

A. Parameters of the Posterior Distribution

Recall that

$$p_n(\boldsymbol{\theta}_n | \mathcal{D}) \propto p^G(\mathbf{Y}_n | \mathbf{X}_n \boldsymbol{\theta}_n, \sigma^2 \mathbf{I}) p^G(\boldsymbol{\theta}_n | \boldsymbol{\alpha}, \boldsymbol{\Sigma}) .$$

We first give a handy proposition for the posterior distribution over the parameters $\boldsymbol{\theta}_n$.

Proposition A.1.

$$\ln p_n(\boldsymbol{\theta}_n) = -\frac{1}{2}(\boldsymbol{\theta}_n - \boldsymbol{\mu})^\top \boldsymbol{\mathcal{T}}^{-1}(\boldsymbol{\theta}_n - \boldsymbol{\mu}) + \text{const}(\boldsymbol{\theta}_n)$$

where covariance is

$$\boldsymbol{\mathcal{T}} = \left(\boldsymbol{\Sigma}^{-1} + \frac{1}{\sigma^2} \mathbf{X}_n^\top \mathbf{X}_n \right)^{-1}$$

and mean is

$$\boldsymbol{\mu} = \boldsymbol{\mathcal{T}} \left(\boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} + \frac{1}{\sigma^2} \mathbf{X}_n^\top \mathbf{Y}_n \right) .$$

Proof. The log-likelihood is the following chain of identities:

$$\begin{aligned} \ln p_n(\boldsymbol{\theta}_n) &= -\frac{\sum_{j=1}^{m_n} (\mathbf{x}_{n,j}^\top \boldsymbol{\theta}_n - Y_{n,j})^2}{2\sigma^2} - \frac{1}{2}(\boldsymbol{\theta}_n - \boldsymbol{\alpha})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\theta}_n - \boldsymbol{\alpha}) + \text{const}(\boldsymbol{\theta}_n) \\ &= -\frac{\sum_{j=1}^{m_n} (\boldsymbol{\theta}_n^\top \mathbf{x}_{n,j} \mathbf{x}_{n,j}^\top \boldsymbol{\theta}_n - 2Y_{n,j} \mathbf{x}_{n,j}^\top \boldsymbol{\theta}_n)}{2\sigma^2} \\ &= -\frac{1}{2}(\boldsymbol{\theta}_n^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\theta}_n - 2\boldsymbol{\alpha}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\theta}_n) + \text{const}(\boldsymbol{\theta}_n) \\ &= -\frac{1}{2} \left(\boldsymbol{\theta}_n^\top \boldsymbol{\mathcal{T}}^{-1} \boldsymbol{\theta}_n - 2 \left(\boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} + \frac{1}{\sigma^2} \mathbf{X}_n^\top \mathbf{Y}_n \right)^\top \boldsymbol{\theta}_n \right) + \text{const}(\boldsymbol{\theta}_n) \\ &= -\frac{1}{2}(\boldsymbol{\theta}_n - \boldsymbol{\mu})^\top \boldsymbol{\mathcal{T}}^{-1}(\boldsymbol{\theta}_n - \boldsymbol{\mu}) + \text{const}(\boldsymbol{\theta}_n) . \end{aligned}$$

□

First note that the first consequence of Proposition A.1 is a MLE for $\boldsymbol{\theta}_n$,

$$\hat{\boldsymbol{\theta}}_n^{\text{MLE}} = \boldsymbol{\mathcal{T}} \left(\boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} + \frac{1}{\sigma^2} \mathbf{X}_n^\top \mathbf{Y}_n \right) .$$

The second consequence is the following corollary which is obtained by taking $Y = \boldsymbol{\theta}_n^\top \mathbf{x} + \varepsilon$ and simply observing that the mean of a $p_n(\boldsymbol{\theta}_n | \mathcal{D})$ is $\boldsymbol{\mu}$, and so $\mathbb{E}[Y | \mathcal{D}] = \mathbf{x}^\top \boldsymbol{\mu}$ while the variance is $\mathbb{V}[Y | \mathcal{D}] = \mathbb{E}[(\mathbf{x}^\top \boldsymbol{\theta}_n + \varepsilon)^2 | \mathcal{D}] - \mathbb{E}[(\mathbf{x}^\top \boldsymbol{\theta}_n + \varepsilon) | \mathcal{D}]^2 = \mathbf{x}^\top \boldsymbol{\mathcal{T}} \mathbf{x} + \sigma^2$.

Proposition 4.6 (restated). *Let $Y = \boldsymbol{\theta}_n^\top \mathbf{x} + \varepsilon$ for $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ and some $\mathbf{x} \in \mathbb{R}^d$. Then,*

$$\begin{aligned} \mathbb{E}[Y | \mathcal{D}] &= \mathbf{x}^\top \boldsymbol{\mathcal{T}} \left(\boldsymbol{\Sigma}^{-1} \boldsymbol{\alpha} + \frac{1}{\sigma^2} \mathbf{X}_n^\top \mathbf{Y}_n \right) \\ \text{and } \mathbb{V}[Y | \mathcal{D}] &= \mathbf{x}^\top \boldsymbol{\mathcal{T}} \mathbf{x} + \sigma^2 . \end{aligned}$$

B. Proof of the Lower Bounds

Our task reduces to establishing lower bounds on

$$\mathbb{E} \left[(\mathbf{x}^\top \boldsymbol{\mathcal{T}} \boldsymbol{\Sigma}^{-1} (\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}))^2 \right] \tag{10}$$

for any choice of estimator $\hat{\boldsymbol{\alpha}}$, which in combination with Lemma 4.5 will prove Theorem 4.1. In the next section we first prove a lower bound for any unbiased estimator relying on the Cramér-Rao inequality. In what follows, in Appendix B.2, we will show a general bound in Lemma B.4 valid for any estimator (possibly biased) using a *hypothesis testing* technique (see, e.g. (Lattimore & Szepesvári, 2018, Chap. 13)). Finally, in Lemma B.5 we prove a high-probability lower bound on Eq. (10).

