)_{i=1}^n \sim \mu^n$ . For any  $\theta \in \Theta$ , let  $w_\theta(Z)$  be the corresponding RERM in Eq. (2) over  $Z$ . Then, the following generalization error bound holds for  $w_\theta(Z)$ :*

$$\mathbb{E}_{Z \sim \mu^n} [\mathcal{R}_\mu(w_\theta(Z)) - \mathcal{R}_Z(w_\theta(Z))] \leq \frac{2L^2}{n} \text{Tr}(\mathbb{E}_{z \sim \mu} \theta x x^\top). \quad (17)$$

*Proof.* During this proof, we need to make explicit the dependency of the RERM (Regularized Empirical Risk Minimizer)  $w_\theta$  in Eq. (2) with respect to the dataset  $Z$ . For any  $i \in \{1, \dots, n\}$ , consider the dataset  $Z^{(i)}$ , a copy of the original dataset  $Z$  in which we exchange the point  $z_i = (x_i, y_i)$  with a new i.i.d. point  $z'_i = (x'_i, y'_i)$ . For a fixed  $\theta \in \Theta$ , we analyze how much this perturbation affects the outputs of the RERM algorithm in Eq. (2). In other words, we study the discrepancy between  $w_\theta(Z)$  and  $w_\theta(Z^{(i)})$ . We start from observing that, since by Asm. 1  $\mathcal{R}_{Z,\theta}$  is 1-strongly convex with respect to  $\|\cdot\|_\theta = \sqrt{\langle \cdot, \theta^\dagger \cdot \rangle}$ , by growth condition and the definition of the RERM algorithm, we can write the following

$$\begin{aligned} \frac{1}{2} \|w_\theta(Z^{(i)}) - w_\theta(Z)\|_\theta^2 &\leq \mathcal{R}_{Z,\theta}(w_\theta(Z^{(i)})) - \mathcal{R}_{Z,\theta}(w_\theta(Z)) \\ \frac{1}{2} \|w_\theta(Z^{(i)}) - w_\theta(Z)\|_\theta^2 &\leq \mathcal{R}_{Z^{(i)},\theta}(w_\theta(Z)) - \mathcal{R}_{Z^{(i)},\theta}(w_\theta(Z^{(i)})). \end{aligned} \quad (18)$$

Hence, summing the two inequalities above, we get

$$\begin{aligned} \|w_\theta(Z^{(i)}) - w_\theta(Z)\|_\theta^2 &\leq \mathcal{R}_{Z,\theta}(w_\theta(Z^{(i)})) - \mathcal{R}_{Z^{(i)},\theta}(w_\theta(Z^{(i)})) + \mathcal{R}_{Z^{(i)},\theta}(w_\theta(Z)) - \mathcal{R}_{Z,\theta}(w_\theta(Z)) \\ &= \frac{A+B}{n}, \end{aligned} \quad (19)$$

where we have introduced the terms

$$\begin{aligned} A &= \ell(\langle x'_i, w_\theta(Z) \rangle, y'_i) - \ell(\langle x'_i, w_\theta(Z^{(i)}) \rangle, y'_i) \\ B &= \ell(\langle x_i, w_\theta(Z^{(i)}) \rangle, y_i) - \ell(\langle x_i, w_\theta(Z) \rangle, y_i). \end{aligned} \quad (20)$$

Now, introducing the subgradients  $s'_{\theta,i} \in \partial \ell(\cdot, y'_i)(\langle x'_i, w_\theta(Z) \rangle)$  and  $s_{\theta,i} \in \partial \ell(\cdot, y_i)(\langle x_i, w_\theta(Z^{(i)}) \rangle)$  and applying Holder's inequality, we can write

$$\begin{aligned} A &\leq \langle x'_i s'_{\theta,i}, w_\theta(Z) - w_\theta(Z^{(i)}) \rangle \leq \|x'_i s'_{\theta,i}\|_{\theta,*} \|w_\theta(Z^{(i)}) - w_\theta(Z)\|_\theta \\ B &\leq \langle x_i s_{\theta,i}, w_\theta(Z^{(i)}) - w_\theta(Z) \rangle \leq \|x_i s_{\theta,i}\|_{\theta,*} \|w_\theta(Z^{(i)}) - w_\theta(Z)\|_\theta, \end{aligned} \quad (21)$$

where  $\|\cdot\|_{\theta,*} = \sqrt{\langle \cdot, \theta \cdot \rangle}$  is the dual norm of  $\|\cdot\|_\theta$ . Combining these last two inequalities with Eq. (19) and simplifying, we get the following

$$\|w_\theta(Z^{(i)}) - w_\theta(Z)\|_\theta \leq \frac{1}{n} (\|x'_i s'_{\theta,i}\|_{\theta,*} + \|x_i s_{\theta,i}\|_{\theta,*}). \quad (22)$$

Hence, combining the first row in Eq. (21) with Eq. (22), we can write

$$\ell(\langle x'_i, w_\theta(Z) \rangle, y'_i) - \ell(\langle x'_i, w_\theta(Z^{(i)}) \rangle, y'_i) \leq \frac{1}{n} \left( \|x'_i s'_{\theta,i}\|_{\theta,*}^2 + \|x'_i s'_{\theta,i}\|_{\theta,*} \|x_i s_{\theta,i}\|_{\theta,*} \right). \quad (23)$$

Now, taking the expectation with respect to  $Z \sim \mu^n$  and  $z'_i \sim \mu$  of the left side member above, according to [8, Lemma 7], we get

$$\mathbb{E}_{Z \sim \mu^n} \mathbb{E}_{z'_i \sim \mu} \left[ \ell(\langle x'_i, w_\theta(Z) \rangle, y'_i) - \ell(\langle x'_i, w_\theta(Z^{(i)}) \rangle, y'_i) \right] = \mathbb{E}_{Z \sim \mu^n} \left[ \mathcal{R}_\mu(w_\theta(Z)) - \mathcal{R}_Z(w_\theta(Z)) \right].$$

Finally, taking the expectation of the right side member, exploiting the fact that the points are i.i.d. according to  $\mu$ , we get

$$\mathbb{E}_{Z \sim \mu^n} \mathbb{E}_{z'_i \sim \mu} \left[ \frac{1}{n} \left( \|x'_i s'_{\theta,i}\|_{\theta,*}^2 + \|x'_i s'_{\theta,i}\|_{\theta,*} \|x_i s_{\theta,i}\|_{\theta,*} \right) \right] \leq \frac{2}{n} \mathbb{E}_{Z \sim \mu^n} \mathbb{E}_{z'_i \sim \mu} \|x'_i s'_{\theta,i}\|_{\theta,*}^2, \quad (24)$$

where we recall that  $s'_{\theta,i} \in \partial \ell(\cdot, y'_i)(\langle x'_i, w_\theta(Z) \rangle)$ . Combining the two last statements above, we get

$$\mathbb{E}_{Z \sim \mu^n} \left[ \mathcal{R}_\mu(w_\theta(Z)) - \mathcal{R}_Z(w_\theta(Z)) \right] \leq \frac{2}{n} \mathbb{E}_{Z \sim \mu^n} \mathbb{E}_{z'_i \sim \mu} \|x'_i s'_{\theta,i}\|_{\theta,*}^2. \quad (25)$$

Finally, substituting the close form of  $\|\cdot\|_{\theta,*}$  and observing that, by Asm. 1 we have  $\|x'_i s'_{\theta,i}\|_{\theta,*}^2 \leq L^2 \|x'_i\|_{\theta,*}^2$ , we get the desired statement:

$$\mathbb{E}_{Z \sim \mu^n} \left[ \mathcal{R}_\mu(w_\theta(Z)) - \mathcal{R}_Z(w_\theta(Z)) \right] \leq \frac{2L^2}{n} \mathbb{E}_{z'_i \sim \mu} \langle x'_i, \theta x'_i \rangle = \frac{2L^2}{n} \text{Tr}(\mathbb{E}_{z \sim \mu} \theta x x^\top). \quad (26)$$

□

## B Proof of Prop. 2

In this section we report the proof to get the closed form of the best conditioning function  $\tau_\rho$  outlined in Prop. 2. In order to do this, we need the following results.

