

Appendix

In this appendix, we report proofs of the results presented in the paper and also additional remarks. We start giving some auxiliary results in [App. A](#) and the proofs of the statements made in [Sec. 2](#) and [Sec. 3](#) in [App. B](#). In [App. C](#) we present the convergence rate of SGD for a general vector-valued LS problem (also known as Least Mean Squares (LMS) algorithm). This rate is then specialized in [App. D](#) to our [Alg. 1](#) minimizing the LS function \mathcal{E}_r . [App. E](#) contains the proof of the main result in the paper ([Prop. 3](#)) and, finally, in [App. F](#), we describe how to tune the hyper-parameters in the LTL setting described in the paper by a cross-validation procedure. In the sequel we will use the notation already introduced in the paper.

A Auxiliary Results

In this section, we give some technical tools that will be used in the sequel of this appendix. We denote by \mathbb{S}^d , \mathbb{S}_{++}^d and \mathbb{S}_+^d the sets of the real $d \times d$ symmetric, symmetric positive definite and symmetric positive semidefinite matrices, respectively. We refer to the book [\[5\]](#) for more details.

Lemma 4. *Let $U \in \mathbb{S}_+^d$ and let $V \in \mathbb{R}^{d \times m}$, then $V^\top UV \in \mathbb{S}_+^m$.*

Proof. For any $a \in \mathbb{R}^m$, defining $b = Va \in \mathbb{R}^d$ and using the assumption $U \in \mathbb{S}_+^d$, we have that $a^\top V^\top UV a = b^\top U b \geq 0$. This proves the desired statement. \blacksquare

A direct consequence of the above remark is the following fact.

Corollary 5. *Let $W \in \mathbb{S}_+^d$, then $W^2 \preceq \|W\|_\infty W$.*

Proof. The statement directly follows from applying [Lemma 4](#) with $V = W^{1/2}$ and $U = \|W\|_\infty I - W$ and using the fact $W \preceq \|W\|_\infty I$. Specifically, $W^2 = W^{1/2} W W^{1/2} \preceq \|W\|_\infty W$. \blacksquare

Lemma 6. *Let $A \in \mathbb{S}_{++}^d$ and $B \in \mathbb{S}_+^d$ such that $AB = BA$, then $AB \in \mathbb{S}_+^d$.*

Proof. We first note that AB is symmetric. Indeed, we have that $(AB)^\top = B^\top A^\top = BA = AB$. Now, we observe that AB is similar to $A^{1/2} B A^{1/2}$, specifically, we have

$$AB = A^{1/2} A^{-1/2} A B A^{1/2} A^{-1/2} \sim A^{-1/2} A B A^{1/2} = A^{1/2} B A^{1/2}.$$

Consequently, the eigenvalues of AB are the same of $A^{1/2} B A^{1/2}$ and the statement follows from [Lemma 4](#), according to which $A^{1/2} B A^{1/2} \in \mathbb{S}_+^d$. \blacksquare

Lemma 7. *The function $\lambda_{\min} : \mathbb{S}^d \rightarrow \mathbb{R}$ is matrix-concave.*

Proof. We recall that for any $A \in \mathbb{S}^d$, we can write

$$\lambda_{\min}(A) = \inf_{v \in \mathbb{R}^d \setminus \{0\}} \frac{v^\top A v}{v^\top v} = \inf_{v \in \mathbb{R}^d \setminus \{0\}} f_v(A),$$

where, for any $v \in \mathbb{R}^d \setminus \{0\}$, we have introduced the linear functions $f_v : \mathbb{S}^d \rightarrow \mathbb{R}$, $A \in \mathbb{S}^d \mapsto \frac{v^\top A v}{v^\top v} \in \mathbb{R}$. Hence, the concavity of $\lambda_{\min}(\cdot)$ derives from the fact that the infimum of linear functions is a concave function. \blacksquare

Lemma 8. *Let $A \in \mathbb{S}_{++}^d$ and $B \in \mathbb{S}_+^d$. Then, $ABA \in \mathbb{S}_+^d$ and the following lower bound on the smallest eigenvalue of ABA holds*

$$\lambda_{\min}(ABA) \geq \lambda_{\min}(B) \lambda_{\min}(A)^2.$$

Proof. We start observing that $ABA \in \mathbb{S}_+^d$, thanks to [Lemma 4](#). Moreover, as already observed in the above lemma, we have that

$$\begin{aligned} \lambda_{\min}(ABA) &= \inf_{v \in \mathbb{R}^d \setminus \{0\}} \frac{v^\top ABAv}{v^\top v} = \inf_{v \in \mathbb{R}^d \setminus \{0\}} \frac{v^\top ABAv}{v^\top A^2v} \frac{v^\top A^2v}{v^\top v} \\ &\geq \inf_{v \in \mathbb{R}^d \setminus \{0\}} \frac{v^\top ABAv}{v^\top A^2v} \inf_{w \in \mathbb{R}^d \setminus \{0\}} \frac{w^\top A^2w}{w^\top w} \\ &= \inf_{u \in \mathbb{R}^d \setminus \{0\}} \frac{u^\top Bu}{u^\top u} \inf_{w \in \mathbb{R}^d \setminus \{0\}} \frac{w^\top A^2w}{w^\top w} \\ &= \lambda_{\min}(B)\lambda_{\min}(A^2), \end{aligned}$$

where the second and the third equality hold thanks to the fact that $A \in \mathbb{S}_{++}^d$, hence $A^2 \in \mathbb{S}_{++}^d$, and consequently, for any $v \in \mathbb{R}^d \setminus \{0\}$, $Av \neq 0$. The inequality above holds since, for any two non-negative functions g and f , we have that $\inf_v g(v)f(v) \geq \inf_v g(v) \inf_w f(w)$; in our case, the two functions are non-negative since $A^2 \in \mathbb{S}_{++}^d$ and, as already observed, $ABA \in \mathbb{S}_+^d$. The statement directly follows observing that $\lambda_{\min}(A^2) = \lambda_{\min}(A)^2$. \blacksquare

The following result is an extension of the Neumann series to matrices.

Lemma 9 (Exercise I.2.6 in [\[5\]](#)). *Let $A \in \mathbb{S}^d$ such that $0 \prec A \prec I$, then*

$$\sum_{k=1}^{\infty} (I - A)^k = A^{-1}(I - A).$$

We conclude this section with a technical lemma which will be used in the following.

Lemma 10. *For the indexes $t, j \in \{i_{\min}, \dots, i_{\max}\}$ consider the sequence of numbers $(a_{t,j})_{t,j}$, where, for any $t, j \in \{i_{\min}, \dots, i_{\max}\}$, $a_{t,j} \in \mathbb{R}$ and $a_{t,j} = a_{j,t}$. Then, we can write*

$$\sum_{t=i_{\min}}^{i_{\max}} \sum_{j=i_{\min}}^{i_{\max}} a_{t,j} = \sum_{t=i_{\min}}^{i_{\max}} a_{t,t} + 2 \sum_{t=i_{\min}}^{i_{\max}-1} \sum_{j=t+1}^{i_{\max}} a_{t,j}.$$

Proof. We can interpret each element $a_{t,j}$ in the sequence as the entry $A_{t,j}$ of a symmetric matrix A with rows and columns indexes in the set $\{i_{\min}, \dots, i_{\max}\}$. The statement follows by first summing the elements in the diagonal and then the elements of the upper and lower triangular parts, observing that, thanks to the symmetry of the matrix, these two last parts bring exactly the same contribution in the sum. Specifically, we have that

$$\sum_{t=i_{\min}}^{i_{\max}} \sum_{j=i_{\min}}^{i_{\max}} a_{t,j} = \sum_{t=i_{\min}}^{i_{\max}} a_{t,t} + \sum_{t=i_{\min}}^{i_{\max}} \sum_{j \neq t} a_{t,j},$$

where, thanks to the symmetry,

$$\sum_{t=i_{\min}}^{i_{\max}} \sum_{j \neq t} a_{t,j} = 2 \sum_{t=i_{\min}}^{i_{\max}-1} \sum_{j=t+1}^{i_{\max}} a_{t,j}.$$

\blacksquare

B Proofs of Results in [Sec. 2](#) and [Sec. 3](#)

In order to prove both [Prop. 1](#) and [Prop. 2](#), we will need the following result, which will be used also in the proof of [Prop. 3](#) in [App. E](#).

Lemma 11 (A Lower Bound on $\lambda_{\min}(\bar{\Sigma}_r)$ and an Upper Bound on $\|\bar{\Sigma}_r\|_{\infty}$). *Let $\mathcal{X} \subseteq \mathcal{B}_1$ and let, for any $r \in \{0\} \cup [n]$, $\bar{\Sigma}_r$ be defined as in [Prop. 1](#) and [Prop. 2](#). Then, for any $r \in \{0\} \cup [n]$, the following inequalities hold.*

$$\frac{\lambda^2 \mathbb{E}_{\mu \sim \rho} \lambda_{\min}(\Sigma_{\mu})}{(r/n + \lambda)^2} \leq \lambda_{\min}(\bar{\Sigma}_r) \leq \|\bar{\Sigma}_r\|_{\infty} \leq \mathbb{E}_{\mu \sim \rho} \|\Sigma_{\mu}\|_{\infty}.$$

Proof. We recall by definition that, for any $r \in [n]$,

$$\bar{\Sigma}_r = \lambda^2 \mathbb{E} \left[C_{\lambda,r}^{-1} \Sigma_\mu C_{\lambda,r}^{-1} \right].$$

We start from proving the lower bound on $\lambda_{\min}(\bar{\Sigma}_r)$. As already observed in [Lemma 7](#), the function $\lambda_{\min} : \mathbb{S}^d \rightarrow \mathbb{R}$ is matrix-concave, hence, by Jensen's inequality, for any $M \in \mathbb{S}^d$, we have that

$$\lambda_{\min}(\mathbb{E}[M]) \geq \mathbb{E}[\lambda_{\min}(M)].$$

Consequently, using the facts $\|X_r\|_\infty^2 \leq \|X_r\|^2 \leq r$ (since $\mathcal{X} \subseteq \mathcal{B}_1$) and $\lambda_{\min}(C_{\lambda,r}^{-1}) = \|C_{\lambda,r}\|_\infty^{-1}$ (since $C_{\lambda,r}^{-1} \in \mathbb{S}_{++}^d$) and applying [Lemma 8](#) to the matrices $A = C_{\lambda,r}^{-1} \in \mathbb{S}_{++}^d$ and $B = \Sigma_\mu \in \mathbb{S}_+^d$, we obtain that

$$\begin{aligned} \lambda_{\min} \left(\mathbb{E} \left[C_{\lambda,r}^{-1} \Sigma_\mu C_{\lambda,r}^{-1} \right] \right) &\geq \mathbb{E} \left[\lambda_{\min} \left(C_{\lambda,r}^{-1} \Sigma_\mu C_{\lambda,r}^{-1} \right) \right] \geq \mathbb{E} \left[\lambda_{\min}(\Sigma_\mu) \lambda_{\min}(C_{\lambda,r}^{-1})^2 \right] \\ &= \mathbb{E} \left[\lambda_{\min}(\Sigma_\mu) \left(\|C_{\lambda,r}\|_\infty \right)^{-2} \right] = \mathbb{E} \left[\lambda_{\min}(\Sigma_\mu) \left(\|X_r\|_\infty^2 / n + \lambda \right)^{-2} \right] \\ &\geq \frac{\mathbb{E}_{\mu \sim \rho} \lambda_{\min}(\Sigma_\mu)}{(r/n + \lambda)^2}. \end{aligned}$$

Multiplying by λ^2 , the result follows. As regards the upper bound on $\|\bar{\Sigma}_r\|_\infty$, we proceed as follows. Exploiting the fact $\|C_{\lambda,r}^{-1}\|_\infty \leq 1/\lambda$, we can write

$$\|\bar{\Sigma}_r\|_\infty \leq \lambda^2 \mathbb{E} \left[\|C_{\lambda,r}^{-1} \Sigma_\mu C_{\lambda,r}^{-1}\|_\infty \right] \leq \lambda^2 \mathbb{E} \left[\|C_{\lambda,r}^{-1}\|_\infty \|\Sigma_\mu\|_\infty \|C_{\lambda,r}^{-1}\|_\infty \right] \leq \mathbb{E}_{\mu \sim \rho} \|\Sigma_\mu\|_\infty. \quad (15)$$

We observe that the previous steps hold also for the extreme case $r = 0$, where $C_{\lambda,0} = \lambda I$ and $\bar{\Sigma}_0 = \mathbb{E}_{\mu \sim \rho} \Sigma_\mu$. However, in such a case, we can also directly get the statement in the proposition by simply applying Jensen's inequality, getting both $\lambda_{\min}(\mathbb{E}_{\mu \sim \rho} \Sigma_\mu) \geq \mathbb{E}_{\mu \sim \rho} \lambda_{\min}(\Sigma_\mu)$ and $\|\mathbb{E}_{\mu \sim \rho} \Sigma_\mu\|_\infty \leq \mathbb{E}_{\mu \sim \rho} \|\Sigma_\mu\|_\infty$. \blacksquare

We now are ready to prove [Prop. 1](#).

