
Towards sample-efficient overparameterized meta-learning

Yue Sun¹ Halil Ibrahim Gulluk² Adhyayan Narang¹ Samet Oymak³ Maryam Fazel¹

Abstract

Meta-learning typically involves two phases. First, one learns a suitable representation from the previously seen tasks. Secondly, this representation is used for learning a new task using only a few samples (i.e., few-shot learning). For both phases, ensuring sample efficient learning is critical. We first observe that existing works do not capture empirical phenomena, that is, typically, both representation learning as well as few-shot learning operate in the overparameterized regime where sample size is less than the degrees of freedom of the problem. To this aim, we study meta-learning with general task and feature covariances. For learning the representation, we investigate multiple spectral approaches showing that, spectral alignment of the feature and tasks covariances can help explain how linear representations can be learned with much fewer samples. Our analysis also leads to refined bounds for achieving optimal sample complexity for subspace-based linear representations. For few-shot learning, we derive the optimal linear representation minimizing the few-shot sample complexity. In contrast to subspace-based representations, this optimal representation can provably be high-dimensional (i.e., downstream task is overparameterized) explaining the empirical observations on few-shot learning in related literature. We then establish favorable robustness properties of this optimal representation to achieve end-to-end learning guarantees for two-phase meta-learning.

1. Introduction

In a multitude of machine learning (ML) tasks with limited data, it is crucial to build accurate models in a sample-efficient way; this goal proves to be especially challenging when the features are high-dimensional, as is often the case in modern ML.

¹University of Washington, Seattle, WA, US ²Bogazici University, Istanbul, Turkey ³University of California, Riverside, CA, US. Correspondence to: Yue Sun <yuesun@uw.edu>.

Constructing a simple yet informative representation of features can be very helpful with learning a model that generalizes well to an unseen test set. Representation learning, which dates back to (Caruana, 1997; Baxter, 2000), aims to jointly utilize information from the training data of multiple related tasks. Such information transfer is the backbone of modern transfer and multitask learning and finds ubiquitous applications in image classification (Deng et al., 2009), machine translation (Bojar et al., 2014) and reinforcement learning (Finn et al., 2017), all of which may involve numerous tasks to be learned¹ with limited data per task.

In this paper, we focus on the fundamental problem of meta-learning with linear features, which consists of two steps. In the meta-training phase, we learn a linear representation from a collection of n_{task} earlier tasks, with n_{spt} training samples per task. Let $(\mathbf{x}_{ij}, y_{ij})$ be the j th training sample from task i and we collect training samples from all tasks in the set $\mathcal{S} = \{(\mathbf{x}_{ij}, y_{ij})\} \subset \mathbb{R}^d \times \mathbb{R}$. Denote $n_{\text{tot}} = n_{\text{task}} \times n_{\text{spt}} = |\mathcal{S}|$. The goal of this stage is to output the matrix

$$\mathbf{\Lambda} := \mathbf{\Lambda}(\mathcal{S}) \in \mathbb{R}^{R \times d}, \quad (1.1)$$

that can be used to map raw input features to a linear feature map of dimension R .

In the second step (few-shot learning), we are given a new task with data $\mathcal{F} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n_{\text{fs}}}$, with possibly very small n_{fs} . We use the representation $\mathbf{\Lambda}$ to transform the input features to obtain a generalizing model $\hat{\alpha}_{\mathbf{\Lambda}}$ by minimizing the empirical risk

$$\hat{\mathcal{L}}_{\mathcal{F}}(\alpha) = \frac{1}{n_{\text{fs}}} \sum_{i=1}^{n_{\text{fs}}} \ell(y_i, \alpha^T \mathbf{\Lambda} \mathbf{x}_i). \quad (1.2)$$

Past analysis of algorithms for both the meta-training and few-shot phases have typically focused on the classical *underparameterized regime*: when the sample size is larger than the degrees of freedom (DoF) which are $\mathcal{O}(Rd)$ for (1.1) and $\mathcal{O}(R)$ for (1.2). These works are also often focused on learning simplistic representations, where $\mathbf{\Lambda}$ is simply a subspace (i.e. orthonormal rows).

However, recent literature in deep learning makes it clear that the high-dimensional (or overparameterized) regime is

¹Each task corresponds to a vector in \mathbb{R}^d and we learn the task vector given task data (features and labels). See Definitions 1 and 2 in Section 2 for details.

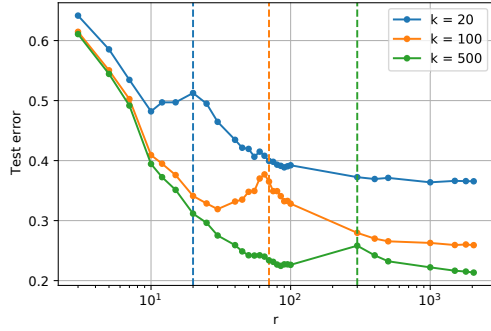


Figure 1. An illustration of the benefit of overparameterization for few-shot learning: A pretrained ResNet50 network on Imagenet (representation learning phase) was utilized to obtain a feature-representation for classification on CIFAR-10. After reducing the dimensionality to r via PCA, the resultant features are to train a logistic regression classifier on CIFAR-10 with k examples from each class (few-shot phase). The figure depicts the test error obtained using projected test features as a function of r . The interpolation threshold (i.e., when accurate fit for data happens) is indicated by the dashed line, and can be seen to coincide exactly with the location of the double-descent peak ($k = r$).

of significant interest for both problems. Deep networks are stellar representation learners despite containing many more parameters than the sample size (Finn et al., 2017). Additionally, the few-shot learning literature demonstrates that very few samples—often much fewer than the representation dimension—can be sufficient for successful adaptation to an unseen task. For example, (Finn et al., 2017) proposes a transfer learning method that adapts a pre-trained model to downstream tasks with only 1-5 new training samples. In Figure 1, the benefits of overparameterization and double-descent (also mentioned in (Mei & Montanari, 2019)) in the few-shot phase are observed in a transfer learning task for image classification; for all choices of k the accuracies are observed to be markedly higher in the overparameterized regime. This motivates the following fundamental questions.

Q1: What is the optimal representation $\mathbf{\Lambda} \in \mathbb{R}^{R \times d}$?

Q2: How to succeed with $n_{\text{tot}} \ll Rd$ and $n_{\text{fs}} \ll R$?

Contributions: Towards answering these, we make several key contributions to the finite-sample understanding of *linear* meta-learning, under assumptions discussed in Section 2. As described below, our results are enabled by studying a general data/task model with *arbitrary task covariance* Σ_{β} and *feature covariance* $\Sigma_{\mathbf{X}}$ which allows for a rich set of observations.

• **Representation learning:** We study representation learning for linear regression tasks using established moment-based algorithms (precisely defined in Section 3.1), and formalize *how feature-task interaction governs learning*.

Specializing our bounds to the case of a rank- r_t task covariance Σ_{β} and a “spiked” (approximately low-rank) feature covariance $\Sigma_{\mathbf{X}}$ reveals that sample complexity for learning the task subspace is governed by the *spectral alignment of the task & feature covariances* and can go well below the DoF $\mathcal{O}(r_t d)$, answering **Q2**. Surprisingly, if this alignment is sufficiently strong, simple PCA estimator on the empirical feature covariance can be preferable to more sophisticated methods (see Fig. 2(a)). When specialized to the case with uninformative features ($\Sigma_{\mathbf{X}} = \mathbf{I}$), we show that $\mathcal{O}(r_t d)$ samples are sufficient when $n_{\text{spt}} \gtrsim \mathcal{O}(r_t)$ achieving *optimal sample size* and improving over prior bounds of $\mathcal{O}(r_t^2 d)$ (Tripuraneni et al., 2020). Importantly, this leads to the critical finding that, for fixed total sample size n_{tot} , *both n_{spt} and n_{task} should be reasonably large* for this near-optimal performance (see Fig. 2(b) for details).

• **Few-shot learning:** With access to infinite representation learning samples, what is the *optimal representation* $\mathbf{\Lambda}_{\text{opt}}$? For least-squares regression, we answer this question by establishing an equivalence between the optimal linear representation and generalized ridge regression over the raw features (see Sec 3.2). This reveals that unlike simplistic subspace projections, the right representation for optimized few-shot performance can in fact be a high-dimensional map with $R \gg n_{\text{fs}}$ addressing **Q1** and **Q2**. In contrast to subspace-based schemes which discard certain features, $\mathbf{\Lambda}_{\text{opt}}$ shapes the feature covariance to emphasize informative features over less informative ones. In practice, $\mathbf{\Lambda}_{\text{opt}}$ is unattainable due to finite sample ($n_{\text{tot}} < \infty$) during the representation learning phase. We establish robustness guarantees under finite n_{tot} to provide an end-to-end guarantee for two phases in terms of $\mathbf{\Lambda}_{\text{opt}}$. This also uncovers a sweet spot between subspace-based approaches and optimal-shaping schemes: As n_{tot} decreases, it becomes more preferable to use a smaller R due to inaccurate estimate of $\mathbf{\Lambda}_{\text{opt}}$. This explains why, in practice, smaller dimensional representations may be preferable (see Fig. 3(b) for details).

Notation: We describe each representation learning task using vector $\beta_i \in \mathbb{R}^d$, data samples are $(\mathbf{x}_{ij}, y_{ij}) \in \mathbb{R}^d \times \mathbb{R}$, and noise $\varepsilon_{ij} \in \mathbb{R}$. The distribution of β_i and \mathbf{x}_{ij} are $\mathcal{N}(0, \Sigma_{\beta})$ and $\mathcal{N}(0, \Sigma_{\mathbf{X}})$. The rank of Σ_{β} and $\Sigma_{\mathbf{X}}$ are r_t and r_f . The few-shot learning (downstream) task is β and samples are (\mathbf{x}_i, y_i) . There are n_{task} tasks with n_{spt} samples per task in the representation learning phase, and a few-shot learning task with n_{fs} samples. Denote $n_{\text{tot}} = n_{\text{task}} n_{\text{spt}}$. Let the representation matrix be $\mathbf{\Lambda} \in \mathbb{R}^{R \times d}$. We use $a \lesssim b$ or $\tilde{\mathcal{O}}(b)$ that means there exists K that depends logarithmically on all parameters such that $a \leq Kb$.

1.1. Prior Art

A commonly studied setup is when tasks are generated from a Gaussian distribution $\mathcal{N}(0, \Sigma_{\beta})$, and Σ_{β} is exactly low

rank or contains a large portion of small eigenvalues. This is studied as mixed linear regression (Zhong et al., 2016; Li & Liang, 2018; Chen et al., 2020), which consists of multiple linear regression problems. The knowledge of the distribution of the tasks, or the subspace where the tasks lie in, enables robust and sample efficient learning for new tasks. (Lounici et al., 2011; Cavallanti et al., 2010; Maurer et al., 2016) propose sample complexity bounds of representation learning for mixed linear regression. There are study of mixed linear regression combined with other structures such as binary task vectors (Balcan et al.) and sparse task vectors (Argyriou et al., 2008).

The recent papers (Kong et al., 2020b;a; Tripuraneni et al., 2020; Du et al., 2020) propose the theoretical bounds on sample complexity and estimation error of representation learning when the tasks lie in an exactly low dimensional subspace, where the first three discuss method of moment estimators and the last two discuss matrix factorized formulations. (Tripuraneni et al., 2020) shows that when features are d dimensional and the task covariance matrix is exactly rank r_t where $r_t \ll d$, the number of samples that enable meaningful representation learning is $\mathcal{O}(dr_t^2)$. However, (Kong et al., 2020b;a; Tripuraneni et al., 2020) all assume that the features follow a standard normal distribution. We study the more general setting of arbitrary feature and task variances, and we show better sample complexity is achieved when the features and tasks are aligned. We also propose another estimator and it succeeds with $\mathcal{O}(dr_t)$ samples when $n_{\text{spt}} \sim r_t$.

After the task covariance is estimated, we will use it to learn a new task generated from the same task distribution. (Kong et al., 2020b;a; Tripuraneni et al., 2020; Du et al., 2020) assume that task covariance is exactly low-rank, so they search for \mathbf{A} in the low dimensional subspace and are able to learn with $\mathcal{O}(r_f)$ samples. We ask whether it is possible to search in the whole d dimensional space with $\mathcal{O}(r_f)$ samples (in which case we have to train in an overparameterized regime). If so, does it yield a smaller error compared to a low dimensional representation? This is meaningful especially when the task covariance Σ_β is approximately but not exactly low rank. (Bartlett et al., 2020) analyzes the generalization guarantee of overparameterized linear regression, (Chang et al., 2020) extends it to non-linear models, and (Nakkinan et al., 2020; Wu & Xu, 2020) analyze ridge regression with weighted regularizer, which also covers overparameterized linear regression. By contrast, we propose learning an (overparameterized) optimally-shaped representation that we then use in the few-shot learning phase, providing an end-to-end performance study of meta-learning.

2. Problem Setup

The meta-learning setup we consider consists of two phases: (i) representation learning, where prior tasks are used to learn a suitable representation, and (ii) few-shot learning, where a new task is learned with only a few samples. To aid in the theoretical analysis and understanding the generic behavior of the sample efficiency of general meta learning algorithms (see Fig. 1), we study linear regression tasks, and focus on a setup where the tasks and the features per task are generated randomly.

We assume that the two phases share the same distributions for features and tasks. In the first phase, there are multiple tasks, each with a batch of available data. The underlying task vectors are generated from the same distribution. We make this setup more precise using the following definitions.

Definition 1 Representation learning: tasks. Suppose the tasks are i.i.d. drawn from the distribution $\mathcal{N}(0, \Sigma_\beta)$ where $\Sigma_\beta \in \mathbf{S}_+^d$. We denote the i th task by $\beta_i \in \mathbb{R}^d$ and the number of tasks by n_{task} .

After we randomly draw tasks from the distribution, we generate data as defined below.

Definition 2 Representation learning: data per task. Fixing β_i , we generate features $\mathbf{x}_{ij} \in \mathbb{R}^d$, labels $y_{ij} \in \mathbb{R}$ for $i = 1, \dots, n_{\text{spt}}$, which follows $y_{ij} = \mathbf{x}_{ij}^\top \beta_i + \varepsilon_{ij}$. \mathbf{x}_{ij} across all tasks are generated from $\mathcal{N}(0, \Sigma_X)$. Without loss of generality, we assume Σ_X is diagonal (the task covariance Σ_β can be any positive semidefinite matrix). ε_{ij} denotes the additive noise with distribution $\mathcal{N}(0, \sigma_\varepsilon)$. For simplicity we assume Σ_X and Σ_β are scaled so that $\|\Sigma_X\|, \|\Sigma_\beta\| \geq 1$.

In Section 3 we study different estimators for the task covariance (or its low-rank approximation). We assume Σ_β is rank r_t and Σ_X is rank r_f . In general r_t, r_f can be as large as d , but we will see the benefit when Σ_β, Σ_X are approximately low rank.

Next we define the few-shot learning phase.

Definition 3 Few shot learning. In the few shot learning phase, suppose the task vector β is generated from $\mathcal{N}(0, \Sigma_\beta)$, the features $\mathbf{x}_i, i = 1, \dots, n_{\text{fs}}$ are independently generated from $\mathcal{N}(0, \Sigma_X)$, and $y_i = \mathbf{x}_i^\top \beta + \varepsilon_i$ where noise ε_i are independently generated from $\mathcal{N}(0, \sigma_\varepsilon)$.

We investigate the training error and the sample efficiency for the few-shot learning phase with respect to the representation of data. In (Kong et al., 2020b;a; Tripuraneni et al., 2020; Du et al., 2020), the task covariance is assumed to be rank r_t (i.e., representation of the features is r_t dimensional). We propose a weighted representation of features as below, which is in general R dimensional and the new task can be learned with $\mathcal{O}(r_t)$ samples even if $R \approx d$.

Definition 4 Shaped representation. We define the weighted representation of features \mathbf{x} with a positive semidefinite shaping matrix $\mathbf{\Lambda} \in \mathbb{R}^{R \times d}$ as $\mathbf{\Lambda}\mathbf{x}$.

In the next section, we will learn the downstream task vector by finding the least norm solution with the shaped representation. The optimal shaping depends on the covariance of task and feature. Since we do not know their covariance, we use the covariance estimators in representation learning phase.

3. Main Results

Our meta-learning scheme involves two phases: representation learning and few-shot learning. The Algorithm 1 frames our linear meta-learning approach. We first compute a moment estimator to obtain a representation $\mathbf{\Lambda}$ in the representation learning phase. This representation can apply a simple subspace-based dimension reduction (PCA option), as in (Kong et al., 2020b;a; Tripuraneni et al., 2020), or can be high-dimensional which allows for a *softer* refinement of raw features by *shaping* their covariance (Shaping option). We then use this feature representation in the few-shot learning phase on a new task to obtain a model $\hat{\beta}_{\mathbf{\Lambda}}$ with small risk.

Algorithm 1 Meta-learning

Require: Representation learning data: As in Def. 1, 2

Few-shot learning data: As in Def. 3

Value of R if using PCA

Representation learning: Calculate any moment estimator from Section 3.1.

if PCA then

Let $\mathbf{\Lambda} \in \mathbb{R}^{R \times d}$ be the projection to the span of top R singular values of moment estimator.

else if Shaping then

Calculate $\mathbf{\Lambda}$ as defined in Def. 5 and (3.12).

end if

Few-shot learning:

Calculate $\hat{\beta}_{\mathbf{\Lambda}}$ from (3.4).

Return: $\hat{\beta}_{\mathbf{\Lambda}}$.

In this section, we will present the analyses of Algorithm 1 in two aspects: In Section 3.1, we discuss the method of moments (MoM), which learns the task distribution from the dataset collected from multiple tasks. In Section 3.2, we propose the optimal overparameterized representation of features, and show how to use the representation for few-shot learning. Finally we integrate the two parts and give an analysis of the overall meta-learning algorithm.

3.1. Representation learning

This section investigates the representation learning via different types of moment-based spectral estimators. In gen-

eral, recovery guarantees for arbitrary covariance matrices $\Sigma_{\mathbf{X}}$ and Σ_{β} suppose we have sufficient samples. We are particularly interested in the case when Σ_{β} is approximately $r_t \ll d$ rank, so that we can estimate the moments with fewer samples. Estimating the moments is our primary goal, which can give full information about task covariance (if $\Sigma_{\mathbf{X}} = \mathbf{I}$, as we will show later). We can also obtain other key variables such as the principal eigenspace. We are further interested in the case when the feature covariance $\Sigma_{\mathbf{X}}$ is approximately low rank and its principal eigenspace covers the top eigenvectors of Σ_{β} . Based on this, we consider the following estimators. Suppose the rank of $\Sigma_{\mathbf{X}}$ is approximately r_f ; we are interested in the $r_t < r_f \ll d$ case but hope that the estimators work well with general r_f, r_t .

