

Towards sample-efficient overparameterized meta-learning

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.

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 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

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.

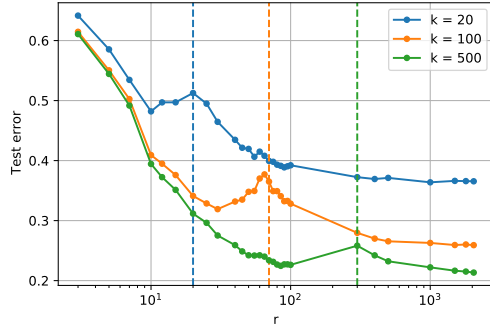


Figure 1. An illustration of the benefit of overparameterization for few-shot learning: A pretrained ResNet50 network on Imagenet (representation learning phase) was utilized to obtain a feature-representation for classification on CIFAR-10. After reducing the dimensionality to r via PCA, the resultant 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 ($k = r$).

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 $\Lambda \in \mathbb{R}^{R \times d}$?

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* Σ_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 Σ_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 ($\Sigma_X = I$), 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* Λ_{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, Λ_{opt} shapes the feature covariance to emphasize informative features over less informative ones. In practice, Λ_{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 Λ_{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 Λ_{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 $(x_{ij}, y_{ij}) \in \mathbb{R}^d \times \mathbb{R}$, and noise $\varepsilon_{ij} \in \mathbb{R}$. The distribution of β_i and x_{ij} are $\mathcal{N}(0, \Sigma_\beta)$ and $\mathcal{N}(0, \Sigma_X)$. The rank of Σ_β and Σ_X are r_t and r_f . The few-shot learning (downstream) task is β and samples are (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 $\Lambda \in \mathbb{R}^{R \times d}$. 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 Λ 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 $\Lambda \in \mathbb{R}^{R \times d}$ as $\Lambda \mathbf{x}$.

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 Λ 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 $\hat{\beta}_\Lambda$ 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 $\Lambda \in \mathbb{R}^{R \times d}$ be the projection to the span of top R singular values of moment estimator.
else if Shaping then
 Calculate Λ as defined in Def. 5 and (3.12).
end if
Few-shot learning:
 Calculate $\hat{\beta}_\Lambda$ from (3.4).
Return: $\hat{\beta}_\Lambda$.

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 general, recovery guarantees for arbitrary covariance matrices Σ_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 $\Sigma_X = I$, 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 Σ_X is approximately low rank and its principal eigenspace cov-

ers the top eigenvectors of Σ_β . Based on this, we consider the following estimators. Suppose the rank of Σ_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).

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

The mean of this estimator is $M = \Sigma_X \Sigma_\beta \Sigma_X$ 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).

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

The mean of this estimator is $Q = 2\Sigma_X \Sigma_\beta \Sigma_X + \text{tr}(\Sigma_\beta \Sigma_X) \Sigma_X$ 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 $\Sigma_X = I$, 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 Σ_X and Σ_β are approximately rank r_t and the span of top r_t eigenvectors of Σ_X and Σ_β are closely aligned, we can directly estimate the covariance of features by

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

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.

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

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

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 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 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

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

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 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 $\Sigma_X = \text{diag}(\mathbf{I}_{r_f}, \iota \mathbf{I}_{d-r_f})$, and $\Sigma_\beta = \text{diag}(\mathbf{I}_{r_t}, 0)$, $\sigma_\varepsilon = 0$, 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., $\Sigma_X = \mathbf{I}$, we need $\tilde{\mathcal{O}}(dr_t^2)$ samples to estimate 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 $\tilde{\mathcal{O}}(r_f r_t^2)$ samples to estimate 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).

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.

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 $x \sim \mathcal{N}(0, \mathbf{I})$, and each task contains $n_{\text{spt}} \gtrsim r_t$ samples, then we can get the span of Σ_β with $\tilde{\mathcal{O}}(dr_t)$ samples.

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

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

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

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

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

In the appendix, we will propose a theorem with general feature covariance Σ_X and noisy data, as a generalization of Theorem 2.