Proposition 1 (LS Problem Around a Common Mean for \mathcal{E}_n). *For any $\lambda > 0$ and $h \in \mathbb{R}^d$, the transfer risk \mathcal{E}_n in Eq. (1) of the learning algorithm w_h in Eqs. (3)-(4), can be rewritten as*

$$\mathcal{E}_n(h) = \frac{1}{2} \mathbb{E}_{\bar{x}_n, \bar{y}_n} \left[(\langle \bar{x}_n, h \rangle - \bar{y}_n)^2 \right] \quad (5)$$

where the meta-data are given by

$$\bar{x}_n = \lambda C_{\lambda,n}^{-1} x, \quad \text{and} \quad \bar{y}_n = y - \left\langle \bar{x}_n, \frac{X_n^\top \mathbf{y}_n}{\lambda n} \right\rangle.$$

Moreover, under [Asm. 1](#), the meta-covariance matrix $\bar{\Sigma}_n = \mathbb{E} \bar{x}_n \bar{x}_n^\top = \lambda^2 \mathbb{E} [C_{\lambda,n}^{-1} \Sigma_\mu C_{\lambda,n}^{-1}]$ is invertible and $h_n^* = \bar{\Sigma}_n^{-1} \mathbb{E} [\bar{y}_n \bar{x}_n]$ is the unique minimizer of the LS function in Eq. (5). In such a case, letting $v = w - h_n^*$, we have that $\bar{y}_n = \langle \bar{x}_n, h_n^* \rangle + \bar{\epsilon}_n$, with $\bar{\epsilon}_n = \epsilon + \left\langle \bar{x}_n, v - \frac{X_n^\top \epsilon_n}{\lambda n} \right\rangle$.

Proof. The rewriting of the transfer risk in Eq. (5) is a direct consequence of the closed form of the algorithm in Eq. (4). The invertibility of the meta-covariance matrix is a direct consequence of [Lemma 11](#) for $r = n$ and the requirement $\mathbb{E}_{\mu \sim \rho} \lambda_{\min}(\Sigma_\mu) > 0$ in [Asm. 1](#). This implies that the LS function is strictly convex and, consequently, it admits a unique minimizer h_n^* . The closed form of this minimizer directly follows from the optimality conditions (normal equations) of the problem. Now, thanks to [Asm. 1](#), using the linear model equations

$$y = \langle x, w \rangle + \epsilon \quad \mathbf{y}_n = X_n w + \epsilon_n$$

and letting $w = h_n^* + v$, for some $v \in \mathbb{R}^d$, we can rewrite

$$\begin{aligned}
\bar{y}_n &= y - \left\langle \bar{x}_n, \frac{X_n^\top \mathbf{y}_n}{\lambda n} \right\rangle = y - \left\langle x, C_{\lambda,n}^{-1} \frac{X_n^\top \mathbf{y}_n}{n} \right\rangle \\
&= \langle x, w \rangle + \epsilon - \left\langle x, C_{\lambda,n}^{-1} \frac{X_n^\top (X_n w + \epsilon_n)}{n} \right\rangle \\
&= \left\langle x, C_{\lambda,n}^{-1} \left(C_{\lambda,n} - \frac{X_n^\top X_n}{n} \right) w \right\rangle + \epsilon - \left\langle x, C_{\lambda,n}^{-1} \frac{X_n^\top \epsilon_n}{n} \right\rangle \\
&= \langle x, \lambda C_{\lambda,n}^{-1} w \rangle + \epsilon - \left\langle x, C_{\lambda,n}^{-1} \frac{X_n^\top \epsilon_n}{n} \right\rangle \\
&= \langle \bar{x}_n, w \rangle + \epsilon - \left\langle \bar{x}_n, \frac{X_n^\top \epsilon_n}{\lambda n} \right\rangle \\
&= \langle \bar{x}_n, h_n^* \rangle + \epsilon + \left\langle \bar{x}_n, v - \frac{X_n^\top \epsilon_n}{\lambda n} \right\rangle.
\end{aligned} \tag{16}$$

We conclude that the noise on the meta-labels is heteroscedastic, since it is given by

$$\bar{\epsilon}_n = \bar{y}_n - \langle \bar{x}_n, h_n^* \rangle = \epsilon + \left\langle \bar{x}_n, v - \frac{X_n^\top \epsilon_n}{\lambda n} \right\rangle. \quad \blacksquare$$

We now proceed with the proof of [Ex. 1](#) and the remarks associated to it. Before doing this, we point out the following aspect which will be used throughout the appendix.

Remark 8. According to the data-generation procedure described in [Asm. 1](#), the samplings of the quantities x, w, ϵ are independent one each other, conditioning with respect to the marginal distribution p .

Example 1. Let $\mathcal{X} \subseteq \mathcal{B}_1$ and let the environment p satisfy [Asm. 1](#). Furthermore, assume for almost every p that: i) $\eta \mid p$ has variance bounded by σ_ϵ^2 , for $\sigma_\epsilon \geq 0$, ii) $\mathbb{E}[w \mid p] = \bar{w}$, and iii) $\mathbb{E}[(w - \bar{w})(w - \bar{w})^\top \mid p] \preceq \sigma_w^2 I$, for $\sigma_w \geq 0$. Then, for any $\lambda > 0$ and $h \in \mathbb{R}^d$, $h_n^* = \bar{w}$ and, letting $A_n = C_{\lambda,n}^{-1} x x^\top C_{\lambda,n}^{-1}$, we have

$$\mathcal{E}_n(h) \leq \frac{1}{2} \|\bar{\Sigma}_n^{1/2}(h - \bar{w})\|^2 + \frac{\sigma_w^2}{2} \text{tr}(\bar{\Sigma}_n) + \frac{\sigma_\epsilon^2}{2} \left(1 + \text{tr} \left(\mathbb{E} \left[\frac{X_n^\top X_n A_n}{n^2} \right] \right) \right). \tag{6}$$

Proof. In the following, because of readability, we condense all the expectations (the one according the sampling of the task $\mu = (w, p, \eta) \sim \rho$, the one referring to the sampling of the training dataset $Z_n \sim \mu^n$ and the one related to the sampling of the test point $z \sim \mu$) in only one symbol. In this way, the transfer risk \mathcal{E}_n of the algorithm w_h on any environment satisfying [Asm. 1](#) can be rewritten as

$$\begin{aligned}
\mathcal{E}_n(h) &= \frac{1}{2} \mathbb{E} \left[(\langle x, w_h(Z_n) \rangle - y)^2 \right] \\
&= \frac{1}{2} \mathbb{E} \left[(\langle x, w_h(Z_n) \rangle - \langle x, w \rangle - \epsilon)^2 \right] \\
&= \frac{1}{2} \mathbb{E} \left[(w_h(Z_n) - w)^\top x x^\top (w_h(Z_n) - w) \right] + \frac{\mathbb{E}[\epsilon^2]}{2},
\end{aligned}$$

where in the second equality we have used the linear model equation $y = \langle x, w \rangle + \epsilon$ and in the third equality we have exploited the fact that the noise is zero-mean, more precisely, thanks to [Rem. 8](#),

$$\mathbb{E} \left[\epsilon \langle x, w_h(Z_n) - w \rangle \right] = \mathbb{E} \left[\mathbb{E}[\epsilon \langle x, w_h(Z_n) - w \rangle \mid p] \right] = \mathbb{E} \left[\mathbb{E}[\epsilon \mid p] \mathbb{E}[\langle x, w_h(Z_n) - w \rangle \mid p] \right] = 0.$$

Using the closed form of the algorithm in [Eq. \(4\)](#) and the equation $\mathbf{y}_n = X_n w + \epsilon_n$ deriving from [Asm. 1](#), a direct computation gives that

$$w_h(Z_n) - w = \lambda C_{\lambda,n}^{-1} (h - w) + C_{\lambda,n}^{-1} \frac{X_n^\top \epsilon_n}{n}.$$

Consequently, we can write

$$\begin{aligned}
(w_h(Z_n) - w)^\top x x^\top (w_h(Z_n) - w) &= \lambda^2 (h - w)^\top C_{\lambda,n}^{-1} x x^\top C_{\lambda,n}^{-1} (h - w) \\
&\quad + \frac{\epsilon_n^\top X_n}{n} C_{\lambda,n}^{-1} x x^\top C_{\lambda,n}^{-1} \frac{X_n^\top \epsilon_n}{n} \\
&\quad + 2\lambda (h - w)^\top C_{\lambda,n}^{-1} x x^\top C_{\lambda,n}^{-1} \frac{X_n^\top \epsilon_n}{n}.
\end{aligned} \tag{17}$$

Hence, recalling the definition of the matrix $A_n = C_{\lambda,n}^{-1} x x^\top C_{\lambda,n}^{-1}$, we have that

$$\begin{aligned}
(w_h(Z_n) - w)^\top x x^\top (w_h(Z_n) - w) &= \lambda^2 (h - w)^\top A_n (h - w) + \frac{\epsilon_n^\top X_n A_n X_n^\top \epsilon_n}{n^2} \\
&\quad + 2\lambda (h - w)^\top \frac{A_n X_n^\top \epsilon_n}{n}.
\end{aligned} \tag{18}$$

Consequently, taking the expectation of Eq. (18) with respect to the sampling of the task $\mu = (w, p, \eta) \sim \rho$ and with respect to the sampling of the data $Z_n \sim \mu^n$ and $z \sim \mu$, we obtain that

$$\mathcal{E}_n(h) = \frac{1}{2} \mathbb{E} \left[\lambda^2 (h - w)^\top A_n (h - w) + \frac{\epsilon_n^\top X_n A_n X_n^\top \epsilon_n}{n^2} \right] + \frac{\mathbb{E}[\epsilon^2]}{2}, \tag{19}$$

where we have exploited again the fact that the noise distribution has zero-mean, more precisely, using [Rem. 8](#), we have that

$$\mathbb{E} \left[(h - w)^\top \frac{A_n X_n^\top \epsilon_n}{n} \right] = \mathbb{E} \left[\mathbb{E} \left[(h - w)^\top \frac{A_n X_n^\top \epsilon_n}{n} \middle| p \right] \right] = \mathbb{E} \left[\mathbb{E} \left[(h - w)^\top \frac{A_n X_n^\top}{n} \middle| p \right] \mathbb{E}[\epsilon_n | p] \right] = 0.$$

Hence, letting $w = \bar{w} + v$, with $v \in \mathbb{R}^d$, we can rewrite

$$\mathcal{E}_n(h) = \frac{1}{2} \mathbb{E} \left[\lambda^2 (h - \bar{w})^\top A_n (h - \bar{w}) - 2\lambda^2 (h - \bar{w})^\top A_n v + \lambda^2 v^\top A_n v + \frac{\epsilon_n^\top X_n A_n X_n^\top \epsilon_n}{n^2} \right] + \frac{\mathbb{E}[\epsilon^2]}{2}.$$

We now observe that, thanks to condition *ii*), $\mathbb{E}[v|p] = 0$ and, consequently, by [Rem. 8](#), we have that

$$\mathbb{E}[(h - \bar{w})^\top A_n v] = (h - \bar{w})^\top \mathbb{E}[\mathbb{E}[A_n v | p]] = (h - \bar{w})^\top \mathbb{E}[\mathbb{E}[A_n | p] \mathbb{E}[v | p]] = 0.$$

Hence, observing that

$$\begin{aligned}
\lambda^2 v^\top A_n v &= \text{tr}(v^\top \lambda^2 A_n v) = \text{tr}(v v^\top \lambda^2 A_n) \\
\frac{\epsilon_n^\top X_n A_n X_n^\top \epsilon_n}{n^2} &= \text{tr} \left(\epsilon_n \epsilon_n^\top \frac{X_n A_n X_n^\top}{n^2} \right),
\end{aligned} \tag{20}$$

and exploiting the relation $\bar{\Sigma}_n = \mathbb{E}[\lambda^2 A_n]$, we can conclude that

$$\mathcal{E}_n(h) = \frac{1}{2} (h - \bar{w})^\top \bar{\Sigma}_n (h - \bar{w}) + \frac{1}{2} \text{tr}(\mathbb{E}[v v^\top \lambda^2 A_n]) + \frac{1}{2} \text{tr} \left(\mathbb{E} \left[\epsilon_n \epsilon_n^\top \frac{X_n A_n X_n^\top}{n^2} \right] \right) + \frac{\mathbb{E}[\epsilon^2]}{2}.$$

From this last equation, taking the derivative with respect to h and exploiting the fact that the covariance matrix $\bar{\Sigma}_n$ is invertible (see [Prop. 1](#)), we conclude that the unique minimizer of $\mathcal{E}_n(h)$ coincides with $h_n^* = \bar{w}$. The upper bound on $\mathcal{E}_n(h)$ given in the last statement of the example directly follows from the following steps. We start observing that, by [Rem. 8](#), we can rewrite

$$\mathbb{E}[v v^\top \lambda^2 A_n] = \mathbb{E}[\mathbb{E}[v v^\top \lambda^2 A_n | p]] = \mathbb{E}[\mathbb{E}[v v^\top | p] \mathbb{E}[\lambda^2 A_n | p]]$$

and, consequently,

$$\begin{aligned}
\text{tr}(\mathbb{E}[v v^\top \lambda^2 A_n]) &= \mathbb{E} \left[\text{tr} \left(\mathbb{E}[v v^\top | p] \mathbb{E}[\lambda^2 A_n | p] \right) \right] \\
&= \mathbb{E} \left[\text{tr} \left(\mathbb{E}[\lambda^2 A_n | p]^{1/2} \mathbb{E}[v v^\top | p] \mathbb{E}[\lambda^2 A_n | p]^{1/2} \right) \right].
\end{aligned} \tag{21}$$

Now, thanks to assumption *iii*), we have that $\mathbb{E}[v v^\top | p] \preceq \sigma_w^2 I$, hence, applying twice [Lemma 4](#) with

$$U = \begin{cases} \mathbb{E}[v v^\top | p] \\ \sigma_w^2 I - \mathbb{E}[v v^\top | p] \end{cases} \quad V = \mathbb{E}[\lambda^2 A_n | p]^{1/2},$$

we have that

$$0 \preceq \mathbb{E}[\lambda^2 A_n | p]^{1/2} \mathbb{E}[v v^\top | p] \mathbb{E}[\lambda^2 A_n | p]^{1/2} \preceq \sigma_w^2 \mathbb{E}[\lambda^2 A_n | p].$$

