

Sample Efficient Linear Meta-Learning by Alternating Minimization

Kiran Koshy Thekumparampil[†], Prateek Jain[‡], Praneeth Netrapalli[‡], Sewoong Oh[±] *

[†]University of Illinois at Urbana-Champaign, [‡]Google Research, India,
[±]University of Washington, Seattle

Abstract

Meta-learning synthesizes and leverages the knowledge from a given set of tasks to rapidly learn new tasks using very little data. Meta-learning of linear regression tasks, where the regressors lie in a low-dimensional subspace, is an extensively-studied fundamental problem in this domain. However, existing results either guarantee highly suboptimal estimation errors, or require $\Omega(d)$ samples per task (where d is the data dimensionality) thus providing little gain over separately learning each task. In this work, we study a simple alternating minimization method (MLLAM), which alternately learns the low-dimensional subspace and the regressors. We show that, for a constant subspace dimension MLLAM obtains nearly-optimal estimation error, despite requiring only $\Omega(\log d)$ samples per task. However, the number of samples required per task grows logarithmically with the number of tasks. To remedy this in the low-noise regime, we propose a novel task subset selection scheme that ensures the same strong statistical guarantee as MLLAM, even with *bounded* number of samples per task for arbitrarily large number of tasks.

1 Introduction

Common real world tasks follow a long tailed distribution where most of the tasks only have a small number of labelled examples [WRH17]. Collecting more clean labels is often costly (e.g., medical imaging). As each task does not have enough examples to be learned in isolation, meta-learning attempts to meta-learn across a large number of tasks by exploiting some structural similarities among those tasks. One popular approach is to learn a shared representation, where all of those tasks can be solved accurately [SSSG17]. Once such a representation has been learnt, we can rapidly adapt to new arriving tasks, learning a model with only a few examples. Empirical evidences suggest that this might also explain recent successes in few-shot supervised learning with optimization based methods like MAML [FAL17; RRBV19a].

In this paper, we study the problem of linear meta-representation learning [Du+20; TJJ20], where the goal is to learn a r -dimensional linear representation/subspace that is shared by a collection of t linear regression tasks in d dimensions. Each task has m labelled examples.

We investigate a fundamental question: as the number of tasks grow, can we learn the underlying r -dimensional shared representation (subspace) more accurately, and consequently learn more accurate regressors per task? The question is important because in general, the number of tasks can be large while a lot of tasks are data starved. Furthermore, in several settings like crowdsourcing or bioinformatics, it might be easier to collect more data for new tasks, instead of collecting more data for the existing tasks.

Most of the existing work do not provide a satisfactory solution for this fundamental problem. In particular, Du et al. [Du+20] require $m = \Omega(d)$ samples per task, which is prohibitively large, and in fact with so many samples, one can solve each task in isolation. The Burer-Monteiro factorization approach of Tripuraneni, Jin, and Jordan [TJJ20] is not able to provide any improvement by increasing the number of tasks, it needs to increase samples per task to $m = \Omega(1/\varepsilon^2)$ to achieve ε accuracy. While the method-of-moments approach proposed in [TJJ20; Kon+20] does provide more accurate representation learning with a larger number of tasks, the method has a highly sub-optimal dependence on the noise variance σ^2 associated with each task. For example, even when each regression task can be solved exactly with 0 error, this method will incur a significant error.

*Author emails are thekump2@illinois.edu, prajain@google.com, pnetrapalli@google.com, and sewoong@cs.washington.edu.

Table 1: Comparison of high-probability error bounds for the distance between the learned (U) and the true (U^*) subspaces, for linear low-rank meta-learning in d dimensions with t tasks, m samples per task, and noise variance σ^2 . Note that \tilde{O} and $\tilde{\Omega}$ hides polylog factors in d and $\log \log$ factors in t . We assume a constant small subspace rank, constant incoherence of tasks, constant magnitude for regressors, and well-conditioned task diversity matrix. Note that non-convex ERM [Du+20] is a result for general non-linear meta-learning.

| Algorithm | Error-bound: $\ (\mathbf{I} - U^*(U^*)^\top)U\ $ | Required samples per task |
|--------------------------------------|--|--|
| Non-convex ERM [Du+20] | $\tilde{O}(\sigma)\sqrt{\frac{t+d}{m t}}$ | $m \geq \tilde{\Omega}(d + \log(t))$ |
| Burer-Monteiro factorization [TJJ20] | $\tilde{O}(\sigma)\sqrt{\frac{\max(t,d)}{m t}}$ | $m \geq \tilde{\Omega}(\log(t))$ |
| Method-of-Moments [TJJ20] | $\tilde{O}(1 + \sigma^2)\sqrt{\frac{d}{m t}}$ | $\Omega(1)$ |
| MLLAM (ours, Theorem 1) | $\tilde{O}(\sigma)\sqrt{\frac{d}{m t}}$ | $m \geq \tilde{\Omega}((1 + \sigma^2) \log t)$ |
| MLLAMS (ours, Theorem 3) | $\tilde{O}(\sigma)\sqrt{\frac{d}{m t}}$ | $m \geq \tilde{\Omega}(1 + \sigma^2 \log(t))$ |
| Lower-bound [TJJ20] | $\Omega(\sigma)\sqrt{\frac{d}{m t}}$ | $\Omega(1)$ |

Contributions. In this paper, we propose the first efficient approach for linear meta-learning with provable guarantees that achieves nearly optimal error rate. According to a Frobenius norm error metric, our bound matches a fundamental lower bound. Our first algorithm MLLAM is based on alternating minimization, inspired by a long line of successes in matrix completion and matrix sensing [JNS13]. Assuming constant dimensionality of the representation, MLLAM requires $m = \Omega(\log t + \log \log(1/\varepsilon))$ samples per task to achieve an accuracy of ε when we have t tasks. Our method obtains nearly optimal dependence on the noise variance σ^2 and the error in representation learning drops nearly optimally with growing number of tasks, which is a significant improvement over the state-of-the-art.

However, the number of samples per task (m) still grows logarithmically on t . To further improve this dependence, we introduce MLLAMS that applies alternating minimization to only a subset of tasks that are well-behaved. When the noise is sufficiently small with variance $O(1/\log t)$, this further reduces the requirement down to $m = \Omega(\log \log(1/\varepsilon))$. That is, despite fixed m , MLLAMS can estimate each task more accurately. Furthermore, due to our improved rates on estimation of the r -dimensional subspace, the best known rates for prediction error on new tasks also improve significantly. Table 1 compares the error and per-task sample complexity of our method against existing state-of-the-art results; see Section 2 for details of the problem setting.

Broadly, our proof structure follows that of existing alternating minimization results [JNS13; NJS15] by showing iterative refinement of the estimates. However, existing techniques are able to rely on restricted isometry property style properties, which are significantly more difficult to prove in our case. Furthermore, most of the existing works in this literature analyze non-noisy setting where each observation is sampled exactly from an underlying model. However, in this work, we also allow each observation to be corrupted by a white noise, leading to more challenging per-iterate analysis.

Notations: For an whole number n , $[n] = \{1, \dots, n\}$. $\|A\|$ and $\|A\|_F$ denote the spectral and Frobenius norms of the matrix A . $\langle A, B \rangle$ denotes inner produce between two matrices. A^\dagger is the Moore-Penrose pseudoinverse and A^\top is the transpose of the matrix A . $x \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$ means that x is d dimensional standard isotropic Gaussian random vector.

1.1 Related work

Representation learning for meta-learning. There is a large body of work in meta-learning from multiple tasks since the seminal work in learning to learn [TP98], inductive bias learning [Bax00], and multitask learning [Car97]. One popular line of work starting from [IE96; Bax95] is to learn a low-dimensional representation for a set of related tasks and use the representation to efficiently train a model for a new arriving task. Recently, these representation learning approaches are gaining more attraction as recent empirical evidence indicates that the success of other popular meta-learning approaches such as MAML [FAL17] is due to their capability to learn useful low-dimensional representations [RRBV19b].

[AZ05; Ris+08; Orl05] address the problem of recovering linear regression parameters that lie on an unknown r -dimensional subspace U^* , where all tasks can be accurately solved. Nuclear-norm minimization ap-

proaches are proposed in [AEP08; Har+12; AFSU07; PM13] but they do not provide subspace/generalization error guarantees and suffer from large training time.

Closest to our work is [TJJ20] that analyzes the landscape of the empirical risk with Burer-Monteiro factorization. It is shown that $mt = \tilde{\Omega}(\max\{t, d\}r^4 + \max\{t, d\}r^2\sigma^2/\epsilon^2)$ is sufficient to achieve a rescaled error $(1/\sqrt{r})\|(\mathbf{I} - U^*(U^*)^T)U\|_F \leq \epsilon$, where we assumed incoherent and well-conditioned regression parameters to simplify the condition. In particular, more tasks do not give any gain beyond a certain point if m is fixed. Further, it is also required that all tasks are of equal strengths, i.e., $\|v^{(i)}\| = \Theta(1)$ for all $i \in [t]$. Another approach is to find the principal directions of a particular 4th moment matrix [TJJ20; Kon+20]. This only requires $mt = \tilde{\Omega}((1 + \sigma^2)dr^2/\epsilon^2)$, but the algorithm is inexact; the error is bounded away from zero even if there is no noise and sample size is sufficiently large to learn all the parameters. This is in a stark contrast with our approach, as illustrated in Figure 1a. Du et al. [Du+20] studies the global minimizer of a non-convex optimization in Eq. (3) without analyzing an efficient algorithm to find it. It is shown that a small generalization error can be achieved if $m = \tilde{\Omega}(d)$.

We also point out a concurrent and independent work [CHMS21], which proposes and analyzes a slightly different variant (with descent step on U) [CHMS21, Algorithm 2] of our alternating minimization algorithm (Algorithm 1, MLLAM) for a similar linear meta-learning setting. However, this work assumes that the linear meta-learning problem is noiseless, i.e. $\sigma = 0$ (as defined in our Assumption 1), and then it provides a per task sample complexity of $m \geq \tilde{\Omega}((\lambda_1^*/\lambda_r^*)^2 r^3 \log(t))$ and a total sample complexity of $mt \geq \tilde{\Omega}((\lambda_1^*/\lambda_r^*)^2 dr^2)$. In contrast, our results are for a more natural noisy setting, and even for the noiseless setting we obtain a tighter per task sample complexity of $m \geq \tilde{\Omega}(r^2)$ and a total sample complexity of $mt \geq \tilde{\Omega}((\lambda_1^*/\lambda_r^*)(d+r^2)r^2)$ (Corollary 4) in terms of the condition number $(\lambda_1^*/\lambda_r^*)$ and the rank r , and our per task complexity does not scale with the number of tasks t . Collins, Hassani, Mokhtari, and Shakkottai [CHMS21] further show that alternating minimization performs better than other baselines for personalized federated learning of neural network classifiers for some datasets.

Matrix sensing. Starting from matrix sensing and completion problems [CR09; MJD09; JNS13], recovering a low-rank matrix from linear measurements have been a popular topic of research. Linear meta-learning is a special case of matrix sensing, but with special sensing operator of the form $\mathcal{A}(UV^T) = [A_1(UV^T), \dots, A_{mt}(UV^T)]$ where $A_{ij}(UV^T) = \langle x_{ij}e_i^\top, UV^\top \rangle$. This operator cannot satisfy sensing properties like restricted isometry property, in general because of sparse sensing matrix, so existing matrix sensing results do not apply directly. Furthermore, [JD13; ZJD15] studied a similar problem but their results also require $O(d)$ samples per task, which limits its applicability to the meta-learning setting where each task has a small number of samples.