1. Method of moment (MoM).

$$\hat{M} = \frac{1}{n_{\text{task}}} \sum_{i=1}^{n_{\text{task}}} \frac{2}{n_{\text{spt}}^2} \left[\sum_{j=1}^{n_{\text{spt}}/2} y_{ij} y_{i(j+n_{\text{spt}}/2)} \cdot (\mathbf{x}_{ij} \mathbf{x}_{i(j+n_{\text{spt}}/2)}^{\top} + \mathbf{x}_{i(j+n_{\text{spt}}/2)} \mathbf{x}_{ij}^{\top}) \right]. \quad (3.1)$$

The mean of this estimator is $M = \Sigma_{\mathbf{X}} \Sigma_{\beta} \Sigma_{\mathbf{X}}$ and we can estimate the top r_t eigenvector with $\tilde{O}(r_f r_t^2)$ samples.

Remark: While this paper will focus on (3.1), one can also consider the symmetric version of the estimator (as in (Tripuraneni et al., 2020)) above (for which results are in similar spirit).

$$\hat{Q} = \frac{1}{n_{\text{task}}} \sum_{i=1}^{n_{\text{task}}} \frac{1}{n_{\text{spt}}} \sum_{j=1}^{n_{\text{spt}}} y_{ij}^2 \mathbf{x}_{ij} \mathbf{x}_{ij}^{\top}. \quad (3.2)$$

The mean of this estimator is $Q = 2\Sigma_{\mathbf{X}} \Sigma_{\beta} \Sigma_{\mathbf{X}} + \text{tr}(\Sigma_{\beta} \Sigma_{\mathbf{X}}) \Sigma_{\mathbf{X}}$ and its error is similar to \hat{M}^2 .

2. MoM after Task-Averaging (MoM-TA). If we further assume each task contains at least $n_{\text{spt}} = \Omega(r_t)$ samples, then we can estimate each single task by $\hat{\beta}_i = \sum_{j=1}^{n_{\text{spt}}} y_{ij} \mathbf{x}_{ij}$. Let $\hat{B} = [\hat{\beta}_1, \dots, \hat{\beta}_k]$. When $\Sigma_{\mathbf{X}} = \mathbf{I}$, we can estimate the top r_t eigenvectors of Σ_{β} with $\tilde{O}(dr_f)$ samples. This is better than (Kong et al., 2020b; Tripuraneni et al., 2020) which require $\tilde{O}(dr_t^2)$ samples. The MoM-TA estimator obeys the result of MoM estimator as we remark that $\frac{1}{n_{\text{task}}} \mathbb{E}[\hat{B} \hat{B}^{\top}] = Q$.

3. PCA on Feature Covariance (MoM-F). If $\Sigma_{\mathbf{X}}$ and Σ_{β} are approximately rank r_t and the span of top r_t eigenvectors of $\Sigma_{\mathbf{X}}$ and Σ_{β} are closely aligned, we can directly estimate the covariance of features by

$$\hat{\Sigma}_{\mathbf{X}} = \frac{1}{n_{\text{tot}}} \sum_{j=1}^{n_{\text{spt}}} \sum_{i=1}^{n_{\text{task}}} \mathbf{x}_{ij} \mathbf{x}_{ij}^{\top} \quad (3.3)$$

²In the appendix we provide the analysis for both \hat{M} and \hat{Q} .

feature cov	$\Sigma_X = I$		$\Sigma_X = \text{diag}(I_{r_f}, 0)$	
	min sample	error	min sample	error
MoM	dr_t^2	$(dr_t^2/n_{\text{tot}})^{1/2}$	$r_f r_t^2$	$(r_f r_t^2/n_{\text{tot}})^{1/2}$
MoM-TA	dr_t	$(r_t/n_{\text{spt}})^{1/2}$	$r_f r_t$	$(r_t/n_{\text{spt}})^{1/2}$
MoM-F	d	$(d/n_{\text{tot}})^{1/2}$	r_f	$(r_f/n_{\text{tot}})^{1/2}$

Table 1. Sample complexity and error of different method of moment (MoM) estimators.

The mean of this estimator is Σ_X and we can estimate the top r_t eigenvector of Σ_X with $\tilde{O}(r_t)$ samples.

Our result about the MoM estimators is summarized in Table 1.

In the following subsections, let $n_{\text{tot}} = n_{\text{spt}} n_{\text{task}}$, and let $\mathcal{S} = \max\{\|\Sigma_X\|, \|\Sigma_\beta\|\}$, $\mathcal{T}_{\beta, X} = \text{tr}(\Sigma_\beta \Sigma_X)$, $\mathcal{T}_X = \text{tr}(\Sigma_X)$, $\mathcal{T}_\beta = \text{tr}(\Sigma_\beta)$.

3.1.1. METHOD OF MOMENTS

We first study \hat{M} in (3.1). Suppose Σ_β is approximately rank r_t and Σ_X is approximately rank r_f , and $r_f \geq r_t$. If the span of top r_t eigenvectors of Σ_β is covered by the span of top r_f eigenvectors of Σ_X , then the r_t -PCA of \hat{M} gives an estimate of the span of top r_t eigenvector of Σ_β , and for this we need only an accurate estimate of \hat{M} , given by \hat{M} .

Theorem 1 Assume the representation learning dataset is generated as in Def 1, 2. Let $N_0 = (\log(n_{\text{task}}) \mathcal{T}_{\beta, X} + \sigma_\varepsilon^2) \mathcal{T}_\beta$, with probability at least $1 - (n_{\text{tot}} \mathcal{T}_\beta)^{-10} - (n_{\text{tot}} N_0)^{-10}$, the moment estimator \hat{M} above satisfies

$$\|\hat{M} - M\| \lesssim \underbrace{\sqrt{\frac{\mathcal{S} \mathcal{T}_{\beta, X}^2 \mathcal{T}_X}{n_{\text{tot}}}}}_{\mathcal{E}_{\text{mb}}(\Sigma_X, \Sigma_\beta) \text{ MoM bias}} + \underbrace{\sigma_\varepsilon \sqrt{\frac{\mathcal{S}(\mathcal{T}_{\beta, X} + \sigma_\varepsilon^2 \mathcal{T}_\beta)}{n_{\text{tot}}}}}_{\mathcal{E}_{\text{mv}}(\Sigma_X, \Sigma_\beta, \sigma_\varepsilon) \text{ MoM variance}} + \underbrace{\mathcal{S}^2 \sqrt{\frac{\mathcal{S} \mathcal{T}_\beta}{n_{\text{task}}}}}_{\mathcal{E}_{\text{tb}}(\Sigma_X, \Sigma_\beta) \text{ task bias}}$$

We define the three parts of the error as \mathcal{E}_{mb} , \mathcal{E}_{mv} , \mathcal{E}_{tb} , and define $\mathcal{E}(\Sigma_X, \Sigma_\beta, \sigma_\varepsilon)$ as the sum of \mathcal{E}_{mb} , \mathcal{E}_{mv} , \mathcal{E}_{tb} . We make clear their dependence on the covariance of tasks and features, which will be used in Theorem 4. Suppose \hat{M} is approximately rank r_t , which means $\lambda_{r_t} - \lambda_{r_t+1}$ is large, then if we do r_t PCA on \hat{M} and denote the span of top r_t eigenvectors as \hat{W} , its angle with the span of top r_t eigenvectors of M , named W , is bounded by the well known result in (Davis & Kahan, 1970). We will give a concrete example below showing the recovery error of \hat{M} and its top eigenvectors.

Example 1 Suppose $\Sigma_X = \text{diag}(I_{r_f}, \iota I_{d-r_f})$, and $\Sigma_\beta =$

$\text{diag}(I_{r_t}, 0)$, $\sigma_\varepsilon = 0$, then³

$$\|\hat{M} - M\| \lesssim \sqrt{\frac{r_t^2(r_f + \iota(d - r_f))}{n_{\text{tot}}}} + \sqrt{\frac{r_t}{n_{\text{task}}}}.$$

Suppose $n_{\text{tot}} \gtrsim r_t^2(r_f + \iota(d - r_f))$, then we have

$$\sin(\angle W, \hat{W}) \lesssim \sqrt{\frac{r_t^2(r_f + \iota(d - r_f))}{n_{\text{tot}}}}$$

When $\iota = 1$, i.e., $\Sigma_X = I$, we need $\tilde{O}(dr_t^2)$ samples to estimate \hat{M} , which is also addressed in the analysis of related literature (Kong et al., 2020b; Tripuraneni et al., 2020). However, when the feature is approximately low rank and aligns with task, we need $\tilde{O}(r_f r_t^2)$ samples to estimate \hat{M} and its top eigenvectors cover the range of tasks. This alignment of tasks and features is discovered in the related work that study the spectrum of the kernel matrices of Hessians of deep neural networks, such as (Sagun et al., 2017; Arora et al., 2019; Oymak et al., 2019).

3.1.2. TASK AVERAGE ESTIMATORS

In the previous MoM estimators, if the features are standard normal vectors, and the task covariance is rank r_t , then we need to estimate the subspace spanned by tasks with dr_t^2 samples. In this section, we will show that, if there are enough samples in each task, we are able to retrieve the subspace with dr_f samples. Let $x \sim \mathcal{N}(0, I)$, and each task contains $n_{\text{spt}} \gtrsim r_t$ samples, then we can get the span of Σ_β with $\tilde{O}(dr_t)$ samples.

Theorem 2 (Standard normal feature, noiseless) Let data be generated as in Def 1, 2. Suppose $\sigma_\varepsilon = 0$, $\Sigma_X = I$, and suppose the rank of Σ_β is r_t . Define $\hat{\beta}_j = n_{\text{spt}}^{-1} \sum_{i=1}^{n_{\text{spt}}} y_{ij} x_{ij}$, $B = [\beta_1, \dots, \beta_{n_{\text{task}}}]$, and $\hat{B} = [\hat{\beta}_1, \dots, \hat{\beta}_{n_{\text{task}}}]$. Let $n_{\text{spt}} > c_1 \mathcal{T}_\beta \lambda_{r_t}^{-1}(\Sigma_\beta)$, $n_{\text{task}} > c_2 \max\{d, \frac{\mathcal{S} \mathcal{T}_\beta}{\lambda_{r_t}^2(\Sigma_\beta)}\}$, with probability $1 - (n_{\text{task}}^{-c_3} + (n_{\text{task}} \mathcal{T}_\beta)^{-c_4} + \exp(-c_5 n_{\text{task}}^2))$, where c_i are constants, $c_{3,4} > 10$,

$$\sigma_{\max}(\hat{B} - B) \lesssim \sqrt{\frac{n_{\text{task}} \mathcal{T}_\beta}{n_{\text{spt}}}}.$$

³The second term can also be written as $\sqrt{r_t n_{\text{spt}} / n_{\text{tot}}}$. Typically $r_t \gtrsim n_{\text{spt}}$ so the first term is larger.

Denote the span of top r_t singular column vectors of $\hat{\mathbf{B}}$ and Σ_β as $\hat{\mathbf{W}}, \mathbf{W}$, then

$$\sin(\angle \hat{\mathbf{W}}, \mathbf{W}) \lesssim \sqrt{\frac{\mathcal{T}_\beta}{n_{\text{spt}} \lambda_{r_t}(\Sigma_\beta)}}.$$

If $\Sigma_\beta = \text{diag}(\mathbf{I}_{r_t}, 0)$, then $\sin(\angle \hat{\mathbf{W}}, \mathbf{W}) \lesssim \sqrt{r_t/n_{\text{spt}}}$.

In the appendix, we will propose a theorem with general feature covariance Σ_X and noisy data, as a generalization of Theorem 2. When Σ_X and Σ_β are approximately rank r_f and r_t , we need $n_{\text{tot}} \gtrsim r_f r_t$.

Remark 1 Theorem 2 requires $n_{\text{spt}} > c_1 \mathcal{T}_\beta \lambda_{r_t}^{-1}(\Sigma_\beta)$, so n_{spt} is lower bounded by $\mathcal{O}(r_t)$ in Theorem 2. If $n_{\text{spt}} = 1$ and we estimate $\hat{\mathbf{B}}$, then $n_{\text{task}} \hat{\mathbf{Q}} = \hat{\mathbf{B}} \hat{\mathbf{B}}^\top$. This means Theorem 1 can be applied to this estimator when $n_{\text{spt}} = 1$.

3.1.3. ESTIMATING THE COVARIANCE OF FEATURES

As we have defined in Def 2, features x_{ij} are generated from $\mathcal{N}(0, \Sigma_X)$. We aim to estimate the covariance Σ_X . Although there are different kinds of algorithms, such as maximum likelihood estimator (Anderson et al., 1970), to be consistent with the algorithms in the latter sections, we study the sample covariance matrix defined by (3.3).

Lemma 1 Suppose x_{ij} are generated independently from $\mathcal{N}(0, \Sigma_X)$. We estimate (3.3), then when $n_{\text{tot}} \gtrsim \mathcal{T}_X$, with probability at least $1 - (n_{\text{tot}} \text{tr}(\Sigma_X))^{-10}$,

$$\|\hat{\Sigma}_X - \Sigma_X\| \lesssim \sqrt{\frac{\mathcal{T}_X}{n_{\text{tot}}}}$$

Denote the span of top r_f eigenvectors of Σ_X as \mathbf{W} and the span of top r_f eigenvectors of $\hat{\Sigma}_X$ as $\hat{\mathbf{W}}$. Let $\delta_\lambda = \lambda_{r_f}(\Sigma_X) - \lambda_{r_f+1}(\Sigma_X)$. Then if $n_{\text{tot}} \gtrsim \frac{\mathcal{T}_X}{\delta_\lambda^2}$, we have

$$\sin(\angle \mathbf{W}, \hat{\mathbf{W}}) \lesssim \sqrt{\frac{\mathcal{T}_X}{n_{\text{tot}} \delta_\lambda^2}}$$

Example 2 When $\Sigma_X = \text{diag}(\mathbf{I}_{r_f}, 0)$, we have $\sin(\angle \mathbf{W}, \hat{\mathbf{W}}) \lesssim \sqrt{\frac{r_f}{n_{\text{tot}}}}$.

Lemma 1 gives the quality of the estimation of the covariance of features x . When the condition number of the matrix Σ_X is close to 1, we need $n_{\text{tot}} \gtrsim d$ to get an estimation with error $\mathcal{O}(1)$. However, when the matrix Σ_X is close to rank r_f , the amount of samples to achieve the same error is smaller, and we can use $n_{\text{tot}} \gtrsim r_f$ samples to get $\mathcal{O}(1)$ estimation error.

3.2. Few-shot learning

In few-shot learning phase, we hope to learn the feature $\beta \in \mathbb{R}^d$ which is generated from $\mathcal{N}(0, \Sigma_\beta)$. The data is $(x_i, y_i)_{i=1, \dots, n_{\text{fs}}}$ where $x_i \sim \mathcal{N}(0, \Sigma_X)$, and $y_i = x_i^\top \beta +$

ε_i where $\varepsilon_i \sim \mathcal{N}(0, \sigma_\varepsilon)$. We assume $n_{\text{fs}} < d$. Denote $\mathbf{X} \in \mathbb{R}^{n_{\text{fs}} \times d}$ whose i th row is x_i , and $\mathbf{y} = [y_1, \dots, y_m]^\top$. Suppose we select a shaping matrix $\Lambda \in \mathbb{R}^{d \times d}$, we are interested in the least norm solution defined as

$$\hat{\alpha}_\Lambda = \arg \min_{\alpha'} \|\alpha'\|_{\ell_2} \text{ s.t. } \mathbf{y} = \mathbf{X} \Lambda \alpha' \quad (3.4a)$$

$$\hat{\beta}_\Lambda = \Lambda \alpha_\Lambda = \Lambda (\mathbf{X} \Lambda)^\dagger \mathbf{y}. \quad (3.4b)$$

The (excess) risk of $\hat{\beta}_\Lambda$ is given by

$$\text{risk}(\Lambda, \Sigma_\beta) = \mathbb{E}_{x, y, \beta} (y - x^\top \hat{\beta}_\Lambda)^2 \quad (3.5)$$

$$= \mathbb{E}_\beta (\hat{\beta}_\Lambda - \beta)^\top \Sigma_X (\hat{\beta}_\Lambda - \beta) + \sigma_\varepsilon^2. \quad (3.6)$$

We want to solve for the optimal representation with Σ_β as

$$\Lambda^* = \arg \min_{\Lambda' \in \mathcal{S}_{++}^d} \text{risk}(\Lambda', \Sigma_\beta) \quad (3.7)$$

Define $\Lambda = \arg \min_{\Lambda' \in \mathcal{S}_{++}^d} \text{risk}(\Lambda', \hat{\Sigma}_\beta)$.

First we observe that, the shaping in Def 4 is a special case of the weighted ridge regression discussed in (Wu & Xu, 2020). We observe the following equivalence of these two descriptions.

Observation 1 Let $\mathbf{X} \in \mathbb{R}^{n_{\text{fs}} \times d}$ and $\mathbf{y} \in \mathbb{R}^{n_{\text{fs}}}$, and define

$$\hat{\beta}_1 = \Lambda (\mathbf{X} \Lambda)^\dagger \mathbf{y}, \quad (3.8)$$

$$\hat{\beta}_2 = \lim_{t \rightarrow 0} \arg \min_{\beta} \|\mathbf{X}^\top \beta - \mathbf{y}\|^2 + t \beta^\top \Lambda^{-2} \beta, \quad (3.9)$$

then $\hat{\beta}_1 = \hat{\beta}_2$.

In Section 3.2.1, we will derive an expression for the risk with an arbitrary representation Λ . Note we can only use $\hat{\Sigma}_\beta$, not Σ_β to obtain Λ . In Section 3.2.2, we bound the sensitivity of risk in $\hat{\Sigma}_\beta - \Sigma_\beta$. Finally we obtain an end to end guarantee of the risk of the whole meta-learning algorithm, including representation learning and few-shot learning phases.

3.2.1. COMPUTING OPTIMAL REPRESENTATION

Prior work typically (Kong et al., 2020b;a; Tripuraneni et al., 2020) projects the features onto the subspace for few shot learning. In Q1 we ask, what can be said about the performance of a general linear representation with arbitrary dimension? For a given Λ , the following theorem characterizes the exact asymptotic risk that helps us differentiate the performance of different linear representations.

In the following discussion, we assume that $\Sigma_X, \Sigma_\beta, \hat{\Sigma}_\beta$ share the same eigenspace, and suppose they are diagonal. In this case, we search over diagonal representation matrix Λ . Next, we will propose an expression for computing the representation matrix Λ as a function of $\hat{\Sigma}_\beta$.

Definition 5 (Precise few-shot risk) When $d > n_{\text{fs}}$ there exists unique $\xi > 0$ such that

$$n_{\text{fs}} = \sum_{i=1}^d (1 + (\xi \Sigma_{\mathbf{X}_i})^{-1})^{-1}, \quad (3.10)$$

Define $\boldsymbol{\theta} \in \mathbb{R}^d$ to be $\boldsymbol{\theta}_i = \frac{\xi \Lambda_i^2}{1 + \xi \Lambda_i^2}$, and define the risk as

$$\text{risk}(\Lambda, \Sigma_\beta) = \frac{1}{n_{\text{fs}} - \|\boldsymbol{\theta}\|^2} \left(\frac{n_{\text{fs}}}{d} \sum_{i=1}^d (1 - \boldsymbol{\theta}_i)^2 \Sigma_{\beta_i} + \|\boldsymbol{\theta}\|^2 \sigma_\varepsilon^2 \right). \quad (3.11)$$

We denote the right hand side as $f(\boldsymbol{\theta}; \Sigma_\beta)$.