B.1. Lower Bound for Unbiased Estimator $\hat{\alpha}$

Theorem B.1 (Cramér-Rao inequality). *Suppose that $\alpha \in \mathbb{R}^d$ is an unknown deterministic parameter with a probability density function $f(x | \alpha)$ and that $\hat{\alpha}$ is an unbiased estimator of α . Moreover assume that for all $i, j \in [d]$, $x : f(x | \alpha) > 0$, $\frac{\partial^2}{\partial \alpha_i \partial \alpha_j} \ln f(x | \alpha)$ exists and is finite, and $\frac{\partial^2}{\partial \alpha_i \partial \alpha_j} \int \hat{\alpha} f(x | \alpha) dx = \int \hat{\alpha} \left(\frac{\partial^2}{\partial \alpha_i \partial \alpha_j} f(x | \alpha) \right) dx$.*

Then, for the Fisher information matrix defined as

$$\mathbf{F} = -\mathbb{E} [\nabla_{\alpha} \ln f(X | \alpha) \nabla_{\alpha} \ln f(X | \alpha)^{\top}]$$

we have

$$\mathbb{E} [(\hat{\alpha} - \mathbb{E}[\hat{\alpha}])(\hat{\alpha} - \mathbb{E}[\hat{\alpha}])^{\top}] \succeq \mathbf{F}^{-1}.$$

Lemma B.2. *For any unbiased estimator $\hat{\alpha}$ of α in Eq. (2) we have*

$$\mathbb{E} [(x^{\top} \mathcal{T} \Sigma^{-1} (\alpha - \hat{\alpha})^2)] \geq x^{\top} \mathcal{T} \Sigma^{-1} (\Psi^{\top} \mathbf{K} \Psi)^{-1} \Sigma^{-1} \mathcal{T} x. \quad (11)$$

Proof. Recall that according to the equivalence (2) $\mathbf{Y} \sim \mathcal{N}(\Psi \alpha, \mathbf{K})$ and the unknown parameter is α . To compute the Fisher information matrix we first observe that

$$\nabla_{\alpha} \ln p^G(\mathbf{Y}; \Psi \alpha, \mathbf{K}) = \Psi^{\top} \mathbf{K}^{-1} (\mathbf{Y} - \Psi \alpha)$$

and so

$$\begin{aligned} \mathbf{F} &= \mathbb{E} [\nabla_{\alpha} \ln p^G(\mathbf{Y}; \Psi \alpha, \mathbf{K}) \nabla_{\alpha} \ln p^G(\mathbf{Y}; \Psi \alpha, \mathbf{K})^{\top}] \\ &= \Psi^{\top} \mathbf{K}^{-1} \mathbb{E} [(\mathbf{Y} - \Psi \alpha)(\mathbf{Y} - \Psi \alpha)^{\top}] \mathbf{K}^{-1} \Psi \\ &= \Psi^{\top} \mathbf{K}^{-1} \Psi. \end{aligned}$$

Thus, by Theorem B.1 we have

$$\mathbb{E} [(\alpha - \hat{\alpha})(\alpha - \hat{\alpha})^{\top}] \succeq (\Psi^{\top} \mathbf{K}^{-1} \Psi)^{-1}.$$

Finally, left-multiplying by $x^{\top} \mathcal{T} \Sigma^{-1}$ and right-multiplying the above by $\Sigma^{-1} \mathcal{T} x$ gives us the statement. \square

B.2. Lower Bound for Any Estimator $\hat{\alpha}$

The proof of is based on the following lemma.

Lemma B.3 (Bretagnolle & Huber 1979). *Let P and Q be probability measures on the same measurable space (Ω, \mathcal{F}) , and let $A \in \mathcal{F}$ be an arbitrary event. Then,*

$$P(A) + Q(A^c) \geq \frac{1}{2} \exp(-D_{\text{KL}}(P, Q)), \quad (12)$$

where $D_{\text{KL}}(P, Q) = \int_{\Omega} \ln(P(\omega)/Q(\omega)) dP(\omega)$ denotes Kullback-Leibler divergence between P and Q and $A^c = \Omega \setminus A$ is the complement of A .

Lemma B.4. *For any estimator $\hat{\alpha}$ of α in Eq. (2) we have*

$$\mathbb{E} [(x^{\top} \mathcal{T} \Sigma^{-1} (\hat{\alpha} - \alpha))^2] \geq \frac{x^{\top} \mathbf{M} x}{16\sqrt{e}}.$$

Proof. Throughout the proof let $\mathbf{q} = \Sigma^{-1} \mathcal{T} x$. Consider two meta-learning problems with target distributions \mathbb{P} and \mathbb{Q} characterized by two means: $\alpha_{\mathbb{P}} = \mathbf{0}$ and $\alpha_{\mathbb{Q}} = \Delta(\Psi^{\top} \mathbf{K}^{-1} \Psi)^{-1} \mathbf{q}$ where $\Delta > 0$ is a free parameter to be tuned later on. Thus, according to our established equivalence (2), in these two cases targets are generated by respective models $\mathbb{P} = \mathcal{N}(\mathbf{0}, \mathbf{K})$ and $\mathbb{Q} = \mathcal{N}(\Delta \Psi(\Psi^{\top} \mathbf{K}^{-1} \Psi)^{-1} \mathbf{q}, \mathbf{K})$.