2 Problem formulation

Suppose there are t d -dimensional linear regression tasks, and each of them have m samples. That is for the i -th task ($i \in [t]$), we are given m samples $\{(x_j^{(i)} \in \mathbb{R}^d, y_j^{(i)} \in \mathbb{R})\}_{j=1}^m$, where $(x_j^{(i)}, y_j^{(i)})$ is the j -th pair of example and observation. The standard goal is to learn accurate regressors $\tilde{v}^{*(i)}$ for each of the tasks. However, in the meta-learning setting, all the tasks are related and share a common but unknown *low-dimensional representation* parameterized by $U^* \in \mathbb{R}^{d \times r}$ where $r \ll d$. Here, the goal is to learn U^* and the task specific regressors $v^{(i)}$ s.t. $v^{(i)}$'s are accurate regressor for samples $\{((U^*)^\top x_j^{(i)} \in \mathbb{R}^r, y_j^{(i)} \in \mathbb{R})\}_{j=1}^m$, for $i \in [t]$. This is equivalent to finding a set of accurate regressors $\tilde{v}^{(i)}$'s, which lie in a low-dimensional subspace.

A natural requirement of the problem is to then learn the tasks accurately with very small number of samples per task, especially for large t . As a task specific regressor has only r parameters, given U^* , we expect the number of samples per task to depend only on r , instead of d . Furthermore, the total number of samples $m \cdot t$ should scale at most linearly with the data dimension d . However, simultaneously learning the representation U and the regressors $v^{(i)}$ is challenging. In fact, since the NP-hard low-rank matrix completion problem [HMRW14] can be reduced to the linear meta-learning problem, the latter is NP-hard. Therefore, similar to Tripuraneni, Jin, and Jordan [TJJ20], we study the problem in the following tractable random design setting.

Assumptions 1. Let $U^* \in \mathbb{R}^{d \times r}$ be an orthonormal matrix. For a task $i \in [t]$, with task specific regressor $v^{*(i)} \in \mathbb{R}^r$ and j -th example $x_j^{(i)} \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$, its observation is

$$y_j^{(i)} = \langle x_j^{(i)}, U^* v^{*(i)} \rangle + \varepsilon_j^{(i)}, \quad (1)$$

where $\varepsilon_j^{(i)} \sim \mathcal{N}(0, \sigma^2)$ is the measurement noise which is independent of $x_j^{(i)}$. So, the optimal regressor $\tilde{v}^{*(i)}$ for each task is given by: $\tilde{v}^{*(i)} = U^* v^{*(i)}$. We denote the matrix of the optimal regressors as: $\tilde{V}^* = U^* (V^*)^T$ where $(V^*)^T = [v^{*(1)}, \dots, v^{*(t)}]$.

Assumptions 2. Let λ_1^* and λ_r^* denote the largest and smallest eigenvalues of the task diversity matrix $(r/t)(V^*)^T V^* \in \mathbb{R}^{r \times r}$ respectively. We assume that V^* is μ -incoherent, i.e.,

$$\max_{i \in [t]} \|v^{*(i)}\|^2 \leq \mu \lambda_r^*. \quad (2)$$

Our goal is to recover the subspace U^* , up to a nearly optimal error, from a small number of samples per task. Recovering U^* enables the estimation of the regressor of any new task in the same subspace, using only $\approx O(r)$ samples. To this end, we minimize the empirical risk of parameter matrices $U \in \mathbb{R}^{d \times r}$ and $V = [v^{(1)}, \dots, v^{(t)}]^T \in \mathbb{R}^{t \times r}$:

$$\mathcal{L}(U, V) = \sum_{i=1}^t \sum_{j=1}^m \frac{1}{2} \left(y_j^{(i)} - \langle U v^{(i)}, x_j^{(i)} \rangle \right)^2. \quad (3)$$

The problem is non-convex due to bi-linearity of U and V . \tilde{O} and $\tilde{\Omega}$ hide logarithmic terms in d and r .

3 Main results

Alternating minimization: We first present our main result for a standard alternating minimization method (Algorithm 1) when applied to the meta-learning linear regression problem in the problem setting described in Section 2.

Theorem 1. Let there be t linear regression tasks, each with m samples satisfying Assumptions 1 and 2, and

$$m \geq \tilde{\Omega}((1 + r(\sigma/\sqrt{\lambda_r^*})^2)r \log t + r^2), \quad \text{and} \quad mt \geq \tilde{\Omega}((1 + (\sigma/\sqrt{\lambda_r^*})^2)(\lambda_1^*/\lambda_r^*)\mu dr^2).$$

Then MLLAM (Algorithm 1), initialized at U_{init} s.t. $\|(\mathbf{I} - U^*(U^*)^\top)U_{\text{init}}\|_F \leq \min(3/4, O(\sqrt{\lambda_r^*/\lambda_1^*}))$ and run for $K = \lceil \log_2(\lambda_r^*\lambda_r^* m t / \lambda_1^* \sigma^2 \mu d r^2) \rceil$ iterations, outputs U so that the following holds (w.p. $\geq 1 - K/(dr)^{10}$):

$$\frac{\|(\mathbf{I} - U^*(U^*)^\top)U\|_F}{\sqrt{r}} \leq \tilde{O} \left(\left(\frac{\sigma}{\sqrt{\lambda_r^*}} \right) \sqrt{\frac{\mu r d}{m t}} \right). \quad (4)$$

Remark 1: Our error rate is nearly optimal, as it matches best possible rate when V^* is specified a priori. This is made formal in the following lower bound, which follows from [TJJ20, Theorem 5]. The upper and lower bounds match up to polynomial factors in the incoherence μ and the condition number λ_1^*/λ_r^* .

Corollary 2. [TJJ20, Theorem 5] Let $r \leq d/2$ and $mt \geq r(d - r)$, then for all V^* , w.p. $\geq 1/2$

$$\inf_{\hat{U}} \sup_{U \in G_{r,d}} \frac{\|(\mathbf{I} - U^*(U^*)^\top)\hat{U}\|_F}{\sqrt{r}} \geq \Omega \left(\left(\frac{\lambda_r^*}{\lambda_1^*} \frac{\sigma}{\sqrt{\lambda_r^*}} \right) \sqrt{\frac{dr}{m t}} \right),$$

where $G_{r,d}$ is the Grassmannian manifold of r -dimensional subspaces in \mathbb{R}^d , the infimum for \hat{U} is taken over the set of all measurable functions that takes mt samples in total from the model in Section 2 satisfying Assumption 1 and 2.

Remark 2: To the best of our knowledge, Theorem 1 presents the first efficient method for achieving optimal error rate in σ , d and r . Tripuraneni, Jin, and Jordan [TJJ20] propose two approaches. The first one is the Burer-Monteiro factorization approach, which achieves a rescaled Frobenius norm error bound of $O((\sigma/\sqrt{\lambda_r^*})\sqrt{\max\{t, d\}r^2 \log(mt)/(mt)})$ if the sample size is $mt \geq O(\max\{t, d\}r^4(\lambda_1^*/\lambda_r^*)^4 \text{polylog}(mt, d))$ and incoherence is $\mu \leq O(\lambda_1^*/\lambda_r^*)$. Several remarks on its sub-optimality are in order: (a) when $t \geq d$ and for $m = \Theta(\log(t))$, the error does not decrease as we increase the number of tasks t , (b) even when $t < d$ the error rate is sub-optimal by a factor of \sqrt{r} , and (c) each task requires $m \geq O(r^4 \text{polylog}(mt))$ samples. In contrast, error of MLLAM decays at a rate of $1/\sqrt{t}$ when $m = \Theta(\log(t))$, and this rate is optimal as it matches a lower bound, and each task requires only $m = \Omega(r^2 + r \log(t))$ samples.

The second approach, based on the method-of-moments, achieves a rescaled Frobenius norm error bound of $\tilde{O}((\sigma/\sqrt{\lambda_r^*})^2 \sqrt{(\lambda_1^*/\lambda_r^*)}(dr^2/(mt)) + \sqrt{\mu dr^2(\lambda_1^*/\lambda_r^*)/(mt)})$ if $m \geq 2$. The first term is suboptimal by a factor of \sqrt{r} . The second term is more problematic as it does not depend on the noise σ ; as we decrease σ , the error does not vanish even if we have enough samples to learn the parameters exactly. This is illustrated in the simulation result in Figure 1a.

Remark 3: One can study the problem in a stochastic setting where we sample a task i and compute stochastic gradient update for U based only on i -th task's samples. In this case, our proof techniques could be combined with that of Jain et al. [Jai+18] to obtain a nearly optimal and efficient one-pass algorithm. But we leave further investigation into such result for future work.

Remark 4: Our result holds if the initial point U_{init} is reasonably accurate. One choice of initialization is to use the Method-of-Moments (MoM) [TJJ20]. Due to sub-optimality of MoM approach (Theorem 7 in Appendix), we get an additional sample complexity requirement of $mt \geq \tilde{\Omega}((\lambda_1^*/\lambda_r^*)dr^2(\mu(\lambda_1^*/\lambda_r^*) + r(\sigma/\sqrt{\lambda_r^*})^4))$. Note that this does not degrade the error rate, $O(\sqrt{dr/mt})$.

Remark 5: Suppose we run Algorithm 1, under the conditions of Theorem 1 to get an estimated subspace U . Let a new task, whose task specific regressor v^{*+} lie in U^* , be introduced with m^+ samples. Now, we can apply the step 4 of Algorithm 1, with U and the new samples, to meta-learn an estimate v^+ of v^{*+} . Then by Tripuraneni, Jin, and Jordan [TJJ20, Theorem 4], the mean-squared-error (MSE) of the estimated regressor is $\tilde{O}((\sigma/\sqrt{\lambda_r^*})(\mu dr^2/mt + r/m^+))$. Therefore, as long as mt was large enough, we only need $m^+ = \Omega(r)$ additional samples to get an arbitrarily small MSE, as opposed to $m^+ = \Omega(d)$ of trivial baseline. We also improve upon other baselines from [TJJ20] (see Table 1) in terms of dependence on σ and t .

Task subset selection: The downside of our Algorithm 1 is that the requirement on m increase with t (i.e., $m = \Omega(\log t)$), which is not natural as the number of required samples per task should not increase as the number of tasks increase. To remove this dependency, we propose a new algorithm (Algorithm 2) that samples a set of tasks at each iteration to ensure we use only the ‘‘well-behaved’’ tasks.

Theorem 3. *Let there be t linear regression tasks, each with m samples satisfying Assumptions 1 and 2, and*

$$m \geq \tilde{\Omega}((\sigma/\sqrt{\lambda_r^*})^2 r^2 \log t + r^2 + \log(\mu)), \quad t \geq \tilde{\Omega}(\mu^2 r^2), \quad \text{and} \quad mt \geq \tilde{\Omega}((1 + (\sigma/\sqrt{\lambda_r^*})^2)(\lambda_1^*/\lambda_r^*)\mu dr^2).$$

Then MLLAMS (Algorithm 2), initialized at U_{init} s.t. $\|(\mathbf{I} - U^(U^*)^\top)U_{\text{init}}\|_F \leq \min(3/4, O(\sqrt{\lambda_r^*/\lambda_1^*}), O(1/\log t))$ and run for $K = \lceil \log_2(\lambda_r^*\lambda_r^*mt/\lambda_1^*\sigma^2\mu dr^2) \rceil$ iterations, outputs U so that, w.p. $\geq 1 - K/(dr)^{10}$*

$$\frac{\|(\mathbf{I} - U^*(U^*)^\top)U\|_F}{\sqrt{r}} \leq \tilde{O}\left(\left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)\sqrt{\frac{\mu r d}{mt}}\right). \quad (5)$$

Remark 6: Note that when $\sigma/\sqrt{\lambda_r^*} \leq 1/\log^2 t$, MLLAMS only needs $m \geq \Omega(r^2 + \log(\mu))$ samples per task. Since, MLLAMS selects a fraction of tasks to perform updates, the time-complexity of the method is similar to that of MLLAM.

