
Provable Generalization of Overparameterized Meta-learning Trained with SGD

Yu Huang
IIS
Tsinghua University
y-huang20@mails.tsinghua.edu.cn

Yingbin Liang
Department of ECE
The Ohio State University
liang.889@osu.edu

Longbo Huang*
IIS
Tsinghua University
longbohuang@tsinghua.edu.cn

Appendices

A Proof of Proposition 1	2
B Analysis for Upper Bound (Theorem 1)	4
B.1 Preliminaries	4
B.2 Fourth Moment Upper Bound for Meta Data	5
B.3 Bias-Variance Decomposition	7
B.4 Bounding the Bias	9
B.5 Bounding the Variance	13
B.6 Proof of Theorem 1	16
C Analysis for Lower Bound (Theorem 2)	17
C.1 Fourth Moment Lower Bound for Meta Noise	17
C.2 Bias-Variance Decomposition	18
C.3 Bounding the Bias	19
C.4 Bounding the Variance	21
C.5 Proof of Theorem 2	23
D Proofs for Section 4.2	23
D.1 Proof of Lemma 1	23
D.2 Proof of Proposition 2	24

*Corresponding author

D.3	Proof of Proposition 3	25
E	Proofs for Section 4.3	25
E.1	Proof of Proposition 4	25
E.2	Proof of Corollary 1	26
F	Discussions on Assumptions	26
G	Further Related Work	27
G.1	Underparameterized Setting	27
G.2	Overparameterized Setting	28
H	Future Directions	28
H.1	Generalizing to Other Learning Methods	28
H.2	Meta-batch Setting	29
H.3	Longer Inner Loops	29

A Proof of Proposition 1

We first show how to connect the loss function associated with MAML to a Meta Least Square Problem.

Proposition A.1 (Proposition 1 Restated). *Under the mixed linear regression model, the expectation of the meta-training loss taken over task and data distributions can be rewritten as:*

$$\mathbb{E} \left[\widehat{\mathcal{L}}(\mathcal{A}, \boldsymbol{\omega}, \beta^{\text{tr}}; \mathcal{D}) \right] = \mathcal{L}(\mathcal{A}, \boldsymbol{\omega}, \beta^{\text{tr}}) = \mathbb{E}_{\mathbf{B}, \gamma} \frac{1}{2} \left[\|\mathbf{B}\boldsymbol{\omega} - \gamma\|^2 \right]. \quad (1)$$

The meta data are given by

$$\mathbf{B} = \frac{1}{\sqrt{n_2}} \mathbf{X}^{\text{out}} \left(\mathbf{I} - \frac{\beta^{\text{tr}}}{n_1} \mathbf{X}^{\text{in}T} \mathbf{X}^{\text{in}} \right) \quad (2)$$

$$\gamma = \frac{1}{\sqrt{n_2}} \left(\mathbf{X}^{\text{out}} \left(\mathbf{I} - \frac{\beta^{\text{tr}}}{n_1} \mathbf{X}^{\text{in}T} \mathbf{X}^{\text{in}} \right) \boldsymbol{\theta} + \mathbf{z}^{\text{out}} - \frac{\beta^{\text{tr}}}{n_1} \mathbf{X}^{\text{out}} \mathbf{X}^{\text{in}T} \mathbf{z}^{\text{in}} \right) \quad (3)$$

where $\mathbf{X}^{\text{in}} \in \mathbb{R}^{n_1 \times d}$, $\mathbf{z}^{\text{in}} \in \mathbb{R}^{n_1}$, $\mathbf{X}^{\text{out}} \in \mathbb{R}^{n_2 \times d}$ and $\mathbf{z}^{\text{out}} \in \mathbb{R}^{n_2}$ denote the inputs and noise for training and validation. Furthermore, we have

$$\gamma = \mathbf{B}\boldsymbol{\theta}^* + \boldsymbol{\xi} \quad \text{with meta noise } \mathbb{E}[\boldsymbol{\xi} \mid \mathbf{B}] = 0. \quad (4)$$

Proof. We first rewrite $\mathcal{L}(\mathcal{A}, \boldsymbol{\omega}, \beta^{\text{tr}})$ as follows:

$$\begin{aligned} \mathcal{L}(\mathcal{A}, \boldsymbol{\omega}, \beta^{\text{tr}}) &= \mathbb{E} \left[\ell(\mathcal{A}(\boldsymbol{\omega}, \beta^{\text{tr}}; \mathcal{D}^{\text{in}}); \mathcal{D}^{\text{out}}) \right] \\ &= \mathbb{E} \left[\frac{1}{2n_2} \sum_{j=1}^{n_2} \left(\langle \mathbf{x}_j^{\text{out}}, \mathcal{A}(\boldsymbol{\omega}, \beta^{\text{tr}}; \mathcal{D}^{\text{in}}) \rangle - y_j^{\text{out}} \right)^2 \right] \\ &= \mathbb{E} \left[\frac{1}{2n_2} \left\| \mathbf{X}^{\text{out}} \left(\mathbf{I} - \frac{\beta^{\text{tr}}}{n_1} \mathbf{X}^{\text{in}T} \mathbf{X}^{\text{in}} \right) \boldsymbol{\omega} + \frac{\beta^{\text{tr}}}{n_1} \mathbf{X}^{\text{in}T} \mathbf{y}^{\text{in}} - \mathbf{y}^{\text{out}} \right\|^2 \right]. \end{aligned}$$

Using the mixed linear model:

$$\mathbf{y}^{\text{in}} = \mathbf{X}^{\text{in}} \boldsymbol{\theta} + \mathbf{z}^{\text{in}}, \quad \mathbf{y}^{\text{out}} = \mathbf{X}^{\text{out}} \boldsymbol{\theta} + \mathbf{z}^{\text{out}}, \quad (5)$$

we have

$$\begin{aligned}
\mathcal{L}(\mathcal{A}, \boldsymbol{\omega}, \beta^{\text{tr}}) &= \mathbb{E} \left[\frac{1}{2n_2} \|\mathbf{X}^{\text{out}} \left(\mathbf{I} - \frac{\beta^{\text{tr}}}{n_1} \mathbf{X}^{\text{in}T} \mathbf{X}^{\text{in}} \right) \boldsymbol{\omega} \right. \\
&\quad \left. - \left(\mathbf{X}^{\text{out}} \left(\mathbf{I} - \frac{\beta^{\text{tr}}}{n_1} \mathbf{X}^{\text{in}T} \mathbf{X}^{\text{in}} \right) \boldsymbol{\theta} + \mathbf{z}^{\text{out}} - \frac{\beta^{\text{tr}}}{n_1} \mathbf{X}^{\text{out}} \mathbf{X}^{\text{in}T} \mathbf{z}^{\text{in}} \right) \right\|^2 \Big] \\
&= \mathbb{E}_{\mathbf{B}, \gamma} \frac{1}{2} \left[\|\mathbf{B}\boldsymbol{\omega} - \gamma\|^2 \right].
\end{aligned}$$

Moreover, note that $\boldsymbol{\theta} - \boldsymbol{\theta}^*$ has mean zero and is independent of data and noise, and define

$$\xi = \frac{1}{\sqrt{n_2}} \left(\mathbf{X}^{\text{out}} \left(\mathbf{I} - \frac{\beta^{\text{tr}}}{n_1} \mathbf{X}^{\text{in}T} \mathbf{X}^{\text{in}} \right) (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \mathbf{z}^{\text{out}} - \frac{\beta^{\text{tr}}}{n_1} \mathbf{X}^{\text{out}} \mathbf{X}^{\text{in}T} \mathbf{z}^{\text{in}} \right). \quad (6)$$

We call ξ as meta noise, and then we have

$$\gamma = \mathbf{B}\boldsymbol{\theta}^* + \xi \quad \text{and} \quad \mathbb{E}[\xi \mid \mathbf{B}] = 0.$$

□

Lemma A.1 (Meta Excess Risk). *Under the mixed linear regression model, the meta excess risk can be rewritten as follows:*

$$R(\bar{\boldsymbol{\omega}}_T, \beta^{\text{te}}) = \frac{1}{2} \mathbb{E} \|\bar{\boldsymbol{\omega}}_T - \boldsymbol{\theta}^*\|_{\mathbf{H}_{m, \beta^{\text{te}}}}^2$$

where $\|\mathbf{a}\|_{\mathbf{A}}^2 = \mathbf{a}^T \mathbf{A} \mathbf{a}$. Moreover, the Bayes error is given by

$$\mathcal{L}(\mathcal{A}, \boldsymbol{\omega}^*, \beta^{\text{te}}) = \frac{1}{2} \text{tr}(\boldsymbol{\Sigma}_{\boldsymbol{\theta}} \mathbf{H}_{m, \beta^{\text{te}}}) + \frac{\sigma^2 \beta^{\text{te}2}}{2m} + \frac{\sigma^2}{2}.$$

Proof. Recall that

$$R(\bar{\boldsymbol{\omega}}_T, \beta^{\text{te}}) \triangleq \mathbb{E} [\mathcal{L}(\mathcal{A}, \bar{\boldsymbol{\omega}}_T, \beta^{\text{te}})] - \mathcal{L}(\mathcal{A}, \boldsymbol{\omega}^*, \beta^{\text{te}})$$

where $\boldsymbol{\omega}^*$ denotes the optimal solution to the population meta-test error. Under the mixed linear model, such a solution can be directly calculated [11], and we obtain $\boldsymbol{\omega}^* = \mathbb{E}[\boldsymbol{\theta}] = \boldsymbol{\theta}^*$. Hence,

$$R(\bar{\boldsymbol{\omega}}_T, \beta^{\text{te}}) = \mathbb{E}_{\mathbf{B}, \gamma} \frac{1}{2} \left[\|\mathbf{B}\bar{\boldsymbol{\omega}}_T - \gamma\|^2 - \|\mathbf{B}\boldsymbol{\theta}^* - \gamma\|^2 \right],$$

where

$$\begin{aligned}
\mathbf{B} &= \mathbf{x}^{\text{out}T} \left(\mathbf{I} - \frac{\beta^{\text{te}}}{m} \mathbf{X}^{\text{in}T} \mathbf{X}^{\text{in}} \right) \\
\gamma &= \mathbf{x}^{\text{out}T} \left(\mathbf{I} - \frac{\beta^{\text{te}}}{m} \mathbf{X}^{\text{in}T} \mathbf{X}^{\text{in}} \right) \boldsymbol{\theta} + \mathbf{z}^{\text{out}} - \frac{\beta^{\text{te}}}{m} \mathbf{x}^{\text{out}T} \mathbf{X}^{\text{in}T} \mathbf{z}^{\text{in}}, \quad (7)
\end{aligned}$$

and $\mathbf{x}^{\text{out}} \in \mathbb{R}^d$, $\mathbf{z}^{\text{out}} \in \mathbb{R}^d$, $\mathbf{X}^{\text{in}} \in \mathbb{R}^{m \times d}$ and $\mathbf{z}^{\text{in}} \in \mathbb{R}^m$. The forms of \mathbf{B} and γ are slightly different since we allow a new adaptation rate β^{te} and the inner loop has m samples at test stage. Similarly

$$\xi = \left(\underbrace{\mathbf{x}^{\text{out}T} \left(\mathbf{I} - \frac{\beta^{\text{te}}}{m} \mathbf{X}^{\text{in}T} \mathbf{X}^{\text{in}} \right) (\boldsymbol{\theta} - \boldsymbol{\theta}^*)}_{\xi_1} + \underbrace{\mathbf{z}^{\text{out}}}_{\xi_2} - \underbrace{\frac{\beta^{\text{te}}}{m} \mathbf{x}^{\text{out}T} \mathbf{X}^{\text{in}T} \mathbf{z}^{\text{in}}}_{\xi_3} \right). \quad (8)$$

Then we have

$$\begin{aligned}
R(\bar{\boldsymbol{\omega}}_T, \beta^{\text{te}}) &= \mathbb{E}_{\mathbf{B}, \gamma} \frac{1}{2} \left[\|\mathbf{B}\bar{\boldsymbol{\omega}}_T - \gamma\|^2 - \|\mathbf{B}\boldsymbol{\theta}^* - \gamma\|^2 \right] \\
&= \mathbb{E}_{\mathbf{B}, \gamma} \frac{1}{2} \left[\|\mathbf{B}(\bar{\boldsymbol{\omega}}_T - \boldsymbol{\theta}^*)\|^2 \right] \\
&= \frac{1}{2} \mathbb{E} \|\bar{\boldsymbol{\omega}}_T - \boldsymbol{\theta}^*\|_{\mathbf{H}_{m, \beta^{\text{te}}}}^2
\end{aligned}$$

where the last equality follows because $\mathbb{E} [\mathbf{B}^\top \mathbf{B}] = \mathbf{H}_{m, \beta^{\text{te}}}$ at the test stage.

The Bayes error can be calculated as follows:

$$\begin{aligned} \mathcal{L}(\mathcal{A}, \boldsymbol{\omega}^*, \beta^{\text{te}}) &= \mathbb{E}_{\mathbf{B}, \gamma} \frac{1}{2} \left[\|\mathbf{B}\boldsymbol{\theta}^* - \gamma\|^2 \right] = \mathbb{E}_{\mathbf{B}, \gamma} \frac{1}{2} [\xi^2] \\ &\stackrel{(a)}{=} \frac{1}{2} (\mathbb{E} [\xi_1^2] + \mathbb{E} [\xi_2^2] + \mathbb{E} [\xi_3^2]) \\ &= \frac{1}{2} (\text{tr}(\boldsymbol{\Sigma}_\theta \mathbf{H}_{m, \beta^{\text{te}}}) + \frac{\beta^{\text{te}2} \sigma^2}{m} + \sigma^2) \end{aligned}$$

where (a) follows because ξ_1, ξ_2, ξ_3 are independent and have zero mean conditioned on \mathbf{X}^{in} and \mathbf{x}^{out} . \square

B Analysis for Upper Bound (Theorem 1)

B.1 Preliminaries

We first introduce some additional notations.

Definition 1 (Inner product of matrices). For any two matrices \mathbf{C}, \mathbf{D} , the inner product of them is defined as

$$\langle \mathbf{C}, \mathbf{D} \rangle = \text{tr}(\mathbf{C}^\top \mathbf{D}).$$

We will use the following property about the inner product of matrices throughout our proof.

Property B.1. If $\mathbf{C} \succeq 0$ and $\mathbf{D} \succeq \mathbf{D}'$, then we have $\langle \mathbf{C}, \mathbf{D} \rangle \geq \langle \mathbf{C}, \mathbf{D}' \rangle$.

Definition 2 (Linear operator). Let \otimes denote the tensor product. Define the following linear operators on symmetric matrices:

$$\begin{aligned} \mathcal{M} &= \mathbb{E} [\mathbf{B}^\top \otimes \mathbf{B}^\top \otimes \mathbf{B} \otimes \mathbf{B}] \quad \widetilde{\mathcal{M}} := \mathbf{H}_{n_1, \beta^{\text{tr}}} \otimes \mathbf{H}_{n_1, \beta^{\text{tr}}} \quad \mathcal{I} := \mathbf{I} \otimes \mathbf{I} \\ \mathcal{T} &:= \mathbf{H}_{n_1, \beta^{\text{tr}}} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}_{n_1, \beta^{\text{tr}}} - \alpha \mathcal{M}, \quad \widetilde{\mathcal{T}} = \mathbf{H}_{n_1, \beta^{\text{tr}}} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{H}_{n_1, \beta^{\text{tr}}} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}} \otimes \mathbf{H}_{n_1, \beta^{\text{tr}}}. \end{aligned}$$

We next define the operation of the above linear operators on a symmetric matrix \mathbf{A} as follows.

$$\begin{aligned} \mathcal{M} \circ \mathbf{A} &= \mathbb{E} [\mathbf{B}^\top \mathbf{B} \mathbf{A} \mathbf{B}^\top \mathbf{B}], \quad \widetilde{\mathcal{M}} \circ \mathbf{A} = \mathbf{H}_{n_1, \beta^{\text{tr}}} \mathbf{A} \mathbf{H}_{n_1, \beta^{\text{tr}}}, \quad \mathcal{I} \circ \mathbf{A} = \mathbf{A}, \\ \mathcal{T} \circ \mathbf{A} &= \mathbf{H}_{n_1, \beta^{\text{tr}}} \mathbf{A} + \mathbf{A} \mathbf{H}_{n_1, \beta^{\text{tr}}} - \alpha \mathbb{E} [\mathbf{B}^\top \mathbf{B} \mathbf{A} \mathbf{B}^\top \mathbf{B}] \\ \widetilde{\mathcal{T}} \circ \mathbf{A} &= \mathbf{H}_{n_1, \beta^{\text{tr}}} \mathbf{A} + \mathbf{A} \mathbf{H}_{n_1, \beta^{\text{tr}}} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}} \mathbf{A} \mathbf{H}_{n_1, \beta^{\text{tr}}}. \end{aligned}$$

Based on the above definitions, we have the following equations hold.

$$\begin{aligned} (\mathcal{I} - \alpha \mathcal{T}) \circ \mathbf{A} &= \mathbb{E} [(\mathbf{I} - \alpha \mathbf{B}^\top \mathbf{B}) \mathbf{A} (\mathbf{I} - \alpha \mathbf{B}^\top \mathbf{B})] \\ (\mathcal{I} - \alpha \widetilde{\mathcal{T}}) \circ \mathbf{A} &= (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}}) \mathbf{A} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}}). \end{aligned}$$

For the linear operators, we have the following technical lemma.

Lemma B.1. We call the linear operator \mathcal{O} a PSD mapping, if for every symmetric PSD matrix \mathbf{A} , $\mathcal{O} \circ \mathbf{A}$ is also PSD matrix. Then we have:

- (i) \mathcal{M} , $\widetilde{\mathcal{M}}$ and $(\mathcal{M} - \widetilde{\mathcal{M}})$ are all PSD mappings.
- (ii) $\widetilde{\mathcal{T}} - \mathcal{T}$, $\mathcal{I} - \alpha \mathcal{T}$ and $\mathcal{I} - \alpha \widetilde{\mathcal{T}}$ are all PSD mappings.
- (iii) If $0 < \alpha < \frac{1}{\max_i \{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})\}}$, then $\widetilde{\mathcal{T}}^{-1}$ exists, and is a PSD mapping.
- (iv) If $0 < \alpha < \frac{1}{\max_i \{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})\}}$, $\widetilde{\mathcal{T}}^{-1} \circ \mathbf{H}_{n_1, \beta^{\text{tr}}} \preceq \mathbf{I}$.
- (v) If $0 < \alpha < \frac{1}{c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})}$, then $\mathcal{T}^{-1} \circ \mathbf{A}$ exists for PSD matrix \mathbf{A} , and \mathcal{T}^{-1} is a PSD mapping.