**Lemma 7.** *For any  $\mu \sim \rho_{\mathcal{M}}$ , define the inputs' covariance matrix  $C_\mu = \mathbb{E}_{x \sim \eta_\mu} x x^\top$ . Then, for any  $w_\mu \in \arg\min_{w \in \mathbb{R}^d} \mathcal{R}_\mu(w)$ , the projection  $w_{0,\mu} = C_\mu^\dagger C_\mu w_\mu$  of  $w_\mu$  onto the range of  $C_\mu$  is still a minimizer of  $\mathcal{R}_\mu$ .*

*Proof.* Consider the decomposition of  $w_\mu$  with respect to the range of  $C_\mu$ :

$$w_\mu = w_{0,\mu} + w^\perp \quad (27)$$

with  $w_{0,\mu} = C_\mu^\dagger C_\mu w_\mu$  and  $w^\perp \in \mathbb{R}^d$  such that  $C_\mu w^\perp = 0$ . We note that, almost surely with respect to the points  $x \in \mathbb{R}^d$  sampled from  $\mu$ , we have  $\langle w^\perp, x \rangle = 0$ . This follows by noting that by the orthogonality between  $C_\mu$  and  $w^\perp$ , we have

$$0 = \langle w^\perp, C_\mu w^\perp \rangle = \mathbb{E}_{x \sim \eta_\mu} \langle w^\perp, x x^\top w^\perp \rangle = \mathbb{E}_{x \sim \eta_\mu} \langle x, w^\perp \rangle^2, \quad (28)$$

that can hold only if  $\langle x, w^\perp \rangle^2 = 0$  almost surely (a.s.) with respect to  $\eta_\mu$ . We conclude that  $\langle w_\mu, x \rangle = \langle w_{0,\mu}, x \rangle + \langle w^\perp, x \rangle = \langle w_{0,\mu}, x \rangle$  a.s. with respect to  $\mu$  and, consequently,  $\mathcal{R}_\mu(w_\mu) = \mathcal{R}_\mu(w_{0,\mu})$ . □

**Corollary 8.** *For any  $s \in \mathcal{S}$ , recall the conditional covariance matrices in Thm. 1. Then,  $\text{Ran}(W(s)) \subset \text{Ran}(C(s))$ , namely the range of the task-vector conditional covariance  $W(s)$  is always contained in the range of the input conditional covariance  $C(s)$ .*

*Proof.* The corollary is a direct consequence of the previous Lemma 7. The result above guarantees that for any  $\mu \sim \rho_{\mathcal{M}}$ , the rank-one operator  $W_\mu = w_\mu w_\mu^\top$  has range contained in the range of  $C_\mu$ . Taking the conditional expectations  $W(s) = \mathbb{E}_{\mu \sim \rho(\cdot|s)} W_\mu$  and  $C(s) = \mathbb{E}_{\mu \sim \rho(\cdot|s)} C_\mu$  maintains this relation unaltered, giving the desired statement. □

**Lemma 9.** Let  $P \in \mathbb{S}_+^d$  be an orthogonal projector, namely such that  $P = P^2$ . Then, for any positive definite matrix  $\theta \in \mathbb{S}_{++}^d$ , we have  $P\theta^{-1}P \succeq (P\theta P)^\dagger$ .

*Proof.* The proof is essentially a corollary of Schur's complement. Let consider the decomposition

$$\theta = \underbrace{P\theta P}_A + \underbrace{P\theta(I-P)}_B + \underbrace{(I-P)\theta P}_{B^\top} + \underbrace{(I-P)\theta(I-P)}_C \quad (29)$$

where  $A, C \in \mathbb{S}_+^d$ ,  $B \in \mathbb{R}^{d \times d}$  and  $CB = B^\top C = 0$  since  $(I-P)P = P(I-P) = P - P^2 = P - P = 0$ . Additionally, since  $C^\dagger = CC^\dagger C^\dagger = C^\dagger C^\dagger C$ , we have that also  $AC^\dagger = ACC^\dagger C^\dagger = 0$  and analogously  $C^\dagger B = B^\top C^\dagger = 0$ . Note that since  $\theta$  is invertible, both  $A$  and  $C$  are full rank. We now observe a few relevant interactions between the objects above. In particular, we observe that  $CC^\dagger B^\top = B^\top$ . To see this, first note that

$$CC^\dagger B^\top = (I-P)\theta(I-P)((I-P)\theta(I-P))^\dagger(I-P)\theta P. \quad (30)$$

By taking  $D = (I-P)\theta^{1/2}$  and using the properties of the pseudoinverse (e.g.  $D = D^\top(DD^\top)^\dagger$ ), we have

$$\begin{aligned} CC^\dagger B^\top &= DD^\top(DD^\top)^\dagger D\theta^{1/2}P \\ &= DD^\dagger D\theta^{1/2}P \\ &= D\theta^{1/2}P \\ &= B^\top. \end{aligned} \quad (31)$$

We now derive an alternative characterization of  $\theta$  in terms of  $A, B, C$ . By adding and removing a term  $BC^\dagger B$  to  $\theta$ , we have

$$\begin{aligned} \theta &= A + B + B^\top + C \\ &= A - BC^\dagger B^\top + BC^\dagger B^\top + B + B^\top + C \\ &= A - BC^\dagger B^\top + B + C + (B + C)(C^\dagger B^\top) \\ &= A - BC^\dagger B^\top + B + C + (A - BC^\dagger B^\top + B + C)(C^\dagger B) \\ &= (A - BC^\dagger B^\top + B + C)(I + C^\dagger B^\top), \end{aligned} \quad (32)$$

where we have first used the equality  $CC^\dagger B^\top = B^\top$  and then the orthogonality  $AC^\dagger = B^\top C^\dagger = 0$ . Following a similar reasoning

$$\begin{aligned} A - BC^\dagger B^\top + B + C &= A - BC^\dagger B^\top + C + BC^\dagger C \\ &= A - BC^\dagger B^\top + C + BC^\dagger(A - BC^\dagger B^\top + C) \\ &= (I + BC^\dagger)(A - BC^\dagger B^\top + C) \end{aligned} \quad (33)$$

since  $BC^\dagger C = C$  (following the same reasoning used for  $B^\top = CC^\dagger B^\top$ ) and  $AC^\dagger = C^\dagger B = 0$ . We conclude that

$$\theta = (I + BC^\dagger)(A - BC^\dagger B^\top + C)(I + C^\dagger B^\top). \quad (34)$$

We now show that all terms in the equation above are invertible. First note that  $(I + BC^\dagger)^{-1} = (I - BC^\dagger)$  and  $(I + C^\dagger B^\top)^{-1} = (I - C^\dagger B^\top)$ . Moreover, since  $\theta \succ 0$  and  $C(A - BC^\dagger B^\top) = 0$ , then also  $A - BC^\dagger B^\top \succ 0$ . We have

$$\theta^{-1} = (I - C^\dagger B^\top)(A - BC^\dagger B^\top + C)^{-1}(I - BC^\dagger), \quad (35)$$

from which we conclude

$$\begin{aligned} P\theta^{-1}P &= P(A - BC^\dagger B^\top + C)^{-1}P \\ &= P((A - BC^\dagger B^\top)^\dagger + C^\dagger)P \\ &= P(A - BC^\dagger B^\top)^\dagger P \\ &= (A - BC^\dagger B^\top)^\dagger. \end{aligned} \quad (36)$$

Since  $BC^\dagger B^\top \succeq 0$ , we have  $A - BC^\dagger B^\top \preceq A$  and therefore  $(A - BC^\dagger B^\top)^\dagger \succeq A^\dagger$  from which we have

$$P\theta^{-1}P = (A - BC^\dagger B^\top)^\dagger \succeq A^\dagger = (P\theta P)^\dagger, \quad (37)$$

as desired.  $\square$

**Proposition 10.** Consider two matrices  $A, B \in \mathbb{S}_+^d$  such that  $\text{Ran}(A) \subseteq \text{Ran}(B)$  and consider the following associated problem:

$$\min_{\theta \in \mathbb{S}_+^d, \text{Ran}(A) \subseteq \text{Ran}(\theta)} \text{Tr}(\theta^{-1}A) + \text{Tr}(\theta B). \quad (38)$$

Then, a minimizer and the corresponding minimum of the problem above are given by

$$\theta_* = B^{-1/2}(B^{1/2}AB^{1/2})^{1/2}B^{-1/2} \quad 2\|B^{1/2}A^{1/2}\|_*. \quad (39)$$

Moreover  $\theta_*$  is the unique minimizer such that  $\text{Ran}(\theta_*) \subset \text{Ran}(B)$ .