Corollary 4. *Consider t linear regression tasks, each with m samples satisfying Assumptions 1 and 2 with $\sigma = 0$, and*

$$m \geq \tilde{\Omega}(r^2 + \log(\mu)), \quad t \geq \tilde{\Omega}((\mu r)^2), \quad \text{and} \quad mt \geq \tilde{\Omega}((\lambda_1^*/\lambda_r^*)\mu dr^2).$$

Then MLLAMS (Algorithm 2), initialized at U_{init} s.t. $\|(\mathbf{I} - U^(U^*)^\top)U_{\text{init}}\|_F \leq \min(3/4, O(\sqrt{\lambda_r^*/\lambda_1^*}), O(1/\log t))$, and run for K iterations outputs U so that the following holds (w.p. $\geq 1 - K/(dr)^{10}$):*

$$\frac{\|(\mathbf{I} - U^*(U^*)^\top)U\|_F}{\sqrt{r}} \leq \tilde{O}\left(\frac{\sqrt{\lambda_r^*/\lambda_1^*}}{\sqrt{r}2^K}\right). \quad (6)$$

Remark 7: In the above corollary, for the noiseless setting $\sigma = 0$, the number of samples per task does not grow with t , and thus it is nearly optimal. Note that the desired initialization point can be obtained using MoM. We leave the extension to noisy setting for future work.

Proofs of Theorems 1 & 3 are in the Appendices A.1 & B.1.

Algorithm 1 MLLAM: Meta-Learning Linear regressors via Alternating Minimization

Required: Data: $\{(x_j^{(i)} \in \mathbb{R}^d, y_j^{(i)} \in \mathbb{R})\}_{j=1}^m$ for all $1 \leq i \leq t$, K : number of steps.

```
1 Initialize  $U \leftarrow U_{\text{init}}$ 
2 Randomly shuffle the tasks  $\{1, \dots, t\}$ 
for  $1 \leq k \leq K$  do
3    $\mathcal{T}_k \leftarrow [1 + \frac{t(k-1)}{K}, \frac{tk}{K}]$ 
   for  $i \in \mathcal{T}_k$  do
4      $v^{(i)} \leftarrow \arg \min_{\hat{v} \in \mathbb{R}^r} \sum_{j \in [m/2]} (y_j^{(i)} - \langle U\hat{v}, x_j^{(i)} \rangle)^2$ 
   end
5    $\hat{U} \leftarrow \arg \min_{\hat{U} \in \mathbb{R}^{d \times r}} \sum_{i \in \mathcal{T}_k} \sum_{j=1+\frac{m}{2}}^m (y_j^{(i)} - \langle \hat{U}v^{(i)}, x_j^{(i)} \rangle)^2$ 
6    $U \leftarrow \text{QR}(\hat{U})$ 
end
return  $U$ 
```

4 Alternating minimization

In this section we discuss the alternating minimization algorithm we study in this paper. The algorithm follows the standard alternating minimization procedure [JNS13; CT84] where we update the representation matrix U and regressors V alternately. Note that, given U , we can estimate each of the regressor $v^{(i)}$ *separately* using standard least squares regression, i.e.,

$$v^{(i)} = \arg \min_v \sum_j (y_j^{(i)} - \langle x_j^{(i)}, Uv \rangle)^2.$$

Similarly, given the updated regressors $v^{(i)}$'s, we can now update U as:

$$\hat{U} = \arg \min_{\hat{U}} \sum_{i,j} (y_j^{(i)} - \langle x_j^{(i)}, \hat{U}v^{(i)} \rangle)^2.$$

To ensure certain normalization, we analyze a modification of the algorithm where the next iterate for U is the orthonormal subspace containing \hat{U} , which we can obtain using the QR-decomposition of \hat{U} .

Our analysis requires that when we update V using current U , we require U to be independent from the training datapoints. Similarly, during the update for U , we require V to be independent of the datapoints. We ensure the independence using two strategies: a) similar to standard online meta-learning settings [FAL17], we select random (previously unseen) tasks and update U and V , b) within each task, we divide the datapoints into two sets to update V and U separately.

Our update for $v^{(i)}$ require $O(mr^2 + r^3)$ time complexity, which can be brought down to $O(m \cdot r)$ by using gradient descent for solving the least squares. Our analysis shows that under the sample complexity assumptions of Theorem 1, each of the least squares problem has a constant condition number. So, the total number of iterations scale as $\log \frac{1}{\epsilon}$ to achieve ϵ error. If we set $\epsilon = 1/\text{poly}(t, \sigma)$, then using standard error analysis, we should be able to obtain the optimal error rate in Theorem 5. Similarly, *exact* update for U requires $O((dr)^3 + mt \cdot (dr)^2)$ time, that decreases to $O(mt \cdot d \cdot r)$ by using gradient descent updates.

4.1 Subset Selection

Algorithm 1 computes regressors $v^{(i)}$ for each of the task and use that to update U . Now, the Hessian for $v^{(i)}$ is given by: $H^{(i)} = \frac{1}{m} U^\top \sum_j x_j^{(i)} (x_j^{(i)})^\top U$. For sub-Gaussian $x_j^{(i)}$'s, $\|H^{(i)} - I\| \leq \sqrt{r/m} \sqrt{\log 1/\delta}$ with probability $1 - \delta$. This implies that if, m is independent of t , and if $t \gg m$ then the Hessian of some of the tasks can be highly ill-conditioned, leading to large estimation error in some of the regressors, which in turn leads to a large error in estimation of U . In Theorem 1, we avoid this issue by selecting m such that it grows logarithmically with t .

However, intuitively the number of required samples for each task should not increase with the number of tasks, especially in noise-less settings, where $t \geq \tilde{\Omega}(d)$ should be enough to ensure exact recovery of U .

Algorithm 2 MLLAMS: Meta-Learning Linear regressors via Alternating Minimization over task Subsets

Required: Data: $\{(x_j^{(i)} \in \mathbb{R}^d, y_j^{(i)} \in \mathbb{R})\}_{j=1}^m$ for all $1 \leq i \leq t$, K : number of steps.

```
1 Initialize  $U \leftarrow U_{\text{init}}$ 
2 Randomly shuffle the tasks  $\{1, \dots, t\}$ 
  for  $1 \leq i \leq t$  do
3    $S^{(i)} \leftarrow \frac{2}{m} \sum_{j \in [m/2]} x_j^{(i)} (x_j^{(i)})^\top$ 
  end
  for  $1 \leq k \leq K$  do
4    $\mathcal{T}_k \leftarrow \{i \in [1 + \frac{t(k-1)}{K}, \frac{tk}{K}] \mid \sigma_{\max}(U^\top S^{(i)} U) \leq 10; \sigma_{\min}(U^\top S^{(i)} U) \geq \frac{1}{2}\}$ 
     for  $i \in \mathcal{T}_k$  do
5        $v^{(i)} \leftarrow \arg \min_{\hat{v} \in \mathbb{R}^r} \sum_{j \in [m/2]} (y_j^{(i)} - \langle U \hat{v}, x_j^{(i)} \rangle)^2$ 
     end
6      $\hat{U} \leftarrow \arg \min_{\hat{U} \in \mathbb{R}^{d \times r}} \sum_{i \in \mathcal{T}_k} \sum_{j=1+\frac{m}{2}}^m (y_j^{(i)} - \langle \hat{U} v^{(i)}, x_j^{(i)} \rangle)^2$ 
7      $U \leftarrow \text{QR}(\hat{U})$ 
  end
return U
```

Practically, also a few poor tasks should not affect representation of the data significantly. So, in Algorithm 2, we propose a method to ignore the poor ill-conditioned tasks. To ensure this, we compute Hessian $H^{(i)}$ for each task, and ignore tasks whose Hessian's eigenvalue is small (see Line 4 in Algorithm 2). As mentioned in Theorem 3, while we condition on a task being *good*, we are still able to provide a similar result as Theorem 1 but with m which is independent of t in low-noise settings, e.g., when $\sigma \leq 1/\log t$.

5 Proof sketch for noiseless case

Here we provide proof sketches of Theorem 1. To highlight the main ideas behind our analysis, we start with the simplest case when there is no noise ($\sigma^2 = 0$) and all the regressors lie on a single dimensional subspace ($r = 1$). The analysis gets quite challenging as we go to multi-dimensional shared subspace ($r > 1$), and we illustrate these challenges and how to resolve them in Section 5.2.

5.1 Proof sketch for the one-dimensional case

Let $u^* \in \mathbb{R}^d$ be the unit vector of the one-dimensional true subspace, and $v^* \in \mathbb{R}^t$ the vector of the true regressor coefficients of the t tasks. In the noiseless setting ($\varepsilon_j^{(i)} = 0$), the k -th step of MLLAM can be written as follows.

$$\begin{aligned} & \text{For all } i \in \mathcal{T}_k \\ & v^{(i)} \leftarrow (u^\top S_1^{(i)} u)^{-1} u^\top S_1^{(i)} (u^*) v^{*(i)}, \\ & \hat{u} \leftarrow \left(\sum_{i \in \mathcal{T}_k} (v^{(i)})^2 S_2^{(i)} \right)^\dagger \left(\sum_{i \in \mathcal{T}_k} v^{*(i)} v^{(i)} S_2^{(i)} u^* \right), \quad u^+ \leftarrow \frac{\hat{u}}{\|\hat{u}\|}, \end{aligned}$$

where $S_\ell^{(i)} = \frac{2}{m} \sum_{j=(\ell-1)m/2+1}^{\ell m/2} x_j^{(i)} (x_j^{(i)})^\top$ is the data covariance matrix of a half of the dataset $[m]$ of task $i \in [t]$. Our incoherence condition for rank-1 case simplifies to $\|v\|_\infty^2 \leq \frac{\mu}{t} \|v\|^2$. The distance between two unit norm vectors u and u^* is commonly measured by the angular distance defined as $\sin \theta(u, u^*) \triangleq \|(\mathbf{I} - u^*(u^*)^\top)u\|^{1/2}$, where $\mathbf{I} - u^*(u^*)^\top$ is the projection operator to the sub-space orthogonal to u^* . In the following we let $q \triangleq \langle u^*, u \rangle$ and use the relation $\sin \theta(u, u^*) = \|u - u^*q\|$ in the analysis. We use the fact that if we have a good previous iterate u close to u^* , i.e. $\sin \theta(u, u^*) \leq 3/4$, then $1/2 \leq |q| \leq 1$.

Our analysis shows that we get geometrically closer to the true subspace u^* at every iteration in this $\sin \theta$ distance, when initialized sufficiently close to u^* .

Our strategy is to show that the v -update achieves $|v^{(i)} - q^{-1}v^{*(i)}| \leq C\|v^{*(i)}\| \sin\theta(u, u^*)$ for some constant C , and the u -update achieves $\sin\theta(u^+, u^*) \leq (c/\|v^*\|)\|v - q^{-1}v^*\|$ where the constant c can be made as small as we want in the assumed sample regime. Together, they imply the desired theorem.

v -update: We can write $v^{(i)}q^{-1} - v^{*(i)}$ as

$$v^{(i)} - q^{-1}v^{*(i)} = u^\top S_1^{(i)}(qu^* - u)(u^\top S_1^{(i)}u)^{-1}q^{-1}v^{*(i)}.$$

In expectation, $\|\mathbb{E}[u^\top S_1^{(i)}(qu^* - u)]\| = \|u^\top(qu^* - u)\| = 1 - q^2 \leq (\sin\theta(u, u^*))^2$ and $\mathbb{E}[u^\top S_1^{(i)}u] = \|u\|^2 = 1$. Therefore, by Lemma A.1, if $\sin\theta(u, u^*) \leq \frac{1}{32}$ and there is enough samples per task, i.e. $m \geq \Omega(\log(t/K\delta))$, we can bound their deviations in terms of $\sin\theta(u, u^*)$. This implies that, with a probability of at least $1 - \delta/2$,

$$\frac{|v^{(i)} - q^{-1}v^{*(i)}|}{|v^{*(i)}|} \leq \frac{\sin\theta(u, u^*)}{4}, \text{ for all } i \in \mathcal{T}_k, \quad (7)$$

where we used the fact that $|q| \geq 1/2$. This in turn implies that $(1/4)|v^{*(i)}| \leq |v^{(i)}|$ and v is incoherent.