Proof. Items (i) and (iii) directly follow from the proofs in [14, 22]. For (iv), by the existence of $\tilde{\mathcal{T}}^{-1}$, we have

$$\begin{aligned}\tilde{\mathcal{T}}^{-1} \circ \mathbf{H}_{n_1, \beta^{\text{tr}}} &= \sum_{t=0}^{\infty} \alpha (\mathcal{I} - \alpha \tilde{\mathcal{T}})^t \circ \mathbf{H}_{n_1, \beta^{\text{tr}}} \\ &= \sum_{t=0}^{\infty} \alpha (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^t \mathbf{H}_{n_1, \beta^{\text{tr}}} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^t \\ &\preceq \sum_{t=0}^{\infty} \alpha (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^t \mathbf{H}_{n_1, \beta^{\text{tr}}} = \mathbf{I}.\end{aligned}$$

For (v), for any PSD matrix \mathbf{A} , consider

$$\mathcal{T}^{-1} \circ \mathbf{A} = \alpha \sum_{k=0}^{\infty} (\mathcal{I} - \alpha \mathcal{T})^k \circ \mathbf{A}.$$

We first show that $\sum_{k=0}^{\infty} (\mathcal{I} - \alpha \mathcal{T})^k \circ \mathbf{A}$ is finite, and then it suffices to show that the trace is finite, i.e.,

$$\sum_{k=0}^{\infty} \text{tr}((\mathcal{I} - \alpha \mathcal{T})^k \circ \mathbf{A}) < \infty. \quad (9)$$

Let $\mathbf{A}_k = (\mathcal{I} - \alpha \mathcal{T})^k \circ \mathbf{A}$. Combining with the definition of \mathcal{T} , we obtain

$$\text{tr}(\mathbf{A}_k) = \text{tr}(\mathbf{A}_{k-1}) - 2\alpha \text{tr}(\mathbf{H}_{n_1, \beta^{\text{tr}}} \mathbf{A}_{k-1}) + \alpha^2 \text{tr}(\mathbf{A} \mathbb{E}[\mathbf{B}^{\top} \mathbf{B} \mathbf{B}^{\top} \mathbf{B}]).$$

Letting $\mathbf{A} = \mathbf{I}$ in Proposition B.1, we have $\mathbb{E}[\mathbf{B}^{\top} \mathbf{B} \mathbf{B}^{\top} \mathbf{B}] \preceq c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}) \mathbf{H}_{n_1, \beta^{\text{tr}}}$. Hence

$$\begin{aligned}\text{tr}(\mathbf{A}_k) &\leq \text{tr}(\mathbf{A}_{k-1}) - (2\alpha - \alpha^2 c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})) \text{tr}(\mathbf{H}_{n_1, \beta^{\text{tr}}} \mathbf{A}_{k-1}) \\ &\leq \text{tr}((\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}}) \mathbf{A}_{k-1}) \quad \text{by } \alpha < \frac{1}{c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})} \\ &\leq \left(1 - \alpha \min_i \{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})\}\right) \text{tr}(\mathbf{A}_{k-1}).\end{aligned}$$

If $\alpha < \frac{1}{\min_i \{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})\}}$, then we substitute it into eq. (9) and obtain

$$\sum_{k=0}^{\infty} \text{tr}((\mathcal{I} - \alpha \mathcal{T})^k \circ \mathbf{A}) = \sum_{k=0}^{\infty} \text{tr}(\mathbf{A}_k) \leq \frac{\text{tr}(\mathbf{A})}{\alpha \min_i \{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})\}} < \infty$$

which guarantees the existence of \mathcal{T}^{-1} . Moreover, \mathbf{A}_k is a PSD matrix for every k since $\mathcal{I} - \alpha \mathcal{T}$ is a PSD mapping. The $\mathcal{T}^{-1} \circ \mathbf{A} = \alpha \sum_{k=0}^{\infty} \mathbf{A}_k$ must be a PSD matrix, which implies that \mathcal{T}^{-1} is PSD mapping. \square

Property B.2 (Commutity). *Suppose Assumption 2 holds, then for all $n > 0$, $|\beta| < 1/\lambda_1$, $\mathbf{H}_{n, \beta}$ with different n and β commute with each other.*

B.2 Fourth Moment Upper Bound for Meta Data

In this section, we provide a technical result for the fourth moment of meta data \mathbf{B} , which is essential throughout the proof of our upper bound.

Proposition B.1. *Suppose Assumptions 1-3 hold. Given $|\beta| < \frac{1}{\lambda_1}$, for any PSD matrix \mathbf{A} , we have*

$$\mathbb{E}[\mathbf{B}^{\top} \mathbf{B} \mathbf{A} \mathbf{B} \mathbf{B}^{\top} \mathbf{B}] \preceq c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) \mathbb{E}[\text{tr}(\mathbf{A} \boldsymbol{\Sigma})] \mathbf{H}_{n_1, \beta^{\text{tr}}}$$

where $c(\beta, \boldsymbol{\Sigma}) := c_1 \left(1 + 8|\beta| \lambda_1 \sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^2 + 64 \sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^4 \beta^2 \text{tr}(\boldsymbol{\Sigma}^2)\right)$.

Proof. Recall that $\mathbf{B} = \frac{1}{\sqrt{n_2}} \mathbf{X}^{\text{out}} (\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^{\text{in}^\top} \mathbf{X}^{\text{in}})$. With a slight abuse of notations, we write β^{tr} as β , \mathbf{X}^{in} as \mathbf{X} in this proof. First consider the case $\beta \geq 0$. By the definition of \mathbf{B} , we have

$$\begin{aligned} & \mathbb{E} [\mathbf{B}^\top \mathbf{B} \mathbf{A} \mathbf{B}^\top \mathbf{B}] \\ &= \mathbb{E} \left[\left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \frac{1}{n_2} \mathbf{X}^{\text{out}^\top} \mathbf{X}^{\text{out}} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{A} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \frac{1}{n_2} \mathbf{X}^{\text{out}^\top} \mathbf{X}^{\text{out}} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\ &\preceq c_1 \mathbb{E} \left[\text{tr} \left(\left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{A} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \right) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\ &\preceq c_1 \mathbb{E} \left[\text{tr} \left(\mathbf{A} \left(\boldsymbol{\Sigma} + \frac{\beta^2}{n_1^2} \mathbf{X}^\top \mathbf{X} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{X} \right) \right) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \end{aligned}$$

where the second inequality follows from Assumption 1. Let \mathbf{x}_i denote the i -th row of \mathbf{X} . Note that $\mathbf{x}_i = \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{z}_i$, where \mathbf{z}_i is independent σ_x -sub-gaussian vector. For any $\mathbf{x}_{i_1}, \mathbf{x}_{i_2}, \mathbf{x}_{i_3}, \mathbf{x}_{i_4}$, where $1 \leq i_1, i_2, i_3, i_4 \leq n_1$, we have:

$$\begin{aligned} & \mathbb{E} \left[\text{tr} (\mathbf{A} \mathbf{x}_{i_1} \mathbf{x}_{i_2}^\top \boldsymbol{\Sigma} \mathbf{x}_{i_3} \mathbf{x}_{i_4}^\top) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\ &= \mathbb{E} \left[\text{tr} (\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{A} \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{z}_{i_1} \mathbf{z}_{i_2}^\top \boldsymbol{\Sigma}^2 \mathbf{z}_{i_3} \mathbf{z}_{i_4}^\top) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\ &= \sum_{k,j} \mu_k \lambda_j^2 \mathbb{E} \left[(\mathbf{z}_{i_4}^\top \mathbf{u}_k) (\mathbf{z}_{i_1}^\top \mathbf{u}_k) (\mathbf{z}_{i_4}^\top \mathbf{v}_j) (\mathbf{z}_{i_1}^\top \mathbf{v}_j) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \end{aligned}$$

where the SVD of $\boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{A} \boldsymbol{\Sigma}^{\frac{1}{2}}$ is $\sum_j \mu_j \mathbf{u}_j \mathbf{u}_j^\top$, the SVD of $\boldsymbol{\Sigma}$ is $\sum_j \lambda_j \mathbf{v}_j \mathbf{v}_j^\top$. For any unit vector $\mathbf{w} \in \mathbb{R}^d$, we have:

$$\begin{aligned} & \mathbf{w}^\top \mathbb{E} \left[\mathbf{H}_{n_1, \beta}^{-\frac{1}{2}} \text{tr} (\mathbf{A} \mathbf{x}_{i_1} \mathbf{x}_{i_2}^\top \boldsymbol{\Sigma} \mathbf{x}_{i_3} \mathbf{x}_{i_4}^\top) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{H}_{n_1, \beta}^{-\frac{1}{2}} \right] \mathbf{w} \\ &\leq \sum_{k,j} \mu_k \lambda_j^2 \sqrt{\mathbb{E} \left[((\mathbf{z}_{i_4}^\top \mathbf{u}_k) (\mathbf{z}_{i_1}^\top \mathbf{u}_k) (\mathbf{z}_{i_4}^\top \mathbf{v}_j) (\mathbf{z}_{i_1}^\top \mathbf{v}_j))^2 \right]} \\ &\quad \times \sqrt{\mathbb{E} \left[\left\| \mathbf{w}^\top \mathbf{H}_{n_1, \beta}^{-\frac{1}{2}} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{H}_{n_1, \beta}^{-\frac{1}{2}} \mathbf{w} \right\|^2 \right]} \\ &\leq 64 \sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^4 \text{tr}(\mathbf{A} \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}^2) \end{aligned}$$

where the first inequality follows from the Cauchy Schwarz inequality; the last inequality is due to Assumption 3 and the property of sub-Gaussian distributions [19]. Therefore,

$$\begin{aligned} & \mathbb{E} \left[\mathbf{H}_{n_1, \beta}^{-\frac{1}{2}} \text{tr} (\mathbf{A} \mathbf{x}_{i_1} \mathbf{x}_{i_2}^\top \boldsymbol{\Sigma} \mathbf{x}_{i_3} \mathbf{x}_{i_4}^\top) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{H}_{n_1, \beta}^{-\frac{1}{2}} \right] \\ &\preceq 64 \sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^4 \text{tr}(\mathbf{A} \boldsymbol{\Sigma}^2) \mathbf{I} \end{aligned}$$

which implies

$$\mathbb{E} \left[\text{tr} (\mathbf{A} \mathbf{x}_{i_1} \mathbf{x}_{i_2}^\top \boldsymbol{\Sigma} \mathbf{x}_{i_3} \mathbf{x}_{i_4}^\top) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \preceq 64 \sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^4 \text{tr}(\mathbf{A} \boldsymbol{\Sigma}^2) \mathbf{H}_{n_1, \beta}.$$

Hence,

$$\begin{aligned} & \mathbb{E} [\mathbf{B}^\top \mathbf{B} \mathbf{A} \mathbf{B}^\top \mathbf{B}] \\ &\preceq c_1 \mathbb{E} \left[\text{tr} \left(\mathbf{A} \left(\boldsymbol{\Sigma} + 64 \sqrt{C} \sigma_x^4 \beta^2 \boldsymbol{\Sigma} \text{tr}(\boldsymbol{\Sigma}^2) \right) \right) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\ &\preceq c_1 (1 + 64 \sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^4 \beta^2 \text{tr}(\boldsymbol{\Sigma}^2)) \mathbb{E} [\text{tr}(\mathbf{A} \boldsymbol{\Sigma})] \mathbf{H}_{n_1, \beta}. \end{aligned}$$

Now we turn to $\beta < 0$, and derive

$$\begin{aligned}
& \mathbb{E} [\mathbf{B}^\top \mathbf{B} \mathbf{A} \mathbf{B}^\top \mathbf{B}] \\
& \preceq c_1 \mathbb{E} \left[\text{tr} \left(\left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{A} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right) \right] \\
& = c_1 \mathbb{E} \left[\text{tr} \left(\mathbf{A} \left(\boldsymbol{\Sigma} - \underbrace{\frac{\beta}{n_1} (\mathbf{X}^\top \mathbf{X} \boldsymbol{\Sigma} + \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{X})}_{\mathbf{J}_1} + \frac{\beta^2}{n_1^2} \mathbf{X}^\top \mathbf{X} \boldsymbol{\Sigma} \mathbf{X}^\top \mathbf{X} \right) \right) \right. \\
& \quad \left. \cdot \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right].
\end{aligned}$$

We can bound the extra term \mathbf{J}_1 in the similar way as $\beta > 0$. For any \mathbf{x}_i , $1 \leq i \leq n_1$, we have

$$\begin{aligned}
& \mathbb{E} \left[\text{tr} \left(\mathbf{A} \mathbf{x}_i \mathbf{x}_i^\top \boldsymbol{\Sigma} \right) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\
& = \mathbb{E} \left[\text{tr} \left(\mathbf{z}_i^\top \boldsymbol{\Sigma}^{\frac{3}{2}} \mathbf{A} \boldsymbol{\Sigma}^{\frac{1}{2}} \mathbf{z}_i \right) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\
& = \sum_k \iota_k \mathbb{E} \left[(\mathbf{z}_i^\top \boldsymbol{\kappa}_k)^2 \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right]
\end{aligned}$$

where the SVD of $\boldsymbol{\Sigma}^{\frac{3}{2}} \mathbf{A} \boldsymbol{\Sigma}^{\frac{1}{2}}$ is $\sum_k \iota_k \boldsymbol{\kappa}_k \boldsymbol{\kappa}_k^\top$. Similarly, for any unit vector $\mathbf{w} \in \mathbb{R}^d$, we can obtain

$$\begin{aligned}
& \mathbf{w}^\top \mathbb{E} \left[\mathbf{H}_{n_1, \beta}^{-\frac{1}{2}} \text{tr} \left(\mathbf{A} \mathbf{x}_i \mathbf{x}_i^\top \boldsymbol{\Sigma} \right) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{H}_{n_1, \beta}^{-\frac{1}{2}} \right] \mathbf{w} \\
& \leq \sum_k \iota_k \sqrt{\mathbb{E}[(\mathbf{z}_i^\top \boldsymbol{\kappa}_k)^4]} \sqrt{\mathbb{E}[\|\mathbf{w}^\top \mathbf{H}_{n_1, \beta}^{-\frac{1}{2}} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{H}_{n_1, \beta}^{-\frac{1}{2}} \mathbf{w}\|^2]} \\
& \leq 4\sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^2 \text{tr}(\mathbf{A} \boldsymbol{\Sigma}^2)
\end{aligned}$$

which implies:

$$\mathbb{E} \left[\text{tr} \left(\mathbf{A} \mathbf{x}_i \mathbf{x}_i^\top \boldsymbol{\Sigma} \right) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \preceq 4\sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^2 \text{tr}(\mathbf{A} \boldsymbol{\Sigma}^2) \mathbf{H}_{n_1, \beta}.$$

Hence,

$$\begin{aligned}
& \mathbb{E} [\mathbf{B}^\top \mathbf{B} \mathbf{A} \mathbf{B}^\top \mathbf{B}] \\
& \preceq c_1 \mathbb{E} \left[\text{tr} \left(\mathbf{A} \left(\boldsymbol{\Sigma} - 8\beta\sqrt{C} \sigma_x^2 \boldsymbol{\Sigma}^2 + 64\sqrt{C} \sigma_x^4 \beta^2 \boldsymbol{\Sigma} \text{tr}(\boldsymbol{\Sigma}^2) \right) \right) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\
& \preceq c_1 \left(1 - 8\beta\lambda_1 \sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^2 + 64\sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^4 \beta^2 \text{tr}(\boldsymbol{\Sigma}^2) \right) \mathbb{E} [\text{tr}(\mathbf{A} \boldsymbol{\Sigma})] \mathbf{H}_{n_1, \beta}.
\end{aligned}$$

Together with the discussions for $\beta > 0$, we have

$$c(\beta, \boldsymbol{\Sigma}) = c_1 (1 + 8|\beta|\lambda_1 \sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^2 + 64\sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^4 \beta^2 \text{tr}(\boldsymbol{\Sigma}^2)),$$

which completes the proof. \square

B.3 Bias-Variance Decomposition

We will use the bias-variance decomposition similar to theoretical studies of classic linear regression [14, 8, 22]. Consider the error at each iteration: $\boldsymbol{\varrho}_t = \boldsymbol{\omega}_t - \boldsymbol{\theta}^*$, where $\boldsymbol{\omega}_t$ is the SGD output at each iteration t . Then the update rule can be written as:

$$\boldsymbol{\varrho}_t := (\mathbf{I} - \alpha \mathbf{B}_t^\top \mathbf{B}_t) \boldsymbol{\varrho}_{t-1} + \alpha \mathbf{B}_t^\top \boldsymbol{\xi}_t$$

where $\mathbf{B}_t, \boldsymbol{\xi}_t$ are the meta data and noise at iteration t (see eqs. (2) and (6)). It is helpful to consider $\boldsymbol{\varrho}_t$ as the sum of the following two random processes:

- If there is no meta noise, the error comes from the bias:

$$\boldsymbol{\varrho}_t^{\text{bias}} := (\mathbf{I} - \alpha \mathbf{B}_t^\top \mathbf{B}_t) \boldsymbol{\varrho}_{t-1}^{\text{bias}} \quad \boldsymbol{\varrho}_t^{\text{bias}} = \boldsymbol{\varrho}_0.$$

- If the SGD trajectory starts from $\boldsymbol{\theta}^*$, the error originates from the variance:

$$\boldsymbol{\varrho}_t^{\text{var}} := (\mathbf{I} - \alpha \mathbf{B}_t^\top \mathbf{B}_t) \boldsymbol{\varrho}_{t-1}^{\text{var}} + \alpha \mathbf{B}_t^\top \boldsymbol{\xi}_t \quad \boldsymbol{\varrho}^{\text{var}} = \mathbf{0}$$

and $\mathbb{E}[\boldsymbol{\varrho}_t^{\text{var}}] = \mathbf{0}$.

With slightly abused notations, we have:

$$\boldsymbol{\varrho}_t = \boldsymbol{\varrho}_t^{\text{bias}} + \boldsymbol{\varrho}_t^{\text{var}}.$$

Define the averaged output of $\boldsymbol{\varrho}_t^{\text{bias}}$, $\boldsymbol{\varrho}_t^{\text{var}}$ and $\boldsymbol{\varrho}_t$ after T iterations as:

$$\bar{\boldsymbol{\varrho}}_T^{\text{bias}} = \frac{1}{T} \sum_{t=1}^T \boldsymbol{\varrho}_t^{\text{bias}}, \quad \bar{\boldsymbol{\varrho}}_T^{\text{var}} = \frac{1}{T} \sum_{t=1}^T \boldsymbol{\varrho}_t^{\text{var}}, \quad \bar{\boldsymbol{\varrho}}_T = \frac{1}{T} \sum_{t=1}^T \boldsymbol{\varrho}_t. \quad (10)$$

Similarly, we have

$$\bar{\boldsymbol{\varrho}}_T = \bar{\boldsymbol{\varrho}}_T^{\text{bias}} + \bar{\boldsymbol{\varrho}}_T^{\text{var}}.$$

Now we are ready to introduce the bias-variance decomposition for the excess risk.