Consequently, taking the trace of the above inequality, we can continue Eq. (21) as follows

$$\text{tr}(\mathbb{E}[v v^\top \lambda^2 A_n]) \leq \sigma_w^2 \mathbb{E}[\text{tr}(\mathbb{E}[\lambda^2 A_n | p])] = \sigma_w^2 \text{tr}(\bar{\Sigma}_n).$$

In a similar way, exploiting again [Rem. 8](#), we observe that

$$\mathbb{E}\left[\epsilon_n \epsilon_n^\top \frac{X_n A_n X_n^\top}{n^2}\right] = \mathbb{E}\left[\mathbb{E}\left[\epsilon_n \epsilon_n^\top \frac{X_n A_n X_n^\top}{n^2} \middle| p\right]\right] = \mathbb{E}\left[\mathbb{E}[\epsilon_n \epsilon_n^\top | p] \mathbb{E}\left[\frac{X_n A_n X_n^\top}{n^2} \middle| p\right]\right]$$

and, consequently,

$$\begin{aligned} \text{tr}\left(\mathbb{E}\left[\epsilon_n \epsilon_n^\top \frac{X_n A_n X_n^\top}{n^2}\right]\right) &= \mathbb{E}\left[\text{tr}\left(\mathbb{E}[\epsilon_n \epsilon_n^\top | p] \mathbb{E}\left[\frac{X_n A_n X_n^\top}{n^2} \middle| p\right]\right)\right] \\ &= \mathbb{E}\left[\text{tr}\left(\mathbb{E}\left[\frac{X_n A_n X_n^\top}{n^2} \middle| p\right]^{1/2} \mathbb{E}[\epsilon_n \epsilon_n^\top | p] \mathbb{E}\left[\frac{X_n A_n X_n^\top}{n^2} \middle| p\right]^{1/2}\right)\right]. \end{aligned} \quad (22)$$

Now, thanks to assumption *i*) and the independence of the points in the datasets, we have that $\mathbb{E}[\epsilon_n \epsilon_n^\top | p] \preceq \sigma_\epsilon^2 I$, hence, applying twice [Lemma 4](#) with

$$U = \begin{cases} \mathbb{E}[\epsilon_n \epsilon_n^\top | p] \\ \sigma_\epsilon^2 I - \mathbb{E}[\epsilon_n \epsilon_n^\top | p] \end{cases}, \quad V = \mathbb{E}\left[\frac{X_n A_n X_n^\top}{n^2} \middle| p\right]^{1/2},$$

we get

$$0 \preceq \mathbb{E}\left[\frac{X_n A_n X_n^\top}{n^2} \middle| p\right]^{1/2} \mathbb{E}[\epsilon_n \epsilon_n^\top | p] \mathbb{E}\left[\frac{X_n A_n X_n^\top}{n^2} \middle| p\right]^{1/2} \preceq \sigma_\epsilon^2 \mathbb{E}\left[\frac{X_n A_n X_n^\top}{n^2} \middle| p\right].$$

Consequently, taking the trace of the above inequality, we can continue Eq. (22) as follows

$$\text{tr}\left(\mathbb{E}\left[\epsilon_n \epsilon_n^\top \frac{X_n A_n X_n^\top}{n^2}\right]\right) \leq \sigma_\epsilon^2 \mathbb{E}\left[\text{tr}\left(\mathbb{E}\left[\frac{X_n A_n X_n^\top}{n^2} \middle| p\right]\right)\right] = \sigma_\epsilon^2 \text{tr}\left(\mathbb{E}\left[\frac{X_n^\top X_n A_n}{n^2}\right]\right).$$

Finally, we observe that, thanks to assumption *i*), we have $\mathbb{E}[\epsilon^2] = \mathbb{E}[\mathbb{E}[\epsilon^2 | p]] \leq \sigma_\epsilon^2$. The statement derives from combining the upper bounds on all the terms. \blacksquare

Remark 9 (Connection to the Mean Estimation Problem). *In [Ex. 1](#), the minimizer of the transfer risk coincides with the mean \bar{w} of the regression vectors. Hence, our problem appears similar to a mean estimation problem (see e.g. [\[11\]](#)). However, in our setting, we do not receive the regression vectors, but we have indirect observations of them by the corresponding datasets. Moreover, in our case, we aim at minimizing the (excess) transfer risk (and not at estimating its minimizer) and, as already observed in [Prop. 1](#), this quantity does not coincide with the quantity $V_h^2 = \frac{1}{2} \mathbb{E}[\|h - \bar{w}\|^2]$. Specifically, also V_h^2 is minimized at $h = \bar{w}$, and, for any $h \in \mathbb{R}^d$, in the setting of [Ex. 1](#), we have*

$$\lambda_{\min}(\bar{\Sigma}_n) V_h^2 \leq \mathcal{E}_n(h) - \mathcal{E}_n(\bar{w}) = \mathcal{E}_n(h) - \mathcal{E}_n(\bar{w}) = \frac{1}{2} \|\bar{\Sigma}_n^{1/2} (h - \bar{w})\|^2 \leq \|\bar{\Sigma}_n\|_\infty V_h^2.$$

Remark 3 (Advantage of Learning Around the Best Mean over ITL in [Ex. 1](#)). *Consider the setting of [Ex. 1](#). If the noise satisfies $\sigma_\epsilon^2 \ll (n^{-1} \lambda^{-2} + 1)^{-1} \|\bar{\Sigma}_n^{1/2} \bar{w}\|^2$ and the regression vectors are such that $\sigma_w^2 \ll \text{tr}(\bar{\Sigma}_n)^{-1} \|\bar{\Sigma}_n^{1/2} \bar{w}\|^2$, then $\mathcal{E}_n(0) - \mathcal{E}_n(\bar{w}) \gg \mathcal{E}_n(\bar{w})$.*

Proof. As already observed in the paper, thanks to [Prop. 1](#), the difference between the transfer risk of the algorithm with $h = 0$ (ITL) and the best algorithm in our class, can be rewritten as $\mathcal{E}_n(0) - \mathcal{E}_n(\bar{w}) = \frac{1}{2} \bar{w}^\top \bar{\Sigma}_n \bar{w}$. Now, in the setting of [Ex. 1](#), we have that

$$\mathcal{E}_n(h) \leq \frac{1}{2} (h - \bar{w})^\top \bar{\Sigma}_n (h - \bar{w}) + \frac{\sigma_w^2 \text{tr}(\bar{\Sigma}_n)}{2} + \frac{\sigma_\epsilon^2}{2} \text{tr}\left(\mathbb{E}\left[\frac{X_n^\top X_n A_n}{n^2}\right]\right) + \frac{\sigma_\epsilon^2}{2}. \quad (23)$$

Then, the improvement over ITL is significant if $\mathcal{E}_n(0) - \mathcal{E}_n(\bar{w})$ is much greater than the RHS term in Eq. (23) evaluated in $h = \bar{w}$, i.e. if

$$\bar{w}^\top \bar{\Sigma}_n \bar{w} \gg \sigma_w^2 \text{tr}(\bar{\Sigma}_n) + \sigma_\epsilon^2 \text{tr}\left(\mathbb{E}\left[\frac{X_n^\top X_n A_n}{n^2}\right]\right) + \sigma_\epsilon^2. \quad (24)$$

Now, we observe that, thanks to the assumption $\mathcal{X} \subseteq \mathcal{B}_1$, we have that $\|xx^\top\|_\infty \leq 1$ and $\text{tr}(X_n^\top X_n) = \|X_n\|^2 \leq n$. Hence, exploiting the definition of the matrix A_n , the fact that $\|C_{\lambda,n}^{-1}\|_\infty \leq 1/\lambda$ and applying Holder's inequality, we get

$$\text{tr}\left(\frac{X_n^\top X_n A_n}{n^2}\right) = \text{tr}\left(\frac{X_n^\top X_n C_{\lambda,n}^{-1} x x^\top C_{\lambda,n}^{-1}}{n^2}\right) \leq \frac{\text{tr}(X_n^\top X_n)}{n^2} \|C_{\lambda,n}^{-1} x x^\top C_{\lambda,n}^{-1}\|_\infty \leq n^{-1} \lambda^{-2}.$$

Consequently, instead of analysing Eq. (24), we can simply require that the following inequality

$$\bar{w}^\top \bar{\Sigma}_n \bar{w} \gg \sigma_w^2 \text{tr}(\bar{\Sigma}_n) + \sigma_\epsilon^2 (n^{-1} \lambda^{-2} + 1)$$

holds. In turn, this corresponds to requiring that $\sigma_w^2 \ll \text{tr}(\bar{\Sigma}_n)^{-1} \bar{w}^\top \bar{\Sigma}_n \bar{w} \leq \|\bar{w}\|^2$ and $\sigma_\epsilon^2 \ll (n^{-1} \lambda^{-2} + 1)^{-1} \bar{w}^\top \bar{\Sigma}_n \bar{w} \leq (n^{-1} \lambda^{-2} + 1)^{-1} \|\bar{w}\|^2$, where in the last inequality we have applied [Lemma 11](#) and exploited the fact that, for any $\mu \sim \rho$, thanks to the assumption $\mathcal{X} \subseteq \mathcal{B}_1$, $\|\Sigma_\mu\|_\infty \leq 1$. \blacksquare

As described in the following, the proof of [Ex. 2](#) follows the same lines of the proof of [Ex. 1](#).

Example 2. Let $\mathcal{X} \subseteq \mathcal{B}_1$ and consider the environment ρ formed by $K \in \mathbb{N} \setminus \{0\}$ clusters of tasks parametrized by the triplet (w, p, η) as in [Asm. 1](#). Assume that each cluster $k \in [K]$ is associated to a marginal distribution p_k that is sampled with probability $\mathbb{P}(p = p_k) = \nu_k > 0$. For any $k \in [K]$ and $\lambda > 0$, let $\bar{w}_k = \mathbb{E}[w|p = p_k]$, $\bar{\Sigma}_{n,k} = (n\lambda)^2 \mathbb{E}[A_n|p = p_k]$ with $A_n = C_{\lambda,n}^{-1} x x^\top C_{\lambda,n}^{-1}$ and assume that $i) \eta|p = p_k$ has variance bounded by $\sigma_{\epsilon,k}^2$ for $\sigma_{\epsilon,k} \geq 0$, $ii) \mathbb{E}[(w - \bar{w}_k)(w - \bar{w}_k)^\top | p = p_k] \preceq \sigma_{w,k}^2 I$ for $\sigma_{w,k} \geq 0$. Then, for any $\lambda > 0$ and $h \in \mathbb{R}^d$, $h_n^* = (\sum_{k=1}^K \nu_k \bar{\Sigma}_{n,k})^{-1} \sum_{k=1}^K \nu_k \bar{\Sigma}_{n,k} \bar{w}_k$ and $\mathcal{E}_n(h) = \sum_{k=1}^K \nu_k \mathcal{E}_{n,k}(h)$, where

$$\mathcal{E}_{n,k}(h) \leq \frac{1}{2} \|\bar{\Sigma}_{n,k}^{1/2} (h - \bar{w}_k)\|^2 + \frac{\sigma_{w,k}^2}{2} \text{tr}(\bar{\Sigma}_{n,k}) + \frac{\sigma_{\epsilon,k}^2}{2} \left(1 + \text{tr}\left(\mathbb{E}\left[\frac{X_n^\top X_n A_n}{n^2} | p = p_k\right]\right)\right).$$

Proof. We start observing that, since the sampling of the marginal distribution p is drawn from a discrete meta-distribution, exploiting the law of total expectation and using the same notation introduced in the proof of [Ex. 1](#), we can rewrite the transfer risk \mathcal{E}_n of the algorithm w_h on any environment satisfying [Asm. 1](#) as follows

$$\mathcal{E}_n(h) = \frac{1}{2} \mathbb{E}\left[\left(\langle x, w_h(Z_n) \rangle - y\right)^2\right] = \frac{1}{2} \sum_{k=1}^K \nu_k \mathbb{E}\left[\left(\langle x, w_h(Z_n) \rangle - y\right)^2 | p = p_k\right] = \sum_{k=1}^K \nu_k \mathcal{E}_{n,k}(h),$$

where in the last step we have introduced the quantity

$$\mathcal{E}_{n,k}(h) = \frac{1}{2} \mathbb{E}\left[\left(\langle x, w_h(Z_n) \rangle - y\right)^2 | p = p_k\right].$$

The proof of the closed form of the function $\mathcal{E}_{n,k}(h)$ follows along the same lines of [Ex. 1](#), taking into account the conditioning with respect to $p = p_k$. Since the steps are exactly the same described in the proof [Ex. 1](#), we skip the derivation and we report only the sketch of the reasoning. Specifically, exploiting again the linear model equation $y = \langle x, w \rangle + \epsilon$, [Rem. 8](#) and the fact that the noise $\eta|p = p_k$ is zero-mean, we can rewrite

$$\mathcal{E}_{n,k}(h) = \frac{1}{2} \mathbb{E}\left[(w_h(Z_n) - w)^\top x x^\top (w_h(Z_n) - w) | p = p_k\right] + \frac{\mathbb{E}[\epsilon^2 | p = p_k]}{2}.$$

Repeating the same steps in Eqs. (17)-(18)-(19) in the proof of [Ex. 1](#) and denoting by $A_n = C_{\lambda,n}^{-1} x x^\top C_{\lambda,n}^{-1}$, we get

$$\mathcal{E}_{n,k}(h) = \frac{1}{2} \mathbb{E}\left[\lambda^2 (h - w)^\top A_n (h - w) + \frac{\epsilon_n^\top X_n A_n X_n^\top \epsilon_n}{n^2} | p = p_k\right] + \frac{\mathbb{E}[\epsilon^2 | p = p_k]}{2}. \quad (25)$$