u -update: We bound the distance between \hat{u} and u^* :

$$\begin{aligned} \hat{u} - u^*q &= \underbrace{\left(\sum_{i \in \mathcal{T}_k} \frac{(v^{(i)})^2}{\|v\|^2} S_2^{(i)} \right)^\dagger}_{:=A} \underbrace{\left(\sum_{i \in \mathcal{T}_k} \frac{v^{(i)}h^{(i)}}{\|v\|^2} S_2^{(i)} u^*q \right)}_{:=\hat{H}}, \end{aligned} \quad (8)$$

where $h^{(i)} = q^{-1}v^{*(i)} - v^{(i)}$. Notice that, in expectation, $\mathbb{E}[A] = \mathbf{I}$ and $\mathbb{E}[\hat{H}u^*q] = \frac{v^\top h}{\|v\|^2}u^*q \leq \frac{\|h\|}{\|v\|}$. Therefore, by Lemma A.2, when there are enough samples, i.e. $mt \geq K\Omega(\mu d \log(\frac{1}{\delta}))$ deviations from these expected values can be bounded using the distance between v and v^* , $\|h\|$. That is with a probability of at least $1 - \frac{\delta}{2}$, A is invertible and well-conditioned,

$$A^{-1} = \mathbf{I} + E_1, \quad \text{and} \quad Hu^*q = \frac{v^\top h}{\|v\|^2}u^*q + e_2,$$

where $\|E_1\| \leq \frac{1}{16}$ and $\|e_2\| \leq \frac{1}{32} \left(\frac{\|h\|}{\|v\|} + \sqrt{\frac{t}{\mu}} \frac{\|h\|_\infty}{\|v\|} \right)$. Note that we had to critically use incoherence of intermediate v to bound e_2 . Therefore

$$\hat{u} - u^*q = \underbrace{\frac{v^\top h}{\|v\|^2}u^*q}_{:=\hat{u}_\parallel} + \underbrace{q \frac{v^\top h}{\|v\|^2}E_1u^* + (\mathbf{I} + E_1)e_2}_{:=f}.$$

Notice that \hat{u}_\parallel is parallel to u^* . Rest of the terms are grouped together as f . The angle distance $\sin(u^+, u^*)$ only depends on the portion of u^+ which lie in the orthogonal subspace to u^* . Therefore, $\|\hat{u}_\parallel\|$ does not directly contribute to the distance, and this is formalized below. Clearly, $\|(\mathbf{I} - u^*(u^*)^\top)u^+\| = \min_{q^+} \|u^+ - u^*q^+\|$. This follows from the trivial solution of the scalar quadratic problem $\min_{q^+ \in \mathbb{R}} \|u - u^*q^+\|^2$. Thus,

$$\begin{aligned} \sin\theta(u^+, u^*) &= \min_{q^+} \|u^+ - u^*q^+\| \\ &\leq \left\| \frac{\hat{u}}{\|\hat{u}\|} - \left(1 + \frac{h^\top v}{\|v\|^2}\right) u^* \frac{q}{\|\hat{u}\|} \right\| \\ &\leq \frac{\|f\|}{\|\hat{u}\|} \leq \frac{\|f\|}{q\|u^*\| - \|f\| - \|h\|/\|v\|}. \end{aligned} \quad (9)$$

Putting them together: We bound f using definitions of E_1 and e_2 , incoherence, and (7) as

$$\|f\| \leq \frac{1}{16} \frac{\|h\|}{\|v\|} + \frac{1}{32} \left(\frac{\|h\|}{\|v\|} + \sqrt{\frac{t}{\mu}} \frac{\|h\|_\infty}{\|v\|} \right) \leq \frac{1}{8} \sin\theta(u, u^*).$$

Combining this with (9), we see that with a probability of at least $1 - \delta$, the angle distance geometrically decreases at each step, i.e.

$$\sin \theta(u^+, u^*) \leq \frac{1}{2} \sin \theta(u, u^*). \quad (10)$$

Finally, if the initialization is good, i.e. $\sin \theta(u_{\text{init}}, u^*) \leq \frac{1}{16}$, we can unroll the above inequality across iterations. Taking union bound over the iterations we get that, with a probability of at least $1 - K\delta$, the output u after K iterations satisfies

$$\sin \theta(u, u^*) \leq \frac{1}{2^K} \sin \theta(u_{\text{init}}, u^*). \quad (11)$$

To achieve this, we need at least $m \geq \Omega(\log(\frac{t}{K\delta}))$ samples per task and at least $mt \geq \Omega(K\mu d \log(\frac{t}{\delta}))$ total samples.

5.2 Proof sketch for the r -dimensional case

Here we do not use $\sin \theta_1(U, u^*)$ distance, as the analysis of $\sin \theta_1$ gets more complicated in the general r -dimensional case. Therefore we use ℓ_2 norm based error, $\Delta(U, U^*) := (\sum_{r'=1}^r \sin^2 \theta_{r'}(U, U^*))^{1/2} := \|(\mathbf{I} - U^*(U^*)^\top)U\|_F$. Let $Q = (U^*)^\top U$, then $\Delta(U, U^*) = \|U - U^*Q\|_F$, and $1/2 \leq \|Q\| \leq 1$ if $\Delta(U, U^*) \leq 3/4$.

$$\begin{aligned} & \text{For all } i \in \mathcal{T}_k \\ & v^{(i)} \leftarrow (U^\top S_1^{(i)} U)^\dagger U^\top S_1^{(i)} U^* v^{*(i)}, \\ & \widehat{U} \leftarrow \left(\mathcal{A}^\dagger \left(\sum_{i \in \mathcal{T}_k} S_2^{(i)} U^* v^{*(i)} (v^{(i)} W^{-\frac{1}{2}})^\top \right) \right) W^{-\frac{1}{2}}, \\ & U \leftarrow \text{QR}(\widehat{U}), \end{aligned}$$

where $W = V^\top V$, $\mathcal{A} : \mathbb{R}^{d \times r} \rightarrow \mathbb{R}^{d \times r}$ is linear operator such that $\mathcal{A}(U) = \sum_{i \in \mathcal{T}_k} S_2^{(i)} U W^{-\frac{1}{2}} v^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}}$, and $S_\ell^{(i)}$ are defined as in the one-dimensional case.

V-update: We will prove that $\|v^{(i)} - Q^{-1}v^{*(i)}\| = O(\Delta(U, U^*))$. Let $h^{(i)} := v^{(i)}Q^{-1} - v^{*(i)}$, then

$$h^{(i)} = (U^\top S_1^{(i)} U)^\dagger \underbrace{U^\top S_1^{(i)} (U^*Q - U) Q^\dagger v^{*(i)}}_{:=G}.$$

Notice that, in expectation, $\|\mathbb{E}[U^\top S_1^{(i)} U]\| = 1$ and $\|\mathbb{E}[G]\| = \|U^\top (U^*Q - U)\| = \|Q^\top Q - \mathbf{I}\| = \Delta^2(U, U^*)$. Therefore, by Lemma A.1, if $\Delta^2(U, U^*) \leq \frac{1}{32}$ and there is enough samples per task, i.e. $m \geq \Omega(r \log(\frac{t}{K\delta}))$, we can bound their deviations in terms of $\sin \theta(u, u^*)$. This implies that, with a probability of at least $1 - \delta/2$,

$$\|h^{(i)}\| \leq \frac{\|v^{*(i)}\| \Delta^2(U, U^*)}{4}, \text{ for all } i \in \mathcal{T}_k. \quad (12)$$

Furthermore, $\|v^{(i)}\| \leq 4\|v^{*(i)}\|$ and V is incoherent.

U-update: We bound the distance between \widehat{U} and U^* :

$$(\widehat{U} - U^*Q)W^{\frac{1}{2}} = \mathcal{A}^\dagger \left(\underbrace{\sum_{i \in \mathcal{T}_k} S_2^{(i)} U^*Q h^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}}}_{:= -\widehat{\mathcal{H}}(U^*Q)} \right).$$

Notice that, in expectation, $\mathbb{E}[\widehat{\mathcal{H}}(U^*Q)] = \mathcal{H}(U^*Q) := U^*Q \sum_{i \in \mathcal{T}_k} h^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}}$ and $\mathcal{H}(U^*Q) \leq \|H\|_F$ and $\mathbb{E}[\mathcal{A}]$ is the identity map \mathcal{I} . Like in the 1-dimensional case, by Lemma A.2, when there are enough samples, i.e. $mt \geq K\Omega(\mu d r^2 \log(\frac{1}{\delta}))$ deviations from these expected values can be bounded using the distance between V and V^* , $\|H\|$. That is, with a probability of at least $1 - \delta/2$, \mathcal{A} is invertible and well-conditioned in Frobenius operator norm,

$$\mathcal{A}^{-1} = \mathcal{I} + \mathcal{E}_1, \quad \text{and} \quad \widehat{\mathcal{H}}(U^*Q) = \mathcal{H}(U^*Q) - E_2,$$

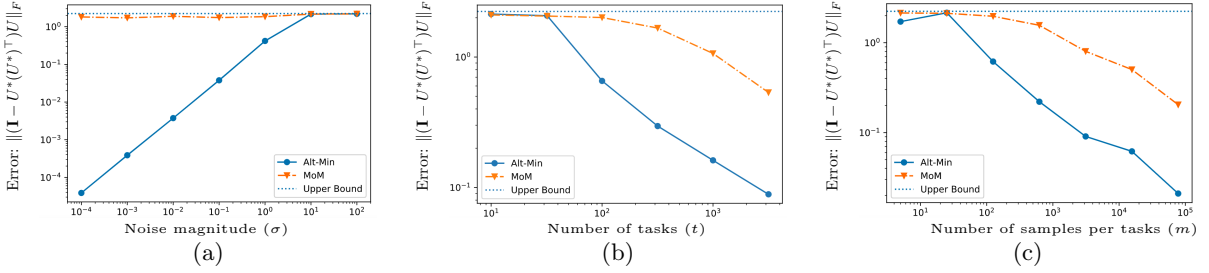


Figure 1: (a): MLLAM achieves vanishing error as noise decreases, whereas the error with Method-of-Moments stay bounded away from zero. (b), (c): MLLAM incurs significantly smaller error in estimation of true subspace U^* than MoM, both for growing number of tasks (t) and for growing number of samples per task (m).

where $\|\mathcal{E}_1\|_F \leq 1/16$ and $\|E_2\|_F \leq 1/32(\|H\|_F + \sqrt{t/\mu}\|H\|_{\infty,2})$. Note that we had to critically use incoherence of intermediate V to bound E_2 . Therefore,

$$(\widehat{U} - U^*Q)W^{\frac{1}{2}} = -\mathcal{H}(U^*Q) - \underbrace{cE_1\mathcal{H}(U^*Q) + (\mathcal{I} + \mathcal{E}_1)E_2}_{:=F}.$$

Now, using similar arguments as in the one-dimensional case, we get

$$\begin{aligned} \Delta(U^+, U^*) &\leq \left\| \widehat{U}R^{-1} - U^*Q + \mathcal{H}(U^*Q) \right\|_F \|W^{-\frac{1}{2}}\| \\ &\leq \frac{\|F\|_F}{\|R^{-1}\|} \leq \frac{\|F\|_F \lambda_r^{-\frac{1}{2}}}{\|QU^*\| - (\|F\|_F + \|H\|_F)\lambda_r^{-\frac{1}{2}}}. \end{aligned}$$

Putting them together: Using similar arguments as in one-dimensional case, if the initialization is good, i.e. $\Delta(U_{\text{init}}, U^*) \leq 1/16$, we can show that with a probability of at least $1 - \delta$, the next iterate U^+ satisfies: $\Delta(U^+, U^*) \leq \frac{1}{2}\Delta(U, U^*)$. To achieve this, we need at least $\Omega(r \log(\frac{t}{K\delta}))$ samples per task (m) and at least $\Omega(K\mu dr^2 \log(\frac{t}{\delta}))$ total samples (mt). Result now follows by applying the above result K times.