Lemma B.2 (Bias-variance decomposition). *Following the notations in eq. (10), then the excess risk can be decomposed as*

$$R(\bar{\boldsymbol{\omega}}_T, \beta^{\text{te}}) \leq 2\mathcal{E}_{\text{bias}} + 2\mathcal{E}_{\text{var}}$$

where

$$\mathcal{E}_{\text{bias}} = \frac{1}{2} \langle \mathbf{H}_{m, \beta^{\text{te}}}, \mathbb{E}[\bar{\boldsymbol{\varrho}}_T^{\text{bias}} \otimes \bar{\boldsymbol{\varrho}}_T^{\text{bias}}] \rangle, \quad \mathcal{E}_{\text{var}} = \frac{1}{2} \langle \mathbf{H}_{m, \beta^{\text{te}}}, \mathbb{E}[\bar{\boldsymbol{\varrho}}_T^{\text{var}} \otimes \bar{\boldsymbol{\varrho}}_T^{\text{var}}] \rangle. \quad (11)$$

Proof. By Lemma A.1, we have

$$\begin{aligned} R(\bar{\boldsymbol{\omega}}_T, \beta^{\text{te}}) &= \frac{1}{2} \langle \mathbf{H}_{m, \beta^{\text{te}}}, \mathbb{E}[\bar{\boldsymbol{\varrho}}_T \otimes \bar{\boldsymbol{\varrho}}_T] \rangle \\ &= \frac{1}{2} \langle \mathbf{H}_{m, \beta^{\text{te}}}, \mathbb{E}[(\bar{\boldsymbol{\varrho}}_T^{\text{bias}} + \bar{\boldsymbol{\varrho}}_T^{\text{var}}) \otimes (\bar{\boldsymbol{\varrho}}_T^{\text{bias}} + \bar{\boldsymbol{\varrho}}_T^{\text{var}})] \rangle \\ &\leq 2 \left(\frac{1}{2} \langle \mathbf{H}_{m, \beta^{\text{te}}}, \mathbb{E}[\bar{\boldsymbol{\varrho}}_T^{\text{bias}} \otimes \bar{\boldsymbol{\varrho}}_T^{\text{bias}}] \rangle + \frac{1}{2} \langle \mathbf{H}_{m, \beta^{\text{te}}}, \mathbb{E}[\bar{\boldsymbol{\varrho}}_T^{\text{var}} \otimes \bar{\boldsymbol{\varrho}}_T^{\text{var}}] \rangle \right) \end{aligned}$$

where the last inequality follows because for vector-valued random variables \mathbf{u} and \mathbf{v} , $\mathbb{E}\|\mathbf{u} + \mathbf{v}\|_H^2 \leq \left(\sqrt{\mathbb{E}\|\mathbf{u}\|_H^2} + \sqrt{\mathbb{E}\|\mathbf{v}\|_H^2} \right)^2$ and from Cauchy-Schwarz inequality. \square

For $t = 0, 1, \dots, T-1$, consider the following bias and variance iterates:

$$\begin{aligned} \mathbf{D}_t &= (\mathcal{I} - \alpha \mathcal{T}) \circ \mathbf{D}_{t-1} \quad \text{and} \quad \mathbf{D}_0 = (\boldsymbol{\omega}_t - \boldsymbol{\theta}^*)(\boldsymbol{\omega}_t - \boldsymbol{\theta}^*)^\top \\ \mathbf{V}_t &= (\mathcal{I} - \alpha \mathcal{T}) \circ \mathbf{V}_{t-1} + \alpha^2 \Pi \quad \text{and} \quad \mathbf{V}_0 = \mathbf{0} \end{aligned} \quad (12)$$

where $\Pi = \mathbb{E}[\mathbf{B}^\top \boldsymbol{\xi} \boldsymbol{\xi}^\top \mathbf{B}]$. One can verify that

$$\mathbf{D}_t = \mathbb{E}[\boldsymbol{\varrho}_t^{\text{bias}} \otimes \boldsymbol{\varrho}_t^{\text{bias}}], \quad \mathbf{V}_t = \mathbb{E}[\boldsymbol{\varrho}_t^{\text{var}} \otimes \boldsymbol{\varrho}_t^{\text{var}}].$$

With such notations, we can further bound the bias and variance terms.

Lemma B.3. *Following the notations in eq. (12), we have*

$$\mathcal{E}_{\text{bias}} \leq \frac{1}{\alpha T^2} \left\langle (\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{te}}})^T) \mathbf{H}_{n_1, \beta^{\text{te}}}^{-1} \mathbf{H}_{m, \beta^{\text{te}}}, \sum_{t=0}^{T-1} \mathbf{D}_t \right\rangle, \quad (13)$$

$$\mathcal{E}_{\text{var}} \leq \frac{1}{T^2} \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \left\langle (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{te}}})^{k-t} \mathbf{H}_{m, \beta^{\text{te}}}, \mathbf{V}_t \right\rangle. \quad (14)$$

Proof. Similar calculations have appeared in the prior works [14, 22]. However, our meta linear model contains additional terms, and hence we provide a proof here for completeness. We first have

$$\begin{aligned}\mathbb{E}[\bar{\boldsymbol{\varrho}}_T^{\text{var}} \otimes \bar{\boldsymbol{\varrho}}_T^{\text{var}}] &= \frac{1}{T^2} \sum_{t=0}^{T-1} \sum_{k=0}^{T-1} \mathbb{E}[\boldsymbol{\varrho}_t^{\text{var}} \otimes \boldsymbol{\varrho}_k^{\text{var}}] \\ &\succeq \frac{1}{T^2} \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \mathbb{E}[\boldsymbol{\varrho}_t^{\text{var}} \otimes \boldsymbol{\varrho}_k^{\text{var}}] + \mathbb{E}[\boldsymbol{\varrho}_k^{\text{var}} \otimes \boldsymbol{\varrho}_t^{\text{var}}]\end{aligned}$$

where the last inequality follows because we double count the diagonal terms $t = k$.

For $t \leq k$, $\mathbb{E}[\boldsymbol{\varrho}_k^{\text{var}} | \boldsymbol{\varrho}_t^{\text{var}}] = (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{vr}}})^{k-t} \boldsymbol{\varrho}_t^{\text{var}}$, since $\mathbb{E}[\mathbf{B}_t^\top \boldsymbol{\xi}_t | \boldsymbol{\varrho}_{t-1}] = \mathbf{0}$. From this, we have

$$\mathbb{E}[\bar{\boldsymbol{\varrho}}_T^{\text{var}} \otimes \bar{\boldsymbol{\varrho}}_T^{\text{var}}] \preceq \frac{1}{T^2} \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \mathbf{V}_t (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{vr}}})^{k-t} + \mathbf{V}_t (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{vr}}})^{k-t}.$$

Substituting the above inequality into $\frac{1}{2} \langle \mathbf{H}_{m, \beta^{\text{ve}}}, \mathbb{E}[\bar{\boldsymbol{\varrho}}_T^{\text{var}} \otimes \bar{\boldsymbol{\varrho}}_T^{\text{var}}] \rangle$, we obtain:

$$\begin{aligned}\mathcal{E}_{\text{var}} &= \frac{1}{2} \langle \mathbf{H}_{m, \beta^{\text{ve}}}, \mathbb{E}[\bar{\boldsymbol{\varrho}}_T^{\text{var}} \otimes \bar{\boldsymbol{\varrho}}_T^{\text{var}}] \rangle \\ &\leq \frac{1}{2T^2} \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \langle \mathbf{H}_{m, \beta^{\text{ve}}}, \mathbf{V}_t (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{vr}}})^{k-t} \rangle + \langle \mathbf{H}_{m, \beta^{\text{ve}}}, \mathbf{V}_t (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{vr}}})^{k-t} \rangle \\ &= \frac{1}{T^2} \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \langle (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{vr}}})^{k-t} \mathbf{H}_{m, \beta^{\text{ve}}}, \mathbf{V}_t \rangle\end{aligned}$$

where the last inequality follows from Assumption 2 that F and $\boldsymbol{\Sigma}$ commute, and hence $\mathbf{H}_{m, \beta^{\text{ve}}}$ and $\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{vr}}}$ commute.

For the bias term, similarly we have:

$$\mathcal{E}_{\text{bias}} \leq \frac{1}{T^2} \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \langle (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{vr}}})^{k-t} \mathbf{H}_{m, \beta^{\text{ve}}}, \mathbf{D}_t \rangle \quad (15)$$

$$= \frac{1}{\alpha T^2} \sum_{t=0}^{T-1} \langle (\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{vr}}})^{T-t}) \mathbf{H}_{n_1, \beta^{\text{vr}}}^{-1} \mathbf{H}_{m, \beta^{\text{ve}}}, \mathbf{D}_t \rangle \quad (16)$$

$$\leq \frac{1}{\alpha T^2} \langle (\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{vr}}})^T) \mathbf{H}_{n_1, \beta^{\text{vr}}}^{-1} \mathbf{H}_{m, \beta^{\text{ve}}}, \sum_{t=0}^{T-1} \mathbf{D}_t \rangle \quad (17)$$

which completes the proof. \square

B.4 Bounding the Bias

Now we start to bound the bias term. By Lemma B.3, we focus on bounding the summation of \mathbf{D}_t , i.e. $\sum_{t=0}^{T-1} \mathbf{D}_t$. Consider $\mathbf{S}_t := \sum_{k=0}^{t-1} \mathbf{D}_k$, and the following lemma shows the properties of \mathbf{S}_t

Lemma B.4. \mathbf{S}_t satisfies the recursion form:

$$\mathbf{S}_t = (\mathcal{I} - \alpha \mathcal{T}) \circ \mathbf{S}_{t-1} + \mathbf{D}_0.$$

Moreover, if $\alpha < \frac{1}{c(\beta^{\text{vr}}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})}$, then we have:

$$\mathbf{D}_0 = \mathbf{S}_0 \preceq \mathbf{S}_1 \preceq \cdots \preceq \mathbf{S}_\infty$$

where $\mathbf{S}_\infty := \sum_{k=0}^{\infty} (\mathcal{I} - \alpha \mathcal{T})^k \circ \mathbf{D}_0 = \alpha^{-1} \mathcal{T}^{-1} \circ \mathbf{D}_0$.

Proof. By eq. (12), we have

$$\begin{aligned}
\mathbf{S}_t &= \sum_{k=0}^{t-1} \mathbf{D}_k = \sum_{k=0}^{t-1} (\mathcal{I} - \alpha\mathcal{T})^k \circ \mathbf{D}_0 \\
&= \mathbf{D}_0 + (\mathcal{I} - \alpha\mathcal{T}) \circ \left(\sum_{k=0}^{t-2} (\mathcal{I} - \alpha\mathcal{T})^k \circ \mathbf{D}_0 \right) \\
&= \mathbf{D}_0 + (\mathcal{I} - \alpha\mathcal{T}) \circ \mathbf{S}_{t-1}.
\end{aligned}$$

By Lemma B.1, $(\mathcal{I} - \alpha\mathcal{T})$ is PSD mapping, and hence $\mathbf{D}_t = (\mathcal{I} - \alpha\mathcal{T}) \circ \mathbf{D}_{t-1}$ is a PSD matrix for every t , which implies $\mathbf{S}_{t-1} \preceq \mathbf{S}_{t-1} + \mathbf{D}_t = \mathbf{S}_t$. The form of \mathbf{S}_∞ can be directly obtained by Lemma B.1. \square

Then we can decompose \mathbf{S}_t as follows:

$$\begin{aligned}
\mathbf{S}_t &= \mathbf{D}_0 + (\mathcal{I} - \alpha\tilde{\mathcal{T}}) \circ \mathbf{S}_{t-1} + \alpha(\tilde{\mathcal{T}} - \mathcal{T}) \circ \mathbf{S}_{t-1} \\
&= \mathbf{D}_0 + (\mathcal{I} - \alpha\tilde{\mathcal{T}}) \circ \mathbf{S}_{t-1} + \alpha^2(\mathcal{M} - \tilde{\mathcal{M}}) \circ \mathbf{S}_{t-1} \\
&\preceq \mathbf{D}_0 + (\mathcal{I} - \alpha\tilde{\mathcal{T}}) \circ \mathbf{S}_{t-1} + \alpha^2\mathcal{M} \circ \mathbf{S}_T \\
&= \sum_{k=0}^{t-1} (\mathcal{I} - \alpha\tilde{\mathcal{T}})^k \circ (\mathbf{D}_0 + \alpha^2\mathcal{M} \circ \mathbf{S}_T)
\end{aligned} \tag{18}$$

where the inequality follows because $\mathbf{S}_t \preceq \mathbf{S}_T$ for any $t \leq T$. Therefore, it is crucial to understand $\mathcal{M} \circ \mathbf{S}_T$.

Lemma B.5. For any symmetric matrix \mathbf{A} , if $\alpha < \frac{1}{c(\beta^{tr}, \Sigma) \text{tr}(\Sigma)}$, it holds that

$$\mathcal{M} \circ \mathcal{T}^{-1} \circ \mathbf{A} \preceq \frac{c(\beta^{tr}, \Sigma) \text{tr}(\Sigma \mathbf{H}_{n_1, \beta^{tr}}^{-1} \mathbf{A})}{1 - \alpha c(\beta^{tr}, \Sigma) \text{tr}(\Sigma)} \cdot \mathbf{H}_{n_1, \beta^{tr}}.$$

Proof. Denote $\mathbf{C} = \mathcal{T}^{-1} \circ \mathbf{A}$. Recalling $\tilde{\mathcal{T}} = \mathcal{T} + \alpha\mathcal{M} - \alpha\tilde{\mathcal{M}}$, we have

$$\begin{aligned}
\tilde{\mathcal{T}} \circ \mathbf{C} &= \mathcal{T} \circ \mathbf{C} + \alpha\mathcal{M} \circ \mathbf{C} - \alpha\tilde{\mathcal{M}} \circ \mathbf{C} \\
&\preceq \mathbf{A} + \alpha\mathcal{M} \circ \mathbf{C}.
\end{aligned}$$

Recalling that $\tilde{\mathcal{T}}^{-1}$ exists and is a PSD mapping, we then have

$$\begin{aligned}
\mathcal{M} \circ \mathbf{C} &\preceq \alpha\mathcal{M} \circ \tilde{\mathcal{T}}^{-1} \circ \mathcal{M} \circ \mathbf{C} + \mathcal{M} \circ \tilde{\mathcal{T}}^{-1} \circ \mathbf{A} \\
&\preceq \sum_{k=0}^{\infty} (\alpha\mathcal{M} \circ \tilde{\mathcal{T}}^{-1})^k \circ (\mathcal{M} \circ \tilde{\mathcal{T}}^{-1} \circ \mathbf{A}).
\end{aligned} \tag{19}$$

By Proposition B.1, we have $\mathcal{M} \circ \tilde{\mathcal{T}}^{-1} \circ \mathbf{A} \preceq \underbrace{c(\beta^{tr}, \Sigma) \text{tr}(\Sigma \tilde{\mathcal{T}}^{-1} \circ \mathbf{A})}_{J_2} \mathbf{H}_{n_1, \beta^{tr}}$. Substituting back

into eq. (19), we obtain:

$$\begin{aligned}
\sum_{k=0}^{\infty} (\alpha\mathcal{M} \circ \tilde{\mathcal{T}}^{-1})^k \circ (\mathcal{M} \circ \tilde{\mathcal{T}}^{-1} \circ \mathbf{A}) &\preceq \sum_{k=0}^{\infty} (\alpha\mathcal{M} \circ \tilde{\mathcal{T}}^{-1})^k \circ (J_2 \mathbf{H}_{n_1, \beta^{tr}}) \\
&\preceq J_2 \sum_{k=0}^{\infty} (\alpha c(\beta^{tr}, \Sigma) \text{tr}(\Sigma))^k \mathbf{H}_{n_1, \beta^{tr}} \preceq \frac{J_2}{1 - \alpha c(\beta^{tr}, \Sigma) \text{tr}(\Sigma)} \mathbf{H}_{n_1, \beta^{tr}}
\end{aligned}$$

where the second inequality follows since $\tilde{\mathcal{T}}^{-1} \circ \mathbf{H}_{n_1, \beta^{tr}} \preceq \mathbf{I}$ (Lemma B.1) and $\mathcal{M} \circ \mathbf{I} \preceq c(\beta^{tr}, \Sigma) \text{tr}(\Sigma) \mathbf{H}_{n_1, \beta^{tr}}$ (Proposition B.1).

Finally, we bound J_2 as follows:

$$\begin{aligned}
\text{tr} \left(\boldsymbol{\Sigma} \tilde{\mathcal{T}}^{-1} \circ \mathbf{A} \right) &= \alpha \text{tr} \left(\sum_{k=0}^{\infty} \boldsymbol{\Sigma} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{tr}})^k \mathbf{A} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{tr}})^k \right) \\
&= \alpha \text{tr} \left(\sum_{k=0}^{\infty} \boldsymbol{\Sigma} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{tr}})^{2k} \mathbf{A} \right) \\
&= \text{tr} \left(\boldsymbol{\Sigma} (2\mathbf{H}_{n_1, \beta^{tr}} - \alpha \mathbf{H}_{n_1, \beta^{tr}}^2)^{-1} \mathbf{A} \right) \\
&\leq \text{tr} \left(\boldsymbol{\Sigma} \mathbf{H}_{n_1, \beta^{tr}}^{-1} \mathbf{A} \right)
\end{aligned}$$

where the second equality follows because $\boldsymbol{\Sigma}$ and $\mathbf{H}_{n_1, \beta^{tr}}$ commute, and the last inequality holds since $\alpha < \frac{1}{\max_i \{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})\}}$. Putting all these results together completes the proof. \square

Lemma B.6 (Bounding $\mathcal{M} \circ \mathbf{S}_T$).

$$\mathcal{M} \circ \mathbf{S}_T \preceq \frac{c(\beta^{tr}, \boldsymbol{\Sigma}) \cdot \text{tr} \left(\boldsymbol{\Sigma} \mathbf{H}_{n_1, \beta^{tr}}^{-1} \left[\mathcal{I} - (\mathcal{I} - \alpha \tilde{\mathcal{T}})^T \right] \circ \mathbf{D}_0 \right)}{\alpha(1 - c(\beta^{tr}, \boldsymbol{\Sigma})) \alpha \text{tr}(\boldsymbol{\Sigma})} \cdot \mathbf{H}_{n_1, \beta^{tr}}.$$

Proof. \mathbf{S}_T can be further derived as follows:

$$\mathbf{S}_T = \sum_{k=0}^{T-1} (\mathcal{I} - \alpha \mathcal{T})^k \circ \mathbf{D}_0 = \alpha^{-1} \mathcal{T}^{-1} \circ \left[\mathcal{I} - (\mathcal{I} - \alpha \mathcal{T})^T \right] \circ \mathbf{D}_0.$$

Since $\tilde{\mathcal{T}} - \mathcal{T}$ is a PSD mapping by Lemma B.1, we have $\mathcal{I} - \alpha \tilde{\mathcal{T}} \preceq \mathcal{I} - \alpha \mathcal{T}$. Hence $\mathcal{I} - (\mathcal{I} - \alpha \mathcal{T})^T \preceq \mathcal{I} - (\mathcal{I} - \alpha \tilde{\mathcal{T}})^T$. Combining with the fact that \mathcal{T}^{-1} is also a PSD mapping, we have:

$$\mathbf{S}_T \preceq \alpha^{-1} \mathcal{T}^{-1} \circ \left[\mathcal{I} - (\mathcal{I} - \alpha \tilde{\mathcal{T}})^T \right] \circ \mathbf{D}_0.$$

Letting $\mathbf{A} = \left[\mathcal{I} - (\mathcal{I} - \alpha \tilde{\mathcal{T}})^T \right] \circ \mathbf{D}_0$ in Lemma B.5, we obtain:

$$\begin{aligned}
\mathcal{M} \circ \mathbf{S}_T &\preceq \alpha^{-1} \mathcal{M} \circ \mathcal{T}^{-1} \circ \left[\mathcal{I} - (\mathcal{I} - \alpha \tilde{\mathcal{T}})^T \right] \circ \mathbf{D}_0 \\
&\preceq \frac{c(\beta^{tr}, \boldsymbol{\Sigma}) \cdot \text{tr} \left(\boldsymbol{\Sigma} \mathbf{H}_{n_1, \beta^{tr}}^{-1} \left[\mathcal{I} - (\mathcal{I} - \alpha \tilde{\mathcal{T}})^T \right] \circ \mathbf{D}_0 \right)}{\alpha(1 - c(\beta^{tr}, \boldsymbol{\Sigma})) \alpha \text{tr}(\boldsymbol{\Sigma})} \cdot \mathbf{H}_{n_1, \beta^{tr}}.
\end{aligned}$$

\square

Now we are ready to derive the upper bound on the bias term.