Now, letting $w = \bar{w}_k + v_k$, with $v_k \in \mathbb{R}^d$, we can rewrite

$$\begin{aligned} \mathcal{E}_{n,k}(h) &= \frac{1}{2} \mathbb{E}\left[\lambda^2 (h - \bar{w}_k)^\top A_n (h - \bar{w}_k) - 2\lambda^2 (h - \bar{w}_k)^\top A_n v_k + \lambda^2 v_k^\top A_n v_k \right. \\ &\quad \left. + \frac{\epsilon_n^\top X_n A_n X_n^\top \epsilon_n}{n^2} | p = p_k\right] + \frac{\mathbb{E}[\epsilon^2 | p = p_k]}{2}. \end{aligned}$$

We now observe that, thanks to the definition of \bar{w}_k , $\mathbb{E}[v_k | p = p_k] = 0$, hence, by [Rem. 8](#), we have

$$\begin{aligned}\mathbb{E}[(h - \bar{w}_k)^\top A_n v_k | p = p_k] &= (h - \bar{w}_k)^\top \mathbb{E}[A_n v_k | p = p_k] \\ &= (h - \bar{w}_k)^\top \mathbb{E}[A_n | p = p_k] \mathbb{E}[v_k | p = p_k] = 0.\end{aligned}$$

Consequently, exploiting again the relations in [Eq. \(20\)](#) and the definition $\bar{\Sigma}_{n,k} = \mathbb{E}[\lambda^2 A_n | p = p_k]$, we can conclude that

$$\begin{aligned}\mathcal{E}_{n,k}(h) &= \frac{1}{2}(h - \bar{w}_k)^\top \bar{\Sigma}_{n,k}(h - \bar{w}_k) + \frac{1}{2}\text{tr}(\mathbb{E}[v_k v_k^\top \lambda^2 A_n | p = p_k]) \\ &\quad + \frac{1}{2}\text{tr}\left(\mathbb{E}\left[\epsilon_n \epsilon_n^\top \frac{X_n A_n X_n^\top}{n^2} \middle| p = p_k\right]\right) + \frac{\mathbb{E}[\epsilon^2 | p = p_k]}{2}.\end{aligned}$$

From this last equation, taking the derivative with respect to h , we get the closed form of the minimizer h_n^* given in the text. Again, the invertibility of the meta-covariance matrix $\sum_{k=1}^K \nu_k \bar{\Sigma}_{n,k}$ is provided by [Prop. 1](#). The upper bound on $\mathcal{E}_{n,k}(h)$ given in the last statement of the example follows by repeating the same steps in the proof of [Ex. 1](#) taking into account the conditioning with respect to $p = p_k$ and exploiting assumptions *i)-ii*). \blacksquare

We conclude this section with the proof of [Prop. 2](#)

Proposition 2 (LS Problem Around a Common Mean for \mathcal{E}_r). *For any $\lambda > 0$, $h \in \mathbb{R}^d$ and $r \in [n-1]$, the transfer risk \mathcal{E}_r in [Eq. \(7\)](#) of the learning algorithm w_h in [Eqs. \(8\)-\(9\)](#), can be rewritten as*

$$\mathcal{E}_r(h) = \frac{1}{2} \mathbb{E}_{\bar{X}_r, \bar{y}_r} \left[\|\bar{X}_r h - \bar{y}_r\|^2 \right] \quad (10)$$

where the meta-data are given by

$$\bar{X}_r = \frac{\lambda}{\sqrt{n-r}} X_{n-r} C_{\lambda,r}^{-1}, \quad \bar{y}_r = \frac{1}{\sqrt{n-r}} \left(\mathbf{y}_{n-r} - \frac{\sqrt{n-r}}{\lambda} \bar{X}_r \frac{X_r^\top \mathbf{y}_r}{n} \right). \quad (11)$$

Moreover, under [Asm. 1](#), the meta-covariance matrix $\bar{\Sigma}_r = \mathbb{E} \bar{X}_r^\top \bar{X}_r = \lambda^2 \mathbb{E} [C_{\lambda,r}^{-1} \Sigma_\mu C_{\lambda,r}^{-1}]$ is invertible and $h_r^* = \bar{\Sigma}_r^{-1} \mathbb{E} [\bar{X}_r^\top \bar{y}_r]$ is the unique minimizer of the LS function in [Eq. \(10\)](#). In such a case, letting $v = w - h_r^*$, we have that $\bar{y}_r = \bar{X}_r h_r^* + \bar{\epsilon}_r$, with

$$\bar{\epsilon}_r = \frac{1}{\sqrt{n-r}} \epsilon_{n-r} + \bar{X}_r \left(v - \frac{X_r^\top \epsilon_r}{\lambda n} \right). \quad (12)$$

Proof. The rewriting of the transfer risk in [Eq. \(10\)](#) is a direct consequence of the closed form of the algorithm in [Eq.\(9\)](#). The invertibility of the meta-covariance matrix is a direct consequence of [Lemma 11](#) for $r \in \{0\} \cup [n-1]$ and the requirement $\mathbb{E}_{\mu \sim \rho} \lambda_{\min}(\Sigma_\mu) > 0$ in [Asm. 1](#). This implies that the LS function is strictly convex and, consequently, it admits a unique minimizer h_r^* . The closed form of this minimizer directly follows from the optimality conditions (normal equations) of the problem. Now, thanks to [Asm. 1](#), using the linear model equations

$$\mathbf{y}_{n-r} = X_{n-r} w + \epsilon_{n-r} \quad \mathbf{y}_r = X_r w + \epsilon_r,$$

and letting $w = h_r^* + v$, for some $v \in \mathbb{R}^d$, the following relations hold

$$\begin{aligned}\bar{y}_r &= \frac{1}{\sqrt{n-r}} \left(\mathbf{y}_{n-r} - \frac{\sqrt{n-r}}{\lambda} \bar{X}_r \frac{X_r^\top \mathbf{y}_r}{n} \right) \\ &= \frac{1}{\sqrt{n-r}} \left(\mathbf{y}_{n-r} - X_{n-r} C_{\lambda,r}^{-1} \frac{X_r^\top \mathbf{y}_r}{n} \right) \\ &= \frac{1}{\sqrt{n-r}} \left(X_{n-r} w + \epsilon_{n-r} - X_{n-r} C_{\lambda,r}^{-1} \frac{X_r^\top (X_r w + \epsilon_r)}{n} \right) \\ &= \frac{1}{\sqrt{n-r}} \left(X_{n-r} - X_{n-r} C_{\lambda,r}^{-1} \frac{X_r^\top X_r}{n} \right) w + \frac{1}{\sqrt{n-r}} \left(\epsilon_{n-r} - X_{n-r} C_{\lambda,r}^{-1} \frac{X_r^\top \epsilon_r}{n} \right) \\ &= \bar{X}_r w + \frac{1}{\sqrt{n-r}} \left(\epsilon_{n-r} - \frac{\sqrt{n-r}}{\lambda} \bar{X}_r \frac{X_r^\top \epsilon_r}{n} \right) \\ &= \bar{X}_r h_r^* + \bar{X}_r v + \frac{1}{\sqrt{n-r}} \left(\epsilon_{n-r} - \frac{\sqrt{n-r}}{\lambda} \bar{X}_r \frac{X_r^\top \epsilon_r}{n} \right) \\ &= \bar{X}_r h_r^* + \frac{1}{\sqrt{n-r}} \epsilon_{n-r} + \bar{X}_r \left(v - \frac{X_r^\top \epsilon_r}{\lambda n} \right),\end{aligned}$$

where in the fifth equality we have used the fact that

$$\begin{aligned} \frac{1}{\sqrt{n-r}} \left(X_{n-r} - X_{n-r} C_{\lambda,r}^{-1} \frac{X_r^\top X_r}{n} \right) &= \frac{1}{\sqrt{n-r}} X_{n-r} C_{\lambda,r}^{-1} \left(C_{\lambda,r} - \frac{X_r^\top X_r}{n} \right) \\ &= \frac{1}{\sqrt{n-r}} X_{n-r} C_{\lambda,r}^{-1} \lambda I = \bar{X}_r. \end{aligned}$$

We conclude that the noise on the meta-labels is heteroscedastic, since it is given by

$$\bar{\epsilon}_r = \bar{y}_r - \bar{X}_r h_r^* = \frac{1}{\sqrt{n-r}} \epsilon_{n-r} + \bar{X}_r \left(v - \frac{X_r^\top \epsilon_r}{\lambda n} \right).$$

■

C Least Mean Squares (LMS) Algorithm

In this section, we extend the convergence rate given in [13] for the Least Mean Squares (LMS) algorithm, i.e. SGD applied to a Least Squares problem, when at each iteration we sample a matrix \bar{X} (instead of a single vector as in the standard case) and a vector \bar{y} (instead of a scalar as in the standard case). The material of this section is essentially based on [13].

C.1 The problem, the algorithm and the convergence rate

Let m be a positive integer, let $\bar{X} \in \mathbb{R}^{m \times d}$, $\bar{y} \in \mathbb{R}^m$, and let \mathcal{D} be a distribution for the pair (\bar{X}, \bar{y}) . We wish to solve the problem

$$\min_{h \in \mathbb{R}^d} \underbrace{\mathbb{E}_{\bar{Z}=(\bar{X}, \bar{y}) \sim \mathcal{D}} \left[\frac{1}{2} \|\bar{X}h - \bar{y}\|^2 \right]}_{f(h)} = \min_{h \in \mathbb{R}^d} f(h). \quad (26)$$

To this end, we assume of having a stream of data $\bar{Z}^{(1)}, \bar{Z}^{(2)}, \dots$ i.i.d sampled from \mathcal{D} and we apply [Alg. 2](#). In the sequel we will omit the subscript $\bar{Z} \sim \mathcal{D}$ in all the expectations. We now introduce some further notation.

1. The exact covariance matrix (which in the following will be assumed to be invertible):

$$\bar{\Sigma} = \mathbb{E}[\bar{X}^{(t)\top} \bar{X}^{(t)}] \in \mathbb{R}^{d \times d}. \quad (27)$$

2. The closed form of the minimizer of the optimization problem in Eq. (26):

$$h^* = \operatorname{argmin}_{h \in \mathbb{R}^d} f(h) = \bar{\Sigma}^{-1} \mathbb{E}[\bar{X}^{(t)\top} \bar{y}^{(t)}]. \quad (28)$$

3. The residual associated to the seen dataset $\bar{Z}^{(t)}$:

$$\bar{\epsilon}^{(t)} = \bar{y}^{(t)} - \bar{X}^{(t)} h^* \in \mathbb{R}^m. \quad (29)$$

It follows that

$$\mathbb{E}[\bar{X}^{(t)\top} \bar{\epsilon}^{(t)}] = 0. \quad (30)$$

Indeed, exploiting the optimality conditions for h^* , i.e. $\bar{\Sigma} h^* = \mathbb{E}[\bar{X}^{(t)\top} \bar{y}^{(t)}]$, we have

$$\mathbb{E}[\bar{X}^{(t)\top} \bar{\epsilon}^{(t)}] = \mathbb{E}[\bar{X}^{(t)\top} (\bar{y}^{(t)} - \bar{X}^{(t)} h^*)] = \mathbb{E}[\bar{X}^{(t)\top} \bar{y}^{(t)}] - \mathbb{E}[\bar{X}^{(t)\top} \bar{X}^{(t)}] h^* = 0.$$

In the sequel we will make the following assumptions.

Assumption 2 (Bounded Inputs). *There exists $R \geq 0$ such that $\|\bar{X}^{(t)}\|_\infty \leq R$ a.s.*

[Asm. 2](#) implies the following points.

1. Applying [Cor. 5](#) to the matrix $W = \bar{X}^{(t)\top} \bar{X}^{(t)}$, we get

$$\mathbb{E}[\bar{X}^{(t)\top} \bar{X}^{(t)} (\bar{X}^{(t)\top} \bar{X}^{(t)})^\top] \preceq \mathbb{E}[\|\bar{X}^{(t)\top} \bar{X}^{(t)}\|_\infty \bar{X}^{(t)\top} \bar{X}^{(t)}] \preceq R^2 \bar{\Sigma}. \quad (31)$$

Algorithm 2 SGD (LMS) applied to the function f in Eq. (26)

Input $\gamma > 0$ (step size)
Initialization $h^{(0)} \in \mathbb{R}^d$
For $t = 1$ to T
 Receive $\bar{Z}^{(t)} = (\bar{X}^{(t)}, \bar{y}^{(t)})$
 Update $h^{(t)} = h^{(t-1)} - \gamma \bar{X}^{(t)\top} (\bar{X}^{(t)} h^{(t-1)} - \bar{y}^{(t)})$
Return $\bar{h}_T = \frac{1}{T+1} \sum_{t=0}^T h^{(t)}$

2. Thanks to the definition of $\bar{\Sigma}$, we have that

$$\|\bar{\Sigma}\|_{\infty} \leq \mathbb{E} \left[\|\bar{X}^{(t)\top} \bar{X}^{(t)}\|_{\infty} \right] \leq R^2. \quad (32)$$

Assumption 3 (Noise on the \bar{y}). *There exists $\sigma \geq 0$ such that*

$$\mathbb{E} \left[\bar{X}^{(t)\top} \bar{\epsilon}^{(t)} (\bar{X}^{(t)\top} \bar{\epsilon}^{(t)})^{\top} \right] = \mathbb{E} \left[\bar{X}^{(t)\top} \bar{\epsilon}^{(t)} \bar{\epsilon}^{(t)\top} \bar{X}^{(t)} \right] \preceq \sigma^2 \bar{\Sigma}.$$

Assumption 4 (Step size γ). *The step size γ is taken such that $\gamma \leq 1/(2R^2)$, where R^2 is the constant in [Asm. 2](#).*

Remark 10. *Eq. (32) implies that $\gamma \bar{\Sigma} \preceq \gamma \|\bar{\Sigma}\|_{\infty} I \preceq \gamma R^2 I$. Hence, under [Asm. 4](#), we have that $\gamma \bar{\Sigma} \preceq (1/2)I \prec I$.*

The following result essentially coincides with a simplified version of [Thm. 2](#) in [\[13\]](#).