Recall our abbreviation $\mathbf{M} = \mathcal{T} \Sigma^{-1} (\Psi^{\top} \mathbf{K}^{-1} \Psi)^{-1} \Sigma^{-1} \mathcal{T}$. Markov's inequality gives

$$\begin{aligned} \mathbb{E}_{\mathbb{P}} [(\hat{\alpha}^{\top} \mathbf{q} - \alpha_{\mathbb{P}}^{\top} \mathbf{q})^2] &= \mathbb{E}_{\mathbb{P}} [(\hat{\alpha}^{\top} \mathbf{q})^2] \geq \frac{\Delta^2}{4} (x^{\top} \mathbf{M} x)^2 \mathbb{P} \left(|\hat{\alpha}^{\top} \mathbf{q}| \geq \frac{\Delta}{2} x^{\top} \mathbf{M} x \right), \quad \text{while} \\ \mathbb{E}_{\mathbb{Q}} [(\hat{\alpha}^{\top} \mathbf{q} - \alpha_{\mathbb{Q}}^{\top} \mathbf{q})^2] &\geq \frac{\Delta^2}{4} (x^{\top} \mathbf{M} x)^2 \mathbb{Q} \left(|\alpha_{\mathbb{Q}}^{\top} \mathbf{q} - \hat{\alpha}^{\top} \mathbf{q}| \geq \frac{\Delta}{2} x^{\top} \mathbf{M} x \right) \\ &\geq \frac{\Delta^2}{4} (x^{\top} \mathbf{M} x)^2 \mathbb{Q} \left(|\hat{\alpha}^{\top} \mathbf{q}| < \frac{\Delta}{2} x^{\top} \mathbf{M} x \right), \end{aligned}$$

where the last inequality comes using the fact that $|a - b| \geq |a| - |b|$ for $a, b \in \mathbb{R}$ and observing that $\alpha_{\mathbb{Q}}^{\top} \mathbf{q} = \mathbf{x}^{\top} \mathbf{M} \mathbf{x}$. Summing both inequalities above and applying Lemma B.3 we get

$$\begin{aligned} \mathbb{E}_{\mathbb{P}} [(\hat{\alpha}^{\top} \mathbf{q} - \alpha_{\mathbb{P}}^{\top} \mathbf{q})^2] + \mathbb{E}_{\mathbb{Q}} [(\hat{\alpha}^{\top} \mathbf{q} - \alpha_{\mathbb{Q}}^{\top} \mathbf{q})^2] &\geq \frac{\Delta^2}{8} (\mathbf{x}^{\top} \mathbf{M} \mathbf{x})^2 \cdot \exp(-D_{\text{KL}}(\mathbb{P}, \mathbb{Q})) \\ &\stackrel{(a)}{=} \frac{\Delta^2}{8} (\mathbf{x}^{\top} \mathbf{M} \mathbf{x})^2 \cdot \exp\left(-\frac{\Delta^2}{2} \mathbf{x}^{\top} \mathbf{M} \mathbf{x}\right), \end{aligned}$$

where step (a) follows from KL-divergence between multivariate Gaussians with the same covariance matrix. Now, using a basic fact that $2 \max\{a, b\} \geq a + b$, we get that for any measure \mathbb{P} given by parameter α we have

$$\mathbb{E} [(\hat{\alpha}^{\top} \mathbf{q} - \alpha^{\top} \mathbf{q})^2] \geq \frac{\Delta^2}{16} (\mathbf{x}^{\top} \mathbf{M} \mathbf{x})^2 \cdot \exp\left(-\frac{\Delta^2}{2} \mathbf{x}^{\top} \mathbf{M} \mathbf{x}\right).$$

The statement then follows by choosing $\Delta^2 = (\mathbf{x}^{\top} \mathbf{M} \mathbf{x})^{-1}$. \square

Now we prove a high-probability version of the just given inequality.

Lemma B.5. *For any estimator $\hat{\alpha}$ of α in Eq. (2) and any $\delta \in (0, 1)$ we have*

$$\mathbb{P}\left((\mathbf{x}^{\top} \mathcal{T} \Sigma^{-1} (\hat{\alpha} - \alpha))^2 \geq \ln\left(\frac{1}{4} \cdot \frac{1}{1 - \delta}\right) \mathbf{x}^{\top} \mathbf{M} \mathbf{x}\right) \geq 1 - \delta.$$

Proof. The proof is very similar to the proof of Lemma B.4 except we will not apply Markov's inequality and focus directly on giving a lower bound the deviation probabilities rather than expectations. Thus, similarly as before introduce mean parameters $\alpha_{\mathbb{P}} = \mathbf{0}$ and $\alpha_{\mathbb{Q}} = \Delta(\Psi^{\top} \mathbf{K}^{-1} \Psi)^{-1} \mathbf{q} / (\mathbf{x}^{\top} \mathbf{M} \mathbf{x})$ and their associated probability measures $\mathbb{P} = \mathcal{N}(\mathbf{0}, \mathbf{K})$ and $\mathbb{Q} = \mathcal{N}\left(\frac{\Delta \Psi(\Psi^{\top} \mathbf{K}^{-1} \Psi)^{-1} \mathbf{q}}{\mathbf{x}^{\top} \mathbf{M} \mathbf{x}}, \mathbf{K}\right)$.

Note that

$$\begin{aligned} \mathbb{P}\left(|\hat{\alpha}^{\top} \mathbf{q}| \geq \frac{\Delta}{2}\right) &= \mathbb{P}\left(|\alpha_{\mathbb{P}}^{\top} \mathbf{q} - \hat{\alpha}^{\top} \mathbf{q}| \geq \frac{\Delta}{2}\right), \\ \mathbb{Q}\left(|\alpha_{\mathbb{Q}}^{\top} \mathbf{q} - \hat{\alpha}^{\top} \mathbf{q}| \geq \frac{\Delta}{2}\right) &\geq \mathbb{Q}\left(|\hat{\alpha}^{\top} \mathbf{q}| < \frac{\Delta}{2}\right) \end{aligned}$$

and so by using Lemma B.3 we obtain an exponential tail bound

$$\mathbb{P}\left(|\alpha_{\mathbb{P}}^{\top} \mathbf{q} - \hat{\alpha}^{\top} \mathbf{q}| \geq \frac{\Delta}{2}\right) + \mathbb{Q}\left(|\alpha_{\mathbb{Q}}^{\top} \mathbf{q} - \hat{\alpha}^{\top} \mathbf{q}| \geq \frac{\Delta}{2}\right) \geq \exp(-D_{\text{KL}}(\mathbb{P} \parallel \mathbb{Q})) = \frac{1}{2} \exp\left(-\frac{\Delta^2}{\mathbf{x}^{\top} \mathbf{M} \mathbf{x}}\right).$$