Lemma B.7 (Bounding the bias). *If $\alpha < \frac{1}{c(\beta^{tr}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})}$, for sufficiently large n_1 , s.t. $\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) > 0$, $\forall i$, then we have*

$$\begin{aligned}
\mathcal{E}_{\text{bias}} &\leq \sum_i \left(\frac{1}{\alpha^2 T^2} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) \geq \frac{1}{\alpha T}} + \mu_i^2(\mathbf{H}_{n_1, \beta^{tr}}) \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) < \frac{1}{\alpha T}} \right) \frac{\omega_i^2 \mu_i(\mathbf{H}_{m, \beta^{tr}})}{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})^2} \\
&\quad + \frac{2c(\beta^{tr}, \boldsymbol{\Sigma})}{T \alpha (1 - c(\beta^{tr}, \boldsymbol{\Sigma})) \alpha \text{tr}(\boldsymbol{\Sigma})} \sum_i \left(\frac{1}{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) \geq \frac{1}{\alpha T}} + T \alpha \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) < \frac{1}{\alpha T}} \right) \cdot \lambda_i \omega_i^2 \\
&\quad \times \sum_i \left(\frac{1}{T} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) \geq \frac{1}{\alpha T}} + T \alpha^2 \mu_i(\mathbf{H}_{n_1, \beta^{tr}})^2 \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) < \frac{1}{\alpha T}} \right) \cdot \frac{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})}{\mu_i(\mathbf{H}_{m, \beta^{tr}})}.
\end{aligned}$$

Proof. Applying Lemma B.6 to eq. (18), we can obtain:

$$\begin{aligned}
\mathbf{S}_t &\preceq \sum_{k=0}^{t-1} (\mathcal{I} - \alpha \tilde{\mathcal{T}})^k \circ \left(\frac{\alpha c(\beta^{\text{tr}}, \Sigma) \cdot \text{tr} \left(\Sigma \mathbf{H}_{n_1, \beta^{\text{tr}}}^{-1} \left[\mathcal{I} - (\mathcal{I} - \alpha \tilde{\mathcal{T}})^T \right] \circ \mathbf{D}_0 \right)}{1 - c(\beta, \Sigma) \alpha \text{tr}(\Sigma)} \cdot \mathbf{H}_{n_1, \beta^{\text{tr}}} + \mathbf{D}_0 \right) \\
&= \sum_{k=0}^{t-1} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^k \cdot \left(\underbrace{\frac{\alpha c(\beta^{\text{tr}}, \Sigma) \cdot \text{tr} \left(\Sigma \mathbf{H}_{n_1, \beta^{\text{tr}}}^{-1} (\mathbf{D}_0 - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^T \mathbf{D}_0 (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^T) \right)}{1 - c(\beta^{\text{tr}}, \Sigma) \alpha \text{tr}(\Sigma)}}_{\mathbf{G}_1} \cdot \mathbf{H}_{n_1, \beta^{\text{tr}}} + \underbrace{\mathbf{D}_0}_{\mathbf{G}_2} \right) \\
&\cdot (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^k.
\end{aligned}$$

Letting $t = T$, and substituting the upper bound of \mathbf{S}_T into the bias term in Lemma B.3, we obtain:

$$\begin{aligned}
\mathcal{E}_{\text{bias}} &\leq \frac{1}{\alpha T^2} \sum_{k=0}^{T-1} \left\langle \left((\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^{2k} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^{T+2k} \right) \mathbf{H}_{n_1, \beta^{\text{tr}}}^{-1} \mathbf{H}_{m, \beta^{\text{tr}}} \mathbf{G}_1 + \mathbf{G}_2 \right\rangle \\
&\leq \frac{1}{\alpha T^2} \sum_{k=0}^{T-1} \left\langle \left((\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^k - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^{T+k} \right) \mathbf{H}_{n_1, \beta^{\text{tr}}}^{-1} \mathbf{H}_{m, \beta^{\text{tr}}} \mathbf{G}_1 + \mathbf{G}_2 \right\rangle.
\end{aligned}$$

We first consider

$$d_1 = \frac{1}{\alpha T^2} \sum_{k=0}^{T-1} \left\langle \left((\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^k - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^{T+k} \right) \mathbf{H}_{n_1, \beta^{\text{tr}}}^{-1} \mathbf{H}_{m, \beta^{\text{tr}}} \mathbf{G}_1 \right\rangle.$$

Since $\mathbf{H}_{n_1, \beta^{\text{tr}}}$, $\mathbf{H}_{m, \beta^{\text{tr}}}$ and $\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}}$ commute, we have

$$\begin{aligned}
d_1 &= \frac{c(\beta^{\text{tr}}, \Sigma) \cdot \text{tr} \left(\Sigma \mathbf{H}_{n_1, \beta^{\text{tr}}}^{-1} (\mathbf{D}_0 - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^T \mathbf{D}_0 (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^T) \right)}{(1 - c(\beta^{\text{tr}}, \Sigma) \alpha \text{tr}(\Sigma)) T^2} \\
&\quad \times \sum_{k=0}^{T-1} \left\langle \left((\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^k - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^{T+k} \right), \mathbf{H}_{m, \beta^{\text{tr}}} \right\rangle.
\end{aligned}$$

For the first term, since Σ , $\mathbf{H}_{n_1, \beta^{\text{tr}}}$ and $\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}}$ can be diagonalized simultaneously, considering the eigen-decompositions under the basis of Σ and recalling $\Sigma = \mathbf{V} \Lambda \mathbf{V}^T$, we have:

$$\begin{aligned}
&\text{tr} \left(\Sigma \mathbf{H}_{n_1, \beta^{\text{tr}}}^{-1} [\mathbf{D}_0 - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^T \mathbf{D}_0 (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^T] \right) \\
&= \sum_i \left(1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}))^{2T} \right) \cdot \langle \mathbf{w}_0 - \mathbf{w}^*, \mathbf{v}_i \rangle^2 \frac{\lambda_i}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \\
&\leq 2 \sum_i \left(\mathbf{1}_{\lambda_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}} + T \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \cdot \langle \mathbf{w}_0 - \mathbf{w}^*, \mathbf{v}_i \rangle^2 \frac{\lambda_i}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})}
\end{aligned}$$

where the last inequality holds since $1 - (1 - \alpha x)^{2T} \leq \min\{2, 2T\alpha x\}$.

For the second term, similarly, $\mathbf{H}_{m, \beta^{\text{tr}}}$ and $\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}}$ can be diagonalized simultaneously. We then have

$$\begin{aligned}
&\sum_{k=0}^{T-1} \left\langle \left((\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^k - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^{T+k} \right), \mathbf{H}_{m, \beta^{\text{tr}}} \right\rangle \\
&\leq \sum_{k=0}^{T-1} \sum_i [(1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}))^k - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}))^{T+k}] \mu_i(\mathbf{H}_{m, \beta^{\text{tr}}}) \\
&= \frac{1}{\alpha} \sum_i [1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}))^T]^2 \frac{\mu_i(\mathbf{H}_{m, \beta^{\text{tr}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \\
&\leq \frac{1}{\alpha} \sum_i \left(\mathbf{1}_{\lambda_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}} + T^2 \alpha^2 \lambda_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \mathbf{1}_{\lambda_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \frac{\mu_i(\mathbf{H}_{m, \beta^{\text{tr}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})}.
\end{aligned}$$

Now we turn to:

$$d_2 = \frac{1}{\alpha T^2} \sum_{k=0}^{T-1} \left\langle ((\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^k - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^{T+k}) \mathbf{H}_{n_1, \beta^{\text{tr}}}^{-1} \mathbf{H}_{m, \beta^{\text{ve}}}, \mathbf{G}_2 \right\rangle.$$

Considering the orthogonal decompositions of $\mathbf{H}_{m, \beta^{\text{ve}}}$ and $\mathbf{H}_{n_1, \beta^{\text{tr}}}$ under \mathbf{V} , $\mathbf{H}_{n_1, \beta^{\text{tr}}} = \mathbf{V} \mathbf{\Lambda}_1 \mathbf{V}^\top$, $\mathbf{H}_{m, \beta^{\text{ve}}} = \mathbf{V} \mathbf{\Lambda}_2 \mathbf{V}^\top$, where the diagonal entries of $\mathbf{\Lambda}_1$ are $\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})$ (and $\mu_i(\mathbf{H}_{m, \beta^{\text{ve}}})$ for $\mathbf{\Lambda}_2$). Then we have:

$$\begin{aligned} d_2 &= \frac{1}{\alpha T^2} \sum_{k=0}^{T-1} \left\langle \underbrace{((\mathbf{I} - \alpha \mathbf{\Lambda}_1)^k - (\mathbf{I} - \alpha \mathbf{\Lambda}_1)^{T+k}) \mathbf{\Lambda}_1^{-1} \mathbf{\Lambda}_2}_{\mathbf{J}_3}, \mathbf{V}^\top \mathbf{D}_0 \mathbf{V} \right\rangle \\ &= \frac{1}{\alpha T^2} \sum_{k=0}^{T-1} \sum_i \left[(1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}))^k - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}))^{T+k} \right] \frac{\omega_i^2 \mu_i(\mathbf{H}_{m, \beta^{\text{ve}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \\ &= \frac{1}{\alpha^2 T^2} \sum_i \left[1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}))^T \right]^2 \frac{\omega_i^2 \mu_i(\mathbf{H}_{m, \beta^{\text{ve}}})}{\mu_i^2(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \\ &\leq \frac{1}{\alpha^2 T^2} \sum_i \left(\mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}} + \alpha^2 T^2 \mu_i^2(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \frac{\omega_i^2 \mu_i(\mathbf{H}_{m, \beta^{\text{ve}}})}{\mu_i^2(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \end{aligned}$$

where $\omega_i = \langle \mathbf{w}_0 - \boldsymbol{\theta}^*, \mathbf{v}_i \rangle$ is the diagonal entry of $\mathbf{V}^\top \mathbf{D}_0 \mathbf{V}$ and the second equality holds since \mathbf{J}_3 is a diagonal matrix. \square

B.5 Bounding the Variance

Note that the noisy part $\Pi = \mathbb{E}[\mathbf{B}^\top \boldsymbol{\xi} \boldsymbol{\xi}^\top \mathbf{B}]$ in eq. (12) is important in the variance iterates. In order to analyze the variance term, we first understand the role of Π by the following lemma.

Lemma B.8 (Bounding the noise).

$$\Pi = \mathbb{E}[\mathbf{B}^\top \boldsymbol{\xi} \boldsymbol{\xi}^\top \mathbf{B}] \preceq f(\beta^{\text{tr}}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta) \mathbf{H}_{n_1, \beta^{\text{tr}}}$$

where $f(\beta, n, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta) = [c(\beta, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}_\theta \boldsymbol{\Sigma}) + 4c_1 \sigma^2 \sigma_x^2 \beta^2 \sqrt{C(\beta, \boldsymbol{\Sigma})} \text{tr}(\boldsymbol{\Sigma}^2) + \sigma^2/n]$.

Proof. With a slight abuse of notations, we write β^{tr} as β in this proof. By definition of meta data and noise, we have

$$\begin{aligned} \Pi &= \mathbb{E}[\mathbf{B}^\top \boldsymbol{\xi} \boldsymbol{\xi}^\top \mathbf{B}] \\ &= \frac{\sigma^2}{n_2} \mathbf{H}_{n_1, \beta} + \mathbb{E}[\mathbf{B}^\top \mathbf{B} \boldsymbol{\Sigma}_\theta \mathbf{B}^\top \mathbf{B}] + \sigma^2 \cdot \frac{\beta^2}{n_2 n_1^2} \mathbb{E}[\mathbf{B}^\top \mathbf{X}^{\text{out}} \mathbf{X}^{\text{in}^\top} \mathbf{X}^{\text{in}} \mathbf{X}^{\text{out}^\top} \mathbf{B}]. \end{aligned}$$

The second term can be directly bounded by Proposition B.1:

$$\mathbb{E}[\mathbf{B}^\top \mathbf{B} \boldsymbol{\Sigma}_\theta \mathbf{B}^\top \mathbf{B}] \preceq c(\beta, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}_\theta \boldsymbol{\Sigma}) \mathbf{H}_{n_1, \beta}.$$

For the third term, we utilize the technique similar to Proposition B.1, and by Assumption 1, we have:

$$\begin{aligned} &\sigma^2 \cdot \frac{\beta^2}{n_2 n_1^2} \mathbb{E} \left[\mathbf{B}^\top \mathbf{X}^{\text{out}} \mathbf{X}^{\text{in}^\top} \mathbf{X}^{\text{in}} \mathbf{X}^{\text{out}^\top} \mathbf{B} \right] \\ &\preceq \sigma^2 c_1 \cdot \frac{\beta^2}{n_1^2} \mathbb{E} \left[\text{tr}(\mathbf{X}^{\text{in}^\top} \mathbf{X}^{\text{in}} \boldsymbol{\Sigma}) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^{\text{in}^\top} \mathbf{X}^{\text{in}} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^{\text{in}^\top} \mathbf{X}^{\text{in}} \right) \right]. \end{aligned}$$

Following the analysis for \mathbf{J}_1 in the proof of Proposition B.1, and letting $\mathbf{A} = \mathbf{I}$, we obtain:

$$\frac{1}{n_1^2} \mathbb{E} \left[\text{tr}(\mathbf{X}^{\text{in}^\top} \mathbf{X}^{\text{in}} \boldsymbol{\Sigma}) \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^{\text{in}^\top} \mathbf{X}^{\text{in}} \right) \boldsymbol{\Sigma} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^{\text{in}^\top} \mathbf{X}^{\text{in}} \right) \right] \preceq 4\sqrt{C(\beta, \boldsymbol{\Sigma})} \sigma_x^2 \text{tr}(\boldsymbol{\Sigma}^2) \mathbf{H}_{n_1, \beta}.$$

Putting all these results together completes the proof. \square

Lemma B.9 (Property of \mathbf{V}_t). *If the stepsize satisfies $\alpha < \frac{1}{c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})}$, it holds that*

$$\mathbf{0} = \mathbf{V}_0 \preceq \mathbf{V}_1 \preceq \dots \preceq \mathbf{V}_\infty \preceq \frac{\alpha f(\beta^{\text{tr}}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta)}{1 - \alpha c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})} \mathbf{I}.$$

Proof. Similar calculations has appeared in prior works [14, 22]. However, our analysis of the meta linear model needs to handle the complicated meta noise, and hence we provide a proof here for completeness.

We first show that $\mathbf{V}_{t-1} \preceq \mathbf{V}_t$. By recursion:

$$\begin{aligned} \mathbf{V}_t &= (\mathcal{I} - \alpha\mathcal{T}) \circ \mathbf{V}_{t-1} + \alpha^2\Pi \\ &\stackrel{(a)}{=} \alpha^2 \sum_{k=0}^{t-1} (\mathcal{I} - \alpha\mathcal{T})^k \circ \Pi \\ &= \mathbf{V}_{t-1} + \alpha^2(\mathcal{I} - \alpha\mathcal{T})^{t-1} \circ \Pi \\ &\stackrel{(b)}{\succeq} \mathbf{V}_{t-1} \end{aligned}$$

where (a) holds by solving the recursion and (b) follows because $\mathcal{I} - \alpha\mathcal{T}$ is a PSD mapping.

The existence of \mathbf{V}_∞ can be shown in the way similar to the proof of Lemma B.1. We first have

$$\mathbf{V}_t = \alpha^2 \sum_{k=0}^{t-1} (\mathcal{I} - \alpha\mathcal{T})^k \circ \Pi \preceq \alpha^2 \sum_{k=0}^{\infty} \underbrace{(\mathcal{I} - \alpha\mathcal{T})^k \circ \Pi}_{\mathbf{A}_k}.$$

By previous analysis in Lemma B.1, if $\alpha < \frac{1}{c(\beta^{\text{tr}}, \Sigma) \text{tr}(\Sigma)}$, we have

$$\text{tr}(\mathbf{A}_k) \leq \left(1 - \alpha \min_i \{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})\}\right) \text{tr}(\mathbf{A}_{t-1}).$$

Therefore,

$$\text{tr}(\mathbf{V}_t) \leq \alpha^2 \sum_{k=0}^{\infty} \text{tr}(\mathbf{A}_k) \leq \frac{\alpha \text{tr}(\Pi)}{\min_i \{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})\}} < \infty.$$

The trace of \mathbf{V}_t is uniformly bounded from above, which indicates that \mathbf{V}_∞ exists.

Finally, we bound \mathbf{V}_∞ . Note that \mathbf{V}_∞ is the solution to:

$$\mathbf{V}_\infty = (\mathcal{I} - \alpha\mathcal{T}) \circ \mathbf{V}_\infty + \alpha^2\Pi.$$

Then we can write \mathbf{V}_∞ as $\mathbf{V}_\infty = \mathcal{T}^{-1} \circ \alpha\Pi$. Following the analysis in the proof of Lemma B.5, we have:

$$\begin{aligned} \tilde{\mathcal{T}} \circ \mathbf{V}_\infty &= \tilde{\mathcal{T}} \circ \mathcal{T}^{-1} \circ \alpha\Pi \\ &\preceq \alpha\Pi + \alpha\mathcal{M} \circ \mathbf{V}_\infty \\ &\preceq \alpha f(\beta^{\text{tr}}, n_2, \sigma, \Sigma, \Sigma_\theta) \mathbf{H}_{n_1, \beta^{\text{tr}}} + \alpha\mathcal{M} \circ \mathbf{V}_\infty \end{aligned}$$

where the last inequality follows from Lemma B.8. Applying $\tilde{\mathcal{T}}^{-1}$, which exists and is a PSD mapping, to the both sides, we have

$$\begin{aligned} \mathbf{V}_\infty &\preceq \alpha f(\beta^{\text{tr}}, n_2, \sigma, \Sigma, \Sigma_\theta) \cdot \tilde{\mathcal{T}}^{-1} \circ \mathbf{H}_{n_1, \beta^{\text{tr}}} + \alpha \tilde{\mathcal{T}}^{-1} \circ \mathcal{M} \circ \mathbf{V}_\infty \\ &\stackrel{(a)}{\preceq} \alpha f(\beta^{\text{tr}}, n_2, \sigma, \Sigma, \Sigma_\theta) \cdot \sum_{t=0}^{\infty} \left(\alpha \tilde{\mathcal{T}}^{-1} \circ \mathcal{M} \right)^t \circ \tilde{\mathcal{T}}^{-1} \circ \mathbf{H}_{n_1, \beta^{\text{tr}}} \\ &\stackrel{(b)}{\preceq} \alpha f(\beta^{\text{tr}}, n_2, \sigma, \Sigma, \Sigma_\theta) \sum_{t=0}^{\infty} (\alpha c(\beta^{\text{tr}}, \Sigma) \text{tr}(\Sigma))^t \mathbf{I} \\ &= \frac{\alpha f(\beta^{\text{tr}}, n_2, \sigma, \Sigma, \Sigma_\theta)}{1 - \alpha c(\beta^{\text{tr}}, \Sigma) \text{tr}(\Sigma)} \mathbf{I} \end{aligned}$$

where (a) holds by directly solving the recursion; (b) follows from the fact that $\tilde{\mathcal{T}}^{-1} \circ \mathbf{H}_{n_1, \beta^{\text{tr}}} \preceq \mathbf{I}$ from Lemma B.1 and $\mathcal{M} \circ \mathbf{I} \preceq c(\beta^{\text{tr}}, \Sigma) \text{tr}(\Sigma) \mathbf{H}_{n_1, \beta^{\text{tr}}}$ by letting $\mathbf{A} = \mathbf{I}$ in Proposition B.1. \square

Now we are ready to provide the upper bound on the variance term.