Theorem 12. *Under [Asm. 2](#), [Asm. 3](#) and [Asm. 4](#), let \bar{h}_T be the output of [Alg. 2](#) applied to solve the stochastic LS problem in Eq. (26). Then, \bar{h}_T satisfies the following convergence rate*

$$\mathbb{E}[f(\bar{h}_T) - f(h^*)] \leq \frac{2\|h^{(0)} - h^*\|^2}{\gamma(T+1)} + \frac{2d\sigma^2}{T+1},$$

where the expectation is over the data $\bar{Z}^{(1)}, \dots, \bar{Z}^{(T)}$.

The bound in the above theorem is the sum of two terms: the first one containing the distance between the initial point and the optimal point is the so-called bias term, the second term is the so-called variance term and it represents the variance on the stochastic gradient at the optimal point. We remark also that it is possible to give a faster convergence rate for the bias term, as described in [Thm. 2](#) in [\[13\]](#), however, in order to keep the presentation simple, we avoid this technical step, which is beyond the scope of this work. The following sub-section is devoted to the proof of [Thm. 12](#).

C.2 Proof of the Convergence Rate

The material of this sub-section is based on [App. B](#) in [\[13\]](#). The starting point for the proof of [Thm. 12](#) is the rewriting of the update rule in [Alg. 2](#) in closed form.

Lemma 13 (Recursive formula of the update in [Alg. 2](#)). *In the setting of [Thm. 12](#), for any positive integers i, k such that $k \leq i + 1$, introduce the operators*

$$M(i, k) = \begin{cases} (I - \gamma \bar{X}^{(i)\top} \bar{X}^{(i)}) \cdots (I - \gamma \bar{X}^{(k)\top} \bar{X}^{(k)}) & \text{if } i \geq k \\ I & \text{if } k = i + 1. \end{cases} \quad (33)$$

Then, for any positive integers i, t such that $i \leq t$, we can rewrite

$$h^{(t)} - h^* = M(t, i+1)(h^{(i)} - h^*) + \gamma \sum_{k=i+1}^t M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)}. \quad (34)$$

In particular, for $i = 0$, we get

$$h^{(t)} - h^* = M(t, 1)(h^{(0)} - h^*) + \gamma \sum_{k=1}^t M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)}. \quad (35)$$

Proof. Looking at the update step in [Alg. 2](#), using the fact $\bar{y}^{(t)} = \bar{\epsilon}^{(t)} + \bar{X}^{(t)}h^*$, we can rewrite

$$h^{(t)} - h^* = (I - \gamma \bar{X}^{(t)\top} \bar{X}^{(t)})(h^{(t-1)} - h^*) + \gamma \bar{X}^{(t)\top} \bar{\epsilon}^{(t)}. \quad (36)$$

Iterating Eq. (36) over t and using the definition in Eq. (33), we get the statement. \blacksquare

Remark 11. We observe that the operator $M(i, k)$ defined in Eq. (33) depends only on $\bar{X}^{(k)}, \dots, \bar{X}^{(i)}$, moreover, since the points are i.i.d., we have that

$$\mathbb{E}[M(i, k)] = (I - \gamma \bar{\Sigma})^{i-k+1}. \quad (37)$$

The proof of the convergence rate in [Thm. 12](#) relies on the bias-variance decomposition described in the following proposition.

Proposition 14 (Bias-Variance Decomposition). *Under the assumptions of [Thm. 12](#), we have that*

$$\mathbb{E}[f(\bar{h}_T) - f(h^*)] \leq \mathbf{V} + \mathbf{B},$$

where the expectation is over the data $\bar{Z}^{(1)}, \dots, \bar{Z}^{(T)}$ and we have introduced the bias and the variance terms, respectively given by

$$\begin{aligned} \mathbf{B} &= \frac{2}{\gamma(T+1)^2} \sum_{t=0}^T \mathbb{E} \left[\left\| M(t, 1)(h^{(0)} - h^*) \right\|^2 \right] \\ \mathbf{V} &= \frac{2\gamma}{(T+1)^2} \sum_{t=0}^T \mathbb{E} \left[\left\| \sum_{k=1}^t M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right\|^2 \right]. \end{aligned}$$

Proof. As well known in Least Squares theory, the starting point of the proof is to exploit the equality

$$f(\bar{h}_T) - f(h^*) = \frac{1}{2} \left\| \bar{\Sigma}^{1/2}(\bar{h}_T - h^*) \right\|^2.$$

We now observe that

$$\begin{aligned} (T+1)^2 \left\| \bar{\Sigma}^{1/2}(\bar{h}_T - h^*) \right\|^2 &= (T+1)^2 \left\langle \bar{h}_T - h^*, \bar{\Sigma}(\bar{h}_T - h^*) \right\rangle \\ &= \sum_{t=0}^T \sum_{j=0}^T \left\langle h^{(t)} - h^*, \bar{\Sigma}(h^{(j)} - h^*) \right\rangle \\ &= \sum_{t=0}^T \left\langle h^{(t)} - h^*, \bar{\Sigma}(h^{(t)} - h^*) \right\rangle + 2 \sum_{t=0}^{T-1} \sum_{j=t+1}^T \left\langle h^{(t)} - h^*, \bar{\Sigma}(h^{(j)} - h^*) \right\rangle \\ &= \sum_{t=0}^T \left\| \bar{\Sigma}^{1/2}(h^{(t)} - h^*) \right\|^2 + 2 \sum_{t=0}^{T-1} \sum_{j=t+1}^T \left\langle h^{(t)} - h^*, \bar{\Sigma}(h^{(j)} - h^*) \right\rangle, \end{aligned}$$

where in the third equality we have applied [Lemma 10](#) to $i_{\min} = 0$, $i_{\max} = T$ and $a_{t,j} = \langle h^{(t)} - h^*, \bar{\Sigma}(h^{(j)} - h^*) \rangle$. Hence, taking the expectation over the data $\bar{Z}^{(1)}, \dots, \bar{Z}^{(T)}$, we get

$$\begin{aligned} (T+1)^2 \mathbb{E} \left[\left\| \bar{\Sigma}^{1/2}(\bar{h}_T - h^*) \right\|^2 \right] &= \sum_{t=0}^T \mathbb{E} \left[\left\| \bar{\Sigma}^{1/2}(h^{(t)} - h^*) \right\|^2 \right] \\ &\quad + 2 \underbrace{\sum_{t=0}^{T-1} \sum_{j=t+1}^T \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma}(h^{(j)} - h^*) \right\rangle}_C. \end{aligned} \quad (38)$$

Using the recursive formula in Eq. (34), more precisely

$$h^{(j)} - h^* = M(j, t+1)(h^{(t)} - h^*) + \gamma \sum_{k=t+1}^j M(j, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)},$$

we can write the term C as follows.

$$\begin{aligned}
C &= \sum_{t=0}^{T-1} \sum_{j=t+1}^T \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma} (h^{(j)} - h^*) \right\rangle \\
&= \sum_{t=0}^{T-1} \sum_{j=t+1}^T \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma} \left(M(j, t+1)(h^{(t)} - h^*) + \gamma \sum_{k=t+1}^j M(j, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right) \right\rangle \\
&= \sum_{t=0}^{T-1} \sum_{j=t+1}^T \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma} M(j, t+1)(h^{(t)} - h^*) \right\rangle \\
&\quad + \gamma \sum_{t=0}^{T-1} \sum_{j=t+1}^T \sum_{k=t+1}^j \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma} M(j, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right\rangle.
\end{aligned}$$

We now observe that, thanks to the constraints we have on the indexes t, j, k and thanks to the independence of the sampled points, the variables $h^{(t)}$, $M(j, k+1)$ and $\bar{X}^{(k)\top} \bar{\epsilon}^{(k)}$ are independent. Consequently, since $\mathbb{E}[\bar{X}^{(k)\top} \bar{\epsilon}^{(k)}] = 0$, the second term vanishes. Hence, exploiting again the independence of $h^{(t)}$ and $M(j, t+1)$, we have that

$$\begin{aligned}
C &= \sum_{t=0}^{T-1} \sum_{j=t+1}^T \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma} M(j, t+1)(h^{(t)} - h^*) \right\rangle \\
&= \sum_{t=0}^{T-1} \sum_{j=t+1}^T \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma} \mathbb{E}[M(j, t+1)](h^{(t)} - h^*) \right\rangle \\
&= \sum_{t=0}^{T-1} \sum_{j=t+1}^T \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma} (I - \gamma \bar{\Sigma})^{j-t} (h^{(t)} - h^*) \right\rangle \\
&= \sum_{t=0}^{T-1} \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma} \sum_{j=t+1}^T (I - \gamma \bar{\Sigma})^{j-t} (h^{(t)} - h^*) \right\rangle \\
&= \sum_{t=0}^{T-1} \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma} \sum_{k=1}^{T-t} (I - \gamma \bar{\Sigma})^k (h^{(t)} - h^*) \right\rangle \\
&\leq \sum_{t=0}^{T-1} \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma} \sum_{k=1}^{\infty} (I - \gamma \bar{\Sigma})^k (h^{(t)} - h^*) \right\rangle,
\end{aligned} \tag{39}$$

where in the third equality we have used Eq. (37), in the last equality we have made the change of indexes $k = j - t$ and in the last inequality, we have exploited the fact that we are adding positive quantities. As a matter of fact, for any k , the matrices $\bar{\Sigma} (I - \gamma \bar{\Sigma})^k \in \mathbb{S}_+^d$ (we apply Lemma 6 to $A = \bar{\Sigma} \in \mathbb{S}_{++}^d$ and $B = I - \gamma \bar{\Sigma} \in \mathbb{S}_{++}^d$, thanks to Rem. 10). Now, thanks to Rem. 10 and the invertibility of $\bar{\Sigma}$, we we can apply Lemma 9 to the matrix $A = \gamma \bar{\Sigma}$, hence, we can write

$$\sum_{k=1}^{\infty} (I - \gamma \bar{\Sigma})^k = \gamma^{-1} \bar{\Sigma}^{-1} (I - \gamma \bar{\Sigma}).$$

Consequently, multiplying by $\bar{\Sigma}$, we obtain

$$\bar{\Sigma} \sum_{k=1}^{\infty} (I - \gamma \bar{\Sigma})^k = \gamma^{-1} (I - \gamma \bar{\Sigma}).$$

Coming back to Eq. (39), since, as already observed, $I - \gamma\bar{\Sigma} \succ 0$, we get

$$\begin{aligned}
C &\leq \gamma^{-1} \sum_{t=0}^{T-1} \mathbb{E} \left\langle h^{(t)} - h^*, (I - \gamma\bar{\Sigma})(h^{(t)} - h^*) \right\rangle \\
&\leq \gamma^{-1} \sum_{t=0}^T \mathbb{E} \left\langle h^{(t)} - h^*, (I - \gamma\bar{\Sigma})(h^{(t)} - h^*) \right\rangle \\
&= \gamma^{-1} \sum_{t=0}^T \mathbb{E} \left[\|h^{(t)} - h^*\|^2 \right] - \sum_{t=0}^T \mathbb{E} \left[\|\bar{\Sigma}^{1/2}(h^{(t)} - h^*)\|^2 \right].
\end{aligned}$$

Continuing with Eq. (38), we get

$$\begin{aligned}
&(T+1)^2 \mathbb{E} \left[\|\bar{\Sigma}^{1/2}(\bar{h}_T - h^*)\|^2 \right] \\
&= \sum_{t=0}^T \mathbb{E} \left[\|\bar{\Sigma}^{1/2}(h^{(t)} - h^*)\|^2 \right] + 2 \underbrace{\sum_{t=0}^{T-1} \sum_{j=t+1}^T \mathbb{E} \left\langle h^{(t)} - h^*, \bar{\Sigma}(h_j - h^*) \right\rangle}_C \\
&\leq \sum_{t=0}^T \mathbb{E} \left[\|\bar{\Sigma}^{1/2}(h^{(t)} - h^*)\|^2 \right] + 2\gamma^{-1} \sum_{t=0}^T \mathbb{E} \left[\|h^{(t)} - h^*\|^2 \right] \\
&\quad - 2 \sum_{t=0}^T \mathbb{E} \left[\|\bar{\Sigma}^{1/2}(h^{(t)} - h^*)\|^2 \right] \\
&= 2\gamma^{-1} \sum_{t=0}^T \mathbb{E} \left[\|h^{(t)} - h^*\|^2 \right] - \sum_{t=0}^T \mathbb{E} \left[\|\bar{\Sigma}^{1/2}(h^{(t)} - h^*)\|^2 \right] \\
&\leq 2\gamma^{-1} \sum_{t=0}^T \mathbb{E} \left[\|h^{(t)} - h^*\|^2 \right].
\end{aligned}$$

To sum up we have obtained that

$$\frac{1}{2} \mathbb{E} \left[\|\bar{\Sigma}^{1/2}(\bar{h}_T - h^*)\|^2 \right] \leq \frac{1}{\gamma(T+1)^2} \sum_{t=0}^T \mathbb{E} \left[\|h^{(t)} - h^*\|^2 \right]. \quad (40)$$

Now, thanks to Eq. (35), we have that

$$h^{(t)} - h^* = \underbrace{M(t, 1)(h^{(0)} - h^*)}_{A_t} + \underbrace{\gamma \sum_{k=1}^t M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)}}_{B_t},$$

hence, using Minkowski's inequality, namely for any two random vectors v_1, v_2 :

$$\mathbb{E} [\|v_1 + v_2\|^2] \leq \left(\sqrt{\mathbb{E}[\|v_1\|^2]} + \sqrt{\mathbb{E}[\|v_2\|^2]} \right)^2 \leq 2 \left(\mathbb{E}[\|v_1\|^2] + \mathbb{E}[\|v_2\|^2] \right),$$

we get

$$\mathbb{E} [\|h^{(t)} - h^*\|^2] = \mathbb{E} [\|A_t + B_t\|^2] \leq 2 \left(\mathbb{E} [\|A_t\|^2] + \mathbb{E} [\|B_t\|^2] \right).$$

Returning to Eq. (40), we get

$$\begin{aligned}
\frac{1}{2}\mathbb{E}\left[\|\bar{\Sigma}^{1/2}(\bar{h}_T - h^*)\|^2\right] &\leq \frac{2}{\gamma(T+1)^2} \sum_{t=0}^T \left(\mathbb{E}\left[\|A_t\|^2\right] + \mathbb{E}\left[\|B_t\|^2\right]\right) \\
&= \underbrace{\frac{2}{\gamma(T+1)^2} \sum_{t=0}^T \mathbb{E}\left[\|M(t,1)(h^{(0)} - h^*)\|^2\right]}_{\mathbf{B}} \\
&\quad + \underbrace{\frac{2\gamma}{(T+1)^2} \sum_{t=0}^T \mathbb{E}\left[\left\|\sum_{k=1}^t M(t,k+1)\bar{X}^{(k)\top}\bar{\epsilon}^{(k)}\right\|^2\right]}_{\mathbf{V}}.
\end{aligned} \tag{41}$$

In the sequel we give a bound on the separate terms \mathbf{B} and \mathbf{V} . We start now from the term \mathbf{B} .