Recall that we defined the few-shot risk as in (3.5). In the appendix, we derive the asymptotical risk applying the CGMT in (Thrapoulidis et al., 2015), and show that (3.5) and (3.11) are asymptotically equivalent when n_{fs}/d is fixed and $n_{\text{fs}}, d \rightarrow \infty$.

Finding optimal representation. Definition 5 grants us access to a closed-form risk for any linear representation. Thus, one can solve for the optimization representation by minimizing this risk using the parameterization $\boldsymbol{\theta}$.⁴

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} f(\boldsymbol{\theta}; \Sigma_\beta), \text{ s.t. } 0 \leq \boldsymbol{\theta} < 1, \sum_{i=1}^d \boldsymbol{\theta}_i = n_{\text{fs}}$$

Recalling $\boldsymbol{\theta}_i = \frac{\xi \Lambda_i^2}{1 + \xi \Lambda_i^2}$, we then find the *optimal representation* via the reverse map $\Lambda_i^* = ((1/\boldsymbol{\theta}_i^* - 1)\xi)^{-2}$.

Ensuring robustness. We use $\hat{\Sigma}_\beta$ instead of Σ_β for computing the optimal representation Λ , thus we need the risk to be robust with respect to $\hat{\Sigma}_\beta - \Sigma_\beta$. Let $0 < \underline{\theta} < n_{\text{fs}}/d$, in implementation we require $\boldsymbol{\theta} \leq \boldsymbol{\theta} \leq 1 - \frac{d - n_{\text{fs}}}{n_{\text{fs}}} \underline{\theta}$ instead of $0 \leq \boldsymbol{\theta} \leq 1$ for robustness concerns. The sensitivity in $\underline{\theta}$ is shown in Theorem 3.

Generally, the optimal representation is a d dimensional matrix and it is not guaranteed to be low rank. When $n_{\text{fs}} < d$, the downstream problem is overparameterized. The overparameterization is often used for complicated models such as neural networks, and we justified it via the linear model. The overparameterized linear regression is also studied in (Wu & Xu, 2020; Bartlett et al., 2020), and we propose the algorithm for arbitrary feature and task covariance.

We give the algorithm computing the optimal shaping with arbitrary dimension R below.

R dimensional optimal representation: Let $\Sigma_\beta = U S U^\top$ be the eigendecomposition of Σ_β . Let $U_1 \in \mathbb{R}^{d \times R}$ be the first R columns block of U and $U_2 \in \mathbb{R}^{d \times (d-R)}$ be the remaining $d - R$ columns block. Let $\mathbf{X}_R = \mathbf{X} U_1 \in \mathbb{R}^{n_{\text{fs}} \times R}$. We again define a shaping matrix $\Lambda \in \mathbb{R}^{R \times R}$ as in 3.4. Λ_i 's are diagonal values of Λ . $\hat{\beta}_R = U_1^\top \hat{\beta}$, where $\hat{\beta}$ is square roots of diagonal values of $\hat{\Sigma}_\beta$, and $\hat{\beta}_{Ri}$ denotes its i th entry. Σ_{Ri} is i th diagonal element of Σ_R . We define the following quantities.

⁴The optimization problem is not convex. In appendix we provide an algorithm solving it in polynomial time following the proof of this theorem.

$$\Sigma_R = U_1^\top \Sigma_{\mathbf{X}} U_1, \quad \Sigma_{\beta_R} = U_1^\top \Sigma_\beta U_1 \quad (3.12a)$$

$$\Lambda_R = \arg \min_{\Lambda' \in \mathbb{R}^{R \times R}} \text{risk}(\Lambda', \Sigma_{\beta_R}) \quad (3.12b)$$

$$\hat{\alpha}_{\Lambda_R} = \arg \min_{\alpha' \in \mathbb{R}^R} \|\alpha'\|_{\ell_2} \text{ s.t. } \mathbf{y} = \mathbf{X}_R \Lambda_R \alpha' \quad (3.12c)$$

$$\hat{\beta}_{\Lambda_R} = \Lambda_R \hat{\alpha}_{\Lambda_R} = \Lambda_R (\mathbf{X}_R \Lambda_R)^\dagger \mathbf{y} \quad (3.12d)$$

Define $\Sigma_\beta^\perp, \Sigma_{\mathbf{X}}^\perp$ as the projection of $\Sigma_\beta, \Sigma_{\mathbf{X}}$ onto U_2 , the noise variance is equivalent to $\sigma_{\varepsilon_R}^2 = \sigma_\varepsilon^2 + \text{tr}(\Sigma_\beta^\perp \Sigma_{\mathbf{X}}^\perp)$. Note we define R dimensional optimal representation as if we know the covariance matrix Σ_β , when we have only $\hat{\Sigma}_\beta$, these definitions can be applied in the same way⁵. We will discuss how truncation level R affects risk in Remark 2.

3.2.2. ROBUSTNESS OF OPTIMAL REPRESENTATION

In this part we focus on the identity feature case $\Sigma_{\mathbf{X}} = \mathbf{I}$, and study the robustness of few-shot learning risk with respect to inaccuracy of representation learning.

In Section 3.1, suppose one uses the estimator⁶ $\hat{\Sigma}_\beta = \hat{M}$ to estimate the task covariance Σ_β , and suffers the error $\|\hat{\Sigma}_\beta - \Sigma_\beta\| \leq \delta_{\Sigma_\beta}$. In few-shot learning phase, one uses $\hat{\Sigma}_\beta$ as the nominal task covariance and get Λ . With the true task distribution, the risk is $\text{risk}(\Lambda, \Sigma_\beta)$. We are interested in its difference from the optimal risk $\text{risk}(\Lambda, \Sigma_\beta) - \text{risk}(\Lambda^*, \Sigma_\beta)$. Note that (Wu & Xu, 2020) Sec.6 gives the exact value of $\text{risk}(\Lambda^*, \Sigma_\beta)$ so we have an end to end error guarantee.

Theorem 3 Suppose the data is generated as Definition 3, Λ and $\underline{\theta}$ are defined in Def. 5 and the estimated task is obtained as (3.4). Suppose $\|\hat{\Sigma}_\beta - \Sigma_\beta\| \leq \delta_{\Sigma_\beta}$. Then the risk of few-shot learning phase suffers at most

$$\text{risk}(\Lambda, \Sigma_\beta) - \text{risk}(\Lambda^*, \Sigma_\beta) \leq \frac{2n_{\text{fs}}^2 \delta_{\Sigma_\beta}}{(d - n_{\text{fs}})(2n_{\text{fs}} - d\underline{\theta})\underline{\theta}}.$$

Theorem 3 shows the robustness of few-shot learning algorithm with respect to the error in representation learning phase. Given Theorem 1 that bounds δ_{Σ_β} , we will propose the bound of Algorithm 1 in the following theorem.

Theorem 4 Suppose $\Sigma_{\mathbf{X}} = \mathbf{I}$. We run Algorithm 1 and set $\hat{\Sigma}_\beta = \hat{M}$. The optimal shaping matrix Λ_R depends on $\hat{\Sigma}_\beta$. Let Σ_{β_R} be the projection of Σ_β onto the top R eigenvector space, $\Sigma_{\beta_R}^\perp = \Sigma_\beta - \Sigma_{\beta_R}$. The asymptotic risk of few-shot learning is upper bounded by

⁵In the definition below, $\Lambda_R \in \mathbb{R}^{R \times R}$, the shaping matrix as defined in Algorithm 1 is $\Lambda = \Lambda_R U_1^\top \in \mathbb{R}^{R \times d}$

⁶In essence we require $\hat{M} - \Sigma_\beta$ being small to compute Λ accurately. If $\Sigma_{\mathbf{X}} \neq \mathbf{I}$, logically $\hat{M} \neq \Sigma_\beta$ and this estimation does not always work. In some special cases, such as $\Sigma_{\mathbf{X}} = \text{diag}(\mathbf{I}_{r_f}, \iota \mathbf{I}_{d-r_f})$ and $\Sigma_\beta = \text{diag}(\mathbf{I}_{r_t}, 0)$ ($r_f > r_t$), \hat{M} is still a good estimator of Σ_β .

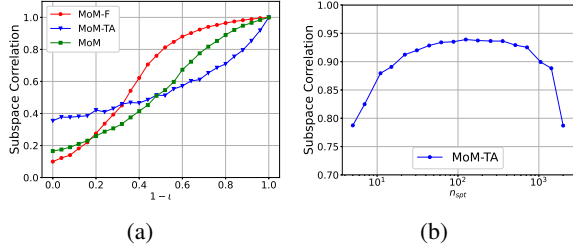


Figure 2. (a) $\Sigma_\beta = (\mathbf{I}_{10}, \mathbf{0}_{90})$. $\sigma_\varepsilon = 0.5$, $\Sigma_X = (\mathbf{I}_{10}, \iota \cdot \mathbf{I}_{90})$, $n_{\text{task}} = 20$, $n_{\text{spt}} = 40$. Learning 10 dimensional top eigenspace of Σ_β with biased features. MoM-F, MoM-TA, MoM stands for Σ_X , $\hat{\mathbf{B}}$, $\hat{\mathbf{Q}}$. Subspace correlation is defined as $\|\hat{\mathbf{U}}^\top \mathbf{U}\|^2 / \|\mathbf{U}\|^2$. When $\iota \rightarrow 0$ the feature and task are more aligned. (b) $\Sigma_\beta = (\mathbf{I}_{10}, \mathbf{0}_{90})$, $\Sigma_X = \mathbf{I}_{100}$, $\sigma_\varepsilon = 0.5$, learning with $n_{\text{tot}} = 20000$ samples with varying n_{task} . High correlation happens when n_{task} and n_{spt} are reasonably large.

$$\text{risk}(\Lambda_R, \Sigma_\beta) - \text{risk}(\Lambda_R^*, \Sigma_\beta) \quad (3.13)$$

$$\lesssim \frac{n_{\text{fs}}^2 \cdot \mathcal{E}(\Sigma_X, \Sigma_{\beta_R}, \sigma_\varepsilon^2 + \text{tr}(\Sigma_X \Sigma_{\beta_R}^\perp))}{(R - n_{\text{fs}})(2n_{\text{fs}} - R\theta)\theta} \quad (3.14)$$

If we do not apply dimension reduction (i.e., $R = d$), then

$$\text{risk}(\Lambda, \Sigma_\beta) \lesssim \text{risk}(\Lambda^*, \Sigma_\beta) + \frac{n_{\text{fs}}^2 \cdot \mathcal{E}(\Sigma_X, \Sigma_\beta, \sigma_\varepsilon)}{(d - n_{\text{fs}})(2n_{\text{fs}} - d\theta)\theta}$$

Remark 2 (Risk with respect to PCA level R) We compare the behavior of (3.14) with different truncation levels R . $\text{risk}(\Lambda_R^*, \Sigma_\beta)$ decrease with R . \mathcal{E} is an increasing function with respect to R and noise. With R increasing, the third term $\sigma_\varepsilon^2 + \text{tr}(\Sigma_X \Sigma_{\beta_R}^\perp)$ decreases. Thus generally the whole end to end error $\text{risk}(\Lambda_R, \Sigma_\beta)$ might not be monotone in R when $R > n_{\text{fs}}$. This is depicted in Fig. 3(b). In essence, this result provides a theoretical justification on the existence of a sweet-spot for the optimal representation strategy. With infinite n_{tot} ($\hat{\Sigma}_\beta = \Sigma_\beta$), it is safe to learn a large representation. As n_{tot} decreases, it becomes difficult to estimate task covariance accurately, especially its small eigenvalues due to finite-sample estimation noise. Thus a large dimensional representation –that uses this noisy covariance– may result in noisier features leading to the excess risk term \mathcal{E} . Thus choosing R adaptively with n_{tot} can strike the right bias-variance tradeoff between the excess risk (variance) and the risk due to suboptimal representation Λ_R^* i.e. $\text{risk}(\Lambda_R^*, \Sigma_\beta) - \text{risk}(\Lambda^*, \Sigma_\beta)$.

4. Numerical Experiments

In this section, we will verify by experiments the three main contributions proposed before: (1) We apply the method of moment estimator for retrieving the covariance of the task or the top eigenspace. We show that when feature space aligns with the task space, we only need $\mathcal{O}(r_f r_t^2)$ instead

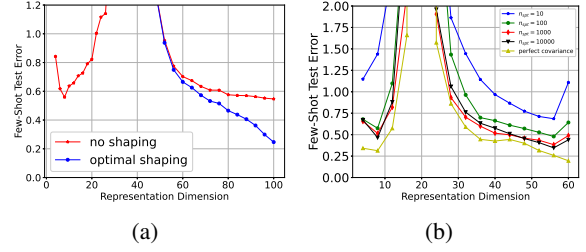


Figure 3. (a) $\Sigma_\beta = (25 \cdot \mathbf{I}_{10}, \mathbf{I}_{90})$, $\Sigma_X = \mathbf{I}_{100}$, $\sigma_\varepsilon = 0.5$, $n_{\text{fs}} = 40$, few shot learning with knowledge of Σ_β , with identity or optimal shaping matrix. For $R < n_{\text{fs}}$, we are not in overparameterized regime, so there is no optimal shaping. (b) Algorithm 1 with $d = 60$, $n_{\text{fs}} = 20$, $n_{\text{task}} = 60$, $\Sigma_\beta = (25 \cdot \mathbf{I}_6, \mathbf{I}_{54})$, $\Sigma_X = \mathbf{I}_{60}$ and varying n_{spt} , R . End to end risk with different representation learning samples and dimensions. Optimal shaping (with respect to $\hat{\Sigma}_\beta$) applied.

of $\mathcal{O}(dr_t^2)$ samples. (2) In Theorem 2, we introduce task average estimator and argue that when each task contains $\Omega(r_t)$ samples, we can retrieve the subspace of tasks with $\mathcal{O}(r_f r_t)$ samples. (3) In Def. 5, we present the optimal representation matrix that minimizes the excess risk and in Theorem 4 we show the overall risk of Algorithm 1.

• **Error of representation learning with respect to task feature alignment.** We depict this in Fig. 2(a). After applying method of moments estimators we do r_t -SVD truncation. We can see from Fig. 2(a) that, MoM’s learn the subspaces accurately when $1 - \iota$ is large, (i.e., with feature-task alignment), whereas behave worse when $\Sigma_X = \mathbf{I}$. Specifically, The MoM-F estimator works better only when ι is small, i.e., high task-feature alignment.

• **Task average helps learning with fewer samples.** We show the sample efficiency of the task average estimator in Fig. 2(b). We fix the total number of samples n_{tot} , and vary the number of samples per task n_{spt} . In Theorem 2, we argue that it estimates the task subspace with $\mathcal{O}(dr_t)$ samples when $n_{\text{spt}} \geq r_t$. We can see from the figure that, when n_{spt} is small, this estimator degenerate to $\hat{\mathbf{Q}}$ which asks for $\mathcal{O}(dr_t^2)$ samples, so the subspace estimation is bad. When n_{task} is small, the sampled task features cannot approximate the true task distribution, which also cause large error. We need n_{spt} and n_{task} both reasonably large.

• **Optimal representation.** In Fig. 3(a), we know Σ_β and do R -SVD truncation and apply optimal shaping for varying R . When the problem becomes overparameterized ($n_{\text{fs}} < R$), the estimation of downstream task depends on the shaping matrix. The risk corresponding to the optimal shaping matrix is smaller than using raw features for solving the downstream task.

• **End to end behavior of meta-learning algorithm.** We run the whole meta-learning process as in Algorithm 1. The

normalized risk is plotted in Fig. 3(b). We use the optimal shaping matrix \mathbf{A} , calculated as a function of $\hat{\Sigma}_\beta$, for solving the downstream task. Both dimension reduction ($n_{fs} > R$) or optimal overparameterized problem ($n_{fs} < R$) end up with small risk, whereas it becomes hard to learn when $n_{fs} \approx R$. The risk of overparameterized case can be smaller than using low dimensional representation (typically applied in (Tripuraneni et al., 2020; Kong et al., 2020b;a; Du et al., 2020)). As discussed in Remark 2, because we cannot learn Σ_β perfectly in Algorithm 1, the smallest risk with finite data happens when the chosen representation dimension $R < d$, unlike the case with known Σ_β .

5. Conclusion

In this paper, we study the sample efficiency of meta-learning with linear representations, motivated by the wide application of meta learning with overparameterized neural networks. On a theoretical level, we show that in representation learning, one can learn the representation with even fewer data if the features follow a spiked distribution and align with task parameters. We also propose an estimator that learns with optimal sample size, improving over prior works. Then we propose the optimal shaped representation which is typically overparameterized. We analyze its robustness while optimal representation is unattainable due to finite samples. Finally we propose an end to end learning guarantee of the overall meta-learning procedures.

References

- Anderson, T. W. et al. Estimation of covariance matrices which are linear combinations or whose inverses are linear combinations of given matrices. *Essays in probability and statistics*, pp. 1–24, 1970.
- Argyriou, A., Evgeniou, T., and Pontil, M. Convex multi-task feature learning. *Machine learning*, 73(3):243–272, 2008.
- Arora, S., Du, S., Hu, W., Li, Z., and Wang, R. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pp. 322–332, 2019.
- Balcan, M.-F., Blum, A., and Vempala, S. Efficient representations for lifelong learning and autoencoding. In *Conference on Learning Theory*.
- Bartlett, P. L., Long, P. M., Lugosi, G., and Tsigler, A. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 2020.
- Baxter, J. A model of inductive bias learning. *Journal of artificial intelligence research*, 12:149–198, 2000.
- Bojar, O., Buck, C., Federmann, C., Haddow, B., Koehn, P., Leveling, J., Monz, C., Pecina, P., Post, M., Saint-Amand, H., et al. Findings of the 2014 workshop on statistical machine translation. In *Proceedings of the ninth workshop on statistical machine translation*, pp. 12–58, 2014.
- Caruana, R. Multitask learning. *Machine learning*, 28(1): 41–75, 1997.
- Cavallanti, G., Cesa-Bianchi, N., and Gentile, C. Linear algorithms for online multitask classification. *The Journal of Machine Learning Research*, 11:2901–2934, 2010.
- Chang, X., Li, Y., Oymak, S., and Thrampoulidis, C. Provable benefits of overparameterization in model compression: From double descent to pruning neural networks. *arXiv preprint arXiv:2012.08749*, 2020.
- Chen, S., Li, J., and Song, Z. Learning mixtures of linear regressions in subexponential time via fourier moments. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pp. 587–600, 2020.
- Davis, C. and Kahan, W. M. The rotation of eigenvectors by a perturbation. iii. *SIAM Journal on Numerical Analysis*, 7(1):1–46, 1970.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Du, S. S., Hu, W., Kakade, S. M., Lee, J. D., and Lei, Q. Few-shot learning via learning the representation, provably. *arXiv preprint arXiv:2002.09434*, 2020.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *International Conference on Machine Learning*, pp. 1126–1135, 2017.
- Kong, W., Somani, R., Kakade, S., and Oh, S. Robust meta-learning for mixed linear regression with small batches. *arXiv preprint arXiv:2006.09702*, 2020a.
- Kong, W., Somani, R., Song, Z., Kakade, S., and Oh, S. Meta-learning for mixed linear regression. *arXiv preprint arXiv:2002.08936*, 2020b.
- Li, Y. and Liang, Y. Learning mixtures of linear regressions with nearly optimal complexity. In *Conference On Learning Theory*, pp. 1125–1144, 2018.
- Lounici, K., Pontil, M., Van De Geer, S., Tsybakov, A. B., et al. Oracle inequalities and optimal inference under group sparsity. *The annals of statistics*, 39(4):2164–2204, 2011.