Lemma B.10 (Bounding the Variance). *If $\alpha < \frac{1}{c(\beta^r, \Sigma) \text{tr}(\Sigma)}$, for sufficiently large n_1 , s.t. $\mu_i(\mathbf{H}_{n_1, \beta^r}) > 0, \forall i$, then we have*

$$\begin{aligned} \mathcal{E}_{\text{var}} &\leq \frac{f(\beta^r, n_2, \sigma, \Sigma, \Sigma_\theta)}{(1 - \alpha c(\beta^r, \Sigma) \text{tr}(\Sigma))} \\ &\quad \times \sum_i \left(\frac{1}{T} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^r}) \geq \frac{1}{\alpha T}} + T \alpha^2 \mu_i^2(\mathbf{H}_{n_1, \beta^r}) \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^r}) < \frac{1}{\alpha T}} \right) \frac{\mu_i(\mathbf{H}_{m, \beta^e})}{\mu_i(\mathbf{H}_{n_1, \beta^r})}. \end{aligned}$$

Proof. Recall

$$\begin{aligned} \mathbf{V}_t &= (\mathcal{I} - \alpha \mathcal{T}) \circ \mathbf{V}_{t-1} + \alpha^2 \Pi \\ &= (\mathcal{I} - \alpha \tilde{\mathcal{T}}) \circ \mathbf{V}_{t-1} + \alpha^2 (\mathcal{M} - \tilde{\mathcal{M}}) \circ \mathbf{V}_{t-1} + \alpha^2 \Pi \\ &\preceq (\mathcal{I} - \alpha \tilde{\mathcal{T}}) \circ \mathbf{V}_{t-1} + \alpha^2 \mathcal{M} \circ \mathbf{V}_{t-1} + \alpha^2 \Pi. \end{aligned} \tag{20}$$

By the uniform bound on \mathbf{V}_t and \mathcal{M} is a PSD mapping, we have:

$$\begin{aligned} \mathcal{M} \circ \mathbf{V}_t &\preceq \mathcal{M} \circ \mathbf{V}_\infty \\ &\stackrel{(a)}{\preceq} \mathcal{M} \circ \frac{\alpha f(\beta^r, n_2, \sigma, \Sigma, \Sigma_\theta)}{1 - \alpha c(\beta^r, \Sigma) \text{tr}(\Sigma)} \mathbf{I} \\ &\stackrel{(b)}{\preceq} \frac{\alpha f(\beta^r, n_2, \sigma, \Sigma, \Sigma_\theta) c(\beta^r, \Sigma) \text{tr}(\Sigma)}{1 - \alpha c(\beta^r, \Sigma) \text{tr}(\Sigma)} \cdot \mathbf{H}_{n_1, \beta^r} \end{aligned}$$

where (a) directly follows from Lemma B.9; (b) holds because $\mathcal{M} \circ \mathbf{I} \preceq c(\beta^r, \Sigma) \text{tr}(\Sigma) \mathbf{H}_{n_1, \beta^r}$ (letting $\mathbf{A} = \mathbf{I}$ in Proposition B.1). Substituting it back into eq. (20), we have:

$$\begin{aligned} \mathbf{V}_t &\preceq (\mathcal{I} - \alpha \tilde{\mathcal{T}}) \circ \mathbf{V}_{t-1} + \alpha^2 \frac{\alpha f c(\beta^r, \Sigma) \text{tr}(\Sigma)}{1 - \alpha c(\beta^r, \Sigma) \text{tr}(\Sigma)} \cdot \mathbf{H}_{n_1, \beta^r} + \alpha^2 f \mathbf{H}_{n_1, \beta^r} \\ &= (\mathcal{I} - \alpha \tilde{\mathcal{T}}) \circ \mathbf{V}_{t-1} + \frac{\alpha^2 f(\beta^r, n_2, \sigma, \Sigma, \Sigma_\theta)}{1 - \alpha c(\beta^r, \Sigma) \text{tr}(\Sigma)} \mathbf{H}_{n_1, \beta^r} \\ &\stackrel{(a)}{=} \frac{\alpha^2 f(\beta^r, n_2, \sigma, \Sigma, \Sigma_\theta)}{1 - \alpha c(\beta^r, \Sigma) \text{tr}(\Sigma)} \sum_{k=0}^{t-1} (\mathbf{I} - \alpha \tilde{\mathcal{T}})^k \circ \mathbf{H}_{n_1, \beta^r} \\ &\stackrel{(b)}{\preceq} \frac{\alpha f(\beta^r, n_2, \sigma, \Sigma, \Sigma_\theta)}{1 - \alpha c(\beta^r, \Sigma) \text{tr}(\Sigma)} (\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n, \beta})^t) \end{aligned}$$

where (a) holds by solving the recursion and (b) is due to the fact that

$$\begin{aligned} \sum_{k=0}^{t-1} (\mathbf{I} - \alpha \tilde{\mathcal{T}})^k \circ \mathbf{H}_{n_1, \beta^r} &= \sum_{k=0}^{t-1} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^k \mathbf{H}_{n_1, \beta^r} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^k \\ &\preceq \sum_{k=0}^{t-1} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^k \mathbf{H}_{n_1, \beta^r} \\ &= \frac{1}{\alpha} [\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^t]. \end{aligned}$$

Substituting the bound for \mathbf{V}_t back into the variance term in Lemma B.3, we have

$$\begin{aligned} \mathcal{E}_{\text{var}} &\leq \frac{1}{T^2} \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \langle (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k-t} \mathbf{H}_{m, \beta^e}, \mathbf{V}_t \rangle \\ &= \frac{1}{\alpha T^2} \sum_{t=0}^{T-1} \langle (\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{T-t}) \mathbf{H}_{n_1, \beta^r}^{-1} \mathbf{H}_{m, \beta^e}, \mathbf{V}_t \rangle \\ &\leq \frac{f(\beta^r, n_2, \sigma, \Sigma, \Sigma_\theta)}{(1 - \alpha c(\beta^r, \Sigma) \text{tr}(\Sigma)) T^2} \sum_{t=0}^{T-1} \langle \mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n, \beta})^{T-t}, (\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n, \beta})^t) \mathbf{H}_{n, \beta}^{-1} \mathbf{H}_{m, \eta} \rangle. \end{aligned}$$

Simultaneously diagonalizing $\mathbf{H}_{n_1, \beta^{tr}}$ and $\mathbf{H}_{m, \beta^{te}}$ as the analysis in Lemma B.7, we have

$$\begin{aligned}
\mathcal{E}_{\text{var}} &\leq \frac{f(\beta^{tr}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})}{(1 - \alpha c(\beta^{tr}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})) T^2} \\
&\quad \cdot \sum_i \sum_{t=0}^{T-1} \left(1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{tr}}))^{T-t}\right) \left(1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{tr}}))^t\right) \frac{\mu_i(\mathbf{H}_{m, \beta^{te}})}{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})} \\
&\leq \frac{f(\beta^{tr}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})}{(1 - \alpha c(\beta^{tr}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})) T^2} \\
&\quad \cdot \sum_i \sum_{t=0}^{T-1} \left(1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{tr}}))^T\right) \left(1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{tr}}))^T\right) \frac{\mu_i(\mathbf{H}_{m, \beta^{te}})}{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})} \\
&= \frac{f(\beta^{tr}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})}{(1 - \alpha c(\beta^{tr}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})) T} \sum_i \left(1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^{tr}}))^T\right)^2 \frac{\mu_i(\mathbf{H}_{m, \beta^{te}})}{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})} \\
&\leq \frac{f(\beta^{tr}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})}{(1 - \alpha c(\beta^{tr}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})) T} \sum_i (\min\{1, \alpha T \mu_i(\mathbf{H}_{n_1, \beta^{tr}})\})^2 \frac{\mu_i(\mathbf{H}_{m, \beta^{te}})}{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})} \\
&\leq \frac{f(\beta^{tr}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_{\boldsymbol{\theta}})}{(1 - \alpha c(\beta^{tr}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}))} \\
&\quad \cdot \sum_i \left(\frac{1}{T} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) \geq \frac{1}{\alpha T}} + T \alpha^2 \mu_i^2(\mathbf{H}_{n_1, \beta^{tr}}) \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) < \frac{1}{\alpha T}} \right) \frac{\mu_i(\mathbf{H}_{m, \beta^{te}})}{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})},
\end{aligned}$$

which completes the proof. \square

B.6 Proof of Theorem 1

Theorem B.3 (Theorem 1 Restated). *Let $\omega_i = \langle \boldsymbol{\omega}_0 - \boldsymbol{\theta}^*, \mathbf{v}_i \rangle$. If $|\beta^{tr}|, |\beta^{te}| < 1/\lambda_1$, n_1 is large ensuring that $\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) > 0, \forall i$ and $\alpha < 1/(c(\beta^{tr}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}))$, then the meta excess risk $R(\bar{\boldsymbol{\omega}}_T, \beta^{te})$ is bounded above as follows*

$$R(\bar{\boldsymbol{\omega}}_T, \beta^{te}) \leq \text{Bias} + \text{Var}$$

where

$$\begin{aligned}
\text{Bias} &= \frac{2}{\alpha^2 T} \sum_i \Xi_i \frac{\omega_i^2}{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})} \\
\text{Var} &= \frac{2}{(1 - \alpha c(\beta^{tr}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}))} \left(\sum_i \Xi_i \right) \\
&\quad \times [f(\beta^{tr}, n_2, \sigma, \boldsymbol{\Sigma}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}) + 2c(\beta^{tr}, \boldsymbol{\Sigma}) \sum_i \left(\frac{\mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) \geq \frac{1}{\alpha T}}}{T \alpha \mu_i(\mathbf{H}_{n_1, \beta^{tr}})} + \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) < \frac{1}{\alpha T}} \right) \lambda_i \omega_i^2].
\end{aligned}$$

Proof. By Lemma B.2, we have

$$R(\bar{\boldsymbol{\omega}}_T, \beta^{te}) \leq 2\mathcal{E}_{\text{bias}} + 2\mathcal{E}_{\text{var}}.$$

Using Lemma B.7 to bound $\mathcal{E}_{\text{bias}}$, and Lemma B.10 to bound \mathcal{E}_{var} , we have

$$\begin{aligned}
& R(\bar{\omega}_T, \beta^{\text{te}}) \\
& \leq \frac{2f(\beta^{\text{tr}}, n_2, \sigma, \Sigma, \Sigma_{\theta})}{(1 - \alpha c(\beta^{\text{tr}}, \Sigma)) \text{tr}(\Sigma)} \\
& \times \sum_i \left(\frac{1}{T} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}} + T\alpha^2 \mu_i^2(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \frac{\mu_i(\mathbf{H}_{m, \beta^{\text{te}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \\
& + \frac{4c(\beta^{\text{tr}}, \Sigma)}{T\alpha(1 - c(\beta^{\text{tr}}, \Sigma))\alpha \text{tr}(\Sigma)} \sum_i \left(\frac{1}{T} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}} + T\alpha^2 \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})^2 \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \\
& \times \sum_i \left(\frac{1}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}} + T\alpha \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \cdot \lambda_i(\langle \omega_0 - \theta^*, \mathbf{v}_i \rangle)^2 \\
& + 2 \sum_i \left(\frac{1}{\alpha^2 T^2} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}} + \mu_i^2(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \frac{\omega_i^2 \mu_i(\mathbf{H}_{m, \beta^{\text{te}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})^2}.
\end{aligned}$$

Incorporating with the definition of effective meta weight

$$\Xi_i(\Sigma, \alpha, T) = \begin{cases} \mu_i(\mathbf{H}_{m, \beta^{\text{te}}}) / (T\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})) & \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}; \\ T\alpha^2 \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \mu_i(\mathbf{H}_{m, \beta^{\text{te}}}) & \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}; \end{cases} \quad (21)$$

we obtain

$$\left(\frac{1}{T} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}} + T\alpha^2 \mu_i^2(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \frac{\mu_i(\mathbf{H}_{m, \beta^{\text{te}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} = \Xi_i(\Sigma, \alpha, T).$$

Therefore,

$$R(\bar{\omega}_T, \beta^{\text{te}}) \leq \text{Bias} + \text{Var}$$

where

$$\begin{aligned}
\text{Bias} &= \frac{2}{\alpha^2 T} \sum_i \Xi_i \frac{\omega_i^2}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \\
\text{Var} &= \frac{2}{(1 - \alpha c(\beta^{\text{tr}}, \Sigma)) \text{tr}(\Sigma)} \left(\sum_i \Xi_i \right) \\
& \times \underbrace{[f(\beta^{\text{tr}}, n_2, \sigma, \Sigma_{\theta}, \Sigma) + 2c(\beta^{\text{tr}}, \Sigma) \sum_i \left(\frac{\mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}}}{T\alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} + \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \lambda_i \omega_i^2]}_{V_2}.
\end{aligned}$$

Note that the term V_2 is obtained by our analysis for $\mathcal{E}_{\text{bias}}$. However, it originates from the stochasticity of SGD, and hence we treat this term as the variance in our final results. \square

C Analysis for Lower Bound (Theorem 2)

C.1 Fourth Moment Lower Bound for Meta Noise

Similarly to upper bound, we need some technical results for the fourth moment of meta data \mathbf{B} and noise ξ to proceed the lower bound analysis.

Lemma C.1. *Suppose Assumption 1-3 hold. Given $|\beta^{\text{tr}}| < \frac{1}{\lambda_1}$, for any PSD matrix \mathbf{A} , we have*

$$\mathbb{E}[\mathbf{B}^{\top} \mathbf{B} \mathbf{A} \mathbf{B}^{\top} \mathbf{B}] \succeq \mathbf{H}_{n_1, \beta^{\text{tr}}} \mathbf{A} \mathbf{H}_{n_1, \beta^{\text{tr}}} + \frac{b_1}{n_2} \text{tr}(\mathbf{H}_{n_1, \beta^{\text{tr}}} \mathbf{A}) \mathbf{H}_{n_1, \beta^{\text{tr}}} \quad (22)$$

$$\Pi \succeq \frac{1}{n_2} g(\beta^{\text{tr}}, n_1, \sigma, \Sigma_{\theta}, \Sigma) \mathbf{H}_{n_1, \beta^{\text{tr}}} \quad (23)$$

where $g(\beta, n, \sigma, \Sigma, \Sigma_{\theta}) := \sigma^2 + b_1 \text{tr}(\Sigma_{\theta} \mathbf{H}_{n, \beta}) + \beta^2 \mathbf{1}_{\beta \leq 0} b_1 \text{tr}(\Sigma^2) / n$.

Proof. With a slight abuse of notations, we write β^{tr} as β , \mathbf{X}^{in} as \mathbf{X} in this proof. Note that $\mathbf{x} \in \mathbb{R}^d \sim \mathcal{P}_{\mathbf{x}}$ is independent of \mathbf{X}^{in} . We first derive

$$\begin{aligned}
& \mathbb{E}[\mathbf{B}^\top \mathbf{B} \mathbf{A} \mathbf{B}^\top \mathbf{B}] \\
&= \frac{1}{n_2} \mathbb{E} \left[\left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{x} \mathbf{x}^\top \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{A} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{x} \mathbf{x}^\top \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\
&+ \frac{n_2 - 1}{n_2} \mathbb{E} \left[\left(\mathbf{I} - \frac{\beta}{n_2} \mathbf{X}^\top \mathbf{X} \right) \Sigma \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{A} \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \Sigma \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\
&\stackrel{(a)}{\succeq} \frac{b_1}{n_2} \mathbb{E} \left[\text{tr} \left(\mathbf{A} \left(\mathbf{I} - \frac{\beta}{n} \mathbf{X}^\top \mathbf{X} \right) \Sigma \left(\mathbf{I} - \frac{\beta}{n} \mathbf{X}^\top \mathbf{X} \right) \right) \left(\mathbf{I} - \frac{\beta}{n} \mathbf{X}^\top \mathbf{X} \right) \Sigma \left(\mathbf{I} - \frac{\beta}{n} \mathbf{X}^\top \mathbf{X} \right) \right] \\
&+ \mathbf{H}_{n_1, \beta} \mathbf{A} \mathbf{H}_{n_1, \beta} \\
&\succeq \frac{b_1}{n_2} \text{tr}(\mathbf{H}_{n_1, \beta} \mathbf{A}) \mathbf{H}_{n_1, \beta} + \mathbf{H}_{n_1, \beta} \mathbf{A} \mathbf{H}_{n_1, \beta}
\end{aligned}$$

where (a) is implied by Assumption 1.

Recall that Π takes the following form:

$$\Pi = \frac{\sigma^2}{n_2} \mathbf{H}_{n_1, \beta} + \mathbb{E}[\mathbf{B}^\top \mathbf{B} \Sigma_\theta \mathbf{B}^\top \mathbf{B}] + \sigma^2 \cdot \frac{\beta^2}{n_2 n_1^2} \mathbb{E}[\mathbf{B}^\top \mathbf{X}^{\text{out}} \mathbf{X}^\top \mathbf{X} \mathbf{X}^{\text{out}^\top} \mathbf{B}].$$

The second term can be directly bounded by letting $\mathbf{A} = \Sigma_\theta$ in eq. (22), and we have:

$$\mathbb{E}[\mathbf{B}^\top \mathbf{B} \Sigma_\theta \mathbf{B}^\top \mathbf{B}] \succeq \frac{b_1}{n_2} \text{tr}(\mathbf{H}_{n_1, \beta} \Sigma_\theta) \mathbf{H}_{n_1, \beta}.$$

For the third term:

$$\begin{aligned}
& \frac{1}{n_2} \mathbb{E}[\mathbf{B}^\top \mathbf{X}^{\text{out}} \mathbf{X}^\top \mathbf{X} \mathbf{X}^{\text{out}^\top} \mathbf{B}] \\
&= \frac{1}{n_2} \mathbb{E} \left[\left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \mathbf{x} \mathbf{x}^\top \mathbf{X}^\top \mathbf{X} \mathbf{x} \mathbf{x}^\top \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\
&+ \frac{n_2 - 1}{n_2} \mathbb{E} \left[\left(\mathbf{I} - \frac{\beta}{n_2} \mathbf{X}^\top \mathbf{X} \right) \Sigma \mathbf{X}^\top \mathbf{X} \Sigma \left(\mathbf{I} - \frac{\beta}{n_1} \mathbf{X}^\top \mathbf{X} \right) \right] \\
&\succeq \frac{n_1 b_1 \text{tr}(\Sigma^2)}{n_2} \mathbf{H}_{n_1, \beta} \mathbf{1}_{\beta \leq 0}
\end{aligned}$$

Putting these results together completes the proof. \square

C.2 Bias-Variance Decomposition

For the lower bound analysis, we also decompose the excess risk into bias and variance terms.