6 Experimental results

In this section we empirically compare the performance of MLLAM (Alt-Min, Algorithm 1) against Method-of-Moments (MoM) [TJJ20]. We generate data samples with dimension $d = 100$ and generate random subspace U^* of rank $r = 5$. We sample the task regressor coefficients as $v^{(i)} \sim \mathcal{N}(0, \mathbf{I})$. In all our experiments, we initialize MLLAM uniformly at random and run it for $K = 20$ iterations. In all the figures, the blue straight line with circular marker denotes the MLLAM algorithm, the orange dashed and dotted line with inverted triangular marker denotes the MoM, and the blue dotted line parallel to x-axis represents the theoretical upper-limit \sqrt{r} of the Frobenius norm distance, $\|(\mathbf{I} - U^*(U^*)^\top)U\|_F$. In all the figures we use log-scaled x and y axes.

Figure 1a plots subspace estimation error ($\|(\mathbf{I} - U^*(U^*)^\top)U\|_F$) against the standard deviation σ of the regression noise, $\varepsilon_j^{(i)} \sim \mathcal{N}(0, \sigma^2)$; see (1). We vary σ from 10^{-4} to 10^2 , while fixing the number of tasks at $t = 200$ and the number of samples per task at $m = 25$. Clearly, as predicted by Theorem 1, our MLLAM (Alt-Min) algorithm achieves a smaller error than MoM over all values of σ . Error of MLLAM is linearly proportional to σ . As predicted by Theorem 7 (in Appendix), the distance of MoM is a constant multiple of $\sqrt{\frac{dr^3}{mt}} = \sqrt{r}$ for all values of σ , and it does not improve when σ decreases.

Figure 1b plots the subspace error against the number of tasks t . We vary t from 10 to 3163, while the number of samples per task is fixed at $m = 25$ and $\sigma = 1$. In Figure 1c, we plot the the error against the number samples per tasks m . We vary m from 5 to 78125, while fixing the number of tasks at $t = 20$ and the standard deviation of the regression noise at $\sigma = 1$. In both of these figures, we observe that, MLLAM (Alt-Min) achieves much smaller subspace error than the MoM. Furthermore, as predicted by Theorems 1 and 7 (in Appendix), the squared error rate for both MLLAM and MoM decreases linearly m and t .

Note that even though we randomly initialize our MLLAM algorithm, it still performs better than the baseline MoM. Similar observations have been made for other non-convex algorithms for solving low-rank problems [CCFM19]. This suggests that the initialization requirement of Theorem 1, $\|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq O(\sqrt{\lambda_r^*/\lambda_1^*})$ may be an artifact of the analysis or may be practically insignificant.

7 Conclusion

In this paper, we analyzed an alternating minimization method for the problem of linear meta-learning, a simple but canonical problem in meta-learning. We showed that Algorithm 1 that alternately learns the shared representation matrix across tasks and the task-specific regressors, can provide nearly optimal error rate along with nearly optimal per-task and overall sample complexities. To the best of our knowledge, we provide the first result with optimal error rate — that scales appropriately with the noise in observations — while still ensuring per-task sample complexity to be nearly independent of d (the dimensionality of data), which is a key requirement in meta-learning as individual tasks are data-starved. We also proposed and analyzed a subset selection based method that further improves per-task sample complexity and ensures that it is *independent* of the number of tasks for noise-less setting.

The work leads to several interesting future directions and questions. For the non-linear version of the problem, ensuring optimal error rate with optimal per-task sample complexity is an interesting open question. Understanding and contrasting standard *MAML* techniques for the linear and non-linear problem is another exciting direction, which is already seeing a fair amount of interest [FMO20; SZKA20]. Finally, analyzing alternating minimization methods with stochastic gradients and streaming tasks is another promising direction.

References

- [AM56] Ali R Amir-Moéz. “Extreme properties of eigenvalues of a Hermitian transformation and singular values of the sum and product of linear transformations”. In: *Duke Mathematical Journal* 23.3 (1956), pp. 463–476.
- [HS81] Harold V Henderson and Shayle R Searle. “On deriving the inverse of a sum of matrices”. In: *Siam Review* 23.1 (1981), pp. 53–60.
- [CT84] I. Csiszár and G. Tusnady. “Information geometry and alternating minimization procedure”. In: *Statistics and Decision* (1984).
- [Bax95] Jonathan Baxter. “Learning internal representations”. In: *Proceedings of the eighth annual conference on Computational learning theory*. 1995, pp. 311–320.
- [IE96] Nathan Intrator and Shimon Edelman. “Making a low-dimensional representation suitable for diverse tasks”. In: *Learning to learn*. Springer, 1996, pp. 135–157.
- [Car97] Rich Caruana. “Multitask learning”. In: *Machine learning* 28.1 (1997), pp. 41–75.
- [TP98] Sebastian Thrun and Lorien Pratt. “Learning to learn: Introduction and overview”. In: *Learning to learn*. Springer, 1998, pp. 3–17.
- [Bax00] Jonathan Baxter. “A model of inductive bias learning”. In: *Journal of artificial intelligence research* 12 (2000), pp. 149–198.
- [AZ05] Rie Kubota Ando and Tong Zhang. “A framework for learning predictive structures from multiple tasks and unlabeled data”. In: *Journal of Machine Learning Research* 6.Nov (2005), pp. 1817–1853.
- [Orl05] Alon Orlitsky. “Supervised dimensionality reduction using mixture models”. In: *Proceedings of the 22nd international conference on Machine learning*. 2005, pp. 768–775.
- [AFSU07] Yonatan Amit, Michael Fink, Nathan Srebro, and Shimon Ullman. “Uncovering shared structures in multiclass classification”. In: *Proceedings of the 24th international conference on Machine learning*. 2007, pp. 17–24.
- [AEP08] Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. “Convex multi-task feature learning”. In: *Machine learning* 73.3 (2008), pp. 243–272.

- [Ris+08] Irina Rish, Genady Grabarnik, Guillermo Cecchi, Francisco Pereira, and Geoffrey J Gordon. “Closed-form supervised dimensionality reduction with generalized linear models”. In: *Proceedings of the 25th international conference on Machine learning*. 2008, pp. 832–839.
- [CR09] Emmanuel J Candès and Benjamin Recht. “Exact matrix completion via convex optimization”. In: *Foundations of Computational mathematics* 9.6 (2009), pp. 717–772.
- [MJD09] Raghu Meka, Prateek Jain, and Inderjit S Dhillon. “Guaranteed rank minimization via singular value projection”. In: *arXiv preprint arXiv:0909.5457* (2009).
- [Ver10] Roman Vershynin. “Introduction to the non-asymptotic analysis of random matrices”. In: *arXiv preprint arXiv:1011.3027* (2010).
- [Har+12] Zaid Harchaoui, Matthijs Douze, Mattis Paulin, Miroslav Dudik, and Jérôme Malick. “Large-scale image classification with trace-norm regularization”. In: *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE. 2012, pp. 3386–3393.
- [JD13] Prateek Jain and Inderjit S Dhillon. “Provable inductive matrix completion”. In: *arXiv preprint arXiv:1306.0626* (2013).
- [JNS13] Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. “Low-rank matrix completion using alternating minimization”. In: *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*. 2013, pp. 665–674.
- [PM13] Massimiliano Pontil and Andreas Maurer. “Excess risk bounds for multitask learning with trace norm regularization”. In: *Conference on Learning Theory*. 2013, pp. 55–76.
- [HMRW14] Moritz Hardt, Raghu Meka, Prasad Raghavendra, and Benjamin Weitz. “Computational limits for matrix completion”. In: *Conference on Learning Theory*. PMLR. 2014, pp. 703–725.
- [NJS15] Praneeth Netrapalli, Prateek Jain, and Sujay Sanghavi. “Phase retrieval using alternating minimization”. In: *IEEE Transactions on Signal Processing* 63.18 (2015), pp. 4814–4826.
- [ZJD15] Kai Zhong, Prateek Jain, and Inderjit S Dhillon. “Efficient matrix sensing using rank-1 gaussian measurements”. In: *International conference on algorithmic learning theory*. Springer. 2015, pp. 3–18.
- [FAL17] Chelsea Finn, Pieter Abbeel, and Sergey Levine. “Model-agnostic meta-learning for fast adaptation of deep networks”. In: *arXiv preprint arXiv:1703.03400* (2017).
- [SSSG17] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. “Revisiting unreasonable effectiveness of data in deep learning era”. In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 843–852.
- [WRH17] Yu-Xiong Wang, Deva Ramanan, and Martial Hebert. “Learning to model the tail”. In: *Advances in Neural Information Processing Systems*. 2017, pp. 7029–7039.
- [Jai+18] Prateek Jain, Sham Kakade, Rahul Kidambi, Praneeth Netrapalli, and Aaron Sidford. “Parallelizing stochastic gradient descent for least squares regression: mini-batching, averaging, and model misspecification”. In: *Journal of Machine Learning Research* 18 (2018).
- [Ver18] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*. Vol. 47. Cambridge university press, 2018.
- [CCFM19] Yuxin Chen, Yuejie Chi, Jianqing Fan, and Cong Ma. “Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval”. In: *Mathematical Programming* 176.1 (2019), pp. 5–37.
- [RRBV19a] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. “Rapid learning or feature reuse? towards understanding the effectiveness of maml”. In: *arXiv preprint arXiv:1909.09157* (2019).
- [RRBV19b] Aniruddh Raghu, Maithra Raghu, Samy Bengio, and Oriol Vinyals. “Rapid learning or feature reuse? towards understanding the effectiveness of maml”. In: *arXiv preprint arXiv:1909.09157* (2019).
- [Du+20] Simon S Du, Wei Hu, Sham M Kakade, Jason D Lee, and Qi Lei. “Few-shot learning via learning the representation, provably”. In: *arXiv preprint arXiv:2002.09434* (2020).

- [FMO20] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. “On the convergence theory of gradient-based model-agnostic meta-learning algorithms”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2020, pp. 1082–1092.
- [Kon+20] Weihao Kong, Raghav Somani, Zhao Song, Sham Kakade, and Sewoong Oh. “Meta-learning for mixed linear regression”. In: *arXiv preprint arXiv:2002.08936* (2020).
- [SZKA20] Nikunj Saunshi, Yi Zhang, Mikhail Khodak, and Sanjeev Arora. “A sample complexity separation between non-convex and convex meta-learning”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 8512–8521.
- [TJJ20] Nilesh Tripuraneni, Chi Jin, and Michael I Jordan. “Provable Meta-Learning of Linear Representations”. In: *arXiv preprint arXiv:2002.11684* (2020).
- [CHMS21] Liam Collins, Hamed Hassani, Aryan Mokhtari, and Sanjay Shakkottai. “Exploiting Shared Representations for Personalized Federated Learning”. In: *arXiv preprint arXiv:2102.07078* (2021).

Appendix

This appendix contains proofs for the claims mentioned main text. Section A and B contain the analyses of Algorithm 1 and 2, respectively. Section C contains corollaries of some known results. Section D contains some general technical lemmas used in this paper.