- Maurer, A., Pontil, M., and Romera-Paredes, B. The benefit of multitask representation learning. *The Journal of Machine Learning Research*, 17(1):2853–2884, 2016.
- Mei, S. and Montanari, A. The generalization error of random features regression: Precise asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355*, 2019.
- Nakkiran, P., Venkat, P., Kakade, S., and Ma, T. Optimal regularization can mitigate double descent. *arXiv preprint arXiv:2003.01897*, 2020.
- Oymak, S., Fabian, Z., Li, M., and Soltanolkotabi, M. Generalization guarantees for neural networks via harnessing the low-rank structure of the jacobian. *arXiv preprint arXiv:1906.05392*, 2019.
- Sagun, L., Evci, U., Guney, V. U., Dauphin, Y., and Bottou, L. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.
- Thrampoulidis, C., Abbasi, E., and Hassibi, B. Lasso with non-linear measurements is equivalent to one with linear measurements. *Advances in Neural Information Processing Systems*, 28:3420–3428, 2015.
- Tripuraneni, N., Jin, C., and Jordan, M. I. Provable meta-learning of linear representations. *arXiv preprint arXiv:2002.11684*, 2020.
- Vershynin, R. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- Wu, D. and Xu, J. On the optimal weighted ℓ_2 regularization in overparameterized linear regression, 2020.
- Zhong, K., Jain, P., and Dhillon, I. S. Mixed linear regression with multiple components. In *Advances in neural information processing systems*, pp. 2190–2198, 2016.

Organization of the appendix

The appendix consists of the proof of our main results including the following parts:

- Representation learning. Sec. A includes the proof for the result about representation learning in Sec. 3.1, including the sample complexity and error guarantee of the three method of moment (MoM) estimators: MoM, MoM-TA, MoM-F.
 - We first analyze MoM-F (feature covariance estimator) in Sec. A.1 which is the most straightforward.
 - we extend the Bernstein type technique for analyzing \hat{M} in Sec. A.2. In Sec. A.3 we analyze \hat{Q} that has similar behavior. We analyze them with general task and feature covariance.
 - Finally we analyze MoM-TA in Sec. A.4. Suppose task covariance is rank r_t . With the assumption that each task has $\Omega(r_t)$ corresponding samples, the sample complexity is **reduced by a factor of r_t compared to MoM**, which meets the *information theoretical lower bound* in (Tripuraneni et al., 2020).
- Optimal representation. The proof for optimal overparameterized representation is in Sec. B. We show that we can use an arbitrary R dimensional representation of feature for few-shot learning, and it can behave better than typical PCA (low dimensional/underparameterized) representation when task is *approximately* low rank.
 - In Sec. B.1 and B.2 we provide the asymptotic analysis of optimal shaping. By asymptotic we refer to the regime where $n_{fs}, d \rightarrow \infty$ and the eigenvalues of task and feature covariance matrices converge to a fixed distribution. We show that $\hat{\beta}_\Lambda$ converges to a Gaussian distribution parameterized by Λ , and use it to express the risk.
 - We extend the asymptotic case (infinite dimensional) to the non-asymptotic (finite dimensional) regime in Sec. B.3. We define the risk function with respect to representation matrix Λ , and solve for the optimal representation by minimizing risk.
 - We prove the robustness of the optimal representation in Sec. B.4, which leads to the overall error guarantee of the proposed meta-learning algorithm.

A. Analysis of MoM estimators

A.1. Covariance estimator

In this part we will first prove the lemma showing the accuracy of covariance estimator, since it is the simplest among the three MoM estimators. We will use Bernstein type concentration results to bound its error, and a similar technique will be used for \hat{M}, \hat{Q} in the next sections.

Lemma 1 Suppose $\mathbf{x}_i, i = 1, \dots, n_{\text{tot}}$ are generated independently from $\mathcal{N}(0, \Sigma_X)$. We estimate (3.3), then when $n_{\text{tot}} \gtrsim \mathcal{T}_X$,

$$\|\hat{\Sigma}_X - \Sigma_X\| \lesssim \sqrt{\frac{\|\Sigma_X\| \text{tr}(\Sigma_X)}{n_{\text{tot}}}}.$$

Denote the span of top r_f eigenvectors of Σ_X as \mathbf{W} and the span of top r_f eigenvectors of $\hat{\Sigma}_X$ as $\hat{\mathbf{W}}$. Let $\delta_\lambda = \lambda_{r_f}(\Sigma_X) - \lambda_{r_f+1}(\Sigma_X)$. Then if $n_{\text{tot}} \gtrsim \frac{\|\Sigma_X\| \text{tr}(\Sigma_X)}{\delta_\lambda^2}$, we have

$$\sin(\angle \mathbf{W}, \hat{\mathbf{W}}) \lesssim \sqrt{\frac{\|\Sigma_X\| \text{tr}(\Sigma_X)}{n_{\text{tot}} \delta_\lambda^2}}$$

Proof First we observe that, the features \mathbf{x}_{ij} (Def. 2) among different tasks are generated i.i.d. from $\mathcal{N}(0, \Sigma_X)$. So we can rewrite (3.3) as

$$\hat{\Sigma}_X = \frac{1}{n_{\text{tot}}} \sum_{i=1}^{n_{\text{tot}}} \mathbf{x}_i \mathbf{x}_i^\top \quad (\text{A.1})$$

where $\mathbf{x}_i \sim \mathcal{N}(0, \Sigma_X)$. The error of $\hat{\Sigma}_X$ depends on n_{tot} regardless of n_{task} and n_{spt} respectively.

First, we know by concentration inequality

$$P(\|\mathbf{x}\mathbf{x}^\top\| - \text{tr}(\Sigma_{\mathbf{X}}) \geq t) = P(\|\mathbf{x}\|^2 - \text{tr}(\Sigma_{\mathbf{X}}) \geq t) \leq \exp(-c \min\{\frac{t^2}{\text{tr}(\Sigma_{\mathbf{X}}^2)}, \frac{t}{\|\Sigma_{\mathbf{X}}\|}\}). \quad (\text{A.2})$$

We will use the fact $\sqrt{\text{tr}(\Sigma_{\mathbf{X}}^2)} \leq \text{tr}(\Sigma_{\mathbf{X}})$. Define $K = C_0 \log(n_{\text{tot}} \text{tr}(\Sigma_{\mathbf{X}})) \text{tr}(\Sigma_{\mathbf{X}})$, $\mathbf{Z} = \mathbf{x}\mathbf{x}^\top$, $\mathbf{Z}' = \mathbf{Z} \cdot \mathbf{1}\{\|\mathbf{Z}\| \leq K\}$ where $\mathbf{1}$ means indicator function ($\mathbf{1}(\text{True}) = 1$, $\mathbf{1}(\text{False}) = 0$), for some positive number C_0 . Then

$$\|\mathbf{E}(\mathbf{Z} - \mathbf{Z}')\| \leq \int_{t=K}^{\infty} (\exp(-c \frac{t^2}{\text{tr}^2(\Sigma_{\mathbf{X}})}) + \exp(-c \frac{t}{\|\Sigma_{\mathbf{X}}\|})) dt \quad (\text{A.3})$$

$$\leq \int_{t=K}^{\infty} (\exp(-c \frac{t}{\text{tr}(\Sigma_{\mathbf{X}})}) + \exp(-c \frac{t}{\|\Sigma_{\mathbf{X}}\|})) dt \quad (\text{A.4})$$

$$\leq 2 \frac{\text{tr}(\Sigma_{\mathbf{X}})}{c} \exp(-c \frac{K}{\text{tr}(\Sigma_{\mathbf{X}})}) \quad (\text{A.5})$$

$$\leq \frac{\sqrt{K \text{tr}^2(\Sigma_{\mathbf{X}})}}{c} \exp(-\frac{cK}{\text{tr}(\Sigma_{\mathbf{X}})}) \quad (\text{A.6})$$

$$\lesssim (n_{\text{tot}} \text{tr}(\Sigma_{\mathbf{X}}))^{-C} \quad (\text{A.7})$$

where C is another constant. Then we compute $(\mathbf{x}\mathbf{x}^\top)^2 = \|\mathbf{x}\|^2 \mathbf{x}\mathbf{x}^\top$. Let $\Sigma_{\mathbf{X}}$ be diagonal. Let $\mathbf{x} = \sqrt{\Sigma_{\mathbf{X}}} \mathbf{z}$. Then

$$\mathbf{E}(\|\mathbf{x}\|^2 \mathbf{x}\mathbf{x}^\top)_{ij} = \begin{cases} \Sigma_{\mathbf{X}ii}(\text{tr}(\Sigma_{\mathbf{X}}) + 2\Sigma_{\mathbf{X}ii}), & i = j, \\ 0, & i \neq j. \end{cases}$$

So $\|\mathbf{E}(\|\mathbf{x}\|^2 \mathbf{x}\mathbf{x}^\top)\| \leq \|\Sigma_{\mathbf{X}}\|(\text{tr}(\Sigma_{\mathbf{X}}) + 2\|\Sigma_{\mathbf{X}}\|) \approx \|\Sigma_{\mathbf{X}}\| \text{tr}(\Sigma_{\mathbf{X}})$.

Using (Tripuraneni et al., 2020) Lemma 29, we get with probability $1 - \mathcal{O}((n_{\text{tot}} \text{tr}(\Sigma_{\mathbf{X}}))^{-C})$,

$$\|\hat{\Sigma}_{\mathbf{X}} - \Sigma_{\mathbf{X}}\| \lesssim \log(n_{\text{tot}} \text{tr}(\Sigma_{\mathbf{X}})) \left(\frac{\log(n_{\text{tot}} \text{tr}(\Sigma_{\mathbf{X}})) \text{tr}(\Sigma_{\mathbf{X}})}{n_{\text{tot}}} + \sqrt{\frac{\|\Sigma_{\mathbf{X}}\| \text{tr}(\Sigma_{\mathbf{X}})}{n_{\text{tot}}}} \right) \quad (\text{A.8})$$

If the number above is smaller than $\lambda_{r_t} - \lambda_{r_t+1}$, we have that

$$n_{\text{tot}} \gtrsim \frac{\|\Sigma_{\mathbf{X}}\| \text{tr}(\Sigma_{\mathbf{X}})}{(\lambda_{r_t} - \lambda_{r_t+1})^2} \quad (\text{A.9})$$

which is $\mathcal{O}(r_t)$ if condition number is 1.

The bound of the angle of top R eigenvector subspace is a direct application of the following lemma.

Lemma 2 (Davis & Kahan, 1970) Let \mathbf{A} be a square matrix. Let $\hat{\mathbf{W}}$, \mathbf{W} denote the span of top r_t singular vectors of $\hat{\mathbf{A}}$ and \mathbf{A} . Suppose $\|\hat{\mathbf{A}} - \mathbf{A}\| \leq \Delta$, and $\sigma_r(\mathbf{A}) - \sigma_{r+1}(\mathbf{A}) \geq \Delta$, then

$$\sin(\angle \mathbf{W}, \hat{\mathbf{W}}) \leq \frac{\Delta}{\sigma_r(\mathbf{A}) - \sigma_{r+1}(\mathbf{A}) - \Delta}.$$

So that the error of principle subspace recovery of feature covariance is upper bounded by $\frac{\|\hat{\Sigma}_{\mathbf{X}} - \Sigma_{\mathbf{X}}\|}{\sigma_{r_f}(\Sigma_{\mathbf{X}}) - \sigma_{r_f+1}(\Sigma_{\mathbf{X}}) - \|\hat{\Sigma}_{\mathbf{X}} - \Sigma_{\mathbf{X}}\|}$,

where $\|\hat{\Sigma}_{\mathbf{X}} - \Sigma_{\mathbf{X}}\|$ is calculated in (A.8). \blacksquare

A.2. Method of moment

We first present the following Bernstein type concentration lemma, also applied in (Tripuraneni et al., 2020):

Lemma 3 Let $\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}$. Choose T_0, σ^2 such that

1. $P(\|Z\| \geq C_0 T_0 + t) \leq \exp(-c\sqrt{t/T_0})$.
2. $\|E(ZZ^\top)\|, \|E(Z^\top Z)\| \leq \sigma^2$.

Then with probability at least $1 - (nT_0)^{-c}$, $c > 10$,

$$\left\| \frac{1}{n} \sum_{i=1}^n Z_i - E(Z_i) \right\| \lesssim \log(nT_0) \left(\frac{T_0 \log(nT_0)}{n} + \frac{\sigma}{\sqrt{n}} \right).$$

Proof Define $K = \log^2(C_K n T_0)$ for $C_K > 0$, $Z' = Z \mathbf{1}(\|Z\| \leq K T_0)$, then

$$\|E(Z - Z')\| \leq \int_{KT_0}^{\infty} \exp(-c\sqrt{t/T_0}) dt \lesssim (1 + \sqrt{K}) \exp(-c\sqrt{K}) T_0 \lesssim (1 + \log(C_K n T_0)) (nT_0)^{-C}. \quad (\text{A.10})$$

We can choose C_K large enough so that $C > 10$. We will use (Tripuraneni et al., 2020) Lemma 29. Set $R = \log^2(C_K n T_0) T_0 + C_0 T_0$, $\Delta = (1 + \log(C_K n T_0)) (nT_0)^{-C}$, $t = C_t \log(nT_0) \left(\frac{T_0 \log(nT_0)}{n} + \frac{\sigma}{\sqrt{n}} \right)$ for some $C_t > 0$, plugging in the last inequality of (Tripuraneni et al., 2020) Lemma 29, the LHS is smaller than $(nT_0)^{-c}$ for some c . We can also check $P(\|Z\| \geq R) \leq (nT_0)^{-c}$ for some c , thus applying (Tripuraneni et al., 2020) Lemma 29 we prove our lemma. ■

Now we will prove Theorem 1 by proving the almost same theorem below.

Theorem 5 We generate tasks $\beta_i \in \mathbb{R}^d \sim \mathcal{N}(0, \Sigma_\beta)$, $i = 1, \dots, n_{\text{task}}$. We generate features and labels with each task by (x_{ij}, y_{ij}) where $y_{ij} = x_{ij}^\top \beta_i + \varepsilon_j$, where $x_{ij} \sim \mathcal{N}(0, \Sigma_X)$, $\varepsilon_j \sim \mathcal{N}(0, \sigma_\varepsilon)$. There are n_{task} tasks and each task has n_{spt} data, let $n_{\text{tot}} = n_{\text{spt}} n_{\text{task}}$. Define $\bar{\Sigma}_\beta = \frac{1}{n_{\text{task}}} \sum_{j=1}^{n_{\text{task}}} \beta_j \beta_j^\top$, $\bar{M} = \Sigma_X \bar{\Sigma}_\beta \Sigma_X$, $\bar{Q} = \Sigma_X \bar{\Sigma}_\beta \Sigma_X + \text{tr}(\bar{\Sigma}_\beta \Sigma_X) \Sigma_X$. The moment estimator above satisfies

$$\begin{aligned} \|\hat{M} - \bar{M}\| &\lesssim \sqrt{\frac{\text{tr}(\Sigma_X) \text{tr}^2(\Sigma_\beta \Sigma_X) \|\Sigma_X\| + (\text{tr}(\Sigma_\beta \Sigma_X) + \sigma_\varepsilon^2 \text{tr}(\Sigma_X)) \cdot \sigma_\varepsilon^2 \|\Sigma_X\|}{n_{\text{tot}}}} \\ &\quad + \frac{\text{tr}^2(\Sigma_X) \text{tr}(\Sigma_\beta \Sigma_X) + (\sigma_\varepsilon + \sqrt{\text{tr}(\Sigma_\beta)}) \sigma_\varepsilon \text{tr}(\Sigma_X)}{n_{\text{tot}}}. \end{aligned}$$

Proof We define $\bar{\Sigma}_\beta = \frac{1}{n_{\text{task}}} \sum_{i=1}^{n_{\text{task}}} \beta_i \beta_i^\top$. The proof contains these steps:

1. First, we show that the norm of $\|\beta_i\|$, $\|\beta_i\|_{\Sigma_X}$, $\|\beta_i\|_{\Sigma_X^2}$, $\forall i = 1, \dots, n_{\text{task}}$ are all around their expected value.
2. For any $i \in \{1, \dots, n_{\text{task}}\}$, we study the covariance of

$$\frac{1}{2} y_{ij} y_{i(j+n_{\text{spt}}/2)} \cdot (x_{ij} x_{i(j+n_{\text{spt}}/2)}^\top + x_{i(j+n_{\text{spt}}/2)} x_{ij}^\top) - \Sigma_X \beta_i \beta_i^\top \Sigma_X$$

Then use Bernstein inequality to analyze the average of them as compared to $\bar{\Sigma}_\beta$.

3. Finally we bound the distance of $\bar{\Sigma}_\beta$ and Σ_β via the same technique in Section A.1.

We first study the property of the tasks $\beta_1, \dots, \beta_{n_{\text{task}}}$. We know that, for any $\beta \sim \mathcal{N}(0, \Sigma_\beta)$,

$$P(\|\beta\|^2 - \text{tr}(\Sigma_\beta) \geq t) \leq \exp(-c \min\{\frac{t^2}{\text{tr}(\Sigma_\beta^2)}, \frac{t}{\|\Sigma_\beta\|}\}). \quad (\text{A.11})$$

So that with probability at least $1 - \delta$, we have

$$\|\beta_i\|^2 \lesssim \text{tr}(\Sigma_\beta) + \sqrt{(\log(1/\delta) + \log(n_{\text{task}})) \text{tr}(\Sigma_\beta^2)} + (\log(1/\delta) + \log(n_{\text{task}})) \|\Sigma_\beta\| \quad (\text{A.12})$$

$$\lesssim \text{tr}(\Sigma_\beta) + \log(n_{\text{task}}/\delta) \sqrt{\text{tr}(\Sigma_\beta^2)}, \quad \forall i = 1, \dots, n_{\text{task}}. \quad (\text{A.13})$$

With similar technique we know that with probability at least $1 - \delta$,

$$\|\Sigma_X \beta_i\|^2 \lesssim \text{tr}(\Sigma_X \Sigma_\beta \Sigma_X) + \log(n_{\text{task}}/\delta) \sqrt{\text{tr}((\Sigma_X \Sigma_\beta \Sigma_X)^2)}, \quad \forall i = 1, \dots, n_{\text{task}}. \quad (\text{A.14})$$

$$\|\sqrt{\Sigma_X} \beta_i\|^2 \lesssim \text{tr}(\sqrt{\Sigma_X} \Sigma_\beta \sqrt{\Sigma_X}) + \log(n_{\text{task}}/\delta) \sqrt{\text{tr}((\sqrt{\Sigma_X} \Sigma_\beta \sqrt{\Sigma_X})^2)}, \quad \forall i = 1, \dots, n_{\text{task}}. \quad (\text{A.15})$$

We will choose $\delta = n_{\text{task}}^{-c}$ for some constant c so that $\log(n_{\text{task}}/\delta) = (c+1) \log(n_{\text{task}}) \approx \log(n_{\text{task}})$.