Lemma C.2 (Bias-variance decomposition, lower bound). *Following the notations in eq. (12), the excess risk can be decomposed as follows:*

$$R(\bar{\omega}_T, \beta^{te}) \geq \mathcal{E}_{\text{bias}} + \mathcal{E}_{\text{var}}$$

where

$$\begin{aligned}
\mathcal{E}_{\text{bias}} &= \frac{1}{2T^2} \cdot \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \langle (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^{k-t} \mathbf{H}_{m, \beta^e}, \mathbf{D}_t \rangle, \\
\mathcal{E}_{\text{var}} &= \frac{1}{2T^2} \cdot \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \langle (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^{\text{tr}}})^{k-t} \mathbf{H}_{m, \beta^e}, \mathbf{V}_t \rangle.
\end{aligned}$$

Proof. The proof is similar to that for Lemma B.3, and the inequality sign is reversed since we only calculate the half of summation. In particular,

$$\begin{aligned}
\mathbb{E}[\bar{\mathbf{e}}_T^{\text{var}} \otimes \bar{\mathbf{e}}_T^{\text{var}}] &= \frac{1}{T^2} \sum_{1 \leq t < k \leq T-1} \mathbb{E}[\mathbf{e}_t^{\text{var}} \otimes \mathbf{e}_k^{\text{var}}] + \frac{1}{T^2} \sum_{1 \leq k < t \leq T-1} \mathbb{E}[\mathbf{e}_t^{\text{var}} \otimes \mathbf{e}_k^{\text{var}}] \\
&\succeq \frac{1}{T^2} \sum_{1 \leq t < k \leq T-1} \mathbb{E}[\mathbf{e}_t^{\text{var}} \otimes \mathbf{e}_k^{\text{var}}].
\end{aligned}$$

For $t \leq k$, $\mathbb{E}[\mathbf{e}_k^{\text{var}} | \mathbf{e}_t^{\text{var}}] = (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k-t} \mathbf{e}_t^{\text{var}}$, since $\mathbb{E}[\mathbf{B}_t^\top \boldsymbol{\xi}_t | \mathbf{e}_{t-1}] = \mathbf{0}$. From this

$$\mathbb{E}[\bar{\mathbf{e}}_T^{\text{var}} \otimes \bar{\mathbf{e}}_T^{\text{var}}] \succeq \frac{1}{T^2} \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \mathbf{V}_t (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k-t}.$$

Plugging this into $\frac{1}{2} \langle \mathbf{H}_{m, \beta^e}, \mathbb{E}[\bar{\mathbf{e}}_T^{\text{var}} \otimes \bar{\mathbf{e}}_T^{\text{var}}] \rangle$, we obtain:

$$\begin{aligned} & \frac{1}{2} \langle \mathbf{H}_{m, \beta^e}, \mathbb{E}[\bar{\mathbf{e}}_T^{\text{var}} \otimes \bar{\mathbf{e}}_T^{\text{var}}] \rangle \\ & \geq \frac{1}{2T^2} \sum_{t=0}^{T-1} \sum_{k=t+1}^{T-1} \langle \mathbf{H}_{m, \beta^e}, \mathbf{V}_t (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k-t} \rangle \\ & = \frac{1}{2T^2} \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \langle (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k-t} \mathbf{H}_{m, \beta^e}, \mathbf{V}_t \rangle \\ & = \underline{\mathcal{E}}_{\text{var}}. \end{aligned}$$

The proof is the same for the term $\underline{\mathcal{E}}_{\text{bias}}$. □

C.3 Bounding the Bias

We first bound the summation of \mathbf{D}_t , i.e. $\mathbf{S}_k = \sum_{t=0}^{k-1} \mathbf{D}_t$.

Lemma C.3 (Bounding \mathbf{S}_t). *If the stepsize satisfies $\alpha < 1/(2 \max_i \{\mu_i(\mathbf{H}_{n_1, \beta^r})\})$, then for any $k \geq 2$, it holds that*

$$\begin{aligned} \mathbf{S}_k & \succeq \frac{b_1}{4n_2} \text{tr} \left(\left(\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k/2} \right) \mathbf{D}_0 \right) \cdot \left(\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k/2} \right) \\ & \quad + \sum_{t=0}^{k-1} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^t \cdot \mathbf{D}_0 \cdot (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^t. \end{aligned}$$

Proof. By eq. (18), since $\widetilde{\mathcal{M}} - \mathcal{M}$ is a PSD mapping, we have

$$\begin{aligned} \mathbf{S}_k & = \mathbf{D}_0 + (\mathcal{I} - \alpha \widetilde{\mathcal{T}}) \circ \mathbf{S}_{k-1} + \alpha^2 (\mathcal{M} - \widetilde{\mathcal{M}}) \circ \mathbf{S}_{k-1} \\ & \succeq \sum_{t=0}^{k-1} (\mathcal{I} - \alpha \widetilde{\mathcal{T}})^t \circ \mathbf{D}_0 \\ & = \sum_{t=0}^{k-1} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^t \cdot \mathbf{D}_0 \cdot (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^t. \end{aligned} \tag{24}$$

Note that for PSD \mathbf{A} ,

$$(\mathcal{M} - \widetilde{\mathcal{M}}) \circ \mathbf{A} = \mathbb{E}[\mathbf{B}^\top \mathbf{B} \mathbf{A} \mathbf{B}^\top \mathbf{B}] - \mathbf{H}_{n_1, \beta^r} \mathbf{A} \mathbf{H}_{n_1, \beta^r}$$

By Lemma C.1, we have

$$\begin{aligned} (\mathcal{M} - \widetilde{\mathcal{M}}) \circ \mathbf{S}_k & \succeq \frac{b_1}{n_2} \text{tr} (\mathbf{H}_{n_1, \beta^r} \mathbf{S}_k) \mathbf{H}_{n_1, \beta^r} \\ & \succeq \frac{b_1}{n_2} \text{tr} \left(\sum_{t=0}^{k-1} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{2t} \mathbf{H}_{n_1, \beta^r} \cdot \mathbf{D}_0 \right) \mathbf{H}_{n_1, \beta^r} \\ & \succeq \frac{b_1}{n_2} \text{tr} \left(\sum_{t=0}^{k-1} (\mathbf{I} - 2\alpha \mathbf{H}_{n_1, \beta^r})^t \mathbf{H}_{n_1, \beta^r} \cdot \mathbf{D}_0 \right) \mathbf{H}_{n_1, \beta^r} \\ & \succeq \frac{b_1}{2n_2 \alpha} \text{tr} \left((\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^k) \mathbf{D}_0 \right) \mathbf{H}_{n_1, \beta^r}. \end{aligned} \tag{25}$$

Substituting eq. (25) back into eq. (24), and solving the recursion, we obtain

$$\begin{aligned}
\mathbf{S}_k &\succeq \sum_{t=0}^{k-1} (\mathcal{I} - \alpha \tilde{\mathcal{T}})^t \circ \left\{ \frac{b_1 \alpha}{2n_2} \text{tr} \left((\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k-1-t}) \mathbf{D}_0 \right) \mathbf{H} + \mathbf{D}_0 \right\} \\
&= \frac{b_1 \alpha}{2n_2} \underbrace{\sum_{t=0}^{k-1} \text{tr} \left((\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k-1-t}) \mathbf{D}_0 \right) \cdot (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{2t} \mathbf{H}_{n_1, \beta^r}}_{\mathbf{J}_4} \\
&\quad + \sum_{t=0}^{k-1} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^t \cdot \mathbf{D}_0 \cdot (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^t.
\end{aligned}$$

The term \mathbf{J}_4 can be further bounded by the following:

$$\begin{aligned}
\mathbf{J}_4 &\succeq \sum_{t=0}^{k-1} \text{tr} \left((\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k-1-t}) \mathbf{D}_0 \right) \cdot (\mathbf{I} - 2\alpha \mathbf{H}_{n_1, \beta^r})^t \mathbf{H}_{n_1, \beta^r} \\
&\succeq \text{tr} \left(\left((\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k/2}) \mathbf{D}_0 \right) \cdot \sum_{t=0}^{k/2-1} (\mathbf{I} - 2\alpha \mathbf{H}_{n_1, \beta^r})^t \mathbf{H}_{n_1, \beta^r} \right) \\
&\succeq \frac{1}{2\alpha} \text{tr} \left(\left((\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k/2}) \mathbf{D}_0 \right) \cdot \left(\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k/2} \right) \right)
\end{aligned}$$

which completes the proof. \square

Then we can bound the bias term.

Lemma C.4 (Bounding the bias). *Let $\omega_i = \langle \omega_0 - \boldsymbol{\theta}^*, \mathbf{v}_i \rangle$. If $\alpha < \frac{1}{c(\beta^r, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})}$, for sufficiently large n_1 , s.t. $\mu_i(\mathbf{H}_{n_1, \beta^r}) > 0, \forall i$, then we have*

$$\begin{aligned}
\underline{\mathcal{E}}_{bias} &\geq \frac{1}{100\alpha^2 T} \sum_i \Xi_i \frac{\omega_i^2}{\mu_i(\mathbf{H}_{n_1, \beta^r})} + \frac{b_1}{1000n_2(1 - \alpha c(\beta^r, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}))} \sum_i \Xi_i \\
&\quad \times \sum_i \left(\frac{\mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^r}) \geq \frac{1}{\alpha T}}}{T\alpha\mu_i(\mathbf{H}_{n_1, \beta^r})} + \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^r}) < \frac{1}{\alpha T}} \right) \lambda_i \omega_i^2.
\end{aligned}$$

Proof. From Lemma C.2, we have

$$\begin{aligned}
\underline{\mathcal{E}}_{bias} &= \frac{1}{2T^2} \cdot \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \langle (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{k-t} \mathbf{H}_{m, \beta^{te}}, \mathbf{D}_t \rangle \\
&= \frac{1}{2\alpha T^2} \cdot \sum_{t=0}^{T-1} \left\langle (\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{T-t}) \mathbf{H}_{n_1, \beta^r}^{-1} \mathbf{H}_{m, \beta^{te}}, \mathbf{D}_t \right\rangle \\
&\geq \frac{1}{2\alpha T^2} \left\langle \left((\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{T/2}) \mathbf{H}_{n_1, \beta^r}^{-1} \mathbf{H}_{m, \beta^{te}}, \sum_{t=0}^{T/2} \mathbf{D}_t \right) \right\rangle \\
&\geq \frac{1}{2\alpha T^2} \left\langle \left((\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{T/2}) \mathbf{H}_{n_1, \beta^r}^{-1} \mathbf{H}_{m, \beta^{te}}, \mathbf{S}_{\frac{T}{2}} \right) \right\rangle.
\end{aligned}$$

Applying Lemma C.3 to $\mathbf{S}_{\frac{T}{2}}$, we obtain:

$$\underline{\mathcal{E}}_{bias} \geq \underline{d}_1 + \underline{d}_2$$

where

$$\begin{aligned} \underline{d}_1 &= \frac{b_1}{8\alpha n_2 T^2} \text{tr} \left(\left(\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{T/4} \right) \mathbf{D}_0 \right) \\ &\quad \times \left\langle \left(\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{T/2} \right) \mathbf{H}_{n_1, \beta^r}^{-1} \mathbf{H}_{m, \beta^e}, \left(\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{T/4} \right) \right\rangle \\ \underline{d}_2 &= \frac{1}{2\alpha T^2} \left\langle \left(\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{T/2} \right) \mathbf{H}_{n_1, \beta^r}^{-1} \mathbf{H}_{m, \beta^e}, \right. \\ &\quad \left. \sum_{t=0}^{T/2-1} (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^t \cdot \mathbf{D}_0 \cdot (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^t \right\rangle. \end{aligned}$$

Moreover,

$$\begin{aligned} \underline{d}_2 &\geq \frac{1}{2\alpha T^2} \left\langle \left(\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{T/2} \right) \mathbf{H}_{n_1, \beta^r}^{-1} \mathbf{H}_{m, \beta^e}, \sum_{t=0}^{T/2-1} (\mathbf{I} - 2\alpha \mathbf{H}_{n_1, \beta^r})^t \mathbf{D}_0 \right\rangle; \\ &\geq \frac{1}{4\alpha^2 T^2} \left\langle \left(\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{T/2} \right)^2 \mathbf{H}_{n_1, \beta^r}^{-2} \mathbf{H}_{m, \beta^e}, \mathbf{D}_0 \right\rangle. \end{aligned}$$

Using the diagonalizing technique similar to the proof for Lemma B.7, we have

$$\underline{d}_1 \geq \frac{b_1}{8\alpha n_2 T^2} \left(\sum_i \left(1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^r}))^{T/4} \right) \omega_i^2 \right) \quad (26)$$

$$\times \left(\sum_i \left(1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^r}))^{T/4} \right)^2 \frac{\mu_i(\mathbf{H}_{m, \beta^e})}{\mu_i(\mathbf{H}_{n_1, \beta^r})} \right), \quad (27)$$

$$\underline{d}_2 \geq \frac{1}{4\alpha^2 T^2} \sum_i \left(1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^r}))^{T/4} \right)^2 \frac{\mu_i(\mathbf{H}_{m, \beta^e})}{\mu_i^2(\mathbf{H}_{n_1, \beta^r})} \omega_i^2. \quad (28)$$

We use the following fact to bound the polynomial term. For $h_1(x) = 1 - (1 - x)^{\frac{T}{4}}$, we have

$$h_1(x) \geq \begin{cases} \frac{1}{5} & x \geq 1/T \\ \frac{T}{5}x & x < 1/T \end{cases}$$

i.e., $1 - (1 - \alpha \mu_i(\mathbf{H}_{n_1, \beta^r}))^{T/4} \geq \left(\frac{1}{5} \mathbf{1}_{\alpha \mu_i(\mathbf{H}_{n_1, \beta^r}) \geq \frac{1}{T}} + \frac{\alpha \mu_i(\mathbf{H}_{n_1, \beta^r})}{5} \mathbf{1}_{\alpha \mu_i(\mathbf{H}_{n_1, \beta^r}) < \frac{1}{T}} \right)$. Substituting this back into eqs. (27) and (28), and using the definition of effective meta weight Ξ_i complete the proof. \square

C.4 Bounding the Variance

We first bound the term \mathbf{V}_t .

Lemma C.5 (Bounding \mathbf{V}_t). *If the stepsize satisfies $\alpha < 1/(\max_i \{\mu_i(\mathbf{H}_{n_1, \beta^r})\})$, it holds that*

$$\mathbf{V}_t \preceq \frac{\alpha g(\beta^r, n_1, \Sigma, \Sigma_\theta)}{2} \cdot (\mathbf{I} - (\mathbf{I} - \alpha \mathbf{H}_{n_1, \beta^r})^{2t}).$$

Proof. With a slight abuse of notations, we write $g(\beta^{\text{tr}}, n_1, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta)$ as g . By definition,

$$\begin{aligned}
\mathbf{V}_t &= (\mathcal{I} - \alpha\mathcal{T}) \circ \mathbf{V}_{t-1} + \alpha^2\Pi \\
&= (\mathcal{I} - \alpha\tilde{\mathcal{T}}) \circ \mathbf{V}_{t-1} + (\mathcal{M} - \tilde{\mathcal{M}}) \circ \mathbf{V}_{t-1} + \alpha^2\Pi \\
&\stackrel{(a)}{\succeq} (\mathcal{I} - \alpha\tilde{\mathcal{T}}) \circ \mathbf{V}_{t-1} + \alpha^2g\mathbf{H}_{n_1, \beta^{\text{tr}}} \\
&\stackrel{(b)}{=} \alpha^2g \cdot \sum_{k=0}^{t-1} (\mathcal{I} - \alpha\tilde{\mathcal{T}})^k \circ \mathbf{H}_{n_1, \beta^{\text{tr}}} \\
&= \alpha^2g \cdot \sum_{k=0}^{t-1} (\mathbf{I} - \alpha\mathbf{H}_{n_1, \beta^{\text{tr}}})^k \mathbf{H}_{n_1, \beta^{\text{tr}}} (\mathbf{I} - \alpha\mathbf{H}_{n_1, \beta^{\text{tr}}})^k \quad (\text{by the definition of } \mathcal{I} - \alpha\tilde{\mathcal{T}}) \\
&= \alpha g \cdot (\mathbf{I} - (\mathbf{I} - \alpha\mathbf{H}_{n_1, \beta^{\text{tr}}})^{2t}) \cdot (2\mathbf{I} - \alpha\mathbf{H}_{n_1, \beta^{\text{tr}}})^{-1} \\
&\stackrel{(c)}{\succeq} \frac{\alpha g}{2} \cdot (\mathbf{I} - (\mathbf{I} - \alpha\mathbf{H}_{n_1, \beta^{\text{tr}}})^{2t})
\end{aligned}$$

where (a) follows from the Lemma C.1, (b) follows by solving the recursion and (c) holds since we directly replace $(2\mathbf{I} - \alpha\mathbf{H}_{n_1, \beta^{\text{tr}}})^{-1}$ by $(2\mathbf{I})^{-1}$. \square

Lemma C.6 (Bounding the variance). *Let $\omega_i = \langle \boldsymbol{\omega}_0 - \boldsymbol{\theta}^*, \mathbf{v}_i \rangle$. If $\alpha < \frac{1}{c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma})}$, for sufficiently large n_1 , s.t. $\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) > 0, \forall i$, for $T > 10$, then we have*

$$\underline{\mathcal{E}}_{\text{var}} \geq \frac{g(\beta^{\text{tr}}, n_1, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta)}{100n_2(1 - \alpha c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}))} \sum_i \Xi_i.$$

Proof. From Lemma C.2, we have

$$\begin{aligned}
\underline{\mathcal{E}}_{\text{var}} &= \frac{1}{2T^2} \cdot \sum_{t=0}^{T-1} \sum_{k=t}^{T-1} \langle (\mathbf{I} - \alpha\mathbf{H}_{n_1, \beta^{\text{tr}}})^{k-t} \mathbf{H}_{m, \beta^{\text{tr}}}, \mathbf{V}_t \rangle \\
&= \frac{1}{2\alpha T^2} \cdot \sum_{t=0}^{T-1} \langle (\mathbf{I} - (\mathbf{I} - \alpha\mathbf{H}_{n_1, \beta^{\text{tr}}})^{T-t}) \mathbf{H}_{n_1, \beta^{\text{tr}}}^{-1} \mathbf{H}_{m, \beta^{\text{tr}}}, \mathbf{V}_t \rangle.
\end{aligned}$$