Lemma 15 (Bound on the Bias Term \mathbf{B}). *Under the assumptions of Thm. 12, we have that*

$$\mathbf{B} = \frac{2}{\gamma(T+1)^2} \sum_{t=0}^T \mathbb{E}\left[\|M(t,1)(h^{(0)} - h^*)\|^2\right] \leq \frac{2\|h^{(0)} - h^*\|^2}{\gamma(T+1)}. \tag{42}$$

Proof. We start observing that, given A, B matrices and v, w vectors, we have that

$$\langle Av, Bw \rangle = \text{tr}((Av)^\top Bw) = \text{tr}(v^\top A^\top Bw) = \text{tr}(A^\top Bwv^\top). \tag{43}$$

Hence, introducing the notation $E^{(0)} = (h^{(0)} - h^*)(h^{(0)} - h^*)^\top$, we can rewrite

$$\begin{aligned}
\sum_{t=0}^T \mathbb{E}\left[\|M(t,1)(h^{(0)} - h^*)\|^2\right] &= \sum_{t=0}^T \mathbb{E}\left\langle \underbrace{M(t,1)}_A \underbrace{(h^{(0)} - h^*)}_v, \underbrace{M(t,1)}_B \underbrace{(h^{(0)} - h^*)}_w \right\rangle \\
&= \mathbb{E}\left[\sum_{t=0}^T \text{tr}\left(M(t,1)^\top M(t,1)E^{(0)}\right)\right] \\
&= \text{tr}\left(\sum_{t=0}^T \mathbb{E}\left[M(t,1)^\top M(t,1)\right]E^{(0)}\right) \\
&= \text{tr}\left(E^{(0)1/2} \sum_{t=0}^T \mathbb{E}\left[M(t,1)^\top M(t,1)\right]E^{(0)1/2}\right).
\end{aligned} \tag{44}$$

Now we observe that

$$\mathbb{E}\left[M(t,1)^\top M(t,1)\right] \prec I, \tag{45}$$

as a matter of fact, exploiting again the independence of the points, according to the definition in Eq. (33), we have that

$$\mathbb{E}\left[M(t,1)^\top M(t,1)\right] = \mathbb{E}\left[M(t-1,1)^\top \mathbb{E}\left[(I - \gamma\bar{X}^{(t)\top}\bar{X}^{(t)})^\top (I - \gamma\bar{X}^{(t)\top}\bar{X}^{(t)})\right] M(t-1,1)\right]$$

and, using Eq. (31) and Asm. 4, according to which $2 - \gamma R^2 \geq 3/2 > 0$, we get

$$\begin{aligned}
\mathbb{E}\left[(I - \gamma\bar{X}^{(t)\top}\bar{X}^{(t)})^\top (I - \gamma\bar{X}^{(t)\top}\bar{X}^{(t)})\right] &= \mathbb{E}\left[I - 2\gamma\bar{X}^{(t)\top}\bar{X}^{(t)} + \gamma^2\bar{X}^{(t)\top}\bar{X}^{(t)}\bar{X}^{(t)}\bar{X}^{(t)\top}\bar{X}^{(t)}\right] \\
&\preceq I - 2\gamma\bar{\Sigma} + \gamma^2 R^2 \bar{\Sigma} \\
&= I - \gamma(2 - \gamma R^2)\bar{\Sigma} \prec I.
\end{aligned}$$

Iterating the previous observation and exploiting Lemma 4 with

$$\begin{aligned}
U &= I - \mathbb{E}\left[(I - \gamma\bar{X}^{(t)\top}\bar{X}^{(t)})^\top (I - \gamma\bar{X}^{(t)\top}\bar{X}^{(t)})\right] \\
V &= M(t-1,1),
\end{aligned}$$

we get Eq. (45). Hence, coming back to Eq. (44), applying Lemma 4 with

$$\begin{aligned} U &= I - \mathbb{E} \left[M(t, 1)^\top M(t, 1) \right] \\ V &= E^{(0)1/2}, \end{aligned}$$

we get

$$\sum_{t=0}^T \mathbb{E} \left[\left\| M(t, 1)(h^{(0)} - h^*) \right\|^2 \right] \leq (T+1) \text{tr}(E^{(0)}) = (T+1) \|h^{(0)} - h^*\|^2.$$

Therefore, for the term \mathbf{B} , we get the following upper bound:

$$\mathbf{B} = \frac{2}{\gamma(T+1)^2} \sum_{t=0}^T \mathbb{E} \left[\left\| M(t, 1)(h^{(0)} - h^*) \right\|^2 \right] \leq \frac{2 \|h^{(0)} - h^*\|^2}{\gamma(T+1)}.$$

■

In order to give a bound on the variance term \mathbf{V} , we will exploit the following lemma.

Lemma 16 (Recursive formula for $M(t, k+1)$). *Under the assumptions of Thm. 12, we have that*

$$\begin{aligned} \mathbb{E} \left[M(t, k+1) \bar{\Sigma} M(t, k+1)^\top \right] &\preceq \\ &\frac{1}{\gamma(2 - \gamma R^2)} \left(\mathbb{E} \left[M(t, k+1) M(t, k+1)^\top \right] - \mathbb{E} \left[M(t, k) M(t, k)^\top \right] \right). \end{aligned}$$

Proof. Exploiting the independence of the points, we have that:

$$\begin{aligned} &\mathbb{E} \left[M(t, k) M(t, k)^\top \right] \\ &= \mathbb{E} \left[M(t, k+1) \mathbb{E} \left[(I - \gamma \bar{X}^{(k)\top} \bar{X}^{(k)}) (I - \gamma \bar{X}^{(k)\top} \bar{X}^{(k)})^\top \right] M(t, k+1)^\top \right] \\ &= \mathbb{E} \left[M(t, k+1) \left(I - 2\gamma \bar{\Sigma} + \gamma^2 \mathbb{E} [\bar{X}^{(k)\top} \bar{X}^{(k)} \bar{X}^{(k)\top} \bar{X}^{(k)}] \right) M(t, k+1)^\top \right] \quad (46) \\ &= \mathbb{E} \left[M(t, k+1) M(t, k+1)^\top \right] \\ &\quad + \gamma \mathbb{E} \left[M(t, k+1) \left(-2\bar{\Sigma} + \gamma \mathbb{E} [\bar{X}^{(k)\top} \bar{X}^{(k)} \bar{X}^{(k)\top} \bar{X}^{(k)}] \right) M(t, k+1)^\top \right]. \end{aligned}$$

But, since, because of Eq. (31), we have that

$$-2\bar{\Sigma} + \gamma \mathbb{E} [\bar{X}^{(k)\top} \bar{X}^{(k)} \bar{X}^{(k)\top} \bar{X}^{(k)}] \preceq (\gamma R^2 - 2) \bar{\Sigma} = -(2 - \gamma R^2) \bar{\Sigma},$$

then, substituting in Eq. (46) and exploiting Lemma 4 with

$$\begin{aligned} U &= -(2 - \gamma R^2) \bar{\Sigma} - \left(-2\bar{\Sigma} + \gamma \mathbb{E} [\bar{X}^{(k)\top} \bar{X}^{(k)} \bar{X}^{(k)\top} \bar{X}^{(k)}] \right) \\ V &= M(t, k+1)^\top, \end{aligned}$$

we get

$$\begin{aligned} \mathbb{E} \left[M(t, k) M(t, k)^\top \right] &\preceq \mathbb{E} \left[M(t, k+1) M(t, k+1)^\top \right] \\ &\quad - \gamma(2 - \gamma R^2) \mathbb{E} \left[M(t, k+1) \bar{\Sigma} M(t, k+1)^\top \right]. \end{aligned}$$

Consequently, since, as already observed, thanks to Asm. 4, we have that $2 - \gamma R^2 \geq 3/2 > 0$, we get

$$\begin{aligned} \mathbb{E} \left[M(t, k+1) \bar{\Sigma} M(t, k+1)^\top \right] &\preceq \\ &\frac{1}{\gamma(2 - \gamma R^2)} \left(\mathbb{E} \left[M(t, k+1) M(t, k+1)^\top \right] - \mathbb{E} \left[M(t, k) M(t, k)^\top \right] \right). \end{aligned}$$

■

We now are ready to give a bound on the term \mathbf{V} .

Lemma 17 (Bound on the Variance Term \mathbf{V}). *Under the assumptions of [Thm. 12](#), we have that*

$$\mathbf{V} = \frac{2\gamma}{(T+1)^2} \sum_{t=0}^T \mathbb{E} \left[\left\| \sum_{k=1}^t M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right\|^2 \right] \leq \frac{2d\sigma^2}{T+1}. \quad (47)$$

Proof. We start observing that applying [Lemma 10](#) to $i_{\min} = 1$, $i_{\max} = t$ and

$$a_{k,j} = (\bar{X}^{(j)\top} \bar{\epsilon}^{(j)})^\top M(t, j+1)^\top M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)},$$

we can rewrite

$$\begin{aligned} \mathbb{E} \left[\left\| \sum_{k=1}^t M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right\|^2 \right] &= \mathbb{E} \left[\sum_{k=1}^t \sum_{j=1}^t (\bar{X}^{(j)\top} \bar{\epsilon}^{(j)})^\top M(t, j+1)^\top M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right] \\ &= \mathbb{E} \left[\sum_{k=1}^t (\bar{X}^{(k)\top} \bar{\epsilon}^{(k)})^\top M(t, k+1)^\top M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right] \\ &\quad + 2\mathbb{E} \left[\sum_{k=1}^{t-1} \sum_{j=k+1}^t (\bar{X}^{(j)\top} \bar{\epsilon}^{(j)})^\top M(t, j+1)^\top M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right]. \end{aligned}$$

But, thanks to the independence of the points and the constraints on the indexes, since $(\bar{X}^{(j)\top} \bar{\epsilon}^{(j)})^\top M(t, j+1)^\top M(t, k+1)$ does not depend on $\bar{X}^{(k)}$ and since $\mathbb{E}[\bar{X}^{(k)\top} \bar{\epsilon}^{(k)}] = 0$, we have that

$$\mathbb{E} \left[\sum_{k=1}^{t-1} \sum_{j=k+1}^t (\bar{X}^{(j)\top} \bar{\epsilon}^{(j)})^\top M(t, j+1)^\top M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right] = 0.$$

Consequently, we can write the following steps.

$$\begin{aligned} \mathbb{E} \left[\left\| \sum_{k=1}^t M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right\|^2 \right] &= \mathbb{E} \left[\sum_{k=1}^t (\bar{X}^{(k)\top} \bar{\epsilon}^{(k)})^\top M(t, k+1)^\top M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right] \\ &= \mathbb{E} \left[\sum_{k=1}^t \text{tr} \left((\bar{X}^{(k)\top} \bar{\epsilon}^{(k)})^\top M(t, k+1)^\top M(t, k+1) \bar{X}^{(k)\top} \bar{\epsilon}^{(k)} \right) \right] \\ &= \mathbb{E} \left[\sum_{k=1}^t \text{tr} \left(\bar{X}^{(k)\top} \bar{\epsilon}^{(k)} (\bar{X}^{(k)\top} \bar{\epsilon}^{(k)})^\top M(t, k+1)^\top M(t, k+1) \right) \right] \\ &= \text{tr} \left(\sum_{k=1}^t \mathbb{E} \left[\bar{X}^{(k)\top} \bar{\epsilon}^{(k)} (\bar{X}^{(k)\top} \bar{\epsilon}^{(k)})^\top \right] \mathbb{E} \left[M(t, k+1)^\top M(t, k+1) \right] \right) \\ &\leq \sigma^2 \text{tr} \left(\sum_{k=1}^t \bar{\Sigma} \mathbb{E} \left[M(t, k+1)^\top M(t, k+1) \right] \right) \\ &= \sigma^2 \text{tr} \left(\sum_{k=1}^t \mathbb{E} \left[\bar{\Sigma} M(t, k+1)^\top M(t, k+1) \right] \right) \\ &= \sigma^2 \text{tr} \left(\sum_{k=1}^t \mathbb{E} \left[M(t, k+1) \bar{\Sigma} M(t, k+1)^\top \right] \right), \end{aligned} \quad (48)$$

where, in the second equality we have used [Eq. \(43\)](#) and in the fourth equality we have exploited the independence of the points and the definition in [Eq. \(33\)](#). Finally, in the above inequality we have applied [Asm. 3](#) and [Lemma 4](#) with

$$\begin{aligned} U &= \sigma^2 \bar{\Sigma} - \mathbb{E} \left[\bar{X}^{(k)\top} \bar{\epsilon}^{(k)} (\bar{X}^{(k)\top} \bar{\epsilon}^{(k)})^\top \right] \\ V &= \mathbb{E} \left[M(t, k+1)^\top M(t, k+1) \right]^{1/2}. \end{aligned}$$