*Proof.* Let  $\Theta = \{\theta \in \mathbb{S}_+^d \mid \text{Ran}(A) \subset \text{Ran}(\theta)\}$  and denote by  $F : \Theta \rightarrow \mathbb{R}$  the objective functional of the problem in Eq. (38), such that for any  $\theta \in \Theta$

$$F(\theta) = \text{Tr}(\theta^{-1}A) + \text{Tr}(\theta B). \quad (40)$$

Note that the sign of inverse is well defined since  $\text{Ran}(A) \subset \text{Ran}(\theta)$ . We begin the proof by showing that the Eq. (38) is equivalent to

$$\min_{\theta \in \mathbb{S}_+^d, \text{Ran}(A) \subset \text{Ran}(\theta) \subset \text{Ran}(B)} \text{Tr}(\theta^{-1}A) + \text{Tr}(\theta B). \quad (41)$$

To see this, let  $P = BB^\dagger$  the orthogonal projector onto the range of  $B$ . By hypothesis,  $A = PAP$  and  $B = PBP$ . Therefore, for any  $\theta \in \mathbb{S}_{++}^d$

$$\begin{aligned} F(\theta) &= \text{Tr}(\theta^{-1}A) + \text{Tr}(\theta B) \\ &= \text{Tr}(P\theta^{-1}PA) + \text{Tr}(P\theta PB) \\ &\geq \text{Tr}((P\theta P)^\dagger A) + \text{Tr}(P\theta PB) \\ &= F(P\theta P), \end{aligned}$$

where we have applied the fact that  $P\theta^{-1}P \geq (P\theta P)^\dagger$  from Lemma 9 and the positive semidefiniteness of  $A$ . The inequality above implies the equivalence between Eq. (38) and Eq. (41). Indeed, let  $\theta_* \in \Theta$  be a minimizer of Eq. (38) and consider a sequence  $(\theta_n)_{n \in \mathbb{N}}$  such that  $\theta_n \in \mathbb{S}_{++}^d$  for any  $n \in \mathbb{N}$  and  $\theta_n \rightarrow \theta_*$ . By continuity of  $F$  we have also that  $F(\theta_n) \rightarrow F(\theta_*)$ . Clearly,  $F(\theta_*) \leq F(P\theta_n P) \leq F(\theta_n)$  and therefore also  $F(P\theta_n P) \rightarrow F(\theta_*)$ . By continuity of  $F$  over  $\Theta$ , this also implies that the limit  $\lim_{n \rightarrow +\infty} P\theta_n P = P\theta_* P$  is a minimizer for Eq. (38) (and one such that  $\text{Ran}(\theta_*) \subset \text{Ran}(B)$ ). We consider now the set  $\Theta_B = \{\theta \in \mathbb{S}_+^d \mid \text{Ran}(\theta) = \text{Ran}(B)\}$  of all positive semidefinite matrices with same range as  $B$ , hence invertible on  $\text{Ran}(B)$ . Note that  $\Theta_B$  is an open subset of  $\Theta$  and its closure in  $\Theta$  corresponds to  $\Theta$  itself. By definition, any  $\theta \in \Theta_B$  is such that  $\theta = B^{\dagger/2}XB^{\dagger/2}$  with  $\text{Ran}(X) = \text{Ran}(B)$ . This implies in particular that  $XB^\dagger B = X$  and  $\theta^\dagger = B^{\dagger/2}X^\dagger B^{\dagger/2}$ . Therefore,

$$F(\theta) = \text{Tr}(\theta^\dagger A) + \text{Tr}(\theta B) \quad (42)$$

$$= \text{Tr}(X^\dagger B^{1/2}AB^{1/2}) + \text{Tr}(X), \quad (43)$$

and  $\text{Ran}(B^{1/2}AB^{1/2}) \subseteq \text{Ran}(B) = \text{Ran}(X)$ . We can now minimize the problem with respect to  $X$ , namely

$$\min_{X \in \mathbb{S}_+^d, \text{Ran}(B^{1/2}AB^{1/2}) \subseteq \text{Ran}(X)} \text{Tr}(X^\dagger B^{1/2}AB^{1/2}) + \text{Tr}(X). \quad (44)$$

The minimization corresponds to the variational form of the trace norm of  $B^{1/2}AB^{1/2}$  [29] and has solution  $X_* = (B^{1/2}AB^{1/2})^{1/2}$ , with minimum corresponding to  $2\text{Tr}((B^{1/2}AB^{1/2})^{1/2}) = 2\|B^{1/2}A^{1/2}\|_*$ . To conclude the proof, let  $G : \{X \in \mathbb{S}_+^d \mid \text{Ran}(B^{1/2}AB^{1/2}) \subseteq \text{Ran}(X)\} \rightarrow \mathbb{R}$  be the objective functional in Eq. (44) such that  $G(X) = \text{Tr}(X^\dagger B^{1/2}AB^{1/2}) + \text{Tr}(X)$ . Let now  $X_* \in \mathbb{S}_+^d$  be a minimizer for  $G$  and  $(X_n)_{n \in \mathbb{N}}$  be a minimizing sequence with  $\text{Ran}(X_n) = \text{Ran}(B)$  for each  $n \in \mathbb{N}$  and  $X_n \rightarrow X_*$ . Let  $(\theta_n)_{n \in \mathbb{N}}$  such that  $\theta_n = B^{\dagger/2}X_n B^{\dagger/2}$  for any  $n \in \mathbb{N}$ . Then we have  $\theta_n \rightarrow B^{\dagger/2}X_* B^{\dagger/2}$  and by continuity  $F(B^{\dagger/2}X_* B^{\dagger/2}) = G(X_*)$ , hence  $\min_X G(X) \leq \min_\theta F(\theta)$ . Note that  $B^{\dagger/2}X_* B^{\dagger/2}$  is a minimizer for  $F$ , since  $F$  and  $G$  have same minimum value. To see this it is sufficient to show that, given a minimizing sequence  $(\theta_n)_{n \in \mathbb{N}}$  such that  $\text{Ran}(\theta_n) = \text{Ran}(B)$  for any  $n \in \mathbb{N}$  and  $\theta_n \rightarrow \theta_*$ , we have  $X_n = B^{1/2}\theta_n B^{1/2} \rightarrow B^{1/2}\theta_* B^{1/2}$

and thus  $F(\theta_*) = G(B^{1/2}\theta_*B^{1/2})$ . We have shown that  $\min_{\theta} F(\theta) \geq \min_X G(X)$ . Therefore  $\theta_* = B^{\dagger/2}X_*B^{\dagger/2} = B^{\dagger/2}(B^{1/2}AB^{1/2})^{1/2}B^{\dagger/2}$  is a minimizer of Eq. (38) as desired. The uniqueness of  $\theta_*$  follows from the uniqueness of  $X_*$  from the standard results on the variational form of the trace norm [29].  $\square$

We now have all the ingredients necessary to prove Prop. 2.

**Proposition 2** (Best conditioning function in hindsight). *The conditioning function minimizer and the minimum of the bound presented in Thm. 1 over the set  $\{\tau \in \mathcal{T} \mid \text{Ran}(W(s)) \subseteq \text{Ran}(\tau(s)), \rho_S\text{-almost surely}\}$ , are respectively*

$$\begin{aligned} \tau_{\rho}(s) &= (2L)^{-1} n^{1/2} C(s)^{\dagger/2} (C(s)^{1/2} W(s) C(s)^{1/2})^{1/2} C(s)^{\dagger/2} \\ \mathcal{E}_{\rho}(\tau_{\rho}) - \mathcal{E}_{\rho}^* &\leq 2L \mathbb{E}_{s \sim \rho_S} \|W(s)^{1/2} C(s)^{1/2}\|_* n^{-1/2}. \end{aligned} \quad (8)$$

*Proof.* We aim to minimize

$$\min_{\substack{\tau: \mathcal{S} \rightarrow \Theta \\ \text{Ran}(W(s)) \subseteq \text{Ran}(\tau(s))}} \mathbb{E}_{s \sim \rho_S} \varphi(s, \tau(s)) \quad \text{with} \quad \varphi(s, \theta) = \frac{\text{Tr}(\theta^{\dagger} W(s))}{2} + \frac{2L^2 \text{Tr}(\theta C(s))}{n}. \quad (45)$$

over the set of all measurable functions  $\tau: \mathcal{S} \rightarrow \Theta$ . Note that from Cor. 8, for any  $s \in \mathcal{S}$  we have  $\text{Ran}(W(s)) \subset \text{Ran}(C(s))$ . Therefore we can apply Prop. 10 to have that for any  $s \in \mathcal{S}$ , the problem

$$\min_{\theta \in \mathbb{S}_+^d, \text{Ran}(W(s)) \subseteq \text{Ran}(\theta)} \varphi(s, \theta) \quad (46)$$

has solution

$$\tau_{\rho}(s) = \frac{\sqrt{n}}{2L} C(s)^{\dagger/2} (C(s)^{1/2} W(s) C(s)^{1/2})^{1/2} C(s)^{\dagger/2}. \quad (47)$$

Therefore, for any  $\tau: \mathcal{S} \rightarrow \Theta$  we have

$$\mathbb{E}_{s \sim \rho_S} \varphi(\tau_{\rho}(s), s) \leq \mathbb{E}_{s \sim \rho_S} \varphi(\tau(s), s), \quad (48)$$

and therefore  $\mathbb{E}_{s \sim \rho_S} \varphi(\tau_{\rho}(s), s) \leq \min_{\tau} \mathbb{E}_{s \sim \rho_S} \varphi(\tau)$ . To conclude the proof we need to show that  $\tau_{\rho}$  is measurable. This follows immediately by applying Aumann's measurable selection principle, see for instance the formulation in [34, Lemma A.3.18]. Under the notation of [34], we can apply the result by taking  $h(s, \theta) = (\theta^{\dagger} - I)W(s)$ , the set  $A = \{0\} \subset Y = \mathbb{S}_+^d$ . This guarantees the existence of a measurable function  $\tau_0: \mathcal{S} \rightarrow \Theta$  such that it minimizes pointwise  $\varphi(s, \cdot)$  for any  $s \in \mathcal{S}$  on the set  $\{\theta \in \mathbb{S}_+^d \mid \text{Ran}(W(s)) \subseteq \text{Ran}(\theta)\}$ . The uniqueness of  $\tau_{\rho}(s)$  for each  $s \in \mathcal{S}$  guarantees that  $\tau_{\rho} = \tau_0$  is measurable as desired.  $\square$

## C Proof of Ex. 1

In this section we report the proof of the statement in Ex. 1.