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{B} , then $n_{\text{task}} \hat{Q} = \hat{B} \hat{B}^\top$. 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 $1 - (n_{\text{tot}} \text{tr}(\Sigma_X))^{-10}$,

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

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

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

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

Lemma 1 gives the quality of the estimation of the covariance of features \mathbf{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 $(\mathbf{x}_i, y_i)_{i=1, \dots, n_{\text{fs}}}$ where $\mathbf{x}_i \sim \mathcal{N}(0, \Sigma_X)$, and $y_i = \mathbf{x}_i^\top \beta + \varepsilon_i$ where $\varepsilon_i \sim \mathcal{N}(0, \sigma_\varepsilon^2)$. We assume $n_{\text{fs}} < d$. Denote $\mathbf{X} \in \mathbb{R}^{n_{\text{fs}} \times d}$ whose i th row is \mathbf{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

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

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

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

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

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

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 θ .⁴

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

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

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 $0 < \underline{\theta} < n_{\text{fs}}/d$,

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

in implementation we require $\underline{\theta} \leq \theta \leq 1 - \frac{d-n_{\text{fs}}}{n_{\text{fs}}} \underline{\theta}$ instead of $0 \leq \theta \leq 1$ 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 $X_R = X U_1 \in \mathbb{R}^{n_{\text{fs}} \times R}$. 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 folloing quantities.

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

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

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

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

Define $\Sigma_\beta^\perp, \Sigma_X^\perp$ as the projection of Σ_β, Σ_X onto U_2 , the noise variance is equivalent to $\sigma_{\varepsilon R}^2 = \sigma_\varepsilon^2 + \text{tr}(\Sigma_\beta^\perp \Sigma_X^\perp)$. 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 $\Sigma_X = I$, 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 .

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 $\Sigma_X \neq I$, logically $\hat{M} \neq \Sigma_\beta$ and this estimation does not always work. In some special cases, such as $\Sigma_X = \text{diag}(I_{r_f}, \iota I_{d-r_f})$ and $\Sigma_\beta = \text{diag}(I_{r_t}, 0)$ ($r_f > r_t$), \hat{M} is still a good estimator of Σ_β .

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 $\Sigma_X = I$. 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

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

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

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

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

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. $\text{risk}(\Lambda_R^*, \Sigma_\beta) - \text{risk}(\Lambda^*, \Sigma_\beta)$.

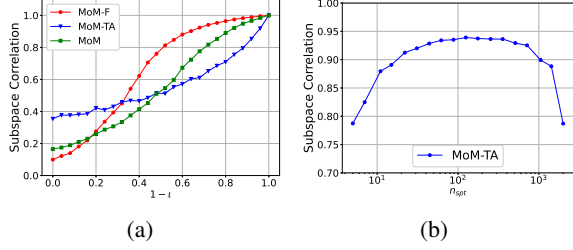


Figure 2. (a) $\Sigma_\beta = (I_{10}, \mathbf{0}_{90})$. $\sigma_\varepsilon = 0.5$, $\Sigma_X = (I_{10}, \iota \cdot I_{90})$, $n_{\text{task}} = 20$, $n_{\text{spt}} = 40$. Learning 10 dimensional top eigenspace of Σ_β with biased features. MoM-F, MoM-TA, MoM stands for Σ_X , \hat{B} , \hat{Q} . Subspace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

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 of $\mathcal{O}(d r_t^2)$ 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 $\Sigma_X = I$. 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 $\mathcal{O}(d r_t)$ 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{Q} which asks for $\mathcal{O}(d r_t^2)$ 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

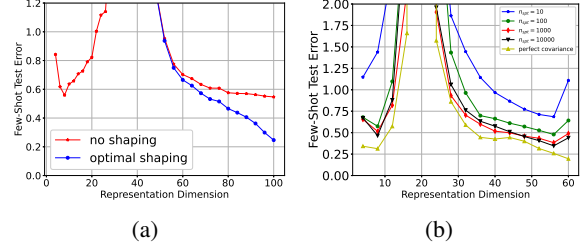


Figure 3. (a) $\Sigma_\beta = (25 \cdot I_{10}, I_{90})$, $\Sigma_X = I_{100}$, $\sigma_\varepsilon = 0.5$, $n_{\text{fs}} = 40$, few shot learning with knowledge of Σ_β , with identity or optimal shaping matrix. For $R < n_{\text{fs}}$, we are not in overparameterized regime, so there is no optimal shaping. (b) Algorithm 1 with $d = 60$, $n_{\text{fs}} = 20$, $n_{\text{task}} = 60$, $\Sigma_\beta = (25 \cdot I_6, I_{54})$, $\Sigma_X = I_{60}$ 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.

($n_{\text{fs}} < R$), 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 Λ , calculated as a function of $\hat{\Sigma}_\beta$, for solving the downstream task. Both dimension reduction ($n_{\text{fs}} > R$) or optimal overparameterized problem ($n_{\text{fs}} < R$) end up with small risk, whereas it becomes hard to learn when $n_{\text{fs}} \approx R$. 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.
- 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.