Setting the r.h.s. in the above to $2(1 - \delta)$ where δ is an error probability, and solving for Δ gives us tuning

$$\Delta^2 = 2 \ln\left(\frac{1}{4} \cdot \frac{1}{1 - \delta}\right) \mathbf{x}^{\top} \mathbf{M} \mathbf{x}.$$

Thus, we get

$$\mathbb{P}\left((\alpha_{\mathbb{P}}^{\top} \mathbf{q} - \hat{\alpha}^{\top} \mathbf{q})^2 \geq \ln\left(\frac{1}{4} \cdot \frac{1}{1 - \delta}\right) \mathbf{x}^{\top} \mathbf{M} \mathbf{x}\right) + \mathbb{Q}\left(|\alpha_{\mathbb{Q}}^{\top} \mathbf{q} - \hat{\alpha}^{\top} \mathbf{q}| \geq \ln\left(\frac{1}{4} \cdot \frac{1}{1 - \delta}\right) \mathbf{x}^{\top} \mathbf{M} \mathbf{x}\right) \geq 2(1 - \delta)$$

and using the fact that $2 \max(a, b) \geq a + b$ we get that for any probability measure \mathbb{P} given by parameter α we have

$$\mathbb{P}\left((\alpha_{\mathbb{P}}^{\top} \mathbf{q} - \hat{\alpha}^{\top} \mathbf{q})^2 \geq \ln\left(\frac{1}{4} \cdot \frac{1}{1 - \delta}\right) \mathbf{x}^{\top} \mathbf{M} \mathbf{x}\right) \geq 1 - \delta.$$

\square

C. Proof of the Upper Bounds

Theorem 4.2 (restated). For the estimator $\hat{\theta}_n(\hat{\alpha}^{\text{MLE}})$ and for any $\mathbf{x} \in \mathbb{R}^d$ we have

$$\mathbb{E}[\mathcal{L}(\mathbf{x})] = \mathbf{x}^\top \mathbf{M} \mathbf{x} + \mathbf{x}^\top \mathcal{T} \mathbf{x} + \sigma^2.$$

Moreover for the same estimator, with probability at least $1 - \delta, \delta \in (0, 1)$ we have

$$\mathcal{L}(\mathbf{x}) \leq 2 \ln \left(\frac{2}{\delta} \right) \mathbf{x}^\top \mathbf{M} \mathbf{x} + \mathbf{x}^\top \mathcal{T} \mathbf{x} + \sigma^2.$$

Proof. Recall that

$$\hat{\alpha}^{\text{MLE}} = (\Psi^\top \mathbf{K}^{-1} \Psi)^{-1} \Psi^\top \mathbf{K}^{-1} \mathbf{Y}.$$

The first result follows from Lemma 4.5 where we have to give an identity for

$$\mathbb{E} \left[(\mathbf{x}^\top \mathcal{T} \Sigma^{-1} (\alpha - \hat{\alpha}^{\text{MLE}}))^2 \right] \quad (13)$$

and the missing piece is a covariance of the estimator $\hat{\alpha}^{\text{MLE}}$

$$\begin{aligned} & \mathbb{E} [(\alpha - \hat{\alpha}^{\text{MLE}})(\alpha - \hat{\alpha}^{\text{MLE}})^\top] \\ &= (\Psi^\top \mathbf{K}^{-1} \Psi)^{-1} \Psi^\top \mathbf{K}^{-1} \text{Cov}(\mathbf{Y}, \mathbf{Y}) \mathbf{K}^{-1} \Psi (\Psi^\top \mathbf{K}^{-1} \Psi)^{-1} \\ &= (\Psi^\top \mathbf{K}^{-1} \Psi)^{-1}. \end{aligned} \quad (14)$$

To prove the second result we have to give a high probability upper bound on Eq. (13).

Let $\mathbf{q} = \Sigma^{-1} \mathcal{T} \mathbf{x}$ and observe that $\mathbf{q}^\top \hat{\alpha}^{\text{MLE}}$ is Gaussian (since \mathbf{Y} is composed of Gaussian entries) with mean $\mathbf{q}^\top \alpha$ by equivalence (2), and covariance $(\Psi^\top \mathbf{K}^{-1} \Psi)^{-1}$ by Eq. (14). Then, by Gaussian concentration for any error probability $\delta \in (0, 1)$ we have

$$\mathbb{P} \left((\mathbf{q}^\top \alpha - \mathbf{q}^\top \hat{\alpha}^{\text{MLE}})^2 \geq \sqrt{2 \mathbf{q}^\top (\Psi^\top \mathbf{K}^{-1} \Psi)^{-1} \mathbf{q} \ln \left(\frac{2}{\delta} \right)} \right) \leq \delta$$

which completes the proof. \square

D. Derivation of EM Steps

Recall that our goal is to solve

$$\max_{\mathcal{E}'} \int \ln(p(\boldsymbol{\vartheta}, \mathcal{D} | \mathcal{E}')) \, \text{d}p(\boldsymbol{\vartheta} | \mathcal{D}, \hat{\mathcal{E}}_t).$$

First, we will focus on the integral. The chain rule readily gives

$$\ln p(\boldsymbol{\Theta}, \mathcal{D} | \mathcal{E}') = \ln p(\boldsymbol{\Theta} | \mathcal{D}, \mathcal{E}') + \ln p(\boldsymbol{\Theta} | \mathcal{E}').$$