**Example 1** (Clusters). *Let  $\mathcal{S} = \mathbb{R}^q$  be the side information space, for some integer  $q > 0$ . Let  $\rho$  be such that the side information marginal distribution  $\rho_S$  is given by a uniform mixture of  $m$  uniform distributions. More precisely, let  $\rho_S = \frac{1}{m} \sum_{i=1}^m \rho_S^{(i)}$ , with  $\rho_S^{(i)} = \mathcal{U}(\mathcal{B}(a_i, 1/2))$  the uniform distribution on the ball of radius  $1/2$  centered at  $a_i \in \mathcal{S}$ , characterizing the cluster  $i$ . For a given side information  $s$ , a task  $\mu \sim \rho(\cdot|s)$  is sampled such that: 1) its inputs' marginal  $\eta_{\mu}$  is a distribution with constant covariance matrix  $C(s) = \mathbb{E}_{\mu \sim \rho(\cdot|s)} \mathbb{E}_{x \sim \eta_{\mu}} x x^{\top} = C$ , for some  $C \in \mathbb{S}_+^d$ , 2)  $w_{\mu}$  is sampled from a distribution with conditional covariance matrix  $W(s) = \mathbb{E}_{\mu \sim \rho(\cdot|s)} w_{\mu} w_{\mu}^{\top}$ , with  $W(s)$  such that  $(C^{1/2} W(s) C^{1/2})(C^{1/2} W(p) C^{1/2}) = 0$  if  $s \neq p$ . Then,  $\mathbb{E}_{s \sim \rho_S} \|C(s)^{1/2} W(s)^{1/2}\|_* = \frac{1}{\sqrt{m}} \|W_{\rho}^{1/2} C_{\rho}^{1/2}\|_*$ .*

*Proof.* According to the setting described in the example, we can rewrite the following:

$$\begin{aligned}
\mathbb{E}_{s \sim \rho_S} \|C(s)^{1/2} W(s)^{1/2}\|_* &= \mathbb{E}_{s \sim \rho_S} \|C^{1/2} W(s)^{1/2}\|_* \\
&= \mathbb{E}_{s \sim \rho_S} \text{Tr} \left( (C^{1/2} W(s) C^{1/2})^{1/2} \right) \\
&= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{s \sim \rho_S^{(i)}} \text{Tr} \left( (C^{1/2} W(s) C^{1/2})^{1/2} \right) \\
&= \frac{1}{m} \sum_{i=1}^m \text{Tr} \left( (C^{1/2} W(a_i) C^{1/2})^{1/2} \right) \\
&= \frac{1}{m} \text{Tr} \left( \sum_{i=1}^m (C^{1/2} W(a_i) C^{1/2})^{1/2} \right) \\
&= \frac{1}{m} \text{Tr} \left( \left( \sum_{i=1}^m C^{1/2} W(a_i) C^{1/2} \right)^{1/2} \right),
\end{aligned} \tag{49}$$

where, in the first equality we have exploited the fact that  $C(s)$  is a constant matrix  $C$ , in the second equality we have applied the definition of the rewriting of the trace norm of a matrix  $A$  as  $\|A\|_* = \text{Tr}((AA^\top)^{1/2})$ , in the third and fourth equality we have exploited the assumption on  $\rho_S$ , and finally, in the last equality, by point 2), we managed to apply the fact that, for two matrices  $A, B \in \mathbb{S}_+^d$  such that  $A^{1/2} B^{1/2} = B^{1/2} A^{1/2} = 0$ , we have

$$(A^{1/2} + B^{1/2})(A^{1/2} + B^{1/2}) = A + B \implies (A + B)^{1/2} = A^{1/2} + B^{1/2}. \tag{50}$$

On the other hand, we observe that we can also write the following:

$$\begin{aligned}
\|C_\rho^{1/2} W_\rho^{1/2}\|_* &= \|C^{1/2} W_\rho^{1/2}\|_* \\
&= \text{Tr} \left( (C^{1/2} W_\rho C^{1/2})^{1/2} \right) \\
&= \text{Tr} \left( (C^{1/2} \mathbb{E}_{s \sim \rho_S} W(s) C^{1/2})^{1/2} \right) \\
&= \text{Tr} \left( \left( C^{1/2} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{s \sim \rho_S^{(i)}} W(s) C^{1/2} \right)^{1/2} \right) \\
&= \frac{1}{\sqrt{m}} \text{Tr} \left( \left( \sum_{i=1}^m C^{1/2} W(a_i) C^{1/2} \right)^{1/2} \right) \\
&= \frac{1}{\sqrt{m}} \text{Tr} \left( \left( C^{1/2} \sum_{i=1}^m W(a_i) C^{1/2} \right)^{1/2} \right),
\end{aligned} \tag{51}$$

where, in the first equality we have exploited the fact that  $C(s)$  is a constant matrix  $C$ , in the second equality we have applied the definition of the rewriting of the trace norm of a matrix  $A$  as  $\|A\|_* = \text{Tr}((AA^\top)^{1/2})$  and in the fourth and fifth equality we have exploited the assumption on  $\rho_S$ . The desired statement directly derives from combining Eq. (49) and Eq. (51).  $\square$

## D Proofs of the statements in Sec. 4

In this section we report the proofs of the statements we used in Sec. 4 in order to prove the expected excess risk bound for Alg. 1 in Thm. 5. We start from proving the matricial rewriting of Prop. 3 in App. D.1. We then prove in App. D.2 the properties of the surrogate functions in Prop. 4. Then, in App. D.3, we prove the convergence rate of Alg. 1 on the surrogate problem in Eq. (15).

### D.1 Proof of Prop. 3

We start from proving the matricial rewriting of Prop. 3..



**Proposition 3** (Matricial re-formulation of  $\tau_M(s)$ ). *Let  $\tau$  be as in Eq. (11). Then,*

$$\tau(s) = (I_d \otimes \Phi(s)^\top) H_M (I_d \otimes \Phi(s)) + C, \quad (12)$$

where  $I_d$  is the identity in  $\mathbb{R}^{d \times d}$  and  $H_M$  is the matrix in  $\mathbb{R}^{dk \times dk}$  defined by the entries

$$(H_M)_{(i-1)k+h, (j-1)k+z} = \langle M(:, i, h), M(:, j, z) \rangle, \quad i, j = 1, \dots, d, \quad h, z = 1, \dots, k.$$

*Proof.* We start from observing that for any  $i, j = 1, \dots, d$ , we can rewrite the following

$$\begin{aligned} \left( (M\Phi(s))^\top M\Phi(s) \right)_{i,j} &= \langle (M\Phi(s))^\top(i, :), (M\Phi(s))(:, j) \rangle \\ &= \langle (M\Phi(s))(:, i), (M\Phi(s))(:, j) \rangle \\ &= \sum_{q=1}^m (M\Phi(s))(:, i)_q (M\Phi(s))(:, j)_q \\ &= \sum_{q=1}^m \left( \sum_{h=1}^k M_{q,i,h} \Phi(s)_h \right) \left( \sum_{z=1}^k M_{q,j,z} \Phi(s)_z \right) \\ &= \sum_{q=1}^m \sum_{h=1}^k \sum_{z=1}^k M_{q,i,h} M_{q,j,z} \Phi(s)_h \Phi(s)_z \\ &= \sum_{h=1}^k \sum_{z=1}^k \Phi(s)_h \Phi(s)_z \sum_{q=1}^m M_{q,i,h} M_{q,j,z} \\ &= \sum_{h=1}^k \sum_{z=1}^k \Phi(s)_h \Phi(s)_z \left( \sum_{q=1}^m M_{q,i,h} M_{q,j,z} \right) \\ &= \sum_{h=1}^k \sum_{z=1}^k \Phi(s)_h \Phi(s)_z \langle M(:, i, h), M(:, j, z) \rangle. \end{aligned} \quad (52)$$