Then we focus on each term $y_{ij} y_{i(j+n_{\text{spt}}/2)} \mathbf{x}_{ij} \mathbf{x}_{i(j+n_{\text{spt}}/2)}^\top$ or $y_{ij} y_{i(j+n_{\text{spt}}/2)} \mathbf{x}_{i(j+n_{\text{spt}}/2)} \mathbf{x}_{ij}^\top$ and mainly bound their variance. For simplicity of notation, denote a pair of data as (\mathbf{x}, y) , (\mathbf{x}', y') and we use $yy' \mathbf{x}(\mathbf{x}')^\top$ to simplify either $y_{ij} y_{i(j+n_{\text{spt}}/2)} \mathbf{x}_{ij} \mathbf{x}_{i(j+n_{\text{spt}}/2)}^\top$ or $y_{ij} y_{i(j+n_{\text{spt}}/2)} \mathbf{x}_{i(j+n_{\text{spt}}/2)} \mathbf{x}_{ij}^\top$. Then we study the quantity $yy' \mathbf{x}(\mathbf{x}')^\top$. We study a fixed i thus we simplify β_i as β in our proof. In the following paragraphs, we denote $y = \mathbf{x}^\top \beta + \varepsilon$ and $y' = (\mathbf{x}')^\top \beta + \varepsilon'$, and take the expectation only over $\mathbf{x}, \mathbf{x}', \varepsilon, \varepsilon'$. That means we treat β as fixed vectors. In Sec. A.2.2 we will further use the bounds about β above.

First, we compute the expectation of \hat{M} .

$$E(y\mathbf{x}) = E(\mathbf{x}\mathbf{x}^\top \beta) = \Sigma_X \beta.$$

So the expectation of M is

$$E_x(\hat{M}) = E_x(yy' \mathbf{x}(\mathbf{x}')^\top) = [E(y\mathbf{x})][E(y\mathbf{x})]^\top = \Sigma_X \beta \beta^\top \Sigma_X. \quad (\text{A.16})$$

Next we use the Bernstein type inequality to bound $\hat{M} - E_x(\hat{M})$.

A.2.0.1 Noiseless Denote $\hat{y} = \beta^\top \mathbf{x}$, we study the estimator $\hat{y}\mathbf{x}$ and $\hat{y}\hat{y}' \mathbf{x}(\mathbf{x}')^\top$. This is equivalent to the noiseless case $\sigma_\varepsilon = 0$.

First we bound $\|\hat{y}\mathbf{x}\|$. With probability $1 - \delta$, $\|\mathbf{x}\|$ is bounded by (A.2), and $|\hat{y}| = |\beta^\top \mathbf{x}| \lesssim \|\sqrt{\Sigma_X} \beta\| (1 + \sqrt{\log(1/\delta)})$. Then with probability $1 - \delta$

$$\|\hat{y}\hat{y}' \mathbf{x}(\mathbf{x}')^\top\| \leq (1 + \log^2(1/\delta)) \text{tr}^2(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|^2.$$

In other words,

$$P(\|\hat{y}\hat{y}' \mathbf{x}(\mathbf{x}')^\top\| > \text{tr}^2(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|^2 + t) \leq \exp(-c \frac{\sqrt{t}}{\text{tr}(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|}) \quad (\text{A.17})$$

Next we need to compute the expectation of $(\hat{y}\hat{y}' \mathbf{x}(\mathbf{x}')^\top)(\hat{y}\hat{y}' \mathbf{x}' \mathbf{x}^\top)$.

Denote $\sqrt{\Sigma_X} \beta = \mathbf{b}$, let $\mathbf{x} = \sqrt{\Sigma_X} \mathbf{z}$ so that \mathbf{z} is standard normal. Denote any vector $\bar{\mathbf{v}} \in \mathbb{R}^d$, and let $\mathbf{v} = \sqrt{\Sigma_X} \bar{\mathbf{v}}$.

$$E_x(\bar{\mathbf{v}}^\top (\hat{y}\hat{y}' \mathbf{x}(\mathbf{x}')^\top) (\hat{y}\hat{y}' \mathbf{x}' \mathbf{x}^\top) \bar{\mathbf{v}}) \quad (\text{A.18})$$

$$= E_x[(\beta^\top \mathbf{x}_1)^2 (\beta^\top \mathbf{x}_2)^2 (\mathbf{x}_1^\top \bar{\mathbf{v}} \bar{\mathbf{v}}^\top \mathbf{x}_1) (\mathbf{x}_2^\top \mathbf{x}_2)]. \quad (\text{A.19})$$

$$= E_x[(\mathbf{b}^\top \mathbf{z})^2 (\mathbf{z}^\top \mathbf{v} \mathbf{v}^\top \mathbf{z})] E_x[(\mathbf{b}^\top \mathbf{z})^2 \mathbf{z}^\top \Sigma_X \mathbf{z}]. \quad (\text{A.20})$$

We compute $E_x[(\mathbf{b}^\top \mathbf{z})^2 (\mathbf{z}^\top \mathbf{v} \mathbf{v}^\top \mathbf{z})]$ and $E_x[(\mathbf{b}^\top \mathbf{z})^2 \mathbf{z}^\top \Sigma_X \mathbf{z}]$ separately.

$$E_x[(\mathbf{b}^\top \mathbf{z})^2 (\mathbf{z}^\top \mathbf{v} \mathbf{v}^\top \mathbf{z})] = \|\mathbf{b}\|^2 \|\mathbf{v}\|^2 + 2 \sum_{i=1}^d b_i^2 v_i^2. \quad (\text{A.21})$$

To bound $E_x[(\mathbf{b}^\top \mathbf{z})^2 \mathbf{z}^\top \Sigma_X \mathbf{z}]$, we denote $\Sigma_X = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$, then

$$E_x[(\mathbf{b}^\top \mathbf{z})^2 \mathbf{z}^\top \Sigma_X \mathbf{z}] = \sum_{i=1}^d \lambda_i (\|\mathbf{b}\|^2 + 2 \sum_{j=1}^d b_j^2 u_{ij}^2) \quad (\text{A.22})$$

$$\leq 3 \text{tr}(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|^2. \quad (\text{A.23})$$

So that $\mathbf{E}_x(\hat{y}\hat{y}'\mathbf{x}(\mathbf{x}')^\top) = \Sigma_X\beta\beta^\top\Sigma_X$, and

$$\mathbf{E}_x(\bar{\mathbf{v}}^\top(\hat{y}\hat{y}'\mathbf{x}(\mathbf{x}')^\top)(\hat{y}\hat{y}'\mathbf{x}(\mathbf{x}')^\top)^\top\bar{\mathbf{v}}) \quad (\text{A.24})$$

$$= (\|\sqrt{\Sigma_X}\beta\|^2\|\sqrt{\Sigma_X}\bar{\mathbf{v}}\|^2 + 2\sum_{i=1}^d(\sqrt{\Sigma_X}\beta)_i^2(\sqrt{\Sigma_X}\bar{\mathbf{v}})_i^2) \cdot (\sum_{i=1}^d\lambda_i(\|\sqrt{\Sigma_X}\beta\|^2\|\mathbf{u}_i\|^2 + 2\sum_{j=1}^d(\sqrt{\Sigma_X}\beta)_j^2\mathbf{u}_{ij}^2)). \quad (\text{A.25})$$

The right hand side is upper bounded as

$$\mathbf{E}_x(\|(\hat{y}\hat{y}'\mathbf{x}(\mathbf{x}')^\top)(\hat{y}\hat{y}'\mathbf{x}(\mathbf{x}')^\top)^\top\|) \lesssim \text{tr}(\Sigma_X)\|\sqrt{\Sigma_X}\beta\|^4\|\Sigma_X\| \quad (\text{A.26})$$

A.2.1. THE VARIANCE CAUSED BY NOISE.

Suppose $y = \mathbf{x}^\top\beta + \varepsilon$ where ε is a noise with variance σ_ε^2 . Then we have

$$\mathbf{E}_x(y y' \mathbf{x}(\mathbf{x}')^\top) = \mathbf{E}_x((\hat{y} + \varepsilon)(\hat{y}' + \varepsilon')\mathbf{x}(\mathbf{x}')^\top) \quad (\text{A.27})$$

$$= \mathbf{E}_x(\hat{y}\hat{y}'\mathbf{x}(\mathbf{x}')^\top) + 2\mathbf{E}_x(\varepsilon(\mathbf{x}^\top\beta)\mathbf{x}'\mathbf{x}^\top) + \mathbf{E}_x(\varepsilon\varepsilon'\mathbf{x}(\mathbf{x}')^\top). \quad (\text{A.28})$$

And with simple computation, we have the following theorem (compare to Lemma 3)

Lemma 4 *Let $c > 0$ be a constant. The following inequalities hold:*

$$P(\|\varepsilon(\mathbf{x}^\top\beta)\mathbf{x}'\mathbf{x}^\top\| \geq \sigma_\varepsilon\|\sqrt{\Sigma_X}\beta\|\text{tr}(\Sigma_X) + t) \leq \exp(-c\sqrt{\frac{t}{\sigma_\varepsilon\|\beta\|\text{tr}(\Sigma_X)}}), \quad (\text{A.29a})$$

$$\|\mathbf{E}_x(\varepsilon(\mathbf{x}^\top\beta)\mathbf{x}'\mathbf{x}^\top)(\varepsilon(\mathbf{x}^\top\beta)\mathbf{x}(\mathbf{x}')^\top)\| \lesssim \sigma_\varepsilon^2\text{tr}(\Sigma_X)\|\Sigma_X\|\|\sqrt{\Sigma_X}\beta\|^2, \quad (\text{A.29b})$$

$$P(\|\varepsilon\varepsilon'\mathbf{x}'\mathbf{x}^\top\| \geq \sigma_\varepsilon^2\text{tr}(\Sigma_X) + t) \leq \exp(-c\sqrt{\frac{t}{\sigma_\varepsilon^2\text{tr}(\Sigma_X)}}), \quad (\text{A.29c})$$

$$\|\mathbf{E}_x((\varepsilon\varepsilon'\mathbf{x}'\mathbf{x}^\top)(\varepsilon\varepsilon'\mathbf{x}(\mathbf{x}')^\top))\| \lesssim \sigma_\varepsilon^4\text{tr}(\Sigma_X)\|\Sigma_X\|. \quad (\text{A.29d})$$

Proof of Lemma 4. By concentration of Gaussian random vecotr, we know that $|\varepsilon| \leq \sigma_\varepsilon(1 + \sqrt{\log(1/\delta)})$, $|\mathbf{x}^\top\beta| \leq \mathcal{O}(\|\sqrt{\Sigma_X}\beta\|(1 + \sqrt{\log(1/\delta)}))$ and $\|\mathbf{x}\|^2 \leq \mathcal{O}(\text{tr}(\Sigma_X) + \text{tr}(\Sigma_X)\log(1/\delta))$ all with probability $1 - \delta$, so that we get (A.29a) and (A.29c).

For (A.29b),

$$\|\mathbf{E}_x(\varepsilon(\mathbf{x}^\top\beta)\mathbf{x}'\mathbf{x}^\top)(\varepsilon(\mathbf{x}^\top\beta)\mathbf{x}(\mathbf{x}')^\top)\| = \|\mathbf{E}_x\sigma_\varepsilon^2(\mathbf{x}^\top\beta)^2\mathbf{x}^\top\mathbf{x}\Sigma_X\| \quad (\text{A.30})$$

Here we can apply (A.23). At the end we get

$$\|\mathbf{E}_x(\varepsilon(\mathbf{x}^\top\beta)\mathbf{x}'\mathbf{x}^\top)(\varepsilon(\mathbf{x}^\top\beta)\mathbf{x}(\mathbf{x}')^\top)\| \lesssim \sigma_\varepsilon^2\text{tr}(\Sigma_X)\|\Sigma_X\|\|\sqrt{\Sigma_X}\beta\|^2 \quad (\text{A.31})$$

For (A.29d),

$$\|\mathbf{E}_x(\varepsilon\varepsilon'\mathbf{x}'\mathbf{x}^\top)(\varepsilon\varepsilon'\mathbf{x}(\mathbf{x}')^\top)\| = \|\mathbf{E}_x\sigma_\varepsilon^4(\mathbf{x}^\top\mathbf{x})\mathbf{x}'(\mathbf{x}')^\top\| \quad (\text{A.32})$$

which easily ends up with (A.29d).

A.2.2. OVERALL ERROR OF MOM

We will assemble the useful bounds of noiseless case and noise together and apply Lemma 3. We look at the data in the j th batch. From (A.16)

$$\mathbf{E}_x(\frac{1}{2}y_{ij}y_{i(j+n_{\text{spt}}/2)}(\mathbf{x}_{ij}\mathbf{x}_{i(j+n_{\text{spt}}/2)}^\top + \mathbf{x}_{i(j+n_{\text{spt}}/2)}\mathbf{x}_{ij}^\top) - \Sigma_X\beta_i\beta_i^\top\Sigma_X) = 0$$

From (A.17) and (A.29),

$$P(\|\frac{1}{2}y_{ij}y_{i(j+n_{\text{spt}}/2)}(\mathbf{x}_{ij}\mathbf{x}_{i(j+n_{\text{spt}}/2)}^\top + \mathbf{x}_{i(j+n_{\text{spt}}/2)}\mathbf{x}_{ij}^\top) - \Sigma_{\mathbf{X}}\beta_i\beta_i^\top\Sigma_{\mathbf{X}}\| \quad (\text{A.33})$$

$$\geq C \log(n_{\text{task}})(\sqrt{\text{tr}((\Sigma_{\mathbf{X}}\Sigma_{\beta}\Sigma_{\mathbf{X}})^2)} + \text{tr}^2(\Sigma_{\mathbf{X}})\sqrt{\text{tr}((\sqrt{\Sigma_{\mathbf{X}}}\Sigma_{\beta}\sqrt{\Sigma_{\mathbf{X}}})^2)}) + \sigma_{\epsilon}^2\text{tr}(\Sigma_{\mathbf{X}}) + t) \quad (\text{A.34})$$

$$\leq \exp\left(-c \sqrt{\frac{t}{\max\left\{\log(n_{\text{task}})\text{tr}(\Sigma_{\mathbf{X}})\sqrt{\text{tr}((\sqrt{\Sigma_{\mathbf{X}}}\Sigma_{\beta}\sqrt{\Sigma_{\mathbf{X}}})^2)}, \sigma_{\epsilon}^2\text{tr}(\Sigma_{\mathbf{X}})\right\}}}\right) \quad (\text{A.35})$$

Note that

$$\begin{aligned} & \|\frac{1}{2}y_{ij}y_{i(j+n_{\text{spt}}/2)}(\mathbf{x}_{ij}\mathbf{x}_{i(j+n_{\text{spt}}/2)}^\top + \mathbf{x}_{i(j+n_{\text{spt}}/2)}\mathbf{x}_{ij}^\top) - \Sigma_{\mathbf{X}}\beta_i\beta_i^\top\Sigma_{\mathbf{X}}\| \\ & \leq \|\frac{1}{2}y_{ij}y_{i(j+n_{\text{spt}}/2)}(\mathbf{x}_{ij}\mathbf{x}_{i(j+n_{\text{spt}}/2)}^\top + \mathbf{x}_{i(j+n_{\text{spt}}/2)}\mathbf{x}_{ij}^\top)\| + \|\Sigma_{\mathbf{X}}\beta_i\beta_i^\top\Sigma_{\mathbf{X}}\| \\ & \lesssim \|\frac{1}{2}y_{ij}y_{i(j+n_{\text{spt}}/2)}(\mathbf{x}_{ij}\mathbf{x}_{i(j+n_{\text{spt}}/2)}^\top + \mathbf{x}_{i(j+n_{\text{spt}}/2)}\mathbf{x}_{ij}^\top)\| + \log^2(n_{\text{task}})\text{tr}((\Sigma_{\mathbf{X}}\Sigma_{\beta}\Sigma_{\mathbf{X}})^2) \end{aligned}$$

Summing up the upper bound of (A.26) and (A.29) in (A.28)

$$\|\mathbf{E}_{\mathbf{x}}(\frac{1}{2}y_{ij}y_{i(j+n_{\text{spt}}/2)}(\mathbf{x}_{ij}\mathbf{x}_{i(j+n_{\text{spt}}/2)}^\top + \mathbf{x}_{i(j+n_{\text{spt}}/2)}\mathbf{x}_{ij}^\top))^2\| \quad (\text{A.36})$$

$$\lesssim \log^2(n_{\text{task}})\text{tr}(\Sigma_{\mathbf{X}})\text{tr}^2((\sqrt{\Sigma_{\mathbf{X}}}\Sigma_{\beta}\sqrt{\Sigma_{\mathbf{X}}})^2)\|\Sigma_{\mathbf{X}}\| + \sigma_{\epsilon}^4\text{tr}(\Sigma_{\mathbf{X}})\|\Sigma_{\mathbf{X}}\|. \quad (\text{A.37})$$

Thus we can denote

$$T_0 = (\log(n_{\text{task}})\sqrt{\text{tr}((\Sigma_{\beta}\Sigma_{\mathbf{X}})^2)} + \sigma_{\epsilon}^2\text{tr}(\Sigma_{\mathbf{X}}), \quad (\text{A.38a})$$

$$\sigma = (\log^2(n_{\text{task}})\text{tr}(\Sigma_{\mathbf{X}})\text{tr}^2((\Sigma_{\beta}\Sigma_{\mathbf{X}})^2)\|\Sigma_{\mathbf{X}}\| + \sigma_{\epsilon}^4\text{tr}(\Sigma_{\mathbf{X}})\|\Sigma_{\mathbf{X}}\|)^{1/2}, \quad (\text{A.38b})$$

and apply Lemma 3 by plugging in T_0 and σ ,

$$\log^{-1}(n_{\text{tot}}T_0)\|\frac{1}{2n}\sum_{i=1}^{n_{\text{spt}}}(y_i y'_i(\mathbf{x}_i(\mathbf{x}'_i)^\top + \mathbf{x}'_i\mathbf{x}_i^\top)) - \frac{1}{n_{\text{task}}}\Sigma_{\mathbf{X}}(\sum_{i=1}^{n_{\text{task}}}\beta_i\beta_i^\top)\Sigma_{\mathbf{X}}\| \leq \mathcal{O}(\frac{T_0 \log(n_{\text{tot}}T_0)}{n_{\text{spt}}} + \frac{\sigma}{\sqrt{n_{\text{spt}}}}). \quad (\text{A.39})$$

The empirical covariance estimator of β has the same concentration property as Lemma 1, which we will omit the proof of.