A Analysis of MLLAM (Algorithm 1)

Initialized at U , the k -th step of alternating minimization-based MLLAM (Algorithm 1) is:

$$v^{(i)} \leftarrow (U^\top S_1^{(i)} U)^\dagger ((U^\top S_1^{(i)} U^*) v^{*(i)} + U^\top z^{(i)}), \quad \text{for } i \in \mathcal{T}_k = [1 + (k-1)t/K, tk/K] \quad (13)$$

$$\hat{U} \leftarrow \mathcal{A}^\dagger \left(\sum_{i \in [t]} S_2^{(i)} U^* v^{*(i)} (v^{(i)})^\top + z^{(i)} (v^{(i)})^\top \right), \quad (14)$$

$$U^+ \leftarrow \text{QR}(\hat{U}), \quad (15)$$

where U^+ is the next iterate, $S_1^{(i)} = \frac{2}{m} \sum_{j \in [1, m/2]} x_j^{(i)} (x_j^{(i)})^\top$, $S_2^{(i)} = \frac{2}{m} \sum_{j \in [1+m/2, m]} x_j^{(i)} (x_j^{(i)})^\top$, $z^{(i)} \triangleq (1/m) \sum_{j \in [m]} \varepsilon_j^{(i)} x_j^{(i)}$ and $\mathcal{A} : \mathbb{R}^{d \times r} \rightarrow \mathbb{R}^{d \times r}$ is a self-adjoint linear operator such that $\mathcal{A}(U) = \sum_{i \in T} S^{(i)} U v^{(i)} (v^{(i)})^\top$. The self-adjointness of \mathcal{A} follows from the symmetry of $S^{(i)}$ when using cyclic property of trace as follows

$$\begin{aligned} \langle U_2, \mathcal{A}(U_1) \rangle &= \sum_{i \in T} \langle U_2, S^{(i)} U_1 v^{(i)} (v^{(i)})^\top \rangle = \sum_{i \in T} \text{tr}(U_2^\top S^{(i)} U_1 v^{(i)} (v^{(i)})^\top) \\ &= \sum_{i \in T} \text{tr}(v^{(i)} (v^{(i)})^\top U_2^\top S^{(i)} U_1) = \langle \mathcal{A}(U_2), U_1 \rangle \end{aligned} \quad (16)$$

Incoherence. $\max_i \|v^{*(i)}\|^2 \leq (\mu r/t) \lambda_r (\sum_{i \in [t]} v^{*(i)} (v^{*(i)})^\top)$, and we define $\nu = (1/t) \lambda_r (\sum_{i \in [t]} v^{*(i)} (v^{*(i)})^\top)$. Notice that, this non-standard definition of incoherence is related to the standard definition: $W^* = (V^*)^\top V^* = \sum_{i \in [t]} v^{*(i)} (v^{*(i)})^\top$, $V^* = \tilde{V}^* R^*$ (QR-decomposition), $\max_i \|\tilde{v}^{*(i)}\|^2 \leq \tilde{\mu} r/t$, as follows $\mu = \hat{\mu} (\sigma_1^2(R^*) / \sigma_r^2(R^*))$.

Theorem 5. *Let there be t linear regression tasks, each with m samples satisfying Assumptions 1 and 2, and $K = \lceil \log_2(\frac{(\lambda_r^*/\lambda_1^*)mt}{\mu dr^2}) \rceil$, $\|(\mathbf{I} - U^*(U^*)^\top)U_{\text{init}}\|_F \leq \min\left(\frac{3}{4}, O\left(\sqrt{\frac{\lambda_r^*}{\lambda_1^*} \frac{1}{\log(t/K)}}\right)\right)$, $m \geq \Omega\left((1+r\left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2)r \log\left(\frac{t}{\delta}\right) + r^2 \log\left(\frac{K}{\delta}\right)\right)$, $t \geq \Omega(\mu^2 r^3 K \log\left(\frac{K}{\delta}\right))$, and $mt \geq \Omega\left(\mu dr^2 K \frac{\lambda_1^*}{\lambda_r^*} \left(\log\left(\frac{t}{\delta}\right) + \left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 \log^2\left(\frac{t}{\delta}\right) \log\left(\frac{rK}{\delta}\right)\right)\right)$. Then, for any $0 < \delta < 1$, after K iterations, MLLAM (Algorithm 1) returns an orthonormal matrix $U \in \mathbb{R}^{d \times r}$, such that with a probability of at least $1 - \delta$*

$$\frac{1}{\sqrt{r}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq O\left(\frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{\mu dr K \log\left(\frac{t}{\delta}\right) \log\left(\frac{rK}{\delta}\right)}{mt}}\right) \quad (17)$$

and the algorithm uses an additional memory of size $O(d^2 r^2)$.

A proof is in Section A.1.

Initialization. If we initialize MLLAM (Algorithm 1) with Method-of-Moments (Theorem 7), we need at least

$$mt \geq \tilde{\Omega}\left(\frac{\lambda_1^{*2}}{\lambda_r^{*2}} \mu dr^2 + \left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^4 \frac{\lambda_1^*}{\lambda_r^*} dr^3\right) \quad (18)$$

initial number of samples, where $\tilde{\Omega}$ hides polylog factors.

A.1 Proof of Theorem 5

Proof sketch: We first prove that distance between U^* and U decreases at each iteration up to some additional noise terms. Then this per iterate result is unrolled to obtain the final guarantees.

First we focus on the k -th iterate. In this analysis, unless specified $[t]$, represents the k -th K -way partition used for the k -th iterate.

In the analysis of an iterate we denote the current iterate using U and the next iterate using U^+ . First we prove that the distance between the true $v^{*(i)}$ and the current $v^{(i)}$ is approximately upper-bounded by multiple of distance between U and U^* . Next we prove that distance between U^+ and U^* is approximately a fraction of the distance between $v^{*(i)}$ and $v^{(i)}$. Finally, combining the above two results gives us desired result.

Preliminaries: Let $Q = (U^*)^\top U$. Using Lemma D.4, if $\|U - U^*(U^*)^\top U\|_F < 1$, Q is invertible. Let Q^{-1} be the right inverse of Q , i.e. $QQ^{-1} = \mathbf{I}$. Let $W = (V^*)^\top V^* = \sum_{i \in [t]} v^{*(i)}(v^{*(i)})^\top$, and $\lambda_1^* = \max_{\|z\|=1} z^\top W^* z$ and $\lambda_r^* = \min_{\|z\|=1} z^\top W^* z$.

Update on V : Let $h^{(i)} = v^{(i)} - Q^{-1}v^{*(i)}$ and $H^T = [h^{(1)}h^{(2)} \dots h^{(t)}]$. Let $\|H\|_F \triangleq \sqrt{\sum_{i \in [t]} \|h^{(i)}\|^2}$ and $\|H\|_{\infty,2} \triangleq \max_{i \in [t]} \|h^{(i)}\|$. Let $W = V^\top V = \sum_{i \in [t]} v^{(i)}(v^{(i)})^\top$, and $\lambda_1 = \max_{\|z\|=1} z^\top W z$ and $\lambda_r = \min_{\|z\|=1} z^\top W z$.

Lemma A.1. *If $\|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq \min\left(\frac{3}{4}, O\left(\sqrt{\frac{\lambda_1^*}{\lambda_r^*}} \frac{1}{\log(t/K)}\right)\right)$ and $m \geq \Omega\left(\left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 r^2 \log\left(\frac{t}{K\delta}\right) + r \log\left(\frac{t}{K\delta}\right)\right)$, then with a probability of at least $1 - \delta/3$,*

$$\|v^{(i)}\| \leq O\left(\mu \lambda_r\right), \text{ and } \lambda_r^* \leq 2\lambda_r \quad (19)$$

and

$$\sqrt{\frac{rK}{t}} \frac{\|H\|_F}{\sqrt{\lambda_r}} \leq O\left(\sqrt{\frac{\log\left(\frac{t}{K\delta}\right)}{\log\left(\frac{1}{\delta}\right)}} \sqrt{\frac{\lambda_1^*}{\lambda_r^*}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log\left(\frac{t}{K\delta}\right)}{m}}\right) \quad (20)$$

$$\sqrt{\frac{rK}{t}} \frac{\|H\|_{\infty,2}}{\sqrt{\lambda_r}} \leq O\left(\sqrt{\frac{\log\left(\frac{t}{K\delta}\right)}{\log\left(\frac{1}{\delta}\right)}} \|(\mathbf{I} - U^*(U^*)^\top)U\| \sqrt{\frac{\mu r K}{t}} + \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 K \log\left(\frac{t}{K\delta}\right)}{mt}}\right) \quad (21)$$

A proof is in Section A.2.1.

Update on U : Let $W, \mathcal{H}, \widehat{\mathcal{H}} : \mathbb{R}^{d \times r} \rightarrow \mathbb{R}^{d \times r}$ be three linear operators, such that $\mathcal{W}(U) = U \sum_{i \in \mathcal{T}_k} v^{(i)}(v^{(i)})^\top = UW$, $\mathcal{H}(U) = U \sum_{i \in \mathcal{T}_k} h^{(i)}(v^{(i)})^\top$ and $\widehat{\mathcal{H}}(U) = \sum_{i \in \mathcal{T}_k} S_2^{(i)} U h^{(i)}(v^{(i)})^\top$, where $h^{(i)} = v^{(i)} - Q^{-1}v^{*(i)}$. \mathcal{W} is invertible and self-adjoint. Therefore $\mathcal{W}^{-\frac{1}{2}}$ and $\mathcal{W}^{\frac{1}{2}}$ exist. Let $\mathcal{I} : \mathbb{R}^{d \times r} \rightarrow \mathbb{R}^{d \times r}$ be the identity mapping, such that $\mathcal{I}(U) = U$.

$$\widehat{U} - U^*Q = \mathcal{A}^\dagger \left(\sum_{i \in \mathcal{T}_k} S_2^{(i)} U^*Q(Q^{-1}v^{*(i)} - v^{(i)})(v^{(i)})^\top + z^{(i)}(v^{(i)})^\top \right) \quad (22)$$

$$= \mathcal{A}^\dagger \left(-\widehat{\mathcal{H}}(U^*Q) + \sum_{i \in \mathcal{T}_k} z^{(i)}(v^{(i)})^\top \right) \quad (23)$$

$$= \mathcal{W}^{-\frac{1}{2}} (\mathcal{W}^{\frac{1}{2}} \mathcal{A}^\dagger \mathcal{W}^{\frac{1}{2}}) \mathcal{W}^{-\frac{1}{2}} \left(-\widehat{\mathcal{H}}(U^*Q) + \sum_{i \in \mathcal{T}_k} z^{(i)}(v^{(i)})^\top \right) \quad (24)$$

$$= \mathcal{W}^{-\frac{1}{2}} (\mathcal{I} + \mathcal{E}_1) \left(-(\mathcal{W}^{-\frac{1}{2}} \mathcal{H} + \mathcal{E}_2)(U^*Q) + \mathcal{W}^{-\frac{1}{2}} \left(\sum_{i \in \mathcal{T}_k} z^{(i)}(v^{(i)})^\top \right) \right) \quad (25)$$

where $\mathcal{E}_1 = (\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}})^\dagger - \mathcal{I}$ and $\mathcal{E}_2 = \mathcal{W}^{-\frac{1}{2}} \widehat{\mathcal{H}} - \mathcal{W}^{-\frac{1}{2}} \mathcal{H}$, and $F = \widehat{U} - U^*Q + \mathcal{W}^{-1}(\mathcal{H}(U^*Q))$. Let $F = \widehat{U} - U^*Q + \mathcal{W}^{-1}(\mathcal{H}(U^*Q))$

Lemma A.2. *Assume that the large probability event in Lemma A.1 holds true. Then,*

$$\|\mathcal{W}^{-1} \mathcal{H}(U^*Q)\|_F \leq O\left(\sqrt{\frac{\lambda_1^*}{\lambda_r^*}} \log\left(\frac{t}{K}\right) \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log\left(\frac{t}{K\delta}\right)}{m}}\right) \quad (26)$$

and if $mt \geq \Omega(\mu dr^2 K \log(t/K\delta))$, then with probability at least $1 - \delta/3$

$$\|F\|_F \leq O\left(\sqrt{\frac{\lambda_1^*}{\lambda_r^*}} \frac{\mu dr^2 K \log\left(\frac{t}{K\delta}\right)}{mt} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \sqrt{\frac{\mu dr^2 K \log\left(\frac{t}{K\delta}\right) \log\left(\frac{r}{\delta}\right)}{mt}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log\left(\frac{1}{\delta}\right)}{m}}\right) \quad (27)$$

A proof is in Section A.3.1.