We now observe that for any  $i, j = 1, \dots, d$ , we can rewrite the following

$$\begin{aligned} \left( (I_d \otimes \Phi(s)^\top) H_M (I_d \otimes \Phi(s)) \right)_{i,j} &= \langle (I_d \otimes \Phi(s)^\top)(i, :), (H_M(I_d \otimes \Phi(s)))(:, j) \rangle \\ &= \langle (I_d \otimes \Phi(s))(:, i), (H_M(I_d \otimes \Phi(s)))(:, j) \rangle \\ &= \sum_{n=1}^{kd} (I_d \otimes \Phi(s))_{n,i} (H_M(I_d \otimes \Phi(s)))_{n,j} \\ &= \sum_{n=1}^{kd} (I_d \otimes \Phi(s))_{n,i} \langle H_M(n, :), (I_d \otimes \Phi(s))(:, j) \rangle \\ &= \sum_{n=1}^{kd} (I_d \otimes \Phi(s))_{n,i} \sum_{p=1}^{kd} (H_M)_{n,p} (I_d \otimes \Phi(s))_{p,j} \\ &= \sum_{n=1}^{kd} \sum_{p=1}^{kd} (I_d \otimes \Phi(s))_{n,i} (H_M)_{n,p} (I_d \otimes \Phi(s))_{p,j} \\ &= \sum_{n=1}^{kd} \sum_{p=1}^{kd} \Phi(s)_h \delta_{n, (i-1)k+h} (H_M)_{n,p} \Phi(s)_z \delta_{p, (j-1)k+z} \\ &= \sum_{h=1}^k \sum_{z=1}^k \Phi(s)_h \Phi(s)_z (H_M)_{(i-1)k+h, (j-1)k+z}, \end{aligned} \quad (53)$$

where, in the seventh equality we have exploited the fact that, by definition,

$$(I_d \otimes \Phi(s))_{n,i} = \begin{cases} \Phi(s)_r & \text{if } r = n - (i-1)k \\ 0 & \text{otherwise} \end{cases} = \Phi(s)_r \delta_{n,r+(i-1)k}. \quad (54)$$

and in the last equality we have defined the new indexes  $h, z = 1, \dots, k$  as

$$h = n - (i-1)k \quad z = p - (j-1)k \quad (55)$$

and, as consequence, we have rewritten

$$n = (i-1)k + h \quad p = (j-1)k + z. \quad (56)$$

As, a consequence, if we define  $H_M$  as the matrix in  $\mathbb{R}^{dk \times dk}$  with entries

$$(H_M)_{(i-1)k+h, (j-1)k+z} = \langle M(:, i, h), M(:, j, z) \rangle, \quad (57)$$

with  $i, j = 1, \dots, d$  and  $h, z = 1, \dots, k$ , then, Eq. (52):

$$\left( (M\Phi(s))^\top M\Phi(s) \right)_{i,j} = \sum_{h=1}^k \sum_{z=1}^k \Phi(s)_h \Phi(s)_z \langle M(:, i, h), M(:, j, z) \rangle \quad (58)$$

and Eq. (53):

$$\left( (I_d \otimes \Phi(s)^\top) H_M (I_d \otimes \Phi(s)) \right)_{i,j} = \sum_{h=1}^k \sum_{z=1}^k \Phi(s)_h \Phi(s)_z (H_M)_{(i-1)k+h, (j-1)k+z} \quad (59)$$

coincide. This coincides with the first desired statement. In order to prove the statement  $H_M \in \mathbb{S}_+^{dk}$ , we show that  $H_M = A_M^\top A_M$ , where  $A_M$  is the matrix in  $\mathbb{R}^{m \times dk}$  defined as

$$A_M(:, (i-1)k + h) = M(:, i, h). \quad (60)$$

We start from recalling that, by definition of  $H_M$ , we have

$$(H_M)_{(i-1)k+h, (j-1)k+z} = \langle M(:, i, h), M(:, j, z) \rangle. \quad (61)$$

Moreover, we observe that, for any  $p, q = 1, \dots, kd$ ,

$$(A_M^\top A_M)_{p,q} = \langle (A_M^\top)(p, :), A_M(:, q) \rangle_{\mathbb{R}^m} = \langle A_M(:, p), A_M(:, q) \rangle. \quad (62)$$

As a consequence, the desired statement is satisfied if we define

$$(A_M)[:, (i-1)k + h] = M(:, i, h). \quad (63)$$

We now prove the last statement. Let  $(e_i)_{i=1}^d$  be the canonical basis in  $\mathbb{R}^d$ . By the definition of the trace and the rewriting of  $\tau(s)$  in Prop. 3, denoting by  $\text{vec}$  the vectorization operation, we can rewrite

$$\begin{aligned} \text{Tr}(\tau(s)) &= \sum_{i=1}^d \langle e_i, \tau(s) e_i \rangle \\ &= \sum_{i=1}^d \langle e_i, (I_d \otimes \Phi(s)^\top) H_M (I_d \otimes \Phi(s)) e_i \rangle \\ &= \sum_{i=1}^d e_i^\top (I_d \otimes \Phi(s)^\top) H_M (I_d \otimes \Phi(s)) e_i \\ &= \sum_{i=1}^d \left( (I_d \otimes \Phi(s)) e_i \right)^\top H_M (I_d \otimes \Phi(s)) e_i \\ &= \sum_{i=1}^d \left( \text{vec}(\Phi(s) e_i^\top) \right)^\top H_M \text{vec}(\Phi(s) e_i^\top) \\ &= \text{Tr} \left( H_M \sum_{i=1}^d \text{vec}(\Phi(s) e_i^\top) \text{vec}(\Phi(s) e_i^\top)^\top \right) \\ &\leq \text{Tr}(H_M) \left\| \sum_{i=1}^d \text{vec}(\Phi(s) e_i^\top) \text{vec}(\Phi(s) e_i^\top)^\top \right\|_\infty \\ &= \text{Tr}(H_M) \|\Phi(s)\|_{\mathbb{R}^k}^2, \end{aligned} \quad (64)$$

where, in the fifth equality, we have applied the relation

$$(C^\top \otimes A) \text{vec}(B) = \text{vec}(ABC) \quad (65)$$

with  $A = \Phi(s)$ ,  $B = e_i^\top$  and  $C = I_d$ , i.e.

$$(I_d \otimes \Phi(s))e_i = \text{vec}(\Phi(s)e_i^\top), \quad (66)$$

in the inequality we have applied Holder's inequality and in the last equality we have applied the following proposition.  $\square$

**Proposition 11.** For any  $i = 1, \dots, d$ , define

$$v_i = \text{vec}(\Phi(s)e_i^\top) \quad (67)$$

Then,

$$\left\| \sum_{i=1}^d \text{vec}(\Phi(s)e_i^\top) \text{vec}(\Phi(s)e_i^\top)^\top \right\|_\infty = \left\| \sum_{i=1}^d v_i v_i^\top \right\|_\infty = \|\Phi(s)\|^2. \quad (68)$$

*Proof.* We start from observing that, for any  $i, j = 1, \dots, d$ , we have

$$\begin{aligned} v_i^\top v_j &= \text{vec}(\Phi(s)e_i^\top)^\top \text{vec}(\Phi(s)e_j^\top) \\ &= \text{Tr}(e_i \Phi(s)^\top \Phi(s) e_j^\top) \\ &= \text{Tr}(\Phi(s)^\top \Phi(s) e_j^\top e_i) \\ &= \Phi(s)^\top \Phi(s) e_j^\top e_i \\ &= \|\Phi(s)\|^2 \delta_{i,j}, \end{aligned} \quad (69)$$

where, in the second equality, we have used the property of the operator  $\text{vec}$ :

$$\text{vec}(A)^\top \text{vec}(B) = \text{Tr}(A^\top B) \quad (70)$$

with

$$A = \Phi(s)e_i^\top \quad B = \Phi(s)e_j^\top. \quad (71)$$

As a consequence, the vectors

$$\tilde{v}_i = \frac{v_i}{\|v_i\|} = \frac{v_i}{\|\Phi(s)\|} \quad i = 1, \dots, d, \quad (72)$$

form an orthonormal basis of the space. Moreover, we can rewrite the operator above as follows

$$\sum_{i=1}^d \text{vec}(\Phi(s)e_i^\top) \text{vec}(\Phi(s)e_i^\top)^\top = \sum_{i=1}^d v_i v_i^\top = \sum_{i=1}^d \|\Phi(s)\|^2 \tilde{v}_i \tilde{v}_i^\top. \quad (73)$$

The rewriting above coincides with the eigenvalues' decomposition of the operator: the vectors  $\tilde{v}_i$  are the eigenvectors with associated constant eigenvalues  $\|\Phi(s)\|^2$ . As a consequence, we can conclude that

$$\left\| \sum_{i=1}^d \text{vec}(\Phi(s)e_i^\top) \text{vec}(\Phi(s)e_i^\top)^\top \right\|_\infty = \|\Phi(s)\|^2. \quad (74)$$

$\square$

## D.2 Proof of Prop. 4

We now prove the properties of the surrogate functions in Prop. 4.