Now we observe that, exploiting the recursive formula in [Lemma 16](#) and the fact that we obtain a telescopic sum, we have that

$$\begin{aligned}
& \sum_{k=1}^t \mathbb{E} \left[M(t, k+1) \bar{\Sigma} M(t, k+1)^\top \right] \\
& \preceq \frac{1}{\gamma(2-\gamma R^2)} \sum_{k=1}^t \left(\mathbb{E} \left[M(t, k+1) M(t, k+1)^\top \right] - \mathbb{E} \left[M(t, k) M(t, k)^\top \right] \right) \\
& = \frac{1}{\gamma(2-\gamma R^2)} \left(\mathbb{E} \left[M(t, t+1) M(t, t+1)^\top \right] - \mathbb{E} \left[M(t, 1) M(t, 1)^\top \right] \right) \\
& \preceq \frac{1}{\gamma(2-\gamma R^2)} I,
\end{aligned}$$

where in the last step we have used the definition $M(t, t+1) = I$ and the fact $\mathbb{E} \left[M(t, 1) M(t, 1)^\top \right] \succeq 0$. Hence, coming back to [Eq. \(48\)](#), we get

$$\mathbb{E} \left[\left\| \sum_{k=1}^t M(t, k+1) \bar{X}_k^\top \bar{\epsilon}_k \right\|^2 \right] \leq \frac{\sigma^2 \text{tr}(I)}{\gamma(2-\gamma R^2)} = \frac{d\sigma^2}{\gamma(2-\gamma R^2)}.$$

Substituting in the definition of \mathbf{V} , we get the following upper bound.

$$\mathbf{V} \leq \frac{2d\sigma^2}{(2-\gamma R^2)(T+1)}.$$

Since $1/(2-\gamma R^2) < 1$ (as already observed, thanks to [Asm. 4](#), $2-\gamma R^2 \geq 3/2 > 0$), we finally obtain that

$$\mathbf{V} \leq \frac{2d\sigma^2}{T+1}.$$

■

Finally, we have all the ingredients necessary to prove [Thm. 12](#).

Proof of [Thm. 12](#). The statement of the theorem directly follows from combining [Prop. 14](#) with [Lemma 15](#) and [Lemma 17](#). ■

D Adaptation of LMS to the Problem of Minimizing \mathcal{E}_r

In this section we adapt the theory described in the previous section to the stochastic problem of minimizing the transfer risk \mathcal{E}_r ($r \in \{0\} \cup [n-1]$) described in the paper. More precisely, when the environment satisfies the assumptions described in [Ex. 1](#), we are able to explicitly estimate the constant σ introduced in [Asm. 3](#) and to derive a consequent version of the rate in [Thm. 12](#) for [Alg. 1](#), coinciding with SGD applied to \mathcal{E}_r . Such convergence rate will be used in the next [App. E](#), to prove the statement in [Prop. 3](#). Starting from now, we use again the notation introduced in the paper.

For any $r \in \{0\} \cup [n-1]$, the problem of minimizing the function \mathcal{E}_r introduced in [Prop. 2](#) in the paper can be written as in [Eq. \(26\)](#) making the following identifications:

$$\begin{aligned}
m & \mapsto n-r \\
\bar{X} & \mapsto \bar{X}_r = \frac{\lambda}{\sqrt{n-r}} X_{n-r} C_{\lambda,r}^{-1} \in \mathbb{R}^{(n-r) \times d} \\
\bar{\mathbf{y}} & \mapsto \bar{\mathbf{y}}_r = \frac{1}{\sqrt{n-r}} \left(\mathbf{y}_{n-r} - \frac{\sqrt{n-r}}{\lambda} \bar{X}_r \frac{X_r^\top \mathbf{y}_r}{n} \right) \in \mathbb{R}^{n-r} \\
\bar{\epsilon} & \mapsto \bar{\epsilon}_r = \frac{1}{\sqrt{n-r}} \epsilon_{n-r} + \bar{X}_r \left(v - \frac{X_r^\top \epsilon_r}{\lambda n} \right) \in \mathbb{R}^{n-r} \\
\bar{\Sigma} & \mapsto \bar{\Sigma}_r \in \mathbb{R}^{d \times d},
\end{aligned} \tag{49}$$

where we recall that $v = w - h_r^*$ and the remaining quantities are introduced in [Prop. 2](#). Moreover, the sampling of the points from the distribution \mathcal{D} coincides, in our case, with the sampling of the meta-datasets (\bar{X}_r, \bar{y}_r) from the distribution induced by the independent sampling of the original datasets (X_n, y_n) from the environment. Hence, in our setting, the meta-points are independently sampled from this induced distribution. In order to give a convergence rate for [Alg. 1](#), we now specialize all the assumptions introduced in [App. C](#) to the setting described in the paper. We start from [Asm. 2](#).

Lemma 18 ([Asm. 2](#) for \mathcal{E}_r). *Consider the setting described in the paper and outlined in Eq. (49) and assume $\mathcal{X} \subseteq \mathcal{B}_1$. Then, for any $r \in \{0\} \cup [n-1]$, [Asm. 2](#) is satisfied with $R = 1$.*

Proof. Using the facts $\|C_{r,\lambda}^{-1}\|_\infty \leq 1/\lambda$, $\|X_{n-r}\|_\infty \leq \|X_{n-r}\| \leq \sqrt{n-r}$ (since we are assuming $\mathcal{X} \subseteq \mathcal{B}_1$) and the definition of \bar{X}_r , we have that

$$\|\bar{X}_r\|_\infty = \frac{\lambda}{\sqrt{n-r}} \|X_{n-r} C_{\lambda,r}^{-1}\|_\infty = \frac{\lambda}{\sqrt{n-r}} \|C_{\lambda,r}^{-1}\|_\infty \|X_{n-r}\|_\infty \leq 1.$$

We remark that the above steps hold also for the extreme case $r = 0$, according to the associated definitions in that case. \blacksquare

In the case of [Ex. 1](#), we manage to get an estimate of the constant σ introduced in [Asm. 3](#). This is reported in the following lemma.

Lemma 19 ([Asm. 3](#) for \mathcal{E}_r and [Ex. 1](#)). *Consider the setting described in the paper and outlined in Eq. (49). Let $\mathcal{X} \subseteq \mathcal{B}_1$ and let the environment satisfy the assumptions in [Ex. 1](#). Then, for any $r \in \{0\} \cup [n-1]$ and for any $\lambda > 0$, [Asm. 3](#) is satisfied with*

$$\sigma_r^2 = \sigma_w^2 + \left(\frac{1}{n-r} + \frac{r}{(n\lambda)^2} \right) \sigma_\epsilon^2.$$

Proof. We wish to estimate σ_r such that

$$\mathbb{E} \left[\bar{X}_r^\top \bar{\epsilon}_r \bar{\epsilon}_r^\top \bar{X}_r \right] \preceq \sigma_r^2 \mathbb{E} \left[\bar{X}_r^\top \bar{X}_r \right].$$

Since in the setting of [Ex. 1](#) $h_r^* = \bar{w}$ for any $r \in \{0\} \cup [n-1]$, the residuals are given by

$$\bar{\epsilon}_r = \bar{y}_r - \bar{X}_r \bar{w} = \bar{X}_r v + \frac{1}{\sqrt{n-r}} \left(\epsilon_{n-r} - \frac{\sqrt{n-r}}{\lambda} \bar{X}_r \frac{X_r^\top \epsilon_r}{n} \right),$$

with $v = w - \bar{w}$. Hence, we have that

$$\begin{aligned} \bar{X}_r^\top \bar{\epsilon}_r \bar{\epsilon}_r^\top \bar{X}_r &= \bar{X}_r^\top \bar{X}_r v v^\top \bar{X}_r^\top \bar{X}_r + \frac{1}{n-r} \bar{X}_r^\top \epsilon_{n-r} \epsilon_{n-r}^\top \bar{X}_r + \frac{1}{(n\lambda)^2} \bar{X}_r^\top \bar{X}_r X_r^\top \epsilon_r \epsilon_r^\top X_r \bar{X}_r^\top \bar{X}_r \\ &\quad + \underbrace{\frac{1}{\sqrt{n-r}} \bar{X}_r^\top \bar{X}_r v \epsilon_{n-r}^\top \bar{X}_r}_{1} - \underbrace{\frac{1}{n\lambda} \bar{X}_r^\top \bar{X}_r v \epsilon_r^\top X_r \bar{X}_r^\top \bar{X}_r}_{2} \\ &\quad + \underbrace{\frac{1}{\sqrt{n-r}} \bar{X}_r^\top \epsilon_{n-r} v^\top \bar{X}_r^\top \bar{X}_r}_{3} - \underbrace{\frac{1}{n\lambda \sqrt{n-r}} \bar{X}_r^\top \epsilon_{n-r} \epsilon_r^\top X_r \bar{X}_r^\top \bar{X}_r}_{4} \\ &\quad - \underbrace{\frac{1}{n\lambda} \bar{X}_r^\top \bar{X}_r X_r^\top \epsilon_r v^\top \bar{X}_r^\top \bar{X}_r}_{5} - \underbrace{\frac{1}{n\lambda \sqrt{n-r}} \bar{X}_r^\top \bar{X}_r X_r^\top \epsilon_r \epsilon_{n-r}^\top \bar{X}_r}_{6}. \end{aligned}$$

Now, we focus on the expectation of the term 1 and we observe that, exploiting the closed form of \bar{X}_r and the fact that $\eta \mid p$ has zero-mean, we can write

$$\begin{aligned} \mathbb{E} \left[\bar{X}_r^\top \bar{X}_r v \epsilon_{n-r}^\top \bar{X}_r \right] &= \mathbb{E} \left[\mathbb{E} \left[\bar{X}_r^\top \bar{X}_r v \epsilon_{n-r}^\top \bar{X}_r \mid p, w, X_n \right] \right] \\ &= \mathbb{E} \left[\bar{X}_r^\top \bar{X}_r v \mathbb{E} \left[\epsilon_{n-r} \mid p, w, X_n \right]^\top \bar{X}_r \right] \\ &= \mathbb{E} \left[\bar{X}_r^\top \bar{X}_r v \mathbb{E} \left[\epsilon_{n-r} \mid p \right]^\top \bar{X}_r \right] = 0, \end{aligned}$$

where in the third equality we have exploited the independence of ϵ_{n-r} with respect to w, X_n , conditioned with respect to p (see [Rem. 8](#)). Applying a similar reasoning, it is easy to show that also the terms 2 – 6 have zero-mean. Consequently, we have that

$$\begin{aligned} \mathbb{E}\left[\bar{X}_r^\top \bar{\epsilon}_r \bar{\epsilon}_r^\top \bar{X}_r\right] &= \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r v v^\top \bar{X}_r^\top \bar{X}_r\right] + \frac{1}{n-r} \mathbb{E}\left[\bar{X}_r^\top \epsilon_{n-r} \epsilon_{n-r}^\top \bar{X}_r\right] \\ &+ \frac{1}{(n\lambda)^2} \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r X_r^\top \epsilon_r \epsilon_r^\top X_r \bar{X}_r^\top \bar{X}_r\right]. \end{aligned} \quad (50)$$

We now treat the three terms in a separate way. As regards the first term in Eq. (50), using assumption *iii*) in [Ex. 1](#) and the independence of w and X_n conditioned with respect to p (see [Rem. 8](#)), we have that

$$\begin{aligned} \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r v v^\top \bar{X}_r^\top \bar{X}_r\right] &= \mathbb{E}\left[\mathbb{E}\left[\bar{X}_r^\top \bar{X}_r v v^\top \bar{X}_r^\top \bar{X}_r \mid p, X_n\right]\right] \\ &= \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r \mathbb{E}[v v^\top \mid p, X_n] \bar{X}_r^\top \bar{X}_r\right] \\ &= \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r \mathbb{E}[v v^\top \mid p] \bar{X}_r^\top \bar{X}_r\right] \\ &\preceq \sigma_w^2 \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r \bar{X}_r^\top \bar{X}_r\right] \\ &\preceq \sigma_w^2 \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r\right], \end{aligned}$$

where, in the first inequality we have applied [Lemma 4](#) with

$$U = \sigma_w^2 I - \mathbb{E}[v v^\top \mid p] \quad V = \bar{X}_r^\top \bar{X}_r,$$

and in the second inequality we have applied [Cor. 5](#) with $W = \bar{X}_r^\top \bar{X}_r$ and [Lemma 18](#), according to which $\|\bar{X}_r^\top \bar{X}_r\|_\infty = \|\bar{X}_r\|_\infty^2 \leq 1$. As regards the second term in Eq. (50), using assumption *i*) in [Ex. 1](#) and the independence of the points in the datasets, we have that $\mathbb{E}[\epsilon_{n-r} \epsilon_{n-r}^\top \mid p] \preceq \sigma_\epsilon^2 I$, consequently, exploiting the independence of ϵ_{n-r} and X_n conditioned with respect to p (see [Rem. 8](#)), we can write

$$\begin{aligned} \frac{1}{n-r} \mathbb{E}\left[\bar{X}_r^\top \epsilon_{n-r} \epsilon_{n-r}^\top \bar{X}_r\right] &= \frac{1}{n-r} \mathbb{E}\left[\mathbb{E}\left[\bar{X}_r^\top \epsilon_{n-r} \epsilon_{n-r}^\top \bar{X}_r \mid p, X_n\right]\right] \\ &= \frac{1}{n-r} \mathbb{E}\left[\bar{X}_r^\top \mathbb{E}[\epsilon_{n-r} \epsilon_{n-r}^\top \mid p, X_n] \bar{X}_r\right] \\ &= \frac{1}{n-r} \mathbb{E}\left[\bar{X}_r^\top \mathbb{E}[\epsilon_{n-r} \epsilon_{n-r}^\top \mid p] \bar{X}_r\right] \\ &\preceq \frac{\sigma_\epsilon^2}{n-r} \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r\right], \end{aligned}$$

where in the last step we have applied [Lemma 4](#) with

$$U = \sigma_\epsilon^2 I - \mathbb{E}[\epsilon_{n-r} \epsilon_{n-r}^\top \mid p] \quad V = \bar{X}_r.$$