Lemma A.3. *If $\frac{1}{2} \leq \sigma_{\min}(Q)$, $\|F\|_F \leq \frac{1}{8}$ and $\|\mathcal{W}^{-1}(\mathcal{H}(U^*Q))\|_F \leq \frac{1}{8}$, then R is invertible and $\|R^{-1}\| \leq 4$.*

A proof is in Section A.4. Clearly, from (26) and (27), a sufficient condition for the above lemma is

$$O\left(\sqrt{\frac{\lambda_1^*}{\lambda_r^*} \log\left(\frac{t}{K}\right)} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log\left(\frac{t}{K\delta}\right)}{m}}\right) \leq \frac{1}{8}, \text{ and} \quad (28)$$

$$O\left(\sqrt{\frac{\lambda_1^*}{\lambda_r^*} \frac{\mu dr^2 K \log\left(\frac{t}{K\delta}\right)}{mt}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \sqrt{\frac{\mu dr^2 K \log\left(\frac{t}{K\delta}\right) \log\left(\frac{r}{\delta}\right)}{mt}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log\left(\frac{1}{\delta}\right)}{m}}\right) \leq \frac{1}{8} \quad (29)$$

which can be satisfied with

$$\|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq O\left(\sqrt{\frac{\lambda_r^*}{\lambda_1^*} \frac{1}{\log(t/K)}}\right), \quad m \geq \Omega\left(\left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 r^2 \log\left(\frac{t}{K\delta}\right) + r^2 \log\left(\frac{1}{\delta}\right)\right), \text{ and} \quad (30)$$

$$mt \geq \Omega\left(\mu dr^2 K \left(1 + \left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 \log\left(\frac{t}{K\delta}\right) \log\left(\frac{r}{\delta}\right)\right)\right) \quad (31)$$

Finally, we bound the Frobenius norm distance of the next iterate U^+ from the optimal U^* .

$$\|(\mathbf{I} - U^*(U^*)^\top)U^+\|_F \quad (32)$$

$$= \min_{Q^+} \|U^+ - U^*Q^+\|_F \quad (33)$$

$$\leq \|\widehat{U}R^{-1} - U^*QR^{-1} + (\mathcal{W}^{-1}\mathcal{H}(U^*Q))R^{-1}\| \quad (34)$$

$$\leq \|\widehat{U} - U^*Q + \mathcal{W}^{-1}\mathcal{H}(U^*Q)\|_F \|R^{-1}\| \quad (35)$$

$$= \|F\|_F \|R^{-1}\| \quad (36)$$

$$\leq O\left(\sqrt{\frac{\lambda_1^*}{\lambda_r^*} \frac{\mu dr^2 K \log\left(\frac{t}{K\delta}\right)}{mt}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \sqrt{\frac{\mu dr^2 K \log\left(\frac{t}{K\delta}\right) \log\left(\frac{r}{\delta}\right)}{mt}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log\left(\frac{1}{\delta}\right)}{m}}\right) \quad (37)$$

If

$$mt \geq \Omega\left(\mu dr^2 K \frac{\lambda_1^*}{\lambda_r^*} \left(\log\left(\frac{t}{K\delta}\right) + \left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 \log^2\left(\frac{t}{K\delta}\right) \log\left(\frac{r}{\delta}\right)\right)\right), \text{ and } m \geq \Omega\left(r^2 \log\left(\frac{1}{\delta}\right)\right) \quad (38)$$

then,

$$\|(\mathbf{I} - U^*(U^*)^\top)U^+\|_F \leq \frac{1}{2} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \min\left(\frac{3}{8}, O\left(\sqrt{\frac{\lambda_r^*}{\lambda_1^*} \frac{1}{\log(t/K)}}\right)\right) \quad (39)$$

Thus if $\|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq \min\left(\frac{3}{4}, O\left(\sqrt{\frac{\lambda_r^*}{\lambda_1^*} \frac{1}{\log(t/K)}}\right)\right)$, then $\|(\mathbf{I} - U^*(U^*)^\top)U^+\|_F \leq \min\left(\frac{3}{4}, O\left(\sqrt{\frac{\lambda_r^*}{\lambda_1^*} \frac{1}{\log(t/K)}}\right)\right)$.

In the following lemma we prove that tasks subset used for each iteration, satisfy approximate incoherence.

Lemma A.4 (Shuffling and partition of tasks). *Let \mathcal{T}_k be the k -th subset ($k \in [K]$) of the K -way partition of the shuffled set of all t tasks. If $t \geq \Omega(\mu^2 r^3 K \log(1/\delta))$, then with a probability of at least $1 - \delta/3$,*

$$\lambda_1\left(\sum_{i \in \mathcal{T}_k} v^{*(i)}(v^{*(i)})^\top\right) = \frac{1}{K} \Theta(\lambda_1((V^*)^\top V^*)) \quad \text{and} \quad \lambda_r\left(\sum_{i \in \mathcal{T}_k} v^{*(i)}(v^{*(i)})^\top\right) = \frac{1}{K} \Theta(\lambda_r((V^*)^\top V^*)), \quad \text{for all } r' \in [r] \quad (40)$$

where are $\lambda_1(\cdot)$ and $\lambda_r(\cdot)$ are the largest and smallest, respectively, eigenvalue operators of real-symmetric $r \times r$ matrix.

A proof is in Section A.5.

Therefore, using union-bound, we can un-roll the relation, between current iterate U and the next iterate U^+ , over K iterations, starting from U_{init} and ending at some U iterations, to get

$$\|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq \frac{1}{2K} \|(\mathbf{I} - U^*(U^*)^\top)U_{\text{init}}\|_F + O\left(\sqrt{\frac{\mu dr^2 K \log(\frac{t}{K\delta}) \log(\frac{r}{\delta})}{mt}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{1}{\delta})}{m}}\right) \quad (41)$$

with probability at least $1 - K\delta$. Finally setting $K = \lceil \log_2(\frac{(\lambda_r^*/\lambda_1^*)mt}{\mu dr^2}) \rceil$ and using $m \geq \Omega(r^2 \log(\frac{1}{\delta}))$ we get that, with a probability of at least $1 - K\delta$

$$\|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq O\left(\frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{\mu dr^2 K \log(\frac{t}{K\delta}) \log(\frac{r}{\delta})}{mt}}\right) \quad (42)$$

A.2 Analysis of update on V

A.2.1 Proof of Lemma A.1

Proof of Lemma A.1. In this proof for brevity, we will first set that $\mathcal{T}_k \leftarrow [t]$, $|\mathcal{T}_k| = t/K \leftarrow t$, $S_1^{(i)} \leftarrow S^{(i)} = \frac{1}{m} \sum_{j \in [m]} x_j^{(i)} (x_j^{(i)})^\top$. This can be done due to the approximate equivalence of the subset \mathcal{T}_k by Lemma A.4. Finally at the end of the analysis we will reset $\mathcal{T}_k \leftarrow \mathcal{T}_k$, $|\mathcal{T}_k| = t/K \leftarrow t/K$, $S_1^{(i)} \leftarrow S_1^{(i)} = \frac{2}{m} \sum_{j \in [1, m/2]} x_j^{(i)} (x_j^{(i)})^\top$.

Recall the definition of $v^{(i)}$ from the update (13), and that Q^{-1} is right inverse of Q , i.e. $QQ^{-1} = \mathbf{I}$.

$$v^{(i)} - Q^{-1}v^{*(i)} = (U^\top S^{(i)}U)^\dagger (U^\top S^{(i)}(U^*Q - U))Q^{-1}v^{*(i)} + (U^\top S^{(i)}U)^\dagger U^\top z^{(i)} \quad (43)$$

We can use re-write the first term as,

$$(U^\top S^{(i)}U)^\dagger U^\top S^{(i)}(U^*Q - U)Q^{-1} \quad (44)$$

$$= (U^\top S^{(i)}U)^\dagger U^\top S^{(i)}(UU^\top + U_\perp U_\perp^\top)(U^*Q - U)Q^{-1} \quad (45)$$

$$= U^\top (U^*Q - U)Q^{-1} + (U^\top S^{(i)}U)^\dagger U^\top S^{(i)}U_\perp U_\perp^\top (U^*Q - U)Q^{-1} \quad (46)$$

$$= -U^\top (\mathbf{I} - U^*(U^*)^\top)^2 UQ^{-1} + (U^\top S^{(i)}U)^\dagger U^\top S^{(i)}U_\perp U_\perp^\top U^* \quad (47)$$

$$= -(U - U^*Q)^\top (U - U^*Q)Q^{-1} + (U^\top S^{(i)}U)^\dagger U^\top S^{(i)}U_\perp U_\perp^\top U^* \quad (48)$$

where we used the fact that $Q = (U^*)^\top U$. Therefore

$$\begin{aligned} \|v^{(i)} - Q^{-1}v^{*(i)}\| &\leq \\ \|U - U^*Q\| \| (U - U^*Q)Q^{-1}v^{*(i)} \| &+ \|(U^\top S^{(i)}U)^\dagger\| (\|U^\top S^{(i)}U_\perp U_\perp^\top U^*v^{*(i)}\| + \|U^\top z^{(i)}\|) \end{aligned} \quad (49)$$

If $m \geq \Omega(r \log(t/\delta))$, then $\alpha = c\sqrt{\frac{r \log(27t/\delta)}{m}} \leq 1/2$ and by Lemma A.5, with a probability of at least $1 - \delta$,

$$\left. \begin{aligned} \|(U^\top S^{(i)}U)^\dagger\| &\leq (1 + 2\alpha), \\ \|U^\top S^{(i)}U_\perp U_\perp^\top U^*v^{*(i)}\| &\leq \alpha \|U_\perp^\top U^*v^{*(i)}\|, \quad \text{and} \\ \|U^\top z^{(i)}\| &\leq \sigma\alpha, \end{aligned} \right\} \text{for all } i \in [t] \quad (50)$$

Now if $m \geq \Omega(r \log(1/\delta))$ and $\|U^*Q - U\| \leq O\left(\sqrt{\frac{\log(\frac{t}{\delta})}{\log(\frac{1}{\delta})}}\right)$, then

$$\|v^{(i)} - Q^{-1}v^{*(i)}\| \leq O\left(\sqrt{\frac{\log(\frac{t}{\delta})}{\log(\frac{1}{\delta})}} (\|(U^*Q - U)Q^{-1}v^{*(i)}\| + \|U_\perp^\top U^*v^{*(i)}\|) + \sigma\sqrt{\frac{r \log(\frac{t}{\delta})}{m}}\right) \quad (51)$$

Next we bound $\|H\|_F$, which by definition is $\|H\|_F = \sqrt{\sum_{i \in [t]} \|h^{(i)}\|^2} = \sqrt{\sum_{i \in [t]} \|v^{(i)} - Q^{-1}v^{*(i)}\|^2}$. Using (51) and the fact that $(a^2 + b^2) \leq 2(a^2 + b^2)$ we get

$$\|H\|_F^2 \leq \frac{\log(\frac{t}{\delta})}{\log(\frac{1}{\delta})} \left[\sum_{i \in \mathcal{T}} O(\|(U^*Q - U)Q^{-1}v^{*(i)}\|^2 + \|U_\perp^\top U^*v^{*(i)}\|^2) \right] + t(\sigma\sqrt{\frac{r \log(\frac{t}{\delta})}{m}})^2 \quad (52)$$

Clearly $\|Q\| = \|(U^*)^\top U\| \leq \|U^*\| \|U\| \leq 1$. If $\|(\mathbf{I} - U^*(U^*)^\top)U\| \leq \|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq \frac{3}{4}$, then by using Lemma D.4, $\|Q^{-1}\| \leq 2$.

$$\sum_{i \in [t]} \|(U^*Q - U)Q^{-1}v^{*(i)}\|^2 = \sum_{i \in [t]} \text{tr}((v^{*(i)})^\top ((U^*Q - U)Q^{-1})^\top (U^*Q - U)Q^{-1}v^{*(i)}) \quad (53)$$

$$= \text{tr}((U^*Q - U)Q^{-1})^\top (U^*Q - U)Q^{-1} \sum_{i \in [t]} v^{*(i)}(v^{*(i)})^\top \quad (54)$$

$$\leq \|(U^*Q - U)\|_F^2 \|Q^{-1}\|^2 O(\lambda_1^*)(t/r) \quad (55)$$

$$\leq 4\|(U^*Q - U)\|_F^2 O(\lambda_1^*)(t/r) \quad (56)$$

Similarly we can use Lemma D.4, to get

$$\sum_{i \in [t]} \|U_\perp^\top U^* v^{*(i)}\|^2 = \sum_{i \in [t]} \text{tr}((v^{*(i)})^\top (U_\perp^\top U^*)^\top U_\perp^\top U^* v^{*(i)}) \quad (57)$$

$$= \text{tr}((U_\perp^\top U^*)^\top (U_\perp^\top U^*) \sum_{i \in [t]} v^{*(i)}(v^{*(i)})^\top) \quad (58)$$

$$\leq \|U_\perp^\top U^*\|_F^2 O(\lambda_1^*)(t/r) \quad (59)$$

$$\leq \|(U^*Q - U)\|_F^2 O(\lambda_1^*)(t/r) \quad (60)$$

Therefore substituting the above two inequalities into (52) and using the fact that $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ for all $0 \leq a, b$ we get

$$\|H\|_F \leq O\left(\sqrt{\frac{\log(\frac{t}{\delta})}{\log(\frac{1}{\delta})}}\right) \|U^*Q - U\|_F \sqrt{\lambda_1^*(t/r)} + \sqrt{t}\sigma \sqrt{\frac{r \log(\frac{t}{\delta})}{m}} \quad (61)$$

Then as $\|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq O\left(\sqrt{\frac{\lambda_r^*}{\lambda_1^*} \frac{1}{\log(t)}}\right)$ and $m \geq \Omega\left(\left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 r^2 \log(\frac{t}{\delta})\right)$, $\|H\|_F \leq (1 - \frac{1}{\sqrt{2}})\sqrt{(t/r)\lambda_r^*}$.