**Proposition 4** (Properties of the surrogate meta-loss  $\mathcal{L}$ ). For any  $Z \in \mathcal{D}$  and  $s \in \mathcal{S}$ , the function  $\mathcal{L}(\cdot, \cdot, s, Z)$  is convex and one of its subgradients is given, for any  $H \in \mathbb{S}_+^{dk}$  and  $C \in \mathbb{S}_+^d$ , by

$$\nabla \mathcal{L}(H, \cdot, s, Z)(C) = \hat{\nabla}, \quad \nabla \mathcal{L}(\cdot, C, s, Z)(H) = (I_d \otimes \Phi(s)) \hat{\nabla} (I_d \otimes \Phi(s)^\top),$$

where

$$\hat{\nabla} = -\frac{\lambda}{2} \tau_{H,C}(s)^\dagger w_{\tau_{H,C}(s)} w_{\tau_{H,C}(s)}^\top \tau_{H,C}(s)^\dagger + \frac{2L^2 X^\top X}{n^2}.$$

Moreover, by Asm. 1 and Asm. 3,

$$\|\nabla \mathcal{L}(\cdot, \cdot, s, Z)(H, C)\|_F \leq (1 + K^2)(LR)^2(2^{-1} + 2n^{-1}).$$

*Proof.* We are interested in studying the properties of the surrogate function  $\mathcal{L}(\cdot, \cdot, s, Z) : \mathbb{S}_+^{dk} \times \mathbb{S}_+^d \rightarrow \mathbb{R}$  in Eq. (15). We start from observing that, such a function coincides with the composition of the function

$$\theta \in \mathbb{S}_+^d \mapsto \Delta(\theta, Z) = F(\theta, Z) + G(\theta, Z) \in \mathbb{R}$$

$$F(\theta, Z) = \min_{w \in \mathbb{R}^d} \mathcal{R}_{Z, \theta}(w) \quad \mathcal{R}_{Z, \theta}(w) = \frac{1}{n} \sum_{i=1}^n \ell(\langle x_i, w \rangle, y_i) + \frac{\lambda}{2} \langle w, \theta^\dagger w \rangle + \iota_{\text{Ran}(\theta)}(w) \quad (75)$$

$$G(\theta, Z) = \frac{2L^2}{n} \text{Tr} \left( \theta \frac{X^\top X}{n} \right).$$

with the linear transformation

$$s \in \mathcal{S} \mapsto \tau_{H, C}(s) = (I_d \otimes \Phi(s)^\top) H (I_d \otimes \Phi(s)) + C \in \mathbb{S}_+^d. \quad (76)$$

In other words, for any  $H \in \mathbb{S}_+^{dk}$  and  $C \in \mathbb{S}_+^d$ , we can write

$$\mathcal{L}(H, C, s, Z) = \Delta(\tau_{H, C}(s), Z) = F(\tau_{H, C}(s), Z) + G(\tau_{H, C}(s), Z). \quad (77)$$

We now observe that both the functions  $F(\cdot, Z)$  and  $G(\cdot, Z)$  are both convex ( $F(\cdot, Z)$  is convex since it is the minimum of a jointly convex function see [15] and  $G(\cdot, Z)$  is a linear function). As a consequence, the function  $\Delta(\cdot, Z)$  is convex over  $\mathbb{S}_+^d$ . This implies the convexity of the surrogate function  $\mathcal{L}(\cdot, \cdot, s, Z)$  over  $\mathbb{S}_+^{dk} \times \mathbb{S}_+^d$  (composition of a convex function with a linear transformation). In order to get the closed form of the gradient in Prop. 4 we proceed in a similar way as in [14]. More precisely, we start from recalling that, as already observed in [15], thanks to strong duality in the within-task problem, for any  $\theta \in \mathbb{S}_+^d$ , we can rewrite

$$F(\theta, Z) = \min_{w \in \text{Ran}(\theta)} \mathcal{R}_{Z, \theta}(w) = \max_{\alpha \in \mathbb{R}^n} \left\{ -\frac{1}{n} \sum_{i=1}^n \ell_i^*(\alpha_i) - \frac{1}{2n^2} \text{Tr}(\theta X^\top \alpha \alpha^\top X) \right\}, \quad (78)$$

where,  $\ell_i^*(\cdot)$  denotes the Fenchel conjugate of  $\ell_i(\cdot) = \ell(\cdot, y_i)$  and  $\alpha \in \mathbb{R}^n$  coincides with the dual variable. As a consequence, we can rewrite

$$\begin{aligned} \Delta(\theta, Z) &= F(\theta, Z) + G(\theta, Z) \\ &= \max_{\alpha \in \mathbb{R}^n} \left\{ -\frac{1}{n} \sum_{i=1}^n \ell_i^*(\alpha_i) - \frac{1}{2n^2} \text{Tr}(\theta X^\top \alpha \alpha^\top X) \right\} + \frac{2L^2}{n} \text{Tr} \left( \theta \frac{X^\top X}{n} \right) \\ &= \max_{\alpha \in \mathbb{R}^n} \left\{ -\frac{1}{n} \sum_{i=1}^n \ell_i^*(\alpha_i) + \text{Tr} \left( \theta \left( -\frac{X^\top \alpha \alpha^\top X}{2n^2} + \frac{2L^2 X^\top X}{n^2} \right) \right) \right\}. \end{aligned} \quad (79)$$

As a consequence, we have

$$\begin{aligned} \Delta(\tau_{H, C}, Z) &= \max_{\alpha \in \mathbb{R}^n} \left\{ -\frac{1}{n} \sum_{i=1}^n \ell_i^*(\alpha_i) + \text{Tr} \left( (I_d \otimes \Phi(s)^\top) H (I_d \otimes \Phi(s)) \left( -\frac{X^\top \alpha \alpha^\top X}{2n^2} + \frac{2L^2 X^\top X}{n^2} \right) \right) \right. \\ &\quad \left. + \text{Tr} \left( C \left( -\frac{X^\top \alpha \alpha^\top X}{2n^2} + \frac{2L^2 X^\top X}{n^2} \right) \right) \right\} \\ &= \max_{\alpha \in \mathbb{R}^n} \left\{ -\frac{1}{n} \sum_{i=1}^n \ell_i^*(\alpha_i) + \text{Tr} \left( H (I_d \otimes \Phi(s)) \left( -\frac{X^\top \alpha \alpha^\top X}{2n^2} + \frac{2L^2 X^\top X}{n^2} \right) (I_d \otimes \Phi(s)^\top) \right) \right. \\ &\quad \left. + \text{Tr} \left( C \left( -\frac{X^\top \alpha \alpha^\top X}{2n^2} + \frac{2L^2 X^\top X}{n^2} \right) \right) \right\} \\ &= \max_{\alpha \in \mathbb{R}^n} Q(\alpha, H, C, s, Z), \end{aligned} \quad (80)$$

where we have introduced the function

$$\begin{aligned} Q(\alpha, H, C, s, Z) &= -\frac{1}{n} \sum_{i=1}^n \ell_i^*(\alpha_i) + \text{Tr} \left( H (I_d \otimes \Phi(s)) \left( -\frac{X^\top \alpha \alpha^\top X}{2n^2} + \frac{2L^2 X^\top X}{n^2} \right) (I_d \otimes \Phi(s)^\top) \right) \\ &\quad + \text{Tr} \left( C \left( -\frac{X^\top \alpha \alpha^\top X}{2n^2} + \frac{2L^2 X^\top X}{n^2} \right) \right). \end{aligned} \quad (81)$$

Hence, applying [15, Lemma 44], we know that, once computed a maximizer  $\alpha_{\tau_{H,C}(s)}$  of the function above  $\alpha \in \mathbb{R}^n \mapsto Q(\alpha, H, C, s, Z)$ ,

$$\nabla Q(\alpha_{\tau_{H,C}(s)}, \cdot, \cdot, s, Z)(H, C) \in \frac{\partial \Delta(\tau_{H,C}(s), Z)}{\partial(H, C)} = \frac{\partial \mathcal{L}(H, C, s, Z)}{\partial(H, C)}. \quad (82)$$