Corollary 1 (of Lemma 1) Suppose $\beta_i, i = 1, \dots, n_{\text{task}}$ are generated independently from $\mathcal{N}(0, \Sigma_{\beta})$. Then with probability $1 - \mathcal{O}((n_{\text{tot}}\text{tr}(\Sigma_{\beta}))^{-C})$,

$$\|\bar{\Sigma}_{\beta} - \Sigma_{\beta}\| \lesssim \log(n_{\text{task}}\text{tr}(\Sigma_{\beta}))(\frac{\log(n_{\text{tot}}\text{tr}(\Sigma_{\beta}))\text{tr}(\Sigma_{\beta})}{n_{\text{task}}} + \sqrt{\frac{\|\Sigma_{\beta}\|\text{tr}(\Sigma_{\beta})}{n_{\text{task}}}})$$

Ensembling Corollary 1 with Theorem 5 we can prove Theorem 1. ■

A.3. Method of moment - another estimator

In this section, we will define another moment estimator, which as well potentially recovers the subspace we want. Define

$$\hat{Q} = \frac{1}{n_{\text{task}}}\sum_{i=1}^{n_{\text{task}}}\frac{1}{n_{\text{spt}}}\sum_{j=1}^{n_{\text{spt}}}y_{ij}^2\mathbf{x}_{ij}\mathbf{x}_{ij}^\top.$$

Theorem 6 *The mean of estimator \hat{Q} is*

$$E(Q) := Q = \Sigma_X \Sigma_\beta \Sigma_X + \text{tr}(\Sigma_\beta \Sigma_X) \Sigma_X$$

Let

$$\bar{Q} = \Sigma_X \bar{\Sigma}_\beta \Sigma_X + \text{tr}(\bar{\Sigma}_\beta \Sigma_X) \Sigma_X$$

with the same assumption as Theorem 5, we have

$$\begin{aligned} \|\hat{Q} - \bar{Q}\| &\lesssim \sqrt{\frac{\text{tr}(\Sigma_X) \text{tr}^2(\Sigma_\beta \Sigma_X) \|\Sigma_X\| + (\text{tr}(\Sigma_\beta \Sigma_X) + \sigma_\varepsilon^2 \text{tr}(\Sigma_X)) \cdot \sigma_\varepsilon^2 \|\Sigma_X\|}{n_{\text{tot}}}} \\ &+ \frac{\text{tr}^2(\Sigma_X) \text{tr}(\Sigma_\beta \Sigma_X) + (\sigma_\varepsilon + \sqrt{\text{tr}(B)}) \sigma_\varepsilon \text{tr}(\Sigma_X)}{n_{\text{tot}}}. \end{aligned}$$

Proof We first compute the expectation of \hat{Q} . We simplify the notation the same way as before: fix β and let $y = x^\top \beta + \varepsilon$, and study yx and $y^2 x x^\top$. Define $z_i \sim \mathcal{N}(0, 1)$, and $\alpha = \sqrt{\Sigma_X} \beta$.

$$Q = E(\hat{Q}) = E y^2 x x^\top = E((\beta^\top x + \varepsilon)^2 x x^\top) = E \left[\sqrt{\sigma_i \sigma_j} z_i z_j \left(\sum_{k=1}^d z_k \alpha_k \right)^2 \right] + \sigma_\varepsilon^2 \Sigma_X \quad (\text{A.40})$$

$$E(\sigma_i z_i^2 \left(\sum_{k=1}^d z_k \alpha_k \right)^2) = E(\|\alpha\|^2 + 2\alpha_i^2) \sigma_i \quad (\text{A.41})$$

$$E(\sqrt{\sigma_i \sigma_j} z_i z_j \left(\sum_{k=1}^d z_k \alpha_k \right)^2) = 2E\sqrt{\sigma_i \sigma_j} \alpha_i \alpha_j, \quad i \neq j. \quad (\text{A.42})$$

Hence

$$Q = \Sigma_X \Sigma_\beta \Sigma_X + \text{tr}(\sqrt{\Sigma_X} \Sigma_\beta \sqrt{\Sigma_X}) \Sigma_X + \sigma_\varepsilon^2 \Sigma_X \quad (\text{A.43})$$

Now we follow the method in the last section, we will bound

$$P(\|y^2 x x^\top\| \geq t) \text{ and } E(\|y^4 (\|x\|^2 x x^\top)\|). \quad (\text{A.44})$$

1. $\|\hat{y}^2 x x^\top\|$. Let $z \sim \mathcal{N}(0, I_d)$. We first write $\|\hat{y}^2 x x^\top\| = \|\beta^\top \sqrt{\Sigma_X} z\|^2 \|\sqrt{\Sigma_X} z\|^2$. We know that with probability $\mathcal{O}(\delta)$

$$\|\sqrt{\Sigma_X} z\|^2 \geq \text{tr}(\Sigma_X) + \text{tr}(\Sigma_X) \log(1/\delta), \quad (\text{A.45})$$

$$\hat{y}^2 = (\beta^\top \sqrt{\Sigma_X} z)^2 \geq \|\sqrt{\Sigma_X} \beta\|^2 (1 + \log(1/\delta)) \quad (\text{A.46})$$

So with probability $1 - \mathcal{O}(\delta)$ we have

$$\|\hat{y}^2 x x^\top\| \lesssim \text{tr}(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|^2 + \text{tr}(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|^2 \log^2(1/\delta) \quad (\text{A.47})$$

2. $E(\|\hat{y}^4 (\|x\|^2 x x^\top)\|)$. With the same derivation above we have that with probability $1 - \mathcal{O}(\delta)$

$$\|\hat{y}^4 (x x^\top)^2\| = \|\hat{y}^4 \|x\|^2 (x x^\top)\| \leq c \left(\text{tr}(\Sigma_X)^2 \|\sqrt{\Sigma_X} \beta\|^4 + \text{tr}^2(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|^4 \log^4(1/\delta) \right) \quad (\text{A.48})$$

So for all $K \geq 1$, the expectation is smaller than

$$K \text{tr}(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|^4 + \int_{K \text{tr}(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|}^{\infty} \exp(-c(\frac{t}{\text{tr}^2(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|^4})^{1/4}) dt \quad (\text{A.49})$$

$$\lesssim K \text{tr}(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|^4 + K^{3/4} \exp(-cK^{1/4}) \text{tr}^2(\Sigma_X) \|\sqrt{\Sigma_X} \beta\|^4. \quad (\text{A.50})$$

3. Noise. Suppose $y = \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon$. Then we have

$$y^2 \mathbf{x} \mathbf{x}^\top = (\mathbf{x}^\top \boldsymbol{\beta} + \varepsilon)^2 \mathbf{x} \mathbf{x}^\top \preceq (\mathbf{x}^\top \boldsymbol{\beta})^2 \mathbf{x} \mathbf{x}^\top + \varepsilon^2 \mathbf{x} \mathbf{x}^\top$$

The first term was computed in the noiseless case, so in this part we focus on $\varepsilon^2 \mathbf{x} \mathbf{x}^\top$. And

$$y^4 (\|\mathbf{x}\|^2 \mathbf{x} \mathbf{x}^\top) = (\mathbf{x}^\top \boldsymbol{\beta} + \varepsilon)^4 (\|\mathbf{x}\|^2 \mathbf{x} \mathbf{x}^\top) \preceq (\mathbf{x}^\top \boldsymbol{\beta})^4 \|\mathbf{x}\|^2 \mathbf{x} \mathbf{x}^\top + \varepsilon^4 \|\mathbf{x}\|^2 \mathbf{x} \mathbf{x}^\top$$

The first term is bounded by the last part, we focus on $\varepsilon^4 \|\mathbf{x}\|^2 \mathbf{x} \mathbf{x}^\top$.

(a) $\varepsilon^2 \mathbf{x} \mathbf{x}^\top$. With probability at least $1 - \delta$, we know that $\varepsilon^2 \leq \mathcal{O}(\sigma_\varepsilon^2(1 + \log(1/\delta)))$ and $\|\mathbf{x} \mathbf{x}^\top\| \leq \text{tr}(\boldsymbol{\Sigma}_\mathbf{X}) + \mathcal{O}(\text{tr}(\boldsymbol{\Sigma}_\mathbf{X}) \log(1/\delta))$. Thus with probability at least $1 - \delta$ we have

$$\|\varepsilon^2 \mathbf{x} \mathbf{x}^\top\| \lesssim \sigma_\varepsilon^2 \text{tr}(\boldsymbol{\Sigma}_\mathbf{X}) + \sigma_\varepsilon^2 \text{tr}(\boldsymbol{\Sigma}_\mathbf{X}) \log^2(1/\delta)$$

(b) $\varepsilon^4 (\mathbf{x} \mathbf{x}^\top)^2$. We know that

$$\|E(\varepsilon^4 \|\mathbf{x}\|^2 \mathbf{x} \mathbf{x}^\top)\| = 3\sigma_\varepsilon^4 \|E(\|\mathbf{x}\|^2 \mathbf{x} \mathbf{x}^\top)\| \lesssim \sigma_\varepsilon^4 \|\boldsymbol{\Sigma}_\mathbf{X}\| \text{tr}(\boldsymbol{\Sigma}_\mathbf{X})$$

So in the noisy case, we know that

$$\|y^2 \mathbf{x} \mathbf{x}^\top\| \lesssim (\|\sqrt{\boldsymbol{\Sigma}_\mathbf{X}} \boldsymbol{\beta}\|^2 + \sigma_\varepsilon^2) \text{tr}(\boldsymbol{\Sigma}_\mathbf{X}) (\|\sqrt{\boldsymbol{\Sigma}_\mathbf{X}} \boldsymbol{\beta}\|^2 + \sigma_\varepsilon^2) \text{tr}(\boldsymbol{\Sigma}_\mathbf{X}) \log^2(1/\delta), \text{ w.p. } 1 - \delta \quad (\text{A.51})$$

$$E(\|y^4 (\mathbf{x} \mathbf{x}^\top)^2\|) \lesssim (\|\sqrt{\boldsymbol{\Sigma}_\mathbf{X}} \boldsymbol{\beta}\|^4 + \sigma_\varepsilon^4 \|\boldsymbol{\Sigma}_\mathbf{X}\|) \text{tr}(\boldsymbol{\Sigma}_\mathbf{X}) \quad (\text{A.52})$$

We can exactly define the same quantities as in (A.38), and get a similar result to (A.39), let $\bar{\boldsymbol{\Sigma}}_\beta = \frac{1}{n_{\text{task}}} \sum_{i=1}^{n_{\text{task}}} \boldsymbol{\beta}_i \boldsymbol{\beta}_i^\top$, then

$$\log^{-1}(n_{\text{tot}} T_0) \left\| \frac{1}{n_{\text{spt}}} \sum_{i=1}^{n_{\text{spt}}} y_i^2 \mathbf{x}_i \mathbf{x}_i^\top - (\boldsymbol{\Sigma}_\mathbf{X} \bar{\boldsymbol{\Sigma}}_\beta \boldsymbol{\Sigma}_\mathbf{X} + \text{tr}(\boldsymbol{\Sigma}_\mathbf{X} \bar{\boldsymbol{\Sigma}}_\beta) \boldsymbol{\Sigma}_\mathbf{X}) \right\| \lesssim \frac{T_0 \log(n_{\text{tot}} T_0)}{n_{\text{spt}}} + \frac{\sigma}{\sqrt{n_{\text{spt}}}}. \quad (\text{A.53})$$

■

A.4. Estimating with fewer samples when each task contains enough samples

In this part we first prove Theorem 2.

Theorem 2 (Standard normal feature, noiseless) *Let data be generated as in Def 1, 2, let $\mathcal{S} = \max\{\|\boldsymbol{\Sigma}_\mathbf{X}\|, \|\boldsymbol{\Sigma}_\beta\|\}$, $\mathcal{T}_{\beta, \mathbf{X}} = \text{tr}(\boldsymbol{\Sigma}_\beta \boldsymbol{\Sigma}_\mathbf{X})$, $\mathcal{T}_\mathbf{X} = \text{tr}(\boldsymbol{\Sigma}_\mathbf{X})$, $\mathcal{T}_\beta = \text{tr}(\boldsymbol{\Sigma}_\beta)$. Suppose $\sigma_\varepsilon = 0$, $\boldsymbol{\Sigma}_\mathbf{X} = \mathbf{I}$, and suppose the rank of $\boldsymbol{\Sigma}_\beta$ is r_t . Define $\hat{\boldsymbol{\beta}}_i = n_{\text{spt}}^{-1} \sum_{j=1}^{n_{\text{spt}}} y_{ij} \mathbf{x}_{ij}$, $\mathbf{B} = [\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_{n_{\text{task}}}]$, and $\hat{\mathbf{B}} = [\hat{\boldsymbol{\beta}}_1, \dots, \hat{\boldsymbol{\beta}}_{n_{\text{task}}}]$. Let $n_{\text{spt}} > c_1 \mathcal{T}_\beta \lambda_{r_t}^{-1}(\boldsymbol{\Sigma}_\beta)$, with probability $1 - (n_{\text{task}}^{-c_3} + (n_{\text{task}} \mathcal{T}_\beta)^{-c_4} + \exp(-c_5 n_{\text{task}}^2))$, where c_i are constants,*

$$\sigma_{\max}(\hat{\mathbf{B}} - \mathbf{B}) \lesssim \sqrt{\frac{n_{\text{task}} \mathcal{T}_\beta}{n_{\text{spt}}}}.$$

Denote the span of top r_t singular column vectors of $\hat{\mathbf{B}}$ and $\boldsymbol{\Sigma}_\beta$ as $\hat{\mathbf{W}}, \mathbf{W}$, then

$$\sin(\angle \hat{\mathbf{W}}, \mathbf{W}) \lesssim \sqrt{\frac{\mathcal{T}_\beta}{n_{\text{spt}} \lambda_{r_t}(\boldsymbol{\Sigma}_\beta)}}.$$

If $\boldsymbol{\Sigma}_\beta = \text{diag}(\mathbf{I}_{r_t}, 0)$, then $\sin(\angle \hat{\mathbf{W}}, \mathbf{W}) \lesssim \sqrt{r_t/n_{\text{spt}}}$.

Proof We first estimate $\boldsymbol{\beta}_i$ as

$$\hat{\boldsymbol{\beta}}_i = \frac{1}{n_{\text{spt}}} \sum_{j=1}^{n_{\text{spt}}} y_{ij} \mathbf{x}_{ij}.$$

Then we fix β_i and compute the covariance of $y_{ij}\mathbf{x}_{ij}$ (its mean is β_i).

$$\text{Cov}(y_{ij}\mathbf{x}_{ij} - \beta_i) = \mathbf{E}(\mathbf{x}_{ij}\mathbf{x}_{ij}^\top \beta_i \beta_i^\top \mathbf{x}_{ij}\mathbf{x}_{ij}^\top) - \beta_i \beta_i^\top \|\beta_i\|^2 \mathbf{I} + \beta_i \beta_i^\top \preceq 2\|\beta_i\|^2 \mathbf{I}.$$

With matrix concentration we know that

$$\text{Cov}(\hat{\beta}_i - \beta_i) \lesssim \frac{\|\beta_i\|^2}{n_{\text{spt}}} \mathbf{I}. \quad (\text{A.54})$$

Suppose $\mathbf{B} = [\beta_1, \dots, \beta_{n_{\text{task}}}]$, and $\hat{\mathbf{B}} = [\hat{\beta}_1, \dots, \hat{\beta}_{n_{\text{task}}}]$. Then we know the covariance of each column of $\hat{\mathbf{B}} - \mathbf{B}$ is bounded by (A.54). Thus with random matrix theory, we know that with probability $1 - \exp(-cn_{\text{task}}^2)$ for constant c ,

$$\sigma_{\max}^2(\hat{\mathbf{B}} - \mathbf{B}) \lesssim \frac{n_{\text{task}} \|\beta_i\|^2}{n_{\text{spt}}}. \quad (\text{A.55})$$

We have proved that $\|\beta_i\|^2 \leq \log(n_{\text{task}}) \text{tr}(\Sigma_\beta)$ with probability $1 - n_{\text{task}}^{-c}$. The columns of \mathbf{B} is generated from $\mathcal{N}(0, \Sigma_\beta)$, so that

$$\sigma_{\max}(\hat{\mathbf{B}} - \mathbf{B}) \lesssim \mathcal{O}\left(\sqrt{\frac{n_{\text{task}} \log(n_{\text{task}}) \text{tr}(\Sigma_\beta)}{n_{\text{spt}}}}\right).$$

Now we study \mathbf{B} . We know that $\mathbf{E}(\mathbf{B}\mathbf{B}^\top) = \mathbf{E}(\sum_{i=1}^{n_{\text{task}}} \beta_i \beta_i^\top) = n_{\text{task}} \Sigma_\beta$. \mathbf{B} is a matrix with independent columns. Thus with the lemma 1 we proved before, let $n_{\text{spt}} > c_1 \text{tr}(\Sigma_\beta) \lambda_{r_t}^{-1}(\Sigma_\beta)$, $n_{\text{task}} > \max\{c_2 d, \frac{\|\Sigma_\beta\| \text{tr}(\Sigma_\beta)}{\lambda_{r_t}^2(\Sigma_\beta)}\}$, then for Gaussian matrix with independent columns (Vershynin, 2010), with probability $1 - \mathcal{O}(n_{\text{task}}^{-c_3} + (n_{\text{task}} \text{tr}(\Sigma_\beta))^{-c_4} + \exp(-c_5 n_{\text{task}}^2))$, where c_i are constants,

$$\sigma_{r_t}(\mathbf{B}) \geq \sqrt{n_{\text{task}} \lambda_{r_t}(\Sigma_\beta) - \mathcal{O}(\sqrt{n_{\text{task}} \|\Sigma_\beta\| \text{tr}(\Sigma_\beta)})}.$$

Denote the span of top r_t singular vectors of $\hat{\mathbf{B}}$ and Σ_β as $\hat{\mathbf{W}}, \mathbf{W}$, with Lemma 2,

$$\sin(\angle \hat{\mathbf{W}}, \mathbf{W}) \leq \sqrt{\frac{\log(n_{\text{task}}) \text{tr}(\Sigma_\beta)}{n_{\text{spt}} \lambda_{r_t}(\Sigma_\beta)}}.$$

■

Next, we will propose a theorem with general feature covariance and noisy data, which is a generalization of Theorem 2.