Using the same reasoning and notation as in the proof of Proposition A.1 we get

$$\begin{aligned} \int \ln p(\boldsymbol{\vartheta} | \mathcal{D}, \mathcal{E}') \, \text{d}p(\boldsymbol{\vartheta} | \mathcal{D}, \hat{\mathcal{E}}_t) &= \sum_{i=1}^n \sum_{j=1}^{m_i} \left(\frac{1}{2} \ln \left(\frac{1}{\sigma^2} \right) - \frac{1}{2\sigma^2} \int (Y_{i,j} - \mathbf{x}_{i,j}^\top \boldsymbol{\vartheta}_i)^2 \, \text{d}p(\boldsymbol{\vartheta}_i | \mathcal{D}, \hat{\mathcal{E}}_t) \right) + \text{const}(\mathcal{E}') \\ &= \sum_{i=1}^n \sum_{j=1}^{m_i} \left(\frac{1}{2} \ln \left(\frac{1}{\sigma^2} \right) - \frac{1}{2\sigma^2} (Y_{i,j} - \mathbf{x}_{i,j}^\top \boldsymbol{\mu}_i)^2 - \mathbf{x}_{i,j}^\top \mathcal{T}_i \mathbf{x}_{i,j} \right) + \text{const}(\mathcal{E}'), \end{aligned}$$

using the fact that $\int (Y_{i,j} - \mathbf{x}_{i,j}^\top \boldsymbol{\vartheta}_i)^2 \, \text{d}p(\boldsymbol{\vartheta}_i | \mathcal{D}, \hat{\mathcal{E}}_t) = (Y_{i,j} - \mathbf{x}_{i,j}^\top \boldsymbol{\mu}_i)^2 + \mathbf{x}_{i,j}^\top \mathcal{T}_i \mathbf{x}_{i,j}$ where we took $\boldsymbol{\theta}_i \sim \mathcal{N}(\boldsymbol{\mu}_i, \mathcal{T}_i)$ according to Proposition A.1.

Now we compute the expected log-likelihood of the vector of task parameters:

$$\int \ln p(\boldsymbol{\vartheta} | \mathcal{E}') \, \text{d}p(\boldsymbol{\vartheta} | \mathcal{D}, \hat{\mathcal{E}}_t) = \frac{n}{2} \ln \det \Sigma^{-1} - \frac{1}{2} \sum_{i=1}^n \int (\boldsymbol{\vartheta}_i - \alpha)^\top \Sigma^{-1} (\boldsymbol{\vartheta}_i - \alpha) \, \text{d}p(\boldsymbol{\vartheta}_i | \mathcal{D}, \hat{\mathcal{E}}_t) + \text{const}(\mathcal{E}').$$

M-step for σ^2 . Now, note that since the likelihood of the vector of task variables Θ does not depend on the parameter σ^2 we can solve for σ^2 based on the first order condition of the problem above. Differentiating the above equation with respect to σ^{-2} (and ignoring the constant) gives

$$\sum_{i=1}^n \sum_{j=1}^{m_i} (\sigma^2 - ((Y_{i,j} - \mathbf{x}_{i,j}^\top \boldsymbol{\mu}_i)^2 + \mathbf{x}_{i,j}^\top \mathcal{T}_i \mathbf{x}_{i,j})) . \quad (15)$$

while setting the derivative to zero gives

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n \frac{1}{m_i} \sum_{j=1}^{m_i} ((Y_{i,j} - \mathbf{x}_{i,j}^\top \boldsymbol{\mu}_i)^2 + \mathbf{x}_{i,j}^\top \mathcal{T}_i \mathbf{x}_{i,j}) . \quad (16)$$

M-step for α . Differentiating the objective w.r.t. α (and ignoring the constant) gives $\sum_{i=1}^n \Sigma^{-1} (\mathbb{E}[\theta_i] - \alpha)$ from which we get

$$\alpha = \sum_{i=1}^n \boldsymbol{\mu}_i . \quad (17)$$

M-step for Σ . Differentiating the expected log-likelihood of the vector of task parameters with respect to $\mathbf{A} = \Sigma^{-1}$ gives

$$\sum_{i=1}^n \text{tr}(\Sigma d\mathbf{A}) - \text{tr} \int ((\boldsymbol{\vartheta}_i - \alpha)(\boldsymbol{\vartheta}_i - \alpha)^\top d\mathbf{A}) dp(\boldsymbol{\vartheta}_i | \mathcal{D}, \hat{\mathcal{E}}_t) \quad (18)$$

from which we get

$$\Sigma = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(\theta_i - \alpha)(\theta_i - \alpha)^\top] . \quad (19)$$

Finally, computing the expectation

$$\sum_{i=1}^n (\mathbb{E}[\theta_i \theta_i^\top] - 2\boldsymbol{\mu}_i \alpha^\top + \alpha \alpha^\top) = \sum_{i=1}^n ((\boldsymbol{\mu}_i - \alpha)(\boldsymbol{\mu}_i - \alpha)^\top + \mathcal{T}_i) \quad (20)$$

shows the update for Σ .

E. Selecting λ in Biased Regression

The parameter λ is selected via random search in the following way. For each of the 50 samples of λ from log-uniform distribution on interval $[0; 100]$ we perform the following procedure to estimate the risk \hat{L} . Firstly, we split the training tasks into $K = 10$ groups $\mathcal{S}_1, \dots, \mathcal{S}_K$ of (approximately) equal size and compute the estimates $\hat{\alpha}_k$ using the data $\mathcal{S}^{\setminus k}$ from all of the groups excluding the group k : $\mathcal{S}^{\setminus k} := \cup_{i \neq k} \mathcal{S}_i$. For each of the estimated values $\hat{\alpha}_k$ we perform adaptation to and testing on the tasks in the group \mathcal{S}_k using the given value of λ . We split the samples of each task data $D_i \in \mathcal{S}_k$ randomly into adaptation and test sets 10 times each time such that the size of adaptation set is close to the size of adaptation sets used with the actual test data. For each of the splits we compute an estimate of the parameter vector $\hat{\theta}_{k,i,l}$ where k is the index of the group which was not used to estimate $\hat{\alpha}_k$, i is the index of a task data $D_i \in \mathcal{S}_k$, l is the index of a random split of the samples in that task into adaptation and test sets. With this parameter vector and using the test set of the task $D_i \in \mathcal{S}_k$ we can also estimate the loss $\hat{L}_{k,i,l}$ after which all the loss values are averaged:

$$\hat{L} = \frac{1}{K} \sum_{k=1}^K \frac{1}{|\mathcal{S}^{\setminus k}|} \sum_{i: D_i \in \mathcal{S}^{\setminus k}} \frac{1}{10} \sum_{l=1}^{10} \hat{L}_{k,i,l} .$$

At the end we select the value of λ which lead to the smallest value of \hat{L} using this cross-validation procedure.