As regards the third term in Eq. (50), we start observing that, using again assumption *i*) in [Ex. 1](#) and the independence of the points in the datasets, we have that $\mathbb{E}[\epsilon_r \epsilon_r^\top \mid p] \preceq \sigma_\epsilon^2 I$. Hence, exploiting the independence of ϵ_r and X_n conditioned with respect to p (see [Rem. 8](#)), we can write

$$\begin{aligned} \frac{1}{(n\lambda)^2} \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r X_r^\top \epsilon_r \epsilon_r^\top X_r \bar{X}_r^\top \bar{X}_r\right] &= \frac{1}{(n\lambda)^2} \mathbb{E}\left[\mathbb{E}\left[\bar{X}_r^\top \bar{X}_r X_r^\top \epsilon_r \epsilon_r^\top X_r \bar{X}_r^\top \bar{X}_r \mid p, X_n\right]\right] \\ &= \frac{1}{(n\lambda)^2} \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r X_r^\top \mathbb{E}[\epsilon_r \epsilon_r^\top \mid p, X_n] X_r \bar{X}_r^\top \bar{X}_r\right] \\ &= \frac{1}{(n\lambda)^2} \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r X_r^\top \mathbb{E}[\epsilon_r \epsilon_r^\top \mid p] X_r \bar{X}_r^\top \bar{X}_r\right] \\ &\preceq \frac{\sigma_\epsilon^2}{(n\lambda)^2} \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r X_r^\top X_r \bar{X}_r^\top \bar{X}_r\right] \\ &\preceq \frac{\sigma_\epsilon^2 r}{(n\lambda)^2} \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r \bar{X}_r^\top \bar{X}_r\right] \\ &\preceq \frac{\sigma_\epsilon^2 r}{(n\lambda)^2} \mathbb{E}\left[\bar{X}_r^\top \bar{X}_r\right], \end{aligned}$$

where, in the first inequality we have applied [Lemma 4](#) with

$$U = \sigma_\epsilon^2 I - \mathbb{E}[\epsilon_r \epsilon_r^\top | p] \quad V = X_r \bar{X}_r^\top \bar{X}_r$$

and in the second inequality we have applied [Lemma 4](#) with

$$U = \|X_r^\top X_r\|_\infty I - X_r^\top X_r \quad V = \bar{X}_r^\top \bar{X}_r$$

and the fact that $\|X_r^\top X_r\|_\infty \leq \|X_r^\top X_r\| \leq r$ (since we are assuming $\mathcal{X} \subseteq \mathcal{B}_1$). Finally, in the third inequality, we have applied [Cor. 5](#) with $W = \bar{X}_r^\top \bar{X}_r$ and [Lemma 18](#), according to which $\|\bar{X}_r^\top \bar{X}_r\|_\infty = \|\bar{X}_r\|^2 \leq 1$. Putting all together, we conclude that

$$\mathbb{E}[\bar{X}_r^\top \bar{\epsilon}_r \bar{\epsilon}_r^\top \bar{X}_r] \preceq \left(\sigma_w^2 + \left(\frac{1}{n-r} + \frac{r}{(n\lambda)^2} \right) \sigma_\epsilon^2 \right) \mathbb{E}[\bar{X}_r^\top \bar{X}_r].$$

We conclude observing that the above steps hold also for the extreme case $r = 0$, according to the associated definitions in that case. \blacksquare

We observe that the upper bound given in [Lemma 19](#) is decreasing in r . This confirms in some way what we have already observed in [Rem. 5](#) in the paper: in the specific setting of [Ex. 1](#), using more than one test point in the definition of the function \mathcal{E}_r has the same effect of traditional mini-batches, i.e. it reduces the variance on the unbiased estimates of the true gradient computed at the true solution. We now are ready to state the adaptation of [Thm. 12](#) for the output of [Alg. 1](#) introduced in the paper, in the setting of [Ex. 1](#).

Theorem 20 (Convergence rate of SGD to \mathcal{E}_r for [Ex. 1](#)). *Assume $\mathcal{X} \subseteq \mathcal{B}_1$ and let $\bar{h}_{T,r,\lambda}$ be the output of [Alg. 1](#), for any $r \in \{0\} \cup [n-1]$ and $\lambda > 0$. Then, the expected excess transfer risk of the algorithm in [Eqs. \(8\)-\(9\)](#) with parameter $h_{T,r,\lambda}$ trained with r points over the environment in [Ex. 1](#) is bounded by*

$$\mathbb{E}[\mathcal{E}_r(\bar{h}_{T,r,\lambda}) - \mathcal{E}_r(h_r^*)] \leq \frac{2\|h^{(0)} - \bar{w}\|^2}{\gamma(T+1)} + \left(\sigma_w^2 + \left(\frac{1}{n-r} + \frac{r}{(n\lambda)^2} \right) \sigma_\epsilon^2 \right) \frac{2d}{T+1},$$

where the expectation is over the datasets $Z^{(1)}, \dots, Z^{(T)}$.

Proof. The statement immediately follows from applying [Thm. 12](#) to the context outlined in [Eq. \(49\)](#) in the setting of [Ex. 1](#). We choose the upper bound on the step size γ in [Asm. 2](#) according to [Lemma 18](#) and we use the estimate for σ^2 in [Asm. 3](#) obtained in [Lemma 19](#). Finally, for [Ex. 1](#), we know that, for any $r \in \{0\} \cup [n-1]$, h_r^* coincides with the mean \bar{w} of the environment. \blacksquare

E Proof of [Prop. 3](#)

In this section we give the proof of [Prop. 3](#).

Proposition 3. *Assume $\mathcal{X} \subseteq \mathcal{B}_1$ and, for any $r \in \{0\} \cup [n-1]$ and $\lambda > 0$, let $\bar{h}_{T,r,\lambda}$ be the output of [Alg. 1](#). Then, the expected excess transfer risk of the algorithm in [Eqs. \(8\)-\(9\)](#) with parameter $h_{T,r,\lambda}$ trained with n points over the environment in [Ex. 1](#) is bounded by*

$$\mathbb{E}[\mathcal{E}_n(\bar{h}_{T,r,\lambda}) - \mathcal{E}_n(h_n^*)] \leq \frac{(r/n + \lambda)^2}{\lambda^2} \frac{4K_\rho}{T+1} \left(\frac{1}{\gamma} \|h^{(0)} - \bar{w}\|^2 + \left(\sigma_w^2 + \left(\frac{1}{n-r} + \frac{r}{(n\lambda)^2} \right) \sigma_\epsilon^2 \right) d \right),$$

where the expectation is over the datasets $Z_n^{(1)}, \dots, Z_n^{(T)}$ and K_ρ is condition number of the environment defined as

$$K_\rho = \frac{\mathbb{E}_{\mu \sim \rho} \|\Sigma_\mu\|_\infty}{\mathbb{E}_{\mu \sim \rho} \lambda_{\min}(\Sigma_\mu)}. \quad (13)$$

Proof. As described in the proof sketch in [Sec. 4](#), in order to prove the statement, we start from the decomposition

$$\mathbb{E}[\mathcal{E}_n(\bar{h}_{T,r,\lambda}) - \mathcal{E}_n(h_n^*)] \leq \underbrace{\mathbb{E}[\|\bar{\Sigma}_n^{1/2}(\bar{h}_{T,r,\lambda} - h_r^*)\|^2]}_A + \underbrace{\mathbb{E}[\|\bar{\Sigma}_n^{1/2}(h_r^* - h_n^*)\|^2]}_B.$$

The proof of the proposition proceeds by bounding the two separate terms A and B . Thanks to the structure of the environment under consideration, we have that $h_r^* = \bar{w}$ for any r (see [Ex. 1](#)), so, the term B vanishes. To bound the term A we observe that

$$A = \mathbb{E} \left[\left\| \bar{\Sigma}_n^{1/2} (\bar{h}_{T,r,\lambda} - h_r^*) \right\|^2 \right] \leq \|\bar{\Sigma}_n\|_\infty \mathbb{E} \left[\left\| \bar{h}_{T,r,\lambda} - h_r^* \right\|^2 \right] \quad (51)$$

and

$$\lambda_{\min}(\bar{\Sigma}_r) \mathbb{E} \left[\left\| \bar{h}_{T,r,\lambda} - h_r^* \right\|^2 \right] \leq \mathbb{E} \left[\left\| \bar{\Sigma}_r^{1/2} (\bar{h}_{T,r,\lambda} - h_r^*) \right\|^2 \right]. \quad (52)$$

Consequently, since the matrix $\bar{\Sigma}_r$ is invertible (see [Prop. 2](#)), combining [Eq. \(51\)](#) with [Eq. \(52\)](#) and exploiting the LS structure of the function \mathcal{E}_r (see [Prop. 2](#)), we can write

$$A \leq \frac{2\|\bar{\Sigma}_n\|_\infty}{\lambda_{\min}(\bar{\Sigma}_r)} \frac{1}{2} \mathbb{E} \left[\left\| \bar{\Sigma}_r^{1/2} (\bar{h}_{T,r,\lambda} - h_r^*) \right\|^2 \right] = \frac{2\|\bar{\Sigma}_n\|_\infty}{\lambda_{\min}(\bar{\Sigma}_r)} \mathbb{E} [\mathcal{E}_r(\bar{h}_{T,r,\lambda}) - \mathcal{E}_r(h_r^*)]. \quad (53)$$

Now, introducing the condition number of the environment

$$K_\rho = \frac{\mathbb{E}_{\mu \sim \rho} \|\Sigma_\mu\|_\infty}{\mathbb{E}_{\mu \sim \rho} \lambda_{\min}(\Sigma_\mu)}$$

and exploiting [Lemma 11](#), we can write

$$\frac{\|\bar{\Sigma}_n\|_\infty}{\lambda_{\min}(\bar{\Sigma}_r)} \leq \frac{K_\rho(r/n + \lambda)^2}{\lambda^2}.$$

Consequently, coming back to [Eq. \(53\)](#), we have

$$A \leq \frac{2K_\rho(r/n + \lambda)^2}{\lambda^2} \mathbb{E} [\mathcal{E}_r(\bar{h}_{T,r,\lambda}) - \mathcal{E}_r(h_r^*)].$$

Using the bound given in [Thm. 20](#) for the term $\mathbb{E} [\mathcal{E}_r(\bar{h}_{T,r,\lambda}) - \mathcal{E}_r(h_r^*)]$, we get the final statement of the proposition. \blacksquare

F How to tune the hyper-parameters in our LTL setting

Denote by $\bar{h}_{T,r,\lambda}$ the output of [Alg. 1](#) computed with T iterations (hence T tasks) with values r and λ . In all experiments, we obtain this estimator $\bar{h}_{T,r,\lambda}$ by learning it on a dataset \mathbf{Z}_{tr} of T_{tr} training tasks, each comprising a dataset Z_n of n input-output pairs $(x, y) \in \mathcal{X} \times \mathcal{Y}$. During the training phase, each of these datasets is splitted into two parts $Z_n = (Z_r, Z_{n-r})$ and we apply the procedure described in the paper in [Alg. 1](#) in order to compute the estimator. We perform this meta-training for different values of $\lambda \in \{\lambda_1, \dots, \lambda_p\}$ and $r \in \{0\} \cup [n-1]$ and we select the best estimator based on the prediction error measured on a separate set \mathbf{Z}_{va} of T_{va} validation tasks. Once such optimal λ and r values have been selected, we report the generalization performance of the corresponding estimator on a set \mathbf{Z}_{te} of T_{te} test tasks. The tasks in the test and validation sets \mathbf{Z}_{te} and \mathbf{Z}_{va} are all provided with both a training and test dataset both sampled from the same distribution. Note that, since we are interested in measuring the performance of the algorithm trained with n points (as already stressed in the paper we aim at minimizing \mathcal{E}_n and not \mathcal{E}_r), the training datasets have all the same sample size n as those in the meta-training datasets in \mathbf{Z}_{tr} , while the test datasets contain n' points each, for some positive integer n' . Indeed, in order to evaluate the performance of a bias h , we need to first train the corresponding algorithm w_h on the training dataset Z_n , and then test its performance on the test set $Z'_{n'}$, by computing the empirical risk $\frac{1}{2n'} \|X'_{n'} w_h(Z_n) - \mathbf{y}'_{n'}\|^2$. For instance, for the synthetic experiments reported in the paper, we chose $n = 20$ and $n' = 100$. Finally, since we are considering the online setting, the training datasets arrive one at the time, therefore model selection is performed *online*: the system keeps track of all candidate values $\bar{h}_{T_{\text{tr}},r,\lambda_j}$, $r \in \{0\} \cup [n-1]$, $j \in [p]$, and, whenever a new training task is presented, these vectors are all updated by incorporating the corresponding new observations. The best bias h is then returned at each iteration, based on its performance on the validation set \mathbf{Z}_{va} . The previous procedure describes how to tune simultaneously both λ and r . In some experiments reported in the paper, we just tuned one of them; in such a case the procedure is analogous to that described above.

Note on Real Data. In real settings, the assumptions of having the same number n of training points available for each task t is often restrictive. Indeed, on the School dataset considered in this work, datasets can have a very different number n_t of training points. Hence, in order to validate the theoretical analysis reported in this paper, in our experiments we down-sampled the datasets with more training examples in order to maintain the same n across all the tasks. However, it is natural to expect that leveraging more training points when available should be more favorable in terms of overall performance. This is a key question that introduces complications to our analysis and that we plan to investigate in the future. We refer to the code for more details.