Theorem 8 Let data be generated as in Def 1, 2. Suppose $\hat{\mathbf{b}}_i = n_{\text{spt}}^{-1} \sum_{j=1}^{n_{\text{spt}}} y_{ij} \mathbf{x}_{ij}$, $\mathbf{B} = \Sigma_{\mathbf{X}} [\beta_1, \dots, \beta_{n_{\text{task}}}]$, and $\hat{\mathbf{B}} = [\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_{n_{\text{task}}}]$. Let $\delta_\lambda = \lambda_{r_t}(\Sigma_{\mathbf{X}} \Sigma_\beta \Sigma_{\mathbf{X}}) - \lambda_{r_t+1}(\Sigma_{\mathbf{X}} \Sigma_\beta \Sigma_{\mathbf{X}})$, suppose $\Sigma_{\mathbf{X}}$ is approximately rank r_f ,

$$\begin{aligned} n_{\text{spt}} &\gtrsim (\text{tr}(\Sigma_\beta \Sigma_{\mathbf{X}}) + \sigma_\varepsilon^2) \|\Sigma_{\mathbf{X}}\|, \\ n_{\text{task}} &\gtrsim \max\left\{r_f, \frac{d \lambda_{r_f+1}(\Sigma_{\mathbf{X}})}{\|\Sigma_{\mathbf{X}}\|}, \frac{\|\Sigma_{\mathbf{X}} \Sigma_\beta \Sigma_{\mathbf{X}}\| \text{tr}(\Sigma_{\mathbf{X}} \Sigma_\beta \Sigma_{\mathbf{X}})}{\delta_\lambda^2}\right\} \end{aligned}$$

then with probability $1 - \mathcal{O}(n_{\text{task}}^{-C_1} + (n_{\text{task}} \text{tr}(\Sigma_{\mathbf{X}} \Sigma_\beta \Sigma_{\mathbf{X}}))^{-C_2} + \exp(-C_3 n_{\text{task}}^2))$, where C_i are constants,

$$\sigma_{\max}(\hat{\mathbf{B}} - \mathbf{B}) \lesssim \sqrt{\frac{n_{\text{task}} (\text{tr}(\Sigma_\beta \Sigma_{\mathbf{X}}) + \sigma_\varepsilon^2) \|\Sigma_{\mathbf{X}}\|}{n_{\text{spt}}}}.$$

Denote the span of top r_t singular vectors of $\hat{\mathbf{B}}$ and $\Sigma_{\mathbf{X}} \Sigma_\beta \Sigma_{\mathbf{X}}$ as $\hat{\mathbf{W}}, \mathbf{W}$,

$$\sin(\angle \hat{\mathbf{W}}, \mathbf{W}) \lesssim \sqrt{\frac{(\text{tr}(\Sigma_\beta \Sigma_{\mathbf{X}}) + \sigma_\varepsilon^2) \|\Sigma_{\mathbf{X}}\|}{n_{\text{spt}} \delta_\lambda^2}}.$$

Example 3 Suppose $\Sigma_X = \text{diag}(\mathbf{I}_{r_f}, \iota \mathbf{I}_{d-r_f})$, and $\Sigma_\beta = \text{diag}(\mathbf{I}_{r_t}, 0)$, $\sigma_\varepsilon = 0$. Suppose $\iota d < r_f$. Then with $n_{\text{task}} \gtrsim r_f$, $n_{\text{spt}} \gtrsim r_t$ so that $n_{\text{tot}} \gtrsim r_f r_t$,

$$\sin(\angle \hat{\mathbf{W}}, \mathbf{W}) \lesssim \sqrt{r_t/n}$$

Proof We let $\mathbf{x}_{ij} \sim \mathcal{N}(0, \Sigma_X)$, and let Σ_X be diagonal. For the i th task, let

$$\hat{\mathbf{b}}_i = \frac{1}{n_{\text{spt}}} \sum_{j=1}^{n_{\text{spt}}} y_{ij} \mathbf{x}_{ij}.$$

We fix β_i and compute the covariance of $y_{ij} \mathbf{x}_{ij}$

$$\mathbf{E}(y_{ij} \mathbf{x}_{ij}) = \mathbf{E}(\mathbf{x}_{ij} \mathbf{x}_{ij}^\top \beta_i) = \Sigma_X \beta_i.$$

and similar to (A.43),

$$\text{Cov}(y_{ij} \mathbf{x}_{ij} - \Sigma_X \beta_i) = (\beta_i^\top \Sigma_X \beta_i) \Sigma_X + \sigma_\varepsilon^2 \Sigma_X$$

With matrix concentration we know that

$$\text{Cov}(\hat{\mathbf{b}}_i - \Sigma_X \beta_i) \lesssim \frac{\beta_i^\top \Sigma_X \beta_i + \sigma_\varepsilon^2}{n_{\text{spt}}} \Sigma_X. \quad (\text{A.56})$$

Suppose $\mathbf{B} = \Sigma_X [\beta_1, \dots, \beta_{n_{\text{task}}}]$, and $\hat{\mathbf{B}} = [\hat{\mathbf{b}}_1, \dots, \hat{\mathbf{b}}_{n_{\text{task}}}]$. $\hat{\mathbf{B}} - \mathbf{B}$ is a matrix with independent columns. Suppose \mathbf{X} is approximately rank r_f , Let $\mathbf{V}_{r_f} \in \mathbb{R}^{d \times d}$ be the projection onto the top- R singular vector space of Σ_X and $\mathbf{V}_{r_f}^\perp \in \mathbb{R}^{d \times d}$ be the projection onto the $r_f + 1$ to d th singular vector space of Σ_X . With there are $n_{\text{task}} \geq r_f$ columns, we know that

$$\begin{aligned} \sigma_{\max}(\mathbf{V}_{r_f}(\hat{\mathbf{B}} - \mathbf{B})) &\lesssim \frac{n_{\text{task}}(\max_i \beta_i^\top \Sigma_X \beta_i + \sigma_\varepsilon^2) \|\Sigma_X\|}{n_{\text{spt}}} \\ \sigma_{\max}(\mathbf{V}_{r_f}^\perp(\hat{\mathbf{B}} - \mathbf{B})) &\lesssim \frac{\max\{n_{\text{task}}, d\}(\max_i \beta_i^\top \Sigma_X \beta_i + \sigma_\varepsilon^2) \lambda_{r_t+1}(\Sigma_X)}{n_{\text{spt}}} \end{aligned}$$

With similar argument as before, with probability $1 - \exp(-cn_{\text{task}}^2)$ for constant c ,

$$\sigma_{\max}^2(\hat{\mathbf{B}} - \mathbf{B}) \lesssim \frac{\max\{n_{\text{task}} \|\Sigma_X\|, d \lambda_{r_f+1}(\Sigma_X)\}(\max_i \beta_i^\top \Sigma_X \beta_i + \sigma_\varepsilon^2) \|\Sigma_X\|}{n_{\text{spt}}}. \quad (\text{A.57})$$

We know that $\|\sqrt{\Sigma_X} \beta_i\|^2 \leq \mathcal{O}(\log(n_{\text{task}}) \text{tr}(\Sigma_\beta \Sigma_X))$ with probability $1 - n_{\text{task}}^{-c}$ for constant c . So that

$$\sigma_{\max}(\hat{\mathbf{B}} - \mathbf{B}) \lesssim \sqrt{\frac{\max\{n_{\text{task}} \|\Sigma_X\|, d \lambda_{r_f+1}(\Sigma_X)\}(\log(n_{\text{task}}) \text{tr}(\Sigma_\beta \Sigma_X) + \sigma_\varepsilon^2) \|\Sigma_X\|}{n_{\text{spt}}}}. \quad (\text{A.58})$$

We can omit the $\log(n_{\text{task}})$ term since we hide \log in \lesssim . Now we study \mathbf{B} . We know that $\mathbf{E}(\mathbf{B} \mathbf{B}^\top) = \mathbf{E}(\Sigma_X (\sum_{i=1}^{n_{\text{task}}} \beta_i \beta_i^\top) \Sigma_X) = n_{\text{task}} \Sigma_X \Sigma_\beta \Sigma_X$.

Thus with the lemma 1 we proved before, let

$$n_{\text{spt}} > C_1 (\log(n_{\text{task}}) \text{tr}(\Sigma_\beta \Sigma_X) + \sigma_\varepsilon^2) \|\Sigma_X\|.$$

then with probability $1 - \mathcal{O}(n_{\text{task}}^{-C_3} + (n_{\text{task}} \text{tr}(\Sigma_X \Sigma_\beta \Sigma_X))^{-C_4} + \exp(-C_5 n_{\text{task}}^2))$, where C_i are constants,

$$\sigma_{r_t}(\mathbf{B}) \geq \sqrt{n_{\text{task}} (\lambda_{r_t}(\Sigma_X \Sigma_\beta \Sigma_X) - \lambda_{r_t+1}(\Sigma_X \Sigma_\beta \Sigma_X))} - \mathcal{O}(\sqrt{n_{\text{task}} \|\Sigma_X \Sigma_\beta \Sigma_X\| \text{tr}(\Sigma_X \Sigma_\beta \Sigma_X)}).$$

Denote the span of top r_t singular vectors of \hat{B} and $\Sigma_X \Sigma_\beta \Sigma_X$ as \hat{W} , W , let

$$n_{\text{task}} \gtrsim \max\left\{r_f, \frac{d\lambda_{r_f+1}(\Sigma_X)}{\|\Sigma_X\|}, \frac{\|\Sigma_X \Sigma_\beta \Sigma_X\| \text{tr}(\Sigma_X \Sigma_\beta \Sigma_X)}{\delta_\lambda^2}\right\} \quad (\text{A.59})$$

we plug in (A.58) with Lemma 2,

$$\begin{aligned} \sin(\angle \hat{W}, W) &\lesssim \sqrt{\left(\frac{d\lambda_{r_f+1}(\Sigma_X)}{n_{\text{task}}\|\Sigma_X\|} + 1\right) \cdot \frac{(\text{tr}(\Sigma_\beta \Sigma_X) + \sigma_\varepsilon^2)\|\Sigma_X\|}{n_{\text{spt}}(\lambda_{r_t}(\Sigma_X \Sigma_\beta \Sigma_X) - \lambda_{r_t+1}(\Sigma_X \Sigma_\beta \Sigma_X))}} \\ &\approx \sqrt{\frac{(\text{tr}(\Sigma_\beta \Sigma_X) + \sigma_\varepsilon^2)\|\Sigma_X\|}{n_{\text{spt}}(\lambda_{r_t}(\Sigma_X \Sigma_\beta \Sigma_X) - \lambda_{r_t+1}(\Sigma_X \Sigma_\beta \Sigma_X))}}. \end{aligned}$$

■

B. Analysis of optimal representation

We denote $\hat{\beta} = \Lambda^*(X\Lambda^*)^\dagger y$. When $n_{\text{fs}} < d$, and Λ^* is full row rank, then $\hat{\beta} = \hat{\alpha}$ where $\hat{\alpha} = (X\Lambda^*)^\dagger y$ is the smallest ℓ_2 norm solution to the equation $X\Lambda^* \alpha = y$.

We first prove Observation 1.

Observation 1 Let $\Lambda \in \mathbb{R}^{d \times d} \succ 0$, $X \in \mathbb{R}^{d \times n_{\text{fs}}}$ and $y \in \mathbb{R}^{n_{\text{fs}}}$, and define

$$\hat{\beta}_1 = \Lambda(X\Lambda)^\dagger y, \quad (\text{B.1})$$

$$\hat{\beta}_2 = \lim_{t \rightarrow 0} \arg\min_{\beta} \|X^\top \beta - y\|^2 + t\beta^\top \Lambda^{-2} \beta \quad (\text{B.2})$$

Then $\hat{\beta}_1 = \hat{\beta}_2$.

Proof Denote the SVD of $\Lambda X = U\Sigma V^\top$, where $U \in \mathbb{R}^{d \times d}$, $\Sigma \in \mathbb{R}^{d \times n_{\text{fs}}}$, $V \in \mathbb{R}^{n_{\text{fs}} \times n_{\text{fs}}}$. Let $\Sigma^{(2)} \in \mathbb{R}^{n_{\text{fs}} \times n_{\text{fs}}} = \Sigma^\top \Sigma = \text{diag}(\Sigma_{1,1}^2, \dots, \Sigma_{n_{\text{fs}}, n_{\text{fs}}}^2) \succ 0$.

$$\begin{aligned} \hat{\beta}_2 &= \lim_{t \rightarrow 0} \arg\min_{\beta} \|X^\top \beta - y\|^2 + t\beta^\top \Lambda^{-2} \beta \\ &= \lim_{t \rightarrow 0} (X X^\top + t\Lambda^{-2})^{-1} X y \\ &= \lim_{s \rightarrow \infty} s\Lambda(s\Lambda X X^\top \Lambda + I)^{-1} \Lambda X y \\ &= \lim_{s \rightarrow \infty} s\Lambda(sU\Sigma V^\top V\Sigma^\top U^\top + I)^{-1} U\Sigma V^\top y \\ &= \lim_{s \rightarrow \infty} s\Lambda(sU \text{diag}(\Sigma^{(2)} + I_{n_{\text{fs}}}, I_{d-n_{\text{fs}}})U^\top)^{-1} U\Sigma V^\top y \\ &= \lim_{s \rightarrow \infty} \Lambda U(\text{diag}(\Sigma^{(2)}, I_{d-n_{\text{fs}}}/s))^{-1} \Sigma V^\top y. \\ &= \Lambda(X\Lambda)^\dagger y \end{aligned}$$

■

The risk of $\hat{\beta}$ is given by

$$\text{risk}(\hat{\beta}) = E(y - x^\top \hat{\beta}) = E\|\hat{\beta} - \beta\|_{\Sigma_X}^2 = E(\hat{\beta} - \beta)^\top \Sigma_X (\hat{\beta} - \beta).$$

In Sec. B.1, we study the asymptotic optimal representation. Below, we characterize the properties of the problem for fixed β and arbitrary input covariance Σ_X . We first go over this and then discuss how to obtain the optimal representation Λ^* minimizing test risk.

B.1. Distributional characterization of least norm solution

Define $\tilde{\mathbf{X}} \in \mathbb{R}^{n \times p}$, $\tilde{\mathbf{y}} \in \mathbb{R}^n$. We study the following least norm solution of the least squares problem

$$\hat{\beta} = \arg \min_{\beta'} \|\beta'\|, \quad \text{s.t., } \tilde{\mathbf{X}}\beta' = \tilde{\mathbf{y}} \quad (\text{B.3})$$

Assumption 1 Assume the rows of $\tilde{\mathbf{X}}$ are independently drawn from $\mathcal{N}(0, \tilde{\Sigma}_{\mathbf{X}})$. We focus on a double asymptotic regime where $p, n \rightarrow \infty$ at fixed overparameterization ratio $\kappa := p/n > 0$.

Assumption 2 The covariance matrix $\tilde{\Sigma}_{\mathbf{X}}$ is diagonal and there exist constants $C_{\min}, C_{\max} \in (0, \infty)$ such that: $C_{\min} \mathbf{I} \preceq \tilde{\Sigma}_{\mathbf{X}} \preceq C_{\max} \mathbf{I}$.

Assumption 3 The joint empirical distribution of $\{(\lambda_i(\tilde{\Sigma}_{\mathbf{X}}), \sqrt{p}\beta_i)\}_{i \in [p]}$ converges in Wasserstein- k distance to a probability distribution μ on $\mathbb{R}_{>0} \times \mathbb{R}$ for some $n_{\text{task}} \geq 4$. That is $\frac{1}{p} \sum_{i \in [p]} \delta_{(\lambda_i(\tilde{\Sigma}_{\mathbf{X}}), \sqrt{p}\beta_i)} \xrightarrow{W_k} \mu$.

Definition 6 (Asymptotic DC – Overparameterized regime) (Thrampoulidis et al., 2015) Let random variables $(\Sigma, B) \sim \mu$ (where μ is defined in Assumption 3) and fix $\kappa > 1$. Define parameter ξ as the unique positive solution to the following equation

$$\mathbb{E}_{\mu} \left[\left(1 + (\xi \cdot \Sigma)^{-1} \right)^{-1} \right] = \kappa^{-1}. \quad (\text{B.4})$$

Further define the positive parameter γ as follows:

$$\gamma := \left(\sigma^2 + \mathbb{E}_{\mu} \left[\frac{B^2 \Sigma}{(1 + \xi \Sigma)^2} \right] \right) / \left(1 - \kappa \mathbb{E}_{\mu} \left[\frac{1}{(1 + (\xi \Sigma)^{-1})^2} \right] \right). \quad (\text{B.5})$$

With these and $H \sim \mathcal{N}(0, 1)$, define the random variable

$$X_{\kappa, \sigma^2}(\Sigma, B, H) := \left(1 - \frac{1}{1 + \xi \Sigma} \right) B + \sqrt{\kappa} \frac{\sqrt{\gamma} \Sigma^{-1/2}}{1 + (\xi \Sigma)^{-1}} H, \quad (\text{B.6})$$

and let Π_{κ, σ^2} be its distribution.

Theorem 9 (Asymptotic DC – Overparameterized LGP) (Thrampoulidis et al., 2015) Fix $\kappa > 1$ and suppose Assumptions 2 and 3 hold. Recall the solution $\hat{\beta}$ from (3.4) and let

$$\hat{\Pi}_n(\tilde{\mathbf{y}}, \tilde{\mathbf{X}}, \beta, \tilde{\Sigma}_{\mathbf{X}}) := \frac{1}{p} \sum_{i=1}^p \delta_{\sqrt{p}\hat{\beta}_i, \sqrt{p}\beta_i, \tilde{\Sigma}_{\mathbf{X}_{i,i}}}$$

be the joint empirical distribution of $(\sqrt{p}\hat{\beta}, \sqrt{p}\beta, \tilde{\Sigma}_{\mathbf{X}})$. Let $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ is a function in $\text{PL}(2)$. We have that

$$\frac{1}{p} \sum_{i=1}^p f(\sqrt{p}\hat{\beta}_i, \sqrt{p}\beta_i, \tilde{\Sigma}_{\mathbf{X}_{i,i}}) \xrightarrow{P} \mathbb{E} [f(X_{\kappa, \sigma^2}, B, \Sigma)]. \quad (\text{B.7})$$

In particular, the risk is given by

$$\text{risk}(\hat{\beta}_n) \xrightarrow{P} \mathbb{E}[\Sigma(B - X_{\kappa, \sigma^2})] \quad (\text{B.8})$$

$$= \mathbb{E} \left[\frac{\Sigma}{(1 + \xi \Sigma)^2} B^2 + \frac{\kappa \gamma}{(1 + (\xi \Sigma)^{-1})^2} \right]. \quad (\text{B.9})$$

B.2. Finding Optimal Representation Λ^*

Now, for simplicity (and actually without losing generality) assume $\tilde{\Sigma}_{\mathbf{X}} = \mathbf{I}$. This means that empirical measure of $\Sigma_{\mathbf{X}}$ trivially converges to $\Sigma = 1$. With the representation Λ^* with asymptotic distribution Λ , the ML problem has the following mapping

$$\beta \rightarrow \Lambda^{*-1} \beta \quad \text{and} \quad \tilde{\Sigma}_{\mathbf{X}} \rightarrow \Lambda^* \tilde{\Sigma}_{\mathbf{X}} \Lambda^*.$$

This means the empirical measure converges to the following mapped distributions

$$B \rightarrow \bar{B} = \Lambda^{-1} B \quad \text{and} \quad \Sigma = 1 \rightarrow \bar{\Sigma} = \Lambda^2 \Sigma = \Lambda^2.$$

Our question: Craft the optimal distribution Λ to minimize the representation learning risk. Specifically, for a given (B, Λ) pair, we know from the theorem above that

$$\text{risk}^{\Lambda^*}(\hat{\beta}_n) \xrightarrow{P} \mathbb{E}\left[\frac{\bar{\Sigma}}{(1 + \xi \bar{\Sigma})^2} \bar{B}^2 + \frac{\kappa \gamma}{(1 + (\xi \bar{\Sigma})^{-1})^2}\right] \quad (\text{B.10})$$

$$= \mathbb{E}\left[\frac{B^2}{(1 + \xi \Lambda^2)^2} + \frac{\kappa \gamma}{(1 + (\xi \Lambda^2)^{-1})^2}\right]. \quad (\text{B.11})$$

Thus, the optimal weighting strategy (asymptotically) is given by the distribution

$$\Lambda^* = \arg \min_{\Lambda} \mathbb{E}\left[\frac{B^2}{(1 + \xi \Lambda^2)^2} + \frac{\kappa \gamma}{(1 + (\xi \Lambda^2)^{-1})^2}\right],$$

where γ, ξ are strictly positive scalars that are also functions of Λ .