F. Supplementary Statements

Proposition F.1. For M (see Eq. (5)) we have

$$M = \sigma^4 \cdot (\Sigma \mathbf{X}_n^\top \mathbf{X}_n + \sigma^2 \mathbf{I})^{-1} \mathbf{A}^{-1} (\Sigma \mathbf{X}_n^\top \mathbf{X}_n + \sigma^2 \mathbf{I})^{-1} ,$$

where we denote

$$\mathbf{A} = \sum_{i=1}^n \mathbf{X}_i^\top (\mathbf{X}_i \boldsymbol{\Sigma} \mathbf{X}_i^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{X}_i.$$

Proof. Recall that

$$\mathbf{M} = \mathcal{T} \boldsymbol{\Sigma}^{-1} (\boldsymbol{\Psi}^\top \mathbf{K}^{-1} \boldsymbol{\Psi})^{-1} \boldsymbol{\Sigma}^{-1} \mathcal{T}$$

and observe that

$$\mathbf{K}^{-1} = \begin{bmatrix} (\mathbf{X}_1 \boldsymbol{\Sigma} \mathbf{X}_1^\top + \sigma^2 \mathbf{I})^{-1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & (\mathbf{X}_2 \boldsymbol{\Sigma} \mathbf{X}_2^\top + \sigma^2 \mathbf{I})^{-1} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & (\mathbf{X}_n \boldsymbol{\Sigma} \mathbf{X}_n^\top + \sigma^2 \mathbf{I})^{-1} \end{bmatrix},$$

which in turn implies

$$\boldsymbol{\Psi}^\top \mathbf{K}^{-1} \boldsymbol{\Psi} = \sum_{i=1}^n \mathbf{X}_i^\top (\mathbf{X}_i \boldsymbol{\Sigma} \mathbf{X}_i^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{X}_i.$$

On the other hand,

$$\begin{aligned} \mathcal{T} \boldsymbol{\Sigma}^{-1} &= \left(\boldsymbol{\Sigma}^{-1} + \frac{1}{\sigma^2} \mathbf{X}_n^\top \mathbf{X}_n \right)^{-1} \boldsymbol{\Sigma}^{-1} \\ &= \sigma^2 (\sigma^2 \mathbf{I} + \boldsymbol{\Sigma} \mathbf{X}_n^\top \mathbf{X}_n)^{-1}. \end{aligned}$$

Combining the above gives the statement. \square

Lemma F.2. In the following assume that $\mathbf{X}_i^\top \mathbf{X}_i = \frac{m_i}{d} \mathbf{I}$ for all i . Let $\lambda_j(\boldsymbol{\Sigma})$ be the j th eigenvalue of $\boldsymbol{\Sigma}$. Then,

$$\lambda_j(\mathbf{M}) = \sigma^4 \cdot \frac{d^2}{(m_n \lambda_j(\boldsymbol{\Sigma}) + d\sigma^2)^2} \cdot \frac{HM\left(\lambda_j(\boldsymbol{\Sigma}) + \frac{d\sigma^2}{m_i}\right)_{i=1}^n}{n},$$

where $HM(z_i)_{i=1}^n$ denotes the harmonic mean of sequence $(z_i)_{i=1}^n$. Moreover,

$$\lambda_j(\mathcal{T}) = \frac{d\sigma^2 \lambda_j(\boldsymbol{\Sigma})}{d\sigma^2 + m_n \lambda_j(\boldsymbol{\Sigma})}.$$

Finally, the eigenvectors of \mathbf{M} and \mathcal{T} coincide with the eigenvectors of $\boldsymbol{\Sigma}$.

Proof. We first characterize eigenvalues of matrix \mathbf{M} . By Proposition F.1,

$$\mathbf{M} = \sigma^4 \cdot (\boldsymbol{\Sigma} \mathbf{X}_n^\top \mathbf{X}_n + \sigma^2 \mathbf{I})^{-1} \mathbf{A}^{-1} (\boldsymbol{\Sigma} \mathbf{X}_n^\top \mathbf{X}_n + \sigma^2 \mathbf{I})^{-1}.$$

We start with \mathbf{A}^{-1} , and by the spectral theorem, $\boldsymbol{\Sigma} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top$ for some unitary \mathbf{U} and diagonal $\boldsymbol{\Lambda}$:

$$\begin{aligned} \mathbf{A}^{-1} &= \left(\sum_{i=1}^n \mathbf{X}_i^\top (\mathbf{X}_i \boldsymbol{\Sigma} \mathbf{X}_i^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{X}_i \right)^{-1} = \left(\sum_{i=1}^n (\boldsymbol{\Sigma} \mathbf{X}_i^\top \mathbf{X}_i + \sigma^2 \mathbf{I})^{-1} \mathbf{X}_i^\top \mathbf{X}_i \right)^{-1} \\ &= \left(\sum_{i=1}^n \left(\boldsymbol{\Sigma} \cdot \frac{m_i}{d} + \sigma^2 \mathbf{I} \right)^{-1} \frac{m_i}{d} \right)^{-1} \\ &= \left(\sum_{i=1}^n \left(\mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^\top \cdot \frac{m_i}{d} + \sigma^2 \mathbf{I} \right)^{-1} \frac{m_i}{d} \right)^{-1} \\ &= \mathbf{U} \left(\sum_{i=1}^n \left(\boldsymbol{\Lambda} + \frac{d\sigma^2}{m_i} \cdot \mathbf{I} \right)^{-1} \right)^{-1} \mathbf{U}^\top. \end{aligned}$$

Now,

$$\begin{aligned} (\Sigma \mathbf{X}_n^\top \mathbf{X}_n + \sigma^2 \mathbf{I})^{-1} &= \left(\Sigma \cdot \frac{m_n}{d} + \sigma^2 \mathbf{I} \right)^{-1} \\ &= \left(\mathbf{U} \Lambda \mathbf{U}^\top \cdot \frac{m_n}{d} + \sigma^2 \mathbf{I} \right)^{-1} \\ &= \mathbf{U} \left(\Lambda \cdot \frac{m_n}{d} + \sigma^2 \mathbf{I} \right)^{-1} \mathbf{U}^\top. \end{aligned}$$