Then applying Lemma C.5, and writing $g(\beta^{\text{tr}}, n_1, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta)$ as g , we obtain

$$\begin{aligned}
\underline{\mathcal{E}}_{\text{var}} &\geq \frac{g}{4T^2} \cdot \sum_{t=0}^{T-1} \left\langle (\mathbf{I} - (\mathbf{I} - \alpha\mathbf{H}_{n_1, \beta^{\text{tr}}})^{T-t}) \mathbf{H}_{n_1, \beta^{\text{tr}}}^{-1} \mathbf{H}_{m, \beta^{\text{tr}}}, (\mathbf{I} - (\mathbf{I} - \alpha\mathbf{H}_{n_1, \beta^{\text{tr}}})^{2t}) \right\rangle \\
&= \frac{g}{4T^2} \sum_i \frac{\mu_i(\mathbf{H}_{m, \beta^{\text{tr}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \sum_{t=0}^{T-1} (1 - (1 - \alpha\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})^{T-t})) (1 - (1 - \alpha\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})^{2t})) \\
&\geq \frac{g}{4T^2} \sum_i \frac{\mu_i(\mathbf{H}_{m, \beta^{\text{tr}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \sum_{t=0}^{T-1} (1 - (1 - \alpha\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})^{T-t-1})) (1 - (1 - \alpha\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})^t)) \quad (29)
\end{aligned}$$

where the equality holds by applying the diagonalizing technique again. Following the trick similar to that in [22] to lower bound the function $h_2(x) := \sum_{t=0}^{T-1} (1 - (1-x)^{T-t-1}) (1 - (1-x)^t)$ defined on $x \in (0, 1)$, for $T > 10$, we have

$$f(x) \geq \begin{cases} \frac{T}{20}, & \frac{1}{T} \leq x < 1 \\ \frac{3T^3}{50} x^2, & 0 < x < \frac{1}{T} \end{cases}$$

Substituting this back into eq. (29), and using the definition of effective meta weight Ξ_i completes the proof. \square

C.5 Proof of Theorem 2

Theorem C.1 (Theorem 2 Restated). *Let $\omega_i = \langle \boldsymbol{\omega}_0 - \boldsymbol{\theta}^*, \mathbf{v}_i \rangle$. If $|\beta^{tr}|, |\beta^{te}| < 1/\lambda_1$, n_1 is large ensuring that $\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) > 0, \forall i$ and $\alpha < 1/(c(\beta^{tr}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}))$. For $T > 10$, the meta excess risk $R(\bar{\boldsymbol{\omega}}_T, \beta^{te})$ is bounded below as follows*

$$\begin{aligned} R(\bar{\boldsymbol{\omega}}_T, \beta^{te}) &\geq \frac{1}{100\alpha^2 T} \sum_i \Xi_i \frac{\omega_i^2}{\mu_i(\mathbf{H}_{n_1, \beta^{tr}})} + \frac{1}{n_2} \cdot \frac{1}{(1 - \alpha c(\beta^{tr}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}))} \sum_i \Xi_i \\ &\quad \times \left[\frac{1}{100} g(\beta^{tr}, n_1, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta) + \frac{b_1}{1000} \sum_i \left(\frac{\mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) \geq \frac{1}{\alpha T}}}{T\alpha\mu_i(\mathbf{H}_{n_1, \beta^{tr}})} + \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{tr}}) < \frac{1}{\alpha T}} \right) \lambda_i \omega_i^2 \right]. \end{aligned}$$

Proof. The proof can be completed by combining Lemmas C.4 and C.6. \square

D Proofs for Section 4.2

D.1 Proof of Lemma 1

Proof of Lemma 1. For the single task setting, we first simplify our notations in Theorem B.3 as follows.

$$c(0, \boldsymbol{\Sigma}) = c_1, \quad f(0, n_2, \sigma, \boldsymbol{\Sigma}, \mathbf{0}) = \sigma^2/n_2, \quad \mathbf{H}_{n_1, \beta^{tr}} = \boldsymbol{\Sigma}.$$

By Theorem B.3, we have

$$\begin{aligned} \text{Bias} &= \frac{2}{\alpha^2 T} \sum_i \left(\frac{1}{T} \mathbf{1}_{\lambda_i \geq \frac{1}{\alpha T}} + T\alpha^2 \lambda_i^2 \mathbf{1}_{\lambda_i < \frac{1}{\alpha T}} \right) \frac{\omega_i^2 \mu_i(\mathbf{H}_{m, \beta^{te}})}{\lambda_i^2} \\ &\leq \frac{2}{\alpha^2 T} \sum_i (\alpha \lambda_i \mathbf{1}_{\lambda_i \geq \frac{1}{\alpha T}} + \alpha \lambda_i \mathbf{1}_{\lambda_i < \frac{1}{\alpha T}}) \frac{\omega_i^2 \mu_i(\mathbf{H}_{m, \beta^{te}})}{\lambda_i^2}. \end{aligned}$$

For large m , we have $\mu_i(\mathbf{H}_{m, \beta^{te}}) = (1 - \beta^{te} \lambda_i)^2 \lambda_i + o(1)$. Therefore,

$$\text{Bias} \leq \frac{2(1 - \beta^{te} \lambda_d)^2}{\alpha^2 T} \sum_i \omega_i^2 \leq \mathcal{O}\left(\frac{1}{T}\right).$$

For the variance term,

$$\begin{aligned} \text{Var} &= \frac{2}{(1 - \alpha c_1 \text{tr}(\boldsymbol{\Sigma}))} \underbrace{\sum_i \left(\frac{1}{T} \mathbf{1}_{\lambda_i \geq \frac{1}{\alpha T}} + T\alpha^2 \lambda_i^2 \mathbf{1}_{\lambda_i < \frac{1}{\alpha T}} \right) \frac{\mu_i(\mathbf{H}_{m, \beta^{te}})}{\lambda_i}}_{J_5} \\ &\quad \times \left[\frac{\sigma^2}{n_2} + 2c_1 \sum_i \left(\frac{\mathbf{1}_{\lambda_i \geq \frac{1}{\alpha T}}}{T\alpha\lambda_i} + \mathbf{1}_{\lambda_i < \frac{1}{\alpha T}} \right) \lambda_i \omega_i^2 \right]. \end{aligned}$$

It is easy to check that

$$\sum_i \left(\frac{\mathbf{1}_{\lambda_i \geq \frac{1}{\alpha T}}}{T\alpha\lambda_i} + \mathbf{1}_{\lambda_i < \frac{1}{\alpha T}} \right) \lambda_i \omega_i^2 \leq \sum_i \left(\frac{\mathbf{1}_{\lambda_i \geq \frac{1}{\alpha T}}}{T\alpha} + \frac{1}{\alpha T} \mathbf{1}_{\lambda_i < \frac{1}{\alpha T}} \right) \omega_i^2 \leq \mathcal{O}(1/T).$$

Moreover,

$$J_5 \leq (1 - \beta^{te} \lambda_d)^2 \sum_i \left(\frac{1}{T} \mathbf{1}_{\lambda_i \geq \frac{1}{\alpha T}} + T\alpha^2 \lambda_i^2 \mathbf{1}_{\lambda_i < \frac{1}{\alpha T}} \right).$$

The term $\sum_i \left(\frac{1}{T} \mathbf{1}_{\lambda_i \geq \frac{1}{\alpha T}} + T\alpha^2 \lambda_i^2 \mathbf{1}_{\lambda_i < \frac{1}{\alpha T}} \right)$ has the form similar to Corollary 2.3 in [22] and we directly have $J_5 = \mathcal{O}(\log^{-p}(T))$, which implies

$$\text{Var} = \mathcal{O}(\log^{-p}(T)).$$

Thus we complete the proof. \square

D.2 Proof of Proposition 2

Proof of Proposition 2. We first consider the bias term in Theorems B.3 and C.1 (up to absolute constants):

$$\text{Bias} = \frac{2}{\alpha^2 T} \sum_i \left(\frac{1}{T} \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}} + \alpha^2 T \mu_i^2(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \frac{\omega_i^2 \mu_i(\mathbf{H}_{m, \beta^{\text{te}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})^2}.$$

If $\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}$, $\frac{1}{T} \leq \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})$; and if $\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}$, then $\alpha^2 T \mu_i^2(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})$. Hence

$$\text{Bias} \leq \frac{1}{\alpha^2 T} \sum_i \frac{\omega_i^2 \mu_i(\mathbf{H}_{m, \beta^{\text{te}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \leq \frac{2}{\alpha^2 T} \cdot \max_i \frac{\mu_i(\mathbf{H}_{m, \beta^{\text{te}}})}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \|\boldsymbol{\omega}_0 - \boldsymbol{\theta}^*\|^2 = \mathcal{O}\left(\frac{1}{T}\right).$$

Moreover,

$$\begin{aligned} & \sum_i \left(\frac{\mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq \frac{1}{\alpha T}}}{T \alpha \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} + \mathbf{1}_{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}} \right) \lambda_i \omega_i^2 \\ & \stackrel{(a)}{\leq} \frac{1}{\alpha T} \sum_i \frac{\lambda_i}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \omega_i^2 \\ & \leq \frac{1}{\alpha T} \max_i \frac{\lambda_i}{\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})} \|\boldsymbol{\omega}_0 - \boldsymbol{\theta}^*\|^2 = \mathcal{O}\left(\frac{1}{T}\right) \end{aligned}$$

where (a) holds since we directly upper bound $\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}})$ by $\frac{1}{\alpha T}$ when $\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) < \frac{1}{\alpha T}$. Therefore, it is essential to analyze $f(\beta^{\text{tr}}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta)$ ($\sum_i \Xi_i$) and $g(\beta^{\text{tr}}, n_1, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta)$ ($\sum_i \Xi_i$) from variance term in the upper and lower bounds respectively.

Then we calculate some rates of interesting in Theorems B.3 and C.1 under the specific data and task distributions in Proposition 2.

If the spectrum of $\boldsymbol{\Sigma}$ satisfies $\lambda_k = k^{-1} \log^{-p}(k+1)$, then it is easily verified that $\text{tr}(\boldsymbol{\Sigma}^s) = \mathcal{O}(1)$ for $s = 1, \dots, 4$. By discussions on Assumption 3 in Appendix F, we have $C(\beta, \boldsymbol{\Sigma}) = \Theta(1)$ for given β . Hence,

$$\begin{aligned} c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) &= \Theta(1) \\ f(\beta^{\text{tr}}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta) &= c(\beta^{\text{tr}}, \boldsymbol{\Sigma}) \text{tr}(\boldsymbol{\Sigma}_\theta \boldsymbol{\Sigma}) + \Theta(1) \\ g(\beta^{\text{tr}}, n_1, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta) &= b_1 \text{tr}(\boldsymbol{\Sigma}_\theta \mathbf{H}_{n_1, \beta^{\text{tr}}}) + \Theta(1). \end{aligned}$$

If $r \geq 2p - 1$, then we have $g(\beta^{\text{tr}}, n_1, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta) \geq \Omega(\log^{r-p+1}(d)) \geq \Omega(\log^{r-p+1}(T))$.

Let $k^\dagger := \text{card}\{i : \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \geq 1/\alpha T\}$. For large n_1 , we have $\mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) = (1 - \beta^{\text{tr}} \lambda_i)^2 \lambda_i + o(1)$. If $k^\dagger = \mathcal{O}(T/\log^p(T+1))$, then

$$\min_{1 \leq i \leq k^\dagger + 1} \mu_i(\mathbf{H}_{n_1, \beta^{\text{tr}}}) = \omega \left(\frac{\log^p(T)}{T[\log(T) - p \log(\log(T))]^p} \right) = \omega \left(\frac{1}{T} \right)$$

which contradicts the definition of k^\dagger . Hence $k^\dagger = \Omega(T/\log^p(T+1))$. Then

$$\sum_i \Xi_i \geq \Omega \left(k^\dagger \cdot \frac{1}{T} \right) = \Omega \left(\frac{1}{\log^p(T)} \right).$$

Therefore, by Theorem C.1, $R(\bar{\boldsymbol{\omega}}, \beta^{\text{te}}) = \Omega(\log^{r-2p+1}(T))$.

For $r < 2p - 1$, if $d = T^l$, where l can be sufficiently large ($d \gg T$) but still finite, then

- If $p - 1 < r < 2p - 1$, $f(\beta^{\text{tr}}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta) \leq \mathcal{O}(\log^{r-p+1}(T))$;
- If $r \leq p - 1$, $f(\beta^{\text{tr}}, n_2, \sigma, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}_\theta) \leq \mathcal{O}(\log(\log(T)))$.

Following the analysis similar to that for Corollary 2.3 in [22], we have $\sum_i \Xi_i = \mathcal{O}(\frac{1}{\log^p(T)})$. Then by Theorem B.3

$$R(\bar{\omega}_T, \beta^{\text{te}}) = \mathcal{O}\left(\frac{1}{\log^{p-(r-p+1)^+}(T)}\right).$$

□

D.3 Proof of Proposition 3

Proof of Proposition 3. Following the analysis in Appendix D.2, it is essential to analyze $f(\beta^{\text{tr}}, n_2, \sigma, \Sigma, \Sigma_\theta)$ ($\sum_i \Xi_i$). If $d = T^l$, where l can be sufficiently large but still finite, then

$$f(\beta^{\text{tr}}, n_2, \sigma, \Sigma, \Sigma_\theta) = \tilde{\Theta}(1)$$

for $\lambda_k = k^q$ ($q > 1$) or $\lambda_k = e^{-k}$.

Following the analysis similar to that for Corollary 2.3 in [22], we have

- If $\lambda_k = k^q$ ($q > 1$), then $\sum_i \Xi_i = \mathcal{O}\left(\frac{1}{T^{\frac{q-1}{q}}}\right)$;
- If $\lambda_k = e^{-k}$, then $\sum_i \Xi_i = \mathcal{O}\left(\frac{\log(T)}{T}\right)$.

Substituting these results back into Theorem B.3, we obtain

- If $\lambda_k = k^q$ ($q > 1$), then $R(\bar{\omega}_T, \beta^{\text{te}}) = \tilde{\mathcal{O}}\left(\frac{1}{T^{\frac{q-1}{q}}}\right)$;
- If $\lambda_k = e^{-k}$, then $R(\bar{\omega}_T, \beta^{\text{te}}) = \tilde{\mathcal{O}}\left(\frac{1}{T}\right)$.

□

E Proofs for Section 4.3

E.1 Proof of Proposition 4

Proof of Proposition 4. Following the analysis in Appendix D.2, it is crucial to analyze $f(\beta^{\text{tr}}, n_2, \sigma, \Sigma, \Sigma_\theta)$ ($\sum_i \Xi_i$).

Then we calculate the rate of interest in Theorems B.3 and C.1 under some specific data and task distributions in Proposition 4. We have $\text{tr}(\Sigma^2) = \frac{1}{s} + \frac{1}{d-s} = \Theta(\frac{\log^p(T)}{T})$. Moreover, by discussions on Assumption 3 in Appendix F, $C(\beta, \Sigma) = \Theta(1)$. Hence

$$c(\beta, \Sigma) := c_1 + \tilde{\mathcal{O}}\left(\frac{1}{T}\right);$$

$$f(\beta, n, \sigma, \Sigma, \Sigma_\theta) := 2c_1 \mathcal{O}(1) + \frac{\sigma^2}{n} + \tilde{\mathcal{O}}\left(\frac{1}{T}\right).$$

By the definition of Ξ_i , we have

$$\begin{aligned} \sum_i \Xi_i &= \mathcal{O}\left(s \cdot \frac{\mu_1(\mathbf{H}_{m, \beta^{\text{te}}})}{T \mu_1(\mathbf{H}_{n_1, \beta^{\text{tr}}})} + \frac{1}{d-s} \cdot T \frac{\mu_d(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \mu_d(\mathbf{H}_{m, \beta^{\text{te}}})}{\lambda_d^2}\right) \\ &= \mathcal{O}\left(\frac{1}{\log^p(T)}\right) \frac{\mu_1(\mathbf{H}_{m, \beta^{\text{te}}})}{\mu_1(\mathbf{H}_{n_1, \beta^{\text{tr}}})} + \mathcal{O}\left(\frac{1}{\log^q(T)}\right) \frac{\mu_d(\mathbf{H}_{n_1, \beta^{\text{tr}}}) \mu_d(\mathbf{H}_{m, \beta^{\text{te}}})}{\lambda_d^2} \\ &= \mathcal{O}\left(\frac{1}{\log^p(T)}\right) \frac{(1 - \beta^{\text{te}} \lambda_1)^2}{(1 - \beta^{\text{tr}} \lambda_1)^2} + \mathcal{O}\left(\frac{1}{\log^q(T)}\right) (1 - \beta^{\text{te}} \lambda_d)^2 (1 - \beta^{\text{tr}} \lambda_d)^2 \end{aligned}$$

where the last equality follows from the fact that for large n , we have $\mu_i(\mathbf{H}_{n, \beta}) = (1 - \beta \lambda_i)^2 \lambda_i + o(1)$. Combining with the bias term which is $\mathcal{O}(\frac{1}{T})$, and applying Theorem B.3 completes the proof. □

E.2 Proof of Corollary 1

Proof of Corollary 1. For $t \in (s, K]$, by Theorem C.1, one can verify that $t = \tilde{\Theta}(K)$ for diminishing risk. Let $t = K \log^{-l}(K)$, where $p > l > 0$. Following the analysis in Appendix E.1, we have

$$R(\bar{\omega}_t^{\beta^{\text{tr}}}, \beta^{\text{te}}) \lesssim \tilde{\mathcal{O}}\left(\frac{1}{K}\right) + (2c_1\nu^2 + \frac{\sigma^2}{n_2}) \quad (30)$$

$$\times \left[\mathcal{O}\left(\frac{1}{\log^{p-l}(K)}\right) \frac{(1 - \beta^{\text{te}}\lambda_1)^2}{(1 - \beta^{\text{tr}}\lambda_1)^2} + \mathcal{O}\left(\frac{1}{\log^{p+l}(K)}\right) (1 - \beta^{\text{tr}}\lambda_d)^2 (1 - \beta^{\text{te}}\lambda_d)^2 \right]. \quad (31)$$

To clearly illustrate the trade-off in the stopping time, we let $l = 0$ for convenience. If $R(\bar{\omega}_t^{\beta^{\text{tr}}}, \beta^{\text{te}}) < \epsilon$, we have

$$t_\epsilon \leq \exp\left(\epsilon^{-\frac{1}{p}} \left[\frac{U_l}{(1 - \beta^{\text{tr}}\lambda_1)^2} + U_t(1 - \beta^{\text{tr}}\lambda_d)^2 \right]^{\frac{1}{p}}\right)$$

where

$$U_l = \mathcal{O}\left(\left(2c_1\nu^2 + \frac{\sigma^2}{n_2}\right)(1 - \beta^{\text{te}}\lambda_1)^2\right) \quad \text{and} \quad U_t = \mathcal{O}\left(\left(2c_1\nu^2 + \frac{\sigma^2}{n_2}\right)(1 - \beta^{\text{te}}\lambda_d)^2\right).$$

The arguments are similar for the lower bound, and we can obtain:

$$L_l = \mathcal{O}\left(\left(2\frac{b_1\nu^2}{n_2} + \frac{\sigma^2}{n_2}\right)(1 - \beta^{\text{te}}\lambda_1)^2\right) \quad \text{and} \quad L_t = \mathcal{O}\left(\left(2\frac{b_1\nu^2}{n_2} + \frac{\sigma^2}{n_2}\right)(1 - \beta^{\text{te}}\lambda_d)^2\right).$$

□

F Discussions on Assumptions

Discussions on Assumption 2 If $\mathcal{P}_{\mathbf{x}}$ is Gaussian distribution, then we have

$$\mathbf{F} = \mathbb{E}[\mathbf{xx}^\top \boldsymbol{\Sigma} \mathbf{xx}^\top] = 2\boldsymbol{\Sigma}^3 + \boldsymbol{\Sigma} \text{tr}(\boldsymbol{\Sigma}^2).$$

This implies that \mathbf{F} and $\boldsymbol{\Sigma}$ commute because $\boldsymbol{\Sigma}^3$ and $\boldsymbol{\Sigma}$ commute. Moreover, in this case

$$\frac{\beta^2}{n}(\mathbf{F} - \boldsymbol{\Sigma}^3) = \frac{\beta^2}{n}(\boldsymbol{\Sigma}^3 + \boldsymbol{\Sigma} \text{tr}(\boldsymbol{\Sigma}^2)).$$

Therefore, if $n \gg \lambda_1(\lambda_1^2 + \text{tr}(\boldsymbol{\Sigma}^2))$, then the eigen-space of $\mathbf{H}_{n,\beta}$ will be dominated by $(\mathbf{I} - \beta\boldsymbol{\Sigma})^2\boldsymbol{\Sigma}$.