B.3. Non-asymptotic Analysis (for simpler insights)

We apply the discussion in Sec. B.1 non-asymptotically in few-shot learning. Remember we define $\mathbf{X} \in \mathbb{R}^{n_{\text{fs}} \times d}$, $\mathbf{y} \in \mathbb{R}^{n_{\text{fs}}}$, each row of \mathbf{X} is independently drawn from $\mathcal{N}(0, \Sigma_{\mathbf{X}})$. We study the following least norm solution of the least squares problem

$$\hat{\beta} = \arg \min_{\beta'} \|\beta'\|, \quad \text{s.t., } \mathbf{X} \beta' = \mathbf{y}. \quad (\text{B.12})$$

Definition 7 (Non-asymptotic DC) Set $\kappa = d/n_{\text{fs}} > 1$. Given $\sigma > 0$, covariance $\Sigma_{\mathbf{X}}$ and latent vector β and define the unique non-negative terms $\xi, \gamma, z \in \mathbb{R}^d$ and $\phi \in \mathbb{R}^d$ as follows:

$$\begin{aligned} \xi > 0 \quad \text{is the solution of} \quad \kappa^{-1} &= d^{-1} \sum_{i=1}^d (1 + (\xi \Sigma_{\mathbf{X}_i})^{-1})^{-1}, \\ \gamma &= \frac{\sigma_{\varepsilon}^2 + \sum_{i=1}^d \frac{\Sigma_{\mathbf{X}_i} \beta_i^2}{(1 + \xi \Sigma_{\mathbf{X}_i})^2}}{1 - \frac{\kappa}{d} \sum_{i=1}^d (1 + (\xi \Sigma_{\mathbf{X}_i})^{-1})^{-2}} \end{aligned}$$

Let $\mathbf{h} \sim \mathcal{N}(0, \mathbf{I}/d)$. The non-asymptotic distributional prediction is given by the following random vector

$$\hat{\beta}(\Sigma_{\mathbf{X}}) = \frac{1}{1 + (\xi \Sigma_{\mathbf{X}})^{-1}} \odot \beta + \frac{\sqrt{\kappa \gamma} \Sigma_{\mathbf{X}}^{-1/2}}{1 + (\xi \Sigma_{\mathbf{X}})^{-1}} \odot \mathbf{h}.$$

Note that, the above formulas can be slightly simplified to have a cleaner look by introducing an additional variable $z = \frac{1}{1 + (\xi \Sigma_{\mathbf{X}})^{-1}}$.

Also note that, the terms in the non-asymptotic DC and asymptotic DC have one to one correspondence. Non-asymptotic DC is essentially a discretized version of asymptotic DC where instead of expectations (which is integral over pdf) we have summations.

Now, we can use this distribution to predict the test risk by substituting $\hat{\beta}(\Sigma_{\mathbf{X}})$ distribution in (3.6).

B.3.1. FINDING OPTIMAL REPRESENTATION Λ^*

Going back to representation question, without losing generality, assume $\Sigma_X = I$ and let us find optimal Λ^* .

$$\hat{\beta}_{\Lambda^*} = \Lambda^* \left[\frac{1}{1 + (\xi \Lambda^{*2})^{-1}} \odot \Lambda^{*-1} \beta + \frac{\sqrt{\kappa \gamma} \Lambda^{*-1}}{1 + (\xi \Lambda^{*2})^{-1}} \odot \mathbf{h} \right].$$

The risk is given by (using $\mathbf{h} \sim \mathcal{N}(0, I_p)$)

$$\text{risk}^{\Lambda^*}(\hat{\beta}) = \mathbb{E}[(\hat{\beta}^{\Lambda^*} - \beta)^\top \Sigma_X (\hat{\beta}^{\Lambda^*} - \beta)] \quad (\text{B.13})$$

$$= \sum_{i=1}^d \frac{\Sigma_{\beta_i}^2}{(1 + \xi(\Lambda_i^*)^2)^2} + \frac{1}{d} \sum_{i=1}^d \frac{\kappa \gamma}{(1 + (\xi(\Lambda_i^*)^2)^{-1})^2} \quad (\text{B.14})$$

Here, note that ξ is function of Λ^* and γ is function of β, Λ^* . If we don't know Σ_β , we use the estimation from representation learning $\hat{\Sigma}_\beta$ instead.

Non-asymptotic question: Find the optimal $\Lambda^* \in \mathbb{R}^p$ vector minimizing risk.

Optimal Representation: To find the optimal representation, we will solve the following optimization problem that minimizes the risk.

$$\begin{aligned} \min_{\Lambda^*} \quad & \sum_{i=1}^d \frac{\beta_i^2}{(1 + \xi(\Lambda_i^*)^2)^2} + \frac{1}{d} \sum_{i=1}^d \frac{\kappa \gamma}{(1 + (\xi(\Lambda_i^*)^2)^{-1})^2} \\ \text{s.t.} \quad & \kappa^{-1} = \frac{1}{d} \sum_{i=1}^d (1 + (\xi(\Lambda_i^*)^2)^{-1})^{-1} \end{aligned} \quad (\text{B.15})$$

$$\gamma = \frac{\sigma_\varepsilon^2 + \frac{1}{d} \sum_{i=1}^d \frac{\beta_i^2}{(1 + \xi(\Lambda_i^*)^2)^2}}{1 - \frac{\kappa}{d} \sum_{i=1}^d (1 + (\xi(\Lambda_i^*)^2)^{-1})^{-2}}$$

So we plug in the expression of γ and get

$$\kappa \gamma = \frac{\sigma_\varepsilon^2 + \frac{1}{d} \sum_{i=1}^d \frac{\beta_i^2}{(1 + \xi(\Lambda_i^*)^2)^2}}{\kappa^{-1} - \frac{1}{d} \sum_{i=1}^d (1 + (\xi(\Lambda_i^*)^2)^{-1})^{-2}} = \frac{\sigma_\varepsilon^2 + \frac{1}{d} \sum_{i=1}^d \frac{\beta_i^2}{(1 + \xi(\Lambda_i^*)^2)^2}}{\frac{1}{d} \sum_{i=1}^d \frac{\xi(\Lambda_i^*)^2}{(1 + \xi(\Lambda_i^*)^2)^2}} \quad (\text{B.16})$$

Let $\theta_i = \frac{\xi(\Lambda_i^*)^2}{1 + \xi(\Lambda_i^*)^2}$ then the objective function becomes

$$\frac{1}{d} \sum_{i=1}^d \Sigma_{\beta_i} (1 - \theta_i)^2 + \left(\sum_{i=1}^d \theta_i^2 \right) \frac{\sigma_\varepsilon^2 + \frac{1}{d} \sum_{i=1}^d \Sigma_{\beta_i} (1 - \theta_i)^2}{\sum_{i=1}^d \theta_i (1 - \theta_i)} = \frac{\frac{n_{\text{fs}}}{d} (\sum_{i=1}^d \Sigma_{\beta_i} (1 - \theta_i)^2) + \sigma_\varepsilon^2 (\sum_{i=1}^d \theta_i^2)}{n_{\text{fs}} - \sum_{i=1}^d \theta_i^2}$$

such that $0 \leq \theta_i < 1$ and $\sum_{i=1}^d \theta_i = \frac{d}{\kappa} = n_{\text{fs}}$.

Solving the optimization problem. We use $\Sigma_{\beta_{i,i}}$ to replace β_i for computing the risk. The objective function is

$$f(\theta) = \frac{1}{d} \sum_{i=1}^d \Sigma_{\beta_i} (1 - \theta_i)^2 + \left(\sum_{i=1}^d \theta_i^2 \right) \frac{\sigma_\varepsilon^2 + \frac{1}{d} \sum_{i=1}^d \Sigma_{\beta_i} (1 - \theta_i)^2}{\sum_{i=1}^d \theta_i (1 - \theta_i)}. \quad (\text{B.17})$$

Lemma 5 Let $C, S, V \in \mathbb{R}$. Define

$$\phi(\Sigma_{\beta_i}; C, V, S) := \frac{Cp(d - n_{\text{fs}} - S)^2}{2n_{\text{fs}}(V + d\sigma_\varepsilon^2 + (d - n_{\text{fs}} - S)\Sigma_{\beta_i}^2)}$$

and we find the root of the following equations:

$$\begin{aligned}\sum_{i=1}^d \phi(\Sigma_{\beta_i}; C, V, S) &= d - n_{\text{fs}}, \\ \sum_{i=1}^d \phi^2(\Sigma_{\beta_i}; C, V, S) &= S - (2n_{\text{fs}} - d), \\ \sum_{i=1}^d \Sigma_{\beta_i} \phi^2(\Sigma_{\beta_i}; C, V, S) &= V\end{aligned}$$

Let $\theta_i = 1 - \phi(\Sigma_{\beta_i}; C^*, V^*, S^*)$ where C^*, V^*, S^* are the roots, then

$$\theta = \arg \min_{\theta'} f(\theta'), \quad \text{s.t., } 0 \leq \theta' < 1, \quad \sum_{i=1}^d \theta'_i = n_{\text{fs}}.$$

Proof Define $s = \sum_{i=1}^d \theta_i^2$, $\phi_i = 1 - \theta_i$. Define $Q = \frac{1}{d} \sum_{i=1}^d \Sigma_{\beta_i} \phi_i^2$. Then

$$\begin{aligned}f(\phi) &= \frac{1}{d} \sum_{i=1}^d \Sigma_{\beta_i} \phi_i^2 + \frac{s}{n_{\text{fs}} - s} (\sigma_\varepsilon^2 + \frac{1}{d} \sum_{i=1}^d \Sigma_{\beta_i} \phi_i^2) \\ &= Q + \frac{s}{n_{\text{fs}} - s} (\sigma_\varepsilon^2 + Q) \\ &= \frac{n_{\text{fs}}}{d - n_{\text{fs}} - \sum_{i=1}^d \phi_i^2} (Q + \sigma_\varepsilon^2).\end{aligned}$$

The last line uses

$$s = \sum_{i=1}^d (1 - \phi_i^2) = d - 2 \sum_{i=1}^d \phi_i + \sum_{i=1}^d \phi_i^2 = d - 2(d - n_{\text{fs}}) + \sum_{i=1}^d \phi_i^2 = 2n_{\text{fs}} - d + \sum_{i=1}^d \phi_i^2.$$

Now define $\sum_{i=1}^d \phi_i^2 = S$, and we compute the gradient of f , we have

$$\frac{df}{d\phi_i} = \left(\frac{2n_{\text{fs}}}{d} (\sum \Sigma_{\beta_j} \phi_j^2 + (d - n_{\text{fs}} - s) \Sigma_{\beta_i}) + 2n_{\text{fs}} \sigma_\varepsilon^2 \right) \phi_i.$$

Suppose $0 < \phi_i < 1$, then we need $\frac{df}{d\phi_i}$ equal to each other for all i . Suppose $\frac{df}{d\phi_i} = C$, and denote $\sum \Sigma_{\beta_j} \phi_j^2 = V$, we can solve for ϕ_i from $\frac{df}{d\phi_i} = C$ as

$$\phi_i = \frac{Cd(d - n_{\text{fs}} - S)^2}{2n_{\text{fs}}(V + d\sigma_\varepsilon^2 + (d - n_{\text{fs}} - S)\Sigma_{\beta_i}^2)} := \phi(\Sigma_{\beta_i}; C, V, S). \quad (\text{B.18})$$

We define the function $\phi(\Sigma_{\beta_i}; C, V, S)$ as above, and use the fact that

$$\begin{aligned}\sum_{i=1}^d \phi(\Sigma_{\beta_i}; C, V, S) &= d - n_{\text{fs}}, \\ \sum_{i=1}^d \phi^2(\Sigma_{\beta_i}; C, V, S) &= S - (2n_{\text{fs}} - d), \\ \sum_{i=1}^d \Sigma_{\beta_i} \phi^2(\Sigma_{\beta_i}; C, V, S) &= V\end{aligned}$$

We can solve⁷ C, V, S and retrieve ϕ_i by (B.18). $\theta_i = 1 - \phi_i$. ■

⁷For the root of 3-dim problem, the worst case we can grid the space and search with time complexity $\mathcal{O}(\varepsilon^{-3})$.

B.4. Proof of Robustness of Optimal Representation

B.4.0.1 Optimal Shaping after R -SVD.

First we discuss the impact of R -SVD truncation.

Let \mathbf{x}_R be the projection of \mathbf{x} onto the R -dimensional subspace spanned by columns of \mathbf{U}_1 , and \mathbf{x}_{R^\perp} is the projection of \mathbf{x} onto the orthogonal complement. Namely, $\mathbf{x}_R = \mathbf{U}_1^\top \mathbf{x} \in \mathbb{R}^R$ and $\mathbf{x}_{R^\perp} = \mathbf{U}_2^\top \mathbf{x} \in \mathbb{R}^{(d-R)}$. Similarly we can define β_R and β_{R^\perp} . Thus,

$$\mathbf{y} = \mathbf{x}^\top \beta + \varepsilon = \mathbf{x}_R^\top \beta_R + \mathbf{x}_{R^\perp}^\top \beta_{R^\perp} + \varepsilon \quad (\text{B.19})$$

We can treat $\varepsilon_R = \mathbf{x}_{R^\perp}^\top \beta_{R^\perp} + \varepsilon$ as the new noise, and try to solve for β_R . Then our noise variance becomes $\sigma_{\varepsilon_R}^2 = \sigma_\varepsilon^2 + \text{tr}(\Sigma_X \Sigma_{\beta_R}^\perp)$. If we are still overparameterized regime, namely $R > n_{\text{fs}}$, then we define optimal representation on top of it. Let's $\kappa_R = R/n_{\text{fs}} > 1$.

In summary, the R -SVD truncation reduces the search space of Λ into $R \times R$ dimensional PSD matrix, where the covariance of the noise in \mathbf{y} increases from $\sigma_\varepsilon^2 \mathbf{I}$ to $(\sigma_\varepsilon^2 + \text{tr}(\Sigma_X \Sigma_{\beta_R}^\perp)) \mathbf{I}$. In the following theorem, we give the proof for d dimensional case with noise level being $\sigma_\varepsilon^2 \mathbf{I}$, and it directly generalize to R dimensional regime by the equivalence above.

In the following proof, we only study the full dimensional case, where $\Lambda \in \mathbb{R}^{d \times d}$. The computation can be generalized to R dimensional case as discussed above by performing R -SVD on Σ_β and replacing the noise with variance $(\sigma_\varepsilon^2 + \text{tr}(\Sigma_X \Sigma_{\beta_R}^\perp)) \mathbf{I}$.

Theorem 4 Suppose the data is generated as Definition 3, Λ and θ are defined in Lemma 5 and the estimated task is obtained as (3.4). Suppose $\|\hat{\Sigma}_\beta - \Sigma_\beta\| \leq \delta_{\Sigma_\beta}$. Then the risk of few-shot learning phase suffers at most

$$\text{risk}(\Lambda, \Sigma_\beta) - \text{risk}(\Lambda^*, \Sigma_\beta) \leq \frac{2n_{\text{fs}}^2 \delta_{\Sigma_\beta}}{(d - n_{\text{fs}})(2n_{\text{fs}} - d\theta)\underline{\theta}}.$$

Proof We first decompose the risk as

$$\begin{aligned} & \text{risk}(\Lambda, \Sigma_\beta) - \text{risk}(\Lambda^*, \Sigma_\beta) \\ &= \underbrace{\text{risk}(\Lambda, \hat{\Sigma}_\beta) - \text{risk}(\Lambda^*, \hat{\Sigma}_\beta)}_{\leq 0} + [\text{risk}(\Lambda, \Sigma_\beta) - \text{risk}(\Lambda, \hat{\Sigma}_\beta)] + [\text{risk}(\Lambda^*, \hat{\Sigma}_\beta) - \text{risk}(\Lambda^*, \Sigma_\beta)]. \end{aligned}$$

We know $\text{risk}(\Lambda, \hat{\Sigma}_\beta) - \text{risk}(\Lambda^*, \hat{\Sigma}_\beta) \leq 0$ due to the optimality of Λ with task covariance $\hat{\Sigma}_\beta$. Now we will bound $\text{risk}(\Lambda, \Sigma_\beta) - \text{risk}(\Lambda, \hat{\Sigma}_\beta)$ and it automatically works for $\text{risk}(\Lambda^*, \hat{\Sigma}_\beta) - \text{risk}(\Lambda^*, \Sigma_\beta)$. Note that in (3.11) we know that

$$\text{risk}(\Lambda', \Sigma'_\beta) = f(\theta; \Sigma_\beta) := \sum_{i=1}^d \frac{n_{\text{fs}}(1 - \theta_i)^2}{d(n_{\text{fs}} - \|\theta\|^2)} \Sigma_{\beta_i} + \frac{\|\theta\|^2}{n_{\text{fs}} - \|\theta\|^2} \sigma_\varepsilon^2. \quad (\text{B.20})$$

This function is linear in Σ_β thus we know that

$$|\text{risk}(\Lambda^*, \hat{\Sigma}_\beta) - \text{risk}(\Lambda^*, \Sigma_\beta)| \leq \frac{n_{\text{fs}}}{n_{\text{fs}} - \|\theta\|^2} \delta_{\Sigma_\beta}. \quad (\text{B.21})$$

Now we need to bound $\|\theta\|^2$. With the constraint $\underline{\theta} \leq \theta < 1 - \frac{d - n_{\text{fs}}}{n_{\text{fs}}} \underline{\theta}$ and $\sum \theta_i = n_{\text{fs}}$, we know that the maximum of $\|\theta\|^2$ happens when $(d - n_{\text{fs}})$ among θ_i are $\underline{\theta}$ and the others are $1 - \frac{d - n_{\text{fs}}}{n_{\text{fs}}} \underline{\theta}$. With this we have

$$\begin{aligned} \|\theta\|^2 &\leq (d - n_{\text{fs}})\underline{\theta}^2 + n_{\text{fs}}\left(1 - \frac{d - n_{\text{fs}}}{n_{\text{fs}}}\underline{\theta}\right)^2 \\ &= (d - n_{\text{fs}})\underline{\theta}^2 + n_{\text{fs}} - 2(d - n_{\text{fs}})\underline{\theta} + \frac{(d - n_{\text{fs}})^2}{n_{\text{fs}}}\underline{\theta}^2 \\ &= n_{\text{fs}} - 2(d - n_{\text{fs}})\underline{\theta} + \frac{(d - n_{\text{fs}})d}{n_{\text{fs}}}\underline{\theta}^2 \end{aligned}$$

Thus

$$n_{\text{fs}} - \|\theta\|^2 \geq (d - n_{\text{fs}})\underline{\theta}(2n_{\text{fs}} - d\underline{\theta}).$$

Plugging it into (B.21) and (B.20) leads to the theorem. ■