Using $\|Q^{-1}\| \leq 2$ in (51) we also get that

$$\|h^{(i)}\| = \|v^{(i)} - Q^{-1}v^{*(i)}\| \leq O\left(\sqrt{\frac{\log(\frac{t}{\delta})}{\log(\frac{1}{\delta})}}\right) \|(U^*Q - U)\| \|v^{*(i)}\| + \sigma \sqrt{\frac{r \log(\frac{t}{\delta})}{m}} \quad (62)$$

By definition is $\|H\|_{\infty,2} = \max_{i \in [t]} \|h^{(i)}\| = \max_{i \in [t]} \|v^{(i)} - Q^{-1}v^{*(i)}\|$. Then as $\|(\mathbf{I} - U^*(U^*)^\top)U\| \leq \|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq O\left(\sqrt{\frac{\lambda_r^*}{\lambda_1^*} \frac{1}{\log(t)}}\right) \leq O(1)$, $m \geq \Omega\left(\left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 r^2 \log(\frac{t}{\delta})\right) \geq \Omega\left(\left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 r \log(\frac{t}{\delta})\right)$, $\|H\|_{\infty,2} \leq O(\mu\lambda_r^*)$. Now, using $\|H\|_F \leq (1 - \frac{1}{\sqrt{2}})\sqrt{(t/r)\lambda_r^*}$, $\|H\|_{\infty,2} \leq O(\mu\lambda_r^*)$, $\|Q\| \leq 1$ and $\frac{1}{2} \leq \sigma_{\min}(Q)$, by Lemma A.6, we get the approximate incoherence relation for the intermediate V

$$\|v^{(i)}\| \leq O(\mu\lambda_r) , \text{ and } \lambda_r^* \leq 2\lambda_r \quad (63)$$

Using this we bound $\|H\|_{\infty,2}$. Using the above incoherence relation and (62), we get

$$\sqrt{\frac{r}{t}} \frac{\|H\|_{\infty,2}}{\sqrt{\lambda_r}} \leq 2\sqrt{\frac{r}{t}} \frac{\|H\|_{\infty,2}}{\sqrt{\lambda_r^*}} \leq O\left(\sqrt{\frac{r}{t}} \sqrt{\frac{\log(\frac{t}{\delta})}{\log(\frac{1}{\delta})}}\right) \|U^*Q - U\| \max_{i \in [t]} \frac{\|v^{*(i)}\|}{\sqrt{\lambda_r^*}} + 2\sqrt{\frac{r}{t}} \frac{2c\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r \log(\frac{27t}{\delta})}{m}} \quad (64)$$

$$\leq O\left(\sqrt{\frac{\log(\frac{t}{\delta})}{\log(\frac{1}{\delta})}}\right) \sqrt{\frac{\mu r}{t}} \|U^*Q - U\| + \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{\delta})}{mt}} \quad (65)$$

Using (63) in (61), we get

$$\sqrt{\frac{r}{t}} \frac{\|H\|_F}{\sqrt{\lambda_r}} \leq 2\sqrt{\frac{r}{t}} \frac{\|H\|_F}{\sqrt{\lambda_r^*}} \leq O\left(\sqrt{\frac{\log(\frac{t}{\delta})}{\log(\frac{1}{\delta})}}\right) \sqrt{\frac{\lambda_1^*}{\lambda_r^*}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{\delta})}{m}} \quad (66)$$

Finally, by resetting $\mathcal{T}_k \leftarrow \mathcal{T}_k$, $|\mathcal{T}_k| = t/K \leftarrow t/K$, $S_1^{(i)} \leftarrow S_1^{(i)} = \frac{2}{m} \sum_{j \in [1, m/2]} x_j^{(i)}(x_j^{(i)})^\top$, we obtain the desired result. \square

A.2.2 Supporting lemmas for the analysis of update on V

Here we bound the linear operators in the $v^{(i)}$ update.

Lemma A.5. *Let $\alpha = c\sqrt{\frac{r \log(27t/\delta)}{m}}$. With a probability of at least $1 - \delta$, the following are true for all $i \in [t]$*

$$\|(U^\top S^{(i)}U)^\dagger\| \leq (1 + 2\alpha), \quad (67)$$

$$\|(U^\top S^{(i)}(U^*Q - U)Q^{-1}v^{*(i)})\| \leq (\|(\mathbf{I} - U^*(U^*)^\top)U\| + \alpha)\|(U^*Q - U)Q^{-1}v^{*(i)}\| \quad (68)$$

$$\leq (1 + \alpha)\|(U^*Q - U)Q^{-1}v^{*(i)}\| \quad (69)$$

$$\|U^\top S^{(i)}U_\perp U_\perp^\top U^*v^{*(i)}\| \leq \alpha\|U_\perp^\top U^*v^{*(i)}\|, \text{ and} \quad (70)$$

$$\|U^\top z^{(i)}\| \leq \sigma\alpha \quad (71)$$

Proof of Lemma A.5. Let $i \in [t]$.

Let $\mathcal{S} = \{v \in \mathbb{R}^r \mid \|v\| = 1\}$ be the set of all real vectors of dimension r with unit Euclidean norm. For $\epsilon \leq 1$, there exists an ϵ -net, $N_\epsilon \subset \mathcal{S}$, of size $(1 + 2/\epsilon)^r$ with respect to the Euclidean norm [Ver10, Lemma 5.2]. That is for any $v' \in \mathcal{S}$, there exists some $v \in N_\epsilon$ such that $\|v' - v\|_F \leq \epsilon$.

Consider a $v \in N_\epsilon$, such that $\|v\|_F = 1$. Now we will prove with high-probability that $\langle (U^\top S^{(i)}U) - \mathbf{I} \rangle v, v \rangle$ is small. Consider the the following quadratic form

$$v^\top (U^\top S^{(i)}U)v = \frac{1}{m} \sum_{j \in [m]} \text{tr}(v^\top (U^\top x_j^{(i)} (x_j^{(i)})^\top U)v) = \frac{1}{m} \sum_{j \in [m]} \text{tr}((x_j^{(i)})^\top U v v^\top U^\top x_j^{(i)}) \quad (72)$$

$x_j^{(i)} \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$ are i.i.d. standard Gaussian random vectors. We will use Hanson-Wright inequality (Lemma D.5) to prove that the above quadratic form concentrates around its mean. In Lemma D.6 (which is a straightforward Corollary of Hanson-Wright inequality), by setting $a \leftarrow Uv, b \leftarrow Uv$, we get that with a probability of at least $1 - \delta$

$$\left| v^\top ((U^\top S^{(i)}U) - \mathbf{I})v \right| \leq c \max \left(\sqrt{\frac{\log(1/\delta)}{m}}, \frac{\log(1/\delta)}{m} \right) := \Delta_\epsilon \quad (73)$$

For brevity, let $E = (U^\top S^{(i)}U) - \mathbf{I}$. Notice that E is a real symmetric matrix, therefore it has an eigen decomposition. Then, let $v' \in \mathcal{S} \subset \mathbb{R}^r$ be the largest ‘‘eigenvector’’ of E , such that $(v')^\top E v' = \|E\| = \max_{\|\tilde{v}\|=1} \tilde{v}^\top E \tilde{v} = \max_{\|\tilde{v}\|=\|\tilde{v}'\|_F=1} \tilde{v}'^\top E \tilde{v}'$. Then there exists some $v \in N_\epsilon$ such that $\|v' - v\| \leq \epsilon$.

$$\|E\|_F = (v')^\top E v = v^\top E v + (v' - v)^\top E v + (v')^\top E (v' - v) \quad (74)$$

$$\leq v^\top E v + \|v' - v\| \|E\| \|v\| + \|v'\| \|E\| \|v' - v\| \quad (75)$$

$$\leq v^\top E v + 2\epsilon \|E\| \quad (76)$$

Re-arranging and setting $\epsilon = 1/4$, and $c \leftarrow 2c$, we get

$$\|(U^\top S^{(i)}U) - \mathbf{I}\| = \|E\| \leq \Delta_{\frac{1}{4}} = \Delta. \quad (77)$$

where $\Delta = c \max \left(\sqrt{\frac{r \log(9/\delta)}{m}}, \frac{r \log(9/\delta)}{m} \right)$. If $m \geq \max(1, 4c^2)r \log(27t/\delta)$, then $\Delta \leq \alpha \leq 1/2$.

Thus with a probability of at least is also implies that

$$\|(U^\top S^{(i)}U)^\dagger\| = (\sigma_{\min}(U^\top S^{(i)}U))^{-1} \leq \frac{1}{1 - \alpha} \leq 2. \quad (78)$$

Using similar arguments we can also prove that with a probability of at least $1 - \delta$

$$\|(U^\top S^{(i)}(U^*Q - U)Q^{-1}v^{*(i)})\| \leq \|U^\top (U^*Q - U)Q^{-1}v^{*(i)}\| + \alpha\|(U^*Q - U)Q^{-1}v^{*(i)}\| \quad (79)$$

$$\leq \|U^\top (\mathbf{I} - U^*(U^*)^\top)UQ^{-1}v^{*(i)}\| + \alpha\|(U^*Q - U)Q^{-1}v^{*(i)}\| \quad (80)$$

$$\leq \|U^\top (\mathbf{I} - U^*(U^*)^\top)^2UQ^{-1}v^{*(i)}\| + \alpha\|(U^*Q - U)Q^{-1}v^{*(i)}\| \quad (81)$$

$$\leq \|U^\top (\mathbf{I} - U^*(U^*)^\top)(U^*Q - U)Q^{-1}v^{*(i)}\| + \alpha\|(U^*Q - U)Q^{-1}v^{*(i)}\| \quad (82)$$

$$\leq \|(\mathbf{I} - U^*(U^*)^\top)U\| \|(U^*Q - U)Q^{-1}v^{*(i)}\| + \alpha\|(U^*Q - U)Q^{-1}v^{*(i)}\| \quad (83)$$

$$\leq (\|(\mathbf{I} - U^*(U^*)^\top)U\| + \alpha)\|(U^*Q - U)Q^{-1}v^{*(i)}\| \quad (84)$$

$$\leq (1 + \alpha)\|(U^*Q - U)Q^{-1}v^{*(i)}\|, \quad (85)$$

Using similar arguments we can also prove that with a probability of at least $1 - \delta$

$$\|U^\top S^{(i)} U_\perp U_\perp^\top U^* v^{*(i)}\| \leq \alpha \|U_\perp^\top U^* v^{*(i)}\| \quad (86)$$

and with a probability of at least $1 - \delta$

$$\