Discussions on Assumption 3 Assumption 3 is an eighth moment condition for $\mathbf{x} := \boldsymbol{\Sigma}^{\frac{1}{2}}\mathbf{z}$, where \mathbf{z} is a σ_x sub-Gaussian vector. Given β , for sufficiently large n s.t. $\mu_i(\mathbf{H}_{n,b}) > 0, \forall i$, and if $\text{tr}(\boldsymbol{\Sigma}^k)$ are all $O(1)$ for $k = 1, \dots, 4$, then by the quadratic form and the sub-Gaussian property, which has finite higher order moments, we can conclude that $C(\beta, \boldsymbol{\Sigma}) = \Theta(1)$.

The following lemma further shows that if $\mathcal{P}_{\mathbf{x}}$ is a Gaussian distribution, we can derive the analytical form for $C(\beta, \boldsymbol{\Sigma})$.

Lemma F.1. Given $|\beta| < \frac{1}{\lambda_1}$, for sufficiently large n s.t. $\mu_i(\mathbf{H}_{n,b}) > 0, \forall i$, and if $\mathcal{P}_{\mathbf{x}}$ is a Gaussian distribution, assuming $\boldsymbol{\Sigma}$ is diagonal, we have:

$$C(\beta, \boldsymbol{\Sigma}) = 210\left(1 + \frac{\beta^4 \text{tr}(\boldsymbol{\Sigma}^2)^2}{(1 - \beta\lambda_1)^4}\right).$$

Proof. Let $\mathbf{e}_i \in \mathbb{R}^d$ denote the vector that the i -th coordinate is 1, and all other coordinates equal 0. For $\mathbf{x} \sim \mathcal{P}_{\mathbf{x}}$, denote $\mathbf{xx}^\top = [x_{ij}]_{1 \leq i, j \leq d}$. Then we have:

$$\begin{aligned} & \mathbb{E}[\|\mathbf{e}_i^\top \mathbf{H}_{n,\beta}^{-\frac{1}{2}} (\mathbf{I} - \frac{\beta}{n} \mathbf{X}^\top \mathbf{X}) \boldsymbol{\Sigma} (\mathbf{I} - \frac{\beta}{n} \mathbf{X}^\top \mathbf{X}) \mathbf{H}_{n,\beta}^{-\frac{1}{2}} \mathbf{e}_i\|^2] \\ & \leq \mathbb{E}[\|\mathbf{e}_i^\top \mathbf{H}_{n,\beta}^{-\frac{1}{2}} (\mathbf{I} - \beta \mathbf{xx}^\top) \boldsymbol{\Sigma} (\mathbf{I} - \beta \mathbf{xx}^\top) \mathbf{H}_{n,\beta}^{-\frac{1}{2}} \mathbf{e}_i\|^2] \\ & = \mathbb{E} \left[\left(\mathbf{e}_i^\top \mathbf{H}_{n,\beta}^{-1} \mathbf{e}_i \right)^2 \left(\sum_{j \neq i} \beta^2 \lambda_j x_{ij}^2 + \lambda_i (1 - \beta x_{ii})^2 \right)^2 \right] \end{aligned}$$

For Gaussian distributions, we have

$$\mathbb{E}[x_{ij}^2 x_{ik}^2] = \begin{cases} 9\lambda_i^2 \lambda_j^2 & j = k \text{ and } \neq i \\ 105\lambda_i^4 & i = j = k \\ 3\lambda_i^2 \lambda_j \lambda_k & i \neq j \neq k \end{cases}$$

We can further obtain:

$$\begin{aligned} & \mathbb{E}[\|\mathbf{e}_i^\top \mathbf{H}_{n,\beta}^{-\frac{1}{2}} (\mathbf{I} - \frac{\beta}{n} \mathbf{X}^\top \mathbf{X}) \boldsymbol{\Sigma} (\mathbf{I} - \frac{\beta}{n} \mathbf{X}^\top \mathbf{X}) \mathbf{H}_{n,\beta}^{-\frac{1}{2}} \mathbf{e}_i\|^2] \\ & \leq 105 (\mathbf{e}_i^\top \mathbf{H}_{n,\beta}^{-1} \mathbf{e}_i)^2 \left(\sum_{j \neq i} \beta^2 \lambda_j^2 + (1 - \beta \lambda_i)^2 \right)^2 \\ & \stackrel{(a)}{\leq} 210 (\mathbf{e}_i^\top \mathbf{H}_{n,\beta}^{-1} \mathbf{e}_i)^2 [\beta^4 \text{tr}(\boldsymbol{\Sigma}^2)^2 + (1 - \beta \lambda_i)^4] \\ & \stackrel{(b)}{\leq} 210 [(\mathbf{e}_i^\top \mathbf{H}_{n,\beta}^{-1} \mathbf{e}_i)^2 \beta^4 \text{tr}(\boldsymbol{\Sigma}^2)^2 + 1] \end{aligned}$$

where (a) follows from the Cauchy-Schwarz inequality, and (b) follows the fact that $(\mathbf{e}_i^\top \mathbf{H}_{n,\beta}^{-1} \mathbf{e}_i)^2 = \frac{1}{[(1 - \beta \lambda_i) \lambda_i^2 + \frac{\beta^2}{n} (\lambda_i^2 + \text{tr}(\boldsymbol{\Sigma}^2) \lambda_i)]^2} \leq 1/(1 - \beta \lambda_i)^4$.

Therefore, for any unit $\mathbf{v} \in \mathbb{R}^d$, we have

$$\begin{aligned} & \mathbb{E}[\|\mathbf{v}^\top \mathbf{H}_{n,\beta}^{-\frac{1}{2}} (\mathbf{I} - \frac{\beta}{n} \mathbf{X}^\top \mathbf{X}) \boldsymbol{\Sigma} (\mathbf{I} - \frac{\beta}{n} \mathbf{X}^\top \mathbf{X}) \mathbf{H}_{n,\beta}^{-\frac{1}{2}} \mathbf{v}\|^2] \\ & \leq \max_{\mathbf{v}} 210 [(\mathbf{v}^\top \mathbf{H}_{n,\beta}^{-1} \mathbf{v})^2 \beta^4 \text{tr}(\boldsymbol{\Sigma}^2)^2 + 1] \leq 210 \left(1 + \frac{\beta^4 \text{tr}(\boldsymbol{\Sigma}^2)^2}{(1 - \beta \lambda_1)^4} \right). \end{aligned}$$

□

G Further Related Work

G.1 Underparameterized Setting

Provable guarantees of meta-learning have been extensively studied in the underparameterized regime, i.e. the number of tasks or the data size is much larger than the data dimension. Here we highlight some existing related studies and discuss their differences from ours.

[4] provides the generalization error bounds for S/Q meta-learners. Their generalization error bound is $O(\beta\sqrt{n}) + O(\frac{M}{\sqrt{n}})$, where the β is the uniform stability parameter, n is the number of task, and M is the uniform bound on the loss function. Note that with the square loss as considered in our setting, M can be as large as the dimension d of the input. For the overparameterized regime, where $d \gg n$, the bound becomes asymptotically large and not useful.

[9] shows the generalization guarantees of MAML on recurring and unseen tasks respectively, where the excess risk bound is roughly $\tilde{O}(\frac{G^2}{T}) + O(\frac{G^2}{mn})$, depending on the total iterations T , the number of tasks m , the available data size for each task n , and the uniform bound for gradient norm G . Notice that the gradient norm typically scales polynomially with the input dimension d . Therefore, in the overparameterized regime ($d \gg mn$), the bound again becomes vacuous.

[6] aims to theoretically characterize the performance between MAML and the standard Empirical Risk Minimization (ERM). Firstly, they show that the empirical training solutions will converge to their population-optimal values with concentration bounds, which have terms $O(\frac{\sqrt{d}}{\sqrt{n}})$ or $O(\frac{\sqrt{d}}{\sqrt{\tau}})$ where n is the sample size and τ is the number of training episodes per task. Such bounds will be crude in the overparameterized regime ($d \gg n, \tau$). Then from the generalization perspective, they only give excess risks for the population-optimal solutions, whereas we analyze the generalization property of empirical training solutions based on their optimization trajectory. Moreover, the comparison between the excess risk of MAML and ERM for the population-optimal solutions requires $m = \Omega(d)$, where m is the number of task. Therefore, such comparison does not apply to the overparameterized regime ($d \gg m$).

From what has been discussed above, all of these bounds, are useful (i.e., yield small or vanishing error) only for underparameterized regime, where sample size n is much larger than the dimension d of input (i.e., $n \gg d$), but become vacuous (i.e., yield large error bound) for overparameterized regime, where $n \ll d$. Thus, they fail to explain why the overparameterized neural network can still generalize well.

In contrast, our work provides a much more refined data-dependent upper bound by incorporating the data and task roles. Our bound can be written in a concise form as $O(\frac{h(\Sigma, \Sigma_\theta, \alpha, T)}{T})$ (see Theorem 1 in Section 4 for exact form of the bound), where the function $h(\cdot)$ is determined by data and task covariances Σ and Σ_θ , step-size α and the total number T of SGD iterations (here T scales linearly with the sample size). Note that the dimension d is implicitly captured in Σ, Σ_θ . In the overparameterized regime with $d \gg T$, the function $h(\cdot)$ explicitly captures the effect of data and task conditions to guarantee the value of $h(\cdot)$ to be small compared to T , and thus the excess risk can diminish under overparameterization. In Propositions 2 and 3, we further give specific examples about data and task covariances Σ, Σ_θ that yield small excess risk and good generalization.

Besides the above key difference, our paper also provides the following important results that were not studied in [4, 9, 6]. (a) We provide the lower bound to justify the tightness of our results in the overparameterized regime, whereas [4, 9, 6] do not have such a result. (b) Our results capture how the task diversity affects the excess risk in the overparameterized regime. In particular, we give an example (in Proposition 2), for which the excess risk exhibits a phase transition with respect to task diversity. Such an interesting behavior is not captured in [4, 9, 6].

G.2 Overparameterized Setting

Despite the overparameterization is crucial to demystify the remarkable generalization ability of deep meta-learning [21, 13], there are only a few theoretical analysis being developed to study the generalization of MAML in the overparameterized regime. [20] studied the MAML with overparameterized deep neural nets with a generalization gap quantifying the difference between the empirical and population loss functions at their optimal solutions. However, their bound is derived by conventional complexity-based techniques and does not consider the data or task-dependency similarly as [4, 9, 6] discussed in Appendix G.1, which tends to be weak in the high dimensional, especially in the overparameterized regime. Most related to our work are recent studies [2, 24], where they develop more precise generalization bounds for overparameterized setting under a mixed linear regression model. Yet, their empirical training solutions are directly calculated by taking the closed-form of training objective’s minimum, which are not obtained by trained with SGD as ours. More crucially, they consider only the simple isotropic covariance for data and tasks, which are directly scaled by d , i.e. $\Sigma = \frac{1}{d}\mathbf{I}$. Thus, they do not explicitly capture how the generalization performance of MAML depends on the data and task distributions.

H Future Directions

This work takes a step towards understanding the benefits of overparameterization for MAML from the generalization aspect. There are many important future work directions, and we elaborate some interesting directions in this section.

H.1 Generalizing to Other Learning Methods

Our result can directly extend to the random feature(RF) model, adopting the similar analysis for nonlinear model in [12]: the data \mathbf{x} is generated by $\sigma(\mathbf{W}\mathbf{z})$, where \mathbf{W} is the random feature matrix and $\sigma(\cdot)$ is the nonlinear (activation) function. Hence, the RF model can be regarded as training a two-layer neural network where the weights in the first layer are chosen randomly and then fixed and only the output layer is optimized.

Beyond the fixed feature approach, understanding the overparameterization in neural networks is much more challenging. Recently, [3] made important progress towards this aspect. They studied the benign overfitting phenomenon in training a two-layer convolutional neural network, and provided new analysis to tackle the neural network learning process. One possible direction is to generalize our analysis to neural networks by further advancing the techniques in [3].

H.2 Meta-batch Setting

We can consider to incorporate a practical meta-batch setting in our framework, i.e. given a set of tasks in advance, and at each iteration, the task is sampled from the set with replacement uniformly at random and will be visited multiple times during optimization. This is broadly referred as multi-pass SGD [16], while we have studied the single-pass setting, that each task is used only once. For multi-pass SGD, the iterate analysis in Appendix B.3 will be more challenging since we have to consider the dependence on history that leads to a complicated calculation of the expected error matrix. To handle such complication, we may adopt some analysis techniques recently developed for multi-pass SGD [23, 17, 16], and further advance them to analyze meta-batch MAML in meta learning.

Intuitively, we expect in such a setting, meta-batch size (of tasks) will play an important role in determining the excess risk. Specifically, how meta-batch size compares with the effect of data covariance and task diversity (coupled with the number of SGD iterations) will determine whether benign fitting in the overparameterized regime can occur.

H.3 Longer Inner Loops

Another important direction is to study the MAML beyond the one-step gradient update in the inner loop. Two possible cases are discussed as follows.

If the inner loop continually takes gradient updates towards the per-task loss function until converges (additionally regularized by the distance between task-specific parameter and the model parameter), then the algorithm is equivalent to another well-known meta method, iMAML [18]. Under the mixed linear regression model, similarly, we can reformulate the problem as a meta least square problem with modified meta inputs and output responses [7, 1]. Therefore, our analysis can directly generalize to such a setting with modified meta least square problem.

If the inner loop only takes a few steps of gradient updates, the analysis will be more challenging, because the update of meta parameter in the outer loop will involve Hessian, and the analysis of such an algorithm will need to handle the complicated statistical correlation between the gradient and Hessian estimators [15], where new techniques are required. Recently, [5] proposed that one can treat the Hessian as identity operator in theoretical analysis, where the corresponding algorithm is called FO-MAML, and FO-MAML typically achieves a performance similar to MAML in practice [10]. Such relaxation may make the problem more tractable. Then the remaining problem is to extend our technique for overparameterized MAML to FO-MAML, which is an interesting research direction for future study.

References

- [1] Yu Bai, Minshuo Chen, Pan Zhou, Tuo Zhao, Jason Lee, Sham Kakade, Huan Wang, and Caiming Xiong. How important is the train-validation split in meta-learning? In *International Conference on Machine Learning*, pages 543–553. PMLR, 2021.
- [2] Alberto Bernacchia. Meta-learning with negative learning rates. In *ICLR*, 2021.
- [3] Yuan Cao, Zixiang Chen, Mikhail Belkin, and Quanquan Gu. Benign overfitting in two-layer convolutional neural networks. *arXiv preprint arXiv:2202.06526*, 2022.
- [4] Jiaxin Chen, Xiao-Ming Wu, Yanke Li, Qimai Li, Li-Ming Zhan, and Fu-lai Chung. A closer look at the training strategy for modern meta-learning. *Advances in Neural Information Processing Systems*, 33:396–406, 2020.
- [5] Liam Collins, Aryan Mokhtari, Sewoong Oh, and Sanjay Shakkottai. Maml and anil provably learn representations. *arXiv preprint arXiv:2202.03483*, 2022.
- [6] Liam Collins, Aryan Mokhtari, and Sanjay Shakkottai. Why does maml outperform erm? an optimization perspective. *arXiv preprint arXiv:2010.14672*, page 6, 2020.
- [7] Giulia Denevi, Carlo Ciliberto, Dimitris Stamos, and Massimiliano Pontil. Learning to learn around a common mean. *Advances in Neural Information Processing Systems*, 31, 2018.

- [8] Aymeric Dieuleveut, Nicolas Flammarion, and Francis Bach. Harder, better, faster, stronger convergence rates for least-squares regression. *The Journal of Machine Learning Research*, 18(1):3520–3570, 2017.
- [9] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Generalization of model-agnostic meta-learning algorithms: Recurring and unseen tasks. *Advances in Neural Information Processing Systems*, 34, 2021.
- [10] Chelsea Finn and Sergey Levine. Meta-learning and universality: Deep representations and gradient descent can approximate any learning algorithm. *arXiv preprint arXiv:1710.11622*, 2017.
- [11] Katelyn Gao and Ozan Sener. Modeling and optimization trade-off in meta-learning. *Advances in Neural Information Processing Systems*, 33:11154–11165, 2020.
- [12] Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *The Annals of Statistics*, 50(2):949–986, 2022.
- [13] Timothy Hospedales, Antreas Antoniou, Paul Micaelli, and Amos Storkey. Meta-learning in neural networks: A survey. *arXiv preprint arXiv:2004.05439*, 2020.
- [14] Prateek Jain, Sham M Kakade, Rahul Kidambi, Praneeth Netrapalli, Venkata Krishna Pillutla, and Aaron Sidford. A markov chain theory approach to characterizing the minimax optimality of stochastic gradient descent (for least squares). *arXiv preprint arXiv:1710.09430*, 2017.
- [15] Kaiyi Ji, Junjie Yang, and Yingbin Liang. Multi-step model-agnostic meta-learning: Convergence and improved algorithms. 2020.
- [16] Nicole Mücke, Gergely Neu, and Lorenzo Rosasco. Beating sgd saturation with tail-averaging and minibatching. *Advances in Neural Information Processing Systems*, 32, 2019.
- [17] Loucas Pillaud-Vivien, Alessandro Rudi, and Francis Bach. Statistical optimality of stochastic gradient descent on hard learning problems through multiple passes. *Advances in Neural Information Processing Systems*, 31, 2018.
- [18] Aravind Rajeswaran, Chelsea Finn, Sham M Kakade, and Sergey Levine. Meta-learning with implicit gradients. *Advances in neural information processing systems*, 32, 2019.
- [19] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [20] Haoxiang Wang, Ruoyu Sun, and Bo Li. Global convergence and induced kernels of gradient-based meta-learning with neural nets. *arXiv preprint arXiv:2006.14606*, 2020.
- [21] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.
- [22] Difan Zou, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, and Sham Kakade. Benign overfitting of constant-stepsize sgd for linear regression. In *Conference on Learning Theory*, pages 4633–4635. PMLR, 2021.
- [23] Difan Zou, Jingfeng Wu, Vladimir Braverman, Quanquan Gu, and Sham M Kakade. Risk bounds of multi-pass sgd for least squares in the interpolation regime. *arXiv preprint arXiv:2203.03159*, 2022.
- [24] Yingtian Zou, Fusheng Liu, and Qianxiao Li. Unraveling model-agnostic meta-learning via the adaptation learning rate. In *International Conference on Learning Representations*, 2021.