As a consequence, since for a given matrix  $A$ ,  $\nabla \text{Tr}(\cdot A)(H) = A$ , we get that

$$\nabla \mathcal{L}(\cdot, \cdot, s, Z)(H, C) = \left( (I_d \otimes \Phi(s)) \hat{\nabla} (I_d \otimes \Phi(s)^\top), \hat{\nabla} \right) \in \frac{\partial \mathcal{L}(H, C, s, Z)}{\partial(H, C)}, \quad (83)$$

with

$$\hat{\nabla} = -\frac{X^\top \alpha_{\tau_{H,C}(s)} \alpha_{\tau_{H,C}(s)}^\top X}{2n^2} + \frac{2L^2 X^\top X}{n^2}. \quad (84)$$

Finally, in order to get the desired closed form in Prop. 4, we just need to observe that, according to the optimality conditions of the within-task problem in (see [15, Lemma 44]) with  $\theta \in \mathbb{S}_+^d$ , we have that

$$X^\top \alpha_\theta = -n\theta^\dagger w_\theta. \quad (85)$$

As a consequence, we can rewrite Eq. (84) as follows by using the primal solution of the within-task problem:

$$\hat{\nabla} = -\frac{\lambda}{2} \tau_{H,C}(s)^\dagger w_{\tau_{H,C}(s)} w_{\tau_{H,C}(s)}^\top \tau_{H,C}(s)^\dagger + \frac{2L^2 X^\top X}{n^2}. \quad (86)$$

Finally, we observe that, by the closed form in Prop. 4,

$$\|\nabla \mathcal{L}(\cdot, \cdot, s, Z)(H, C)\|_F \leq \|\nabla \mathcal{L}(\cdot, \cdot, s, Z)(H, C)\|_* \leq A + B + C + D \quad (87)$$

with

$$\begin{aligned} A &= \left\| (I_d \otimes \Phi(s)) \frac{X^\top \alpha_{\tau_{H,C}(s)} \alpha_{\tau_{H,C}(s)}^\top X}{2n^2} (I_d \otimes \Phi(s)^\top) \right\|_* \\ B &= \left\| (I_d \otimes \Phi(s)) \frac{2L^2 X^\top X}{n^2} (I_d \otimes \Phi(s)^\top) \right\|_* \\ C &= \left\| \frac{X^\top \alpha_{\tau_{H,C}(s)} \alpha_{\tau_{H,C}(s)}^\top X}{2n^2} \right\|_* \\ D &= \left\| \frac{2L^2 X^\top X}{n^2} \right\|_*. \end{aligned} \quad (88)$$

We now observe that all the matrices inside the trace norms above are positive semidefinite (as a matter of fact, if a matrix  $Q$  is positive semidefinite, then,  $P^\top Q P$  is positive semidefinite for any matrix  $P$ ). As a consequence, all the trace norms above coincide with the trace of the corresponding matrices, namely,

$$\begin{aligned} A &= \text{Tr} \left( (I_d \otimes \Phi(s)) \frac{X^\top \alpha_{\tau_{H,C}(s)} \alpha_{\tau_{H,C}(s)}^\top X}{2n^2} (I_d \otimes \Phi(s)^\top) \right) \\ B &= \text{Tr} \left( (I_d \otimes \Phi(s)) \frac{2L^2 X^\top X}{n^2} (I_d \otimes \Phi(s)^\top) \right) \\ C &= \text{Tr} \left( \frac{X^\top \alpha_{\tau_{H,C}(s)} \alpha_{\tau_{H,C}(s)}^\top X}{2n^2} \right) \\ D &= \text{Tr} \left( \frac{2L^2 X^\top X}{n^2} \right). \end{aligned} \quad (89)$$

We now observe that, proceeding as above in Eq. (64) and exploiting Asm. 3, we can write

$$\begin{aligned} A &\leq \|\Phi(s)\|^2 \text{Tr} \left( \frac{X^\top \alpha_{\tau_{H,C}(s)} \alpha_{\tau_{H,C}(s)}^\top X}{2n^2} \right) = \|\Phi(s)\|^2 C \leq K^2 C \\ B &\leq \|\Phi(s)\|^2 \text{Tr} \left( \frac{2L^2 X^\top X}{n^2} \right) = \|\Phi(s)\|^2 D \leq K^2 D. \end{aligned} \quad (90)$$

Hence, combining everything in Eq. (87), we get

$$\|\nabla \mathcal{L}(\cdot, \cdot, s, Z)(H, C)\|_F \leq (1 + K^2)(C + D). \quad (91)$$

The desired statement derives from observing that, since, by Asm. 1,  $\text{Tr}(X^\top \alpha_{\tau_{H,C}(s)} \alpha_{\tau_{H,C}(s)}^\top X) \leq (nLR)^2$  (see [15, Lemma 44]) and  $\text{Tr}(X^\top X) = \text{Tr}(XX^\top) = \sum_{i=1}^n \|x_i\|^2 \leq nR^2$ , then

$$C \leq \frac{(LR)^2}{2\lambda} \quad D \leq \frac{2(LR)^2}{n}. \quad (92)$$

□

### D.3 Convergence rate of Alg. 1 on the surrogate problem in Eq. (15)

We now give the convergence rate of Alg. 1 on the surrogate problem in Eq. (15).

**Proposition 12** (Convergence rate on the surrogate problem in Eq. (15)). *Let  $\bar{H}$  and  $\bar{C}$  be the average of the iterations obtained from the application of Alg. 1 over the training data  $(Z_t, s_t)_{t=1}^T$  with constant meta-step size  $\gamma > 0$ . Then, under Asm. 1 and Asm. 3, for any  $\tau_{H,C} \in \mathcal{T}_\Phi$ , in expectation with respect to the sampling of  $(Z_t, s_t)_{t=1}^T$ ,*

$$\mathbb{E} \hat{\mathcal{E}}_\rho(\tau_{\bar{H}, \bar{C}}) - \hat{\mathcal{E}}_\rho(\tau_{H,C}) \leq \frac{\gamma(1 + K^2)^2(LR)^4}{2\lambda^2} \left( \frac{1}{2} + \frac{2}{n} \right)^2 + \frac{\|(H - H_0, C - C_0)\|_F^2}{2\gamma T}. \quad (93)$$

*Proof.* We observe that Alg. 1 coincides with projected Stochastic Gradient Descent applied to the convex and Lipschitz (see Prop. 4) surrogate problem in Eq. (15):

$$\min_{H \in \mathbb{S}_+^{d_k}, C \in \mathbb{S}_+^d} \hat{\mathcal{E}}_\rho(\tau_{H,C}) \quad \hat{\mathcal{E}}_\rho(\tau_{H,C}) = \mathbb{E}_{(\mu, s) \sim \rho} \mathbb{E}_{Z \sim \mu^n} \mathcal{L}(H, C, s, Z). \quad (94)$$

As a consequence, by standard arguments (see e.g. [33, Lemma 14.1, Thm. 14.8] and references therein), for any  $\tau_{H,C} \in \mathcal{T}_\Phi$ , we have

$$\mathbb{E} \hat{\mathcal{E}}_\rho(\tau_{\bar{H}, \bar{C}}) - \hat{\mathcal{E}}_\rho(\tau_{H,C}) \leq \frac{\gamma}{2T} \sum_{t=1}^T \mathbb{E} \|\nabla \mathcal{L}(\cdot, \cdot, s, Z_t)(H_t, C_t)\|_F^2 + \frac{\|(H - H_0, C - C_0)\|_F^2}{2\gamma T}. \quad (95)$$

The desired statement derives from combining this bound with the bound on the norm of the meta-subgradients in Prop. 4. □

### D.4 Proof of Thm. 5

We now have all the ingredients necessary to prove Thm. 5.