Thus,

$$\mathbf{M} = \mathbf{U} \left(\left(\Lambda \cdot \frac{m_n}{d} + \sigma^2 \mathbf{I} \right)^2 \sum_{i=1}^n \left(\Lambda + \frac{d\sigma^2}{m_i} \right)^{-1} \right)^{-1} \mathbf{U}^\top.$$

and moreover the j th eigenvalue of \mathbf{M} is

$$\begin{aligned} \lambda_j(\mathbf{M}) &= \frac{1}{\left(\frac{m_n}{d} \lambda_j(\Sigma) + \sigma^2 \right)^2} \cdot \frac{1}{\sum_{i=1}^n \frac{1}{\lambda_j(\Sigma) + \frac{d\sigma^2}{m_i}}} \\ &= \frac{1}{\left(\frac{m_n}{d} \lambda_j(\Sigma) + \sigma^2 \right)^2} \cdot \frac{\text{HM} \left(\lambda_j(\Sigma) + \frac{d\sigma^2}{m_i} \right)_{i=1}^n}{n}. \end{aligned}$$

where recall that $\text{HM}(z_i)_{i=1}^n$ denotes the harmonic mean of sequence $(z_i)_{i=1}^n$.

Using the same arguments as above

$$\begin{aligned} \mathcal{T} &= \left(\Sigma^{-1} + \frac{1}{\sigma^2} \mathbf{X}_n^\top \mathbf{X}_n \right)^{-1} \\ &= \left(\mathbf{U} \Lambda^{-1} \mathbf{U}^\top + \frac{m_n}{d\sigma^2} \right)^{-1} \end{aligned}$$

and so

$$\lambda_j(\mathcal{T}) = \frac{1}{\frac{1}{\lambda_j(\Sigma)} + \frac{m_n}{d\sigma^2}} = \frac{d\sigma^2 \lambda_j(\Sigma)}{d\sigma^2 + m_n \lambda_j(\Sigma)}.$$

Finally, in both cases of \mathbf{M} and \mathcal{T} we observe that their eigenvectors are eigenvectors of Σ . □

Corollaries 4.3 and 4.4 (restated). *In the following assume that $\mathbf{X}_i^\top \mathbf{X}_i = \frac{m_i}{d} \mathbf{I}$ for all i . For $\Sigma = \tau^2 \mathbf{I}$, any $\mathbf{x} \in \mathbb{R}^d$, and any $c > 0$,*

$$c \mathbf{x}^\top \mathbf{M} \mathbf{x} + \mathbf{x}^\top \mathcal{T} \mathbf{x} + \sigma^2 = c \cdot \frac{H_{\tau^2}}{n} \cdot \frac{d^2 \sigma^4}{(\tau^2 m_n + d\sigma^2)^2} \cdot \|\mathbf{x}\|^2 + \frac{d\sigma^2 \tau^2}{\tau^2 m_n + d\sigma^2} \cdot \|\mathbf{x}\|^2 + \sigma^2,$$

where H_{τ^2} is a harmonic mean of the sequence $\left(\tau^2 + \frac{d\sigma^2}{m_i} \right)_{i=1}^n$.

Moreover, let Σ be a PSD matrix of rank $s \leq d$ with eigenvalues $\lambda_1 \geq \dots \geq \lambda_s > 0$. Then for any $\mathbf{x} \in \mathbb{R}^d$ and any $c > 0$,

$$c \mathbf{x}^\top \mathbf{M} \mathbf{x} + \mathbf{x}^\top \mathcal{T} \mathbf{x} + \sigma^2 \geq c \cdot \frac{H_{\lambda_s}}{n} \cdot \frac{d^2 \sigma^4}{(\lambda_1 m_n + d\sigma^2)^2} \cdot \|\mathbf{x}\|_{\mathbf{P}_s^\top \mathbf{P}_s}^2 + \frac{d\sigma^2 \lambda_s}{\lambda_s m_n + d\sigma^2} \cdot \|\mathbf{x}\|_{\mathbf{P}_s^\top \mathbf{P}_s}^2 + \sigma^2$$

where $\mathbf{P}_s = [\mathbf{u}_1, \dots, \mathbf{u}_s]^\top$ and $(\mathbf{u}_j)_{j=1}^s$ are eigenvectors of Σ .

Proof. Recalling that by Proposition F.1,

$$\mathbf{M} = \sigma^4 \cdot (\Sigma \mathbf{X}_n^\top \mathbf{X}_n + \sigma^2 \mathbf{I})^{-1} \mathbf{A}^{-1} (\Sigma \mathbf{X}_n^\top \mathbf{X}_n + \sigma^2 \mathbf{I})^{-1}.$$

and using Lemma F.2 with $\Sigma = \tau^2 \mathbf{I}$ we get the first result.

Now we turn to the low-rank case. We start by considering a PSD matrix Σ_ε with s eigenvalues $\lambda_1 \geq \dots \geq \lambda_s > 0$ and remaining $d - s$ are $\varepsilon > 0$. Denote also by $M_\varepsilon, \mathcal{T}_\varepsilon$ matrices w.r.t. Σ_ε . The idea is to lower bound $\mathbf{x}^\top M_\varepsilon \mathbf{x}$ and $\mathbf{x}^\top \mathcal{T}_\varepsilon \mathbf{x}$ and then analyze a limiting behavior as $\varepsilon \rightarrow 0$.