**Theorem 5** (Excess risk bound for the conditioning function returned by Alg. 1). *Let Asm. 1 and Asm. 3 hold. For any  $s \sim \rho_S$ , recall the conditional covariance matrices  $W(s)$  and  $C(s)$  introduced in Thm. 1. Let  $\tau_{H,C}$  be a fixed function in  $\mathcal{T}_\Phi$  such that  $\text{Ran}(W(s)) \subseteq \text{Ran}(\tau_{H,C}(s))$  for any  $s \sim \rho_S$ . Let  $\bar{H}$  and  $\bar{C}$  be the outputs of Alg. 1 applied to a sequence  $(Z_t, s_t)_{t=1}^T$  of i.i.d. pairs sampled from  $\rho$  with an appropriate meta-step size  $\gamma$ . Then, in expectation with respect to the sampling of  $(Z_t, s_t)_{t=1}^T$ ,*

$$\begin{aligned} \mathbb{E} \mathcal{E}_\rho(\tau_{\bar{H}, \bar{C}}) - \mathcal{E}_\rho^* &\leq \frac{\mathbb{E}_{s \sim \rho_S} \text{Tr}(\tau_{H,C}(s)^\dagger W(s))}{2} + \frac{2L^2 \mathbb{E}_{s \sim \rho_S} \text{Tr}(\tau_{H,C}(s) C(s))}{n} \\ &\quad + \left( \frac{1}{2} + \frac{2}{n} \right) \frac{(1 + K^2)(LR)^2 \|(H - H_0, C - C_0)\|_F}{\sqrt{T}}. \end{aligned}$$



*Proof.* We start from observing that, in expectation with respect to the meta-training set, for any fixed conditioning function  $\tau_{H,C} \in \mathcal{T}_\Phi$ , we can write the following decomposition

$$\begin{aligned}\mathbb{E} \mathcal{E}_\rho(\tau_{\bar{H},\bar{C}}) - \mathcal{E}_\rho^* &\leq \mathbb{E} \hat{\mathcal{E}}_\rho(\tau_{\bar{H},\bar{C}}) - \mathcal{E}_\rho^* \\ &= \mathbb{E} \hat{\mathcal{E}}_\rho(\tau_{\bar{H},\bar{C}}) - \mathcal{E}_\rho^* \pm \hat{\mathcal{E}}_\rho(\tau_{H,C}) \\ &= \underbrace{\mathbb{E} \hat{\mathcal{E}}_\rho(\tau_{\bar{H},\bar{C}}) - \hat{\mathcal{E}}_\rho(\tau_{H,C})}_{A(\tau_{H,C})} + \underbrace{\hat{\mathcal{E}}_\rho(\tau_{H,C}) - \mathcal{E}_\rho^*}_{B(\tau_{H,C})},\end{aligned}\tag{96}$$

where in the inequality above we have exploited the fact that, for any  $\tau \in \mathcal{T}$ ,  $\mathcal{E}_\rho(\tau) \leq \hat{\mathcal{E}}_\rho(\tau)$ . We now observe that the term  $A(\tau_{H,C})$  can be controlled according to the convergence properties of the meta-algorithm in Alg. 1 as described in Prop. 12:

$$\mathbb{E} \hat{\mathcal{E}}_\rho(\tau_{\bar{H},\bar{C}}) - \hat{\mathcal{E}}_\rho(\tau_{H,C}) \leq \frac{\gamma(1+K^2)^2(LR)^4}{2} \left(\frac{1}{2} + \frac{2}{n}\right)^2 + \frac{\|(H - H_0, C - C_0)\|_F^2}{2\gamma T}.\tag{97}$$

Regarding the term  $B(\tau_{H,C})$ , we observe that, for any  $\tau$ , we can rewrite

$$\begin{aligned}B(\tau) &= \hat{\mathcal{E}}_\rho(\tau) - \mathcal{E}_\rho^* \\ &= \mathbb{E}_{(\mu,s) \sim \rho} \mathbb{E}_{Z \sim \mu^n} \left[ \mathcal{R}_{Z,\tau(s)}(A(\tau(s), Z)) - \mathcal{R}_\mu(w_\mu) \right] + \frac{2L^2 \mathbb{E}_{(\mu,s) \sim \rho} \text{Tr}(\tau(s) \mathbb{E}_{x \sim \eta_\mu} x x^\top)}{n} \\ &\leq \frac{\mathbb{E}_{(\mu,s) \sim \rho} \text{Tr}(\tau(s)^\dagger w_\mu w_\mu^\top)}{2} + \frac{2L^2 \mathbb{E}_{(\mu,s) \sim \rho} \text{Tr}(\tau(s) \mathbb{E}_{x \sim \eta_\mu} x x^\top)}{n},\end{aligned}\tag{98}$$

where in the inequality we have exploited the fact that, thanks to the definition of the algorithm, for any  $(\mu, s) \sim \rho$ , we can write

$$\mathbb{E}_{Z \sim \mu^n} \left[ \mathcal{R}_{Z,\tau(s)}(A(\tau(s), Z)) - \mathcal{R}_\mu(w_\mu) \right] \leq \frac{\text{Tr}(\tau(s)^\dagger w_\mu w_\mu^\top)}{2}.\tag{99}$$

Combining the bounds on the two terms above in Eq. (96), we get

$$\begin{aligned}\mathbb{E} \mathcal{E}_\rho(\tau_{\bar{H},\bar{C}}) - \mathcal{E}_\rho^* &\leq \frac{\mathbb{E}_{(\mu,s) \sim \rho} \text{Tr}(\tau_{H,C}(s)^\dagger w_\mu w_\mu^\top)}{2} + \frac{2L^2 \mathbb{E}_{(\mu,s) \sim \rho} \text{Tr}(\tau_{H,C}(s) \mathbb{E}_{x \sim \eta_\mu} x x^\top)}{n} \\ &\quad + \frac{\gamma(1+K^2)^2(LR)^4}{2} \left(\frac{1}{2} + \frac{2}{n}\right)^2 + \frac{\|(H - H_0, C - C_0)\|_F^2}{2\gamma T}.\end{aligned}\tag{100}$$

Optimizing with respect to the hyper-parameter  $\gamma > 0$ , for

$$\gamma = \frac{\|(H - H_0, C - C_0)\|_F}{(1+K^2)(LR)^2} \left(\frac{1}{2} + \frac{2}{n}\right)^{-1} \frac{1}{\sqrt{T}},\tag{101}$$

we get the desired statement.  $\square$

## E Experimental details

We now report the experimental details we missed in the main body. Specifically, we report the details regarding the tuning of the hyper-parameter  $\gamma$  and the characteristics of the machine we used for running our experiments.

**Synthetic clusters.** In order to tune the hyper-parameter  $\gamma$  we applied the procedure above with 14 candidates values for  $\gamma$  in the range  $[10^{-5}, 10^5]$  with logarithmic spacing and we evaluated the performance of the estimated meta-parameters (linear representations) by using  $T = T_{\text{tr}} = 500$ ,  $T_{\text{va}} = 300$ ,  $T_{\text{te}} = 100$  of the available tasks for meta-training, meta-validation and meta-testing, respectively. In order to train and to test the inner algorithm, we splitted each within-task dataset into  $n = n_{\text{tr}} = 50\% n_{\text{tot}}$  for training and  $n_{\text{te}} = 50\% n_{\text{tot}}$  for test.

**Lenk dataset.** In order to tune the hyper-parameter  $\gamma$  we applied the procedure above with 14 candidates values for  $\gamma$  in the range  $[10^{-5}, 10^5]$  with logarithmic spacing and we evaluated the

performance of the estimated meta-parameters (linear representations) by using  $T = T_{\text{tr}} = 100$ ,  $T_{\text{va}} = 40$ ,  $T_{\text{te}} = 30$  of the available tasks for meta-training, meta-validation and meta-testing, respectively. In order to train and to test the inner algorithm, we splitted each within-task dataset into  $n = n_{\text{tr}} = 16$  for training and  $n_{\text{te}} = 4$  for test.

**Movielens-100k dataset.** In order to tune the hyper-parameter  $\gamma$  we applied the procedure above with 14 candidates values for  $\gamma$  in the range  $[10^{-5}, 10^5]$  with logarithmic spacing and we evaluated the performance of the estimated meta-parameters (linear representations) by using  $T = T_{\text{tr}} = 200$ ,  $T_{\text{va}} = 100$ ,  $T_{\text{te}} = 100$  of the available tasks for meta-training, meta-validation and meta-testing, respectively. In order to train and to test the inner algorithm, we splitted each within-task dataset into  $n = n_{\text{tr}} = 15$  for training and  $n_{\text{te}} = 5$  for test.

**Jester-1 dataset.** In order to tune the hyper-parameter  $\gamma$  we applied the procedure above with 14 candidates values for  $\gamma$  in the range  $[10^{-5}, 10^5]$  with logarithmic spacing and we evaluated the performance of the estimated meta-parameters (linear representations) by using  $T = T_{\text{tr}} = 250$ ,  $T_{\text{va}} = 100$ ,  $T_{\text{te}} = 100$  of the available tasks for meta-training, meta-validation and meta-testing, respectively. In order to train and to test the inner algorithm, we splitted each within-task dataset into  $n = n_{\text{tr}} = 15$  for training and  $n_{\text{te}} = 5$  for test.

All the experiments were conducted on a workstation with 4 Intel Xeon E5-2697 V3 2.60Ghz CPUs and 256GB RAM.