By Lemma F.2, $M_\varepsilon, \mathcal{T}_\varepsilon$, and Σ_ε share the same eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_s$, and so

$$\begin{aligned} c\mathbf{x}^\top M_\varepsilon \mathbf{x} + \mathbf{x}^\top \mathcal{T}_\varepsilon \mathbf{x} &= c \sum_{j=1}^d (\mathbf{u}_j^\top \mathbf{x})^2 \lambda_j(M_\varepsilon) + \sum_{j=1}^d (\mathbf{u}_j^\top \mathbf{x})^2 \lambda_j(\mathcal{T}_\varepsilon) \\ &= c \cdot \sum_{j=1}^s \frac{H_{\lambda_j}}{n} \cdot \frac{\sigma^4}{(\lambda_j \frac{m_n}{d} + \sigma^2)^2} (\mathbf{u}_j^\top \mathbf{x})^2 + c \cdot \underbrace{\frac{H_\varepsilon}{n} \cdot \frac{\sigma^4}{(\varepsilon \frac{m_n}{d} + \sigma^2)^2}}_{(a)} \left(\sum_{j=s+1}^d (\mathbf{u}_j^\top \mathbf{x})^2 \right) \\ &\quad + \sum_{j=1}^s \frac{\sigma^2 \lambda_j}{\lambda_j \frac{m_n}{d} + \sigma^2} (\mathbf{u}_j^\top \mathbf{x})^2 + \frac{\sigma^2 \varepsilon}{\varepsilon \frac{m_n}{d} + \sigma^2} \left(\sum_{j=s+1}^d (\mathbf{u}_j^\top \mathbf{x})^2 \right). \end{aligned}$$

Now,

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} (c\mathbf{x}^\top M_\varepsilon \mathbf{x} + \mathbf{x}^\top \mathcal{T}_\varepsilon \mathbf{x}) &= c \cdot \sum_{j=1}^s \frac{H_{\lambda_j}}{n} \cdot \frac{\sigma^4}{(\lambda_j \frac{m_n}{d} + \sigma^2)^2} (\mathbf{u}_j^\top \mathbf{x})^2 + \frac{d\sigma^2}{M} \sum_{j=s+1}^d (\mathbf{u}_j^\top \mathbf{x})^2 + \sum_{j=1}^s \frac{\sigma^2 \lambda_j}{\lambda_j \frac{m_n}{d} + \sigma^2} (\mathbf{u}_j^\top \mathbf{x})^2 \\ &\geq c \cdot \frac{H_{\lambda_s}}{n} \cdot \frac{\sigma^4}{(\lambda_1 \frac{m_n}{d} + \sigma^2)^2} \sum_{j=1}^s (\mathbf{u}_j^\top \mathbf{x})^2 + \frac{\sigma^2 \lambda_s}{\lambda_s \frac{m_n}{d} + \sigma^2} \sum_{j=1}^s (\mathbf{u}_j^\top \mathbf{x})^2, \end{aligned}$$

where we note that the limit of term (a) is handled as

$$\lim_{\varepsilon \rightarrow 0} \frac{1}{\sum_{i=1}^n \frac{1}{\varepsilon + \frac{d\sigma^2}{m_i}}} \cdot \frac{\sigma^4}{(\varepsilon \frac{m_n}{d} + \sigma^2)^2} = \frac{d\sigma^2}{M} \geq 0.$$

□

G. Further Experimental Details and Results

In this section we provide extra figures for our experimental results. Fig. 1 is complemented with Fig. 7, adding a second example in addition to the one shown in the previous figure.

For completeness, the pseudocode of MoM is given in Algorithm 2.

Algorithm 2 MoM Estimator for Learning Linear Features of (Tripuraneni et al., 2020)

Input: $((\mathbf{x}_{1,j}, y_{1,j}))_{j=1}^{m_1}, \dots, ((\mathbf{x}_{m_{n-1},j}, y_{m_{n-1},j}))_{j=1}^{m_{n-1}}$ — training examples from $n - 1$ past tasks, s — problem rank.
 $UDV^\top \leftarrow \text{SVD} \left(\frac{1}{M - m_n} \sum_{i=1}^{n-1} \sum_{j=1}^{m_n} y_{i,j}^2 \mathbf{x}_{i,j} \mathbf{x}_{i,j}^\top \right)$
 $\hat{B} \leftarrow [D_{1,1} \mathbf{u}_1, \dots, D_{s,s} \mathbf{u}_s]$
return \hat{B}

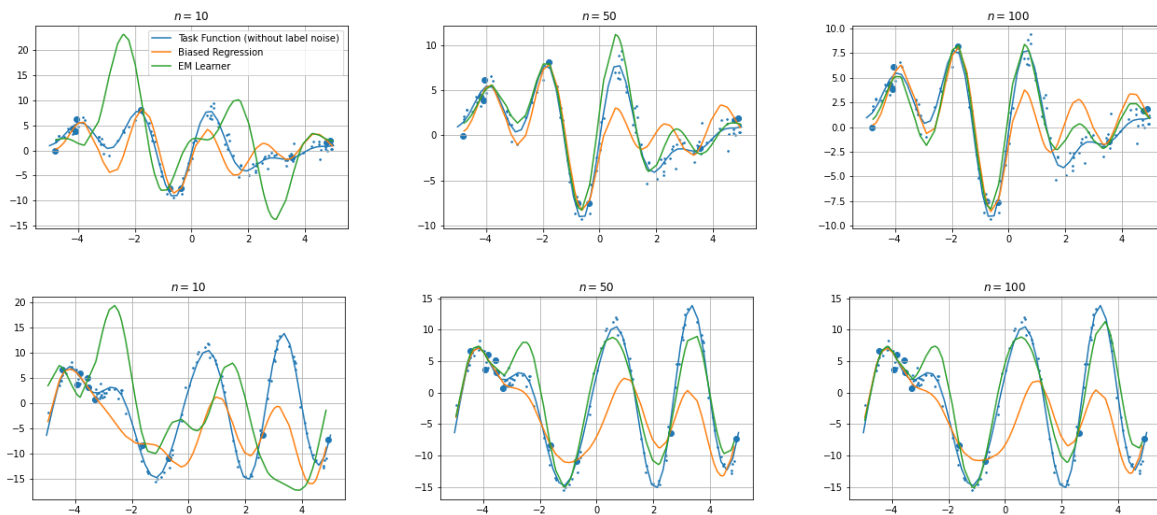


Figure 7. Two sets of examples of predictions on the synthetic, ‘Fourier’ meta-learning problem. Top and bottom rows correspond to different (random) instances; the top row in fact replicates Fig. 1. Training data is shown in bold, small dots show test data. We also show the predictions for two learners (at every input) and the target function. The column correspond to outputs obtained training on $n \in \{10, 50, 100\}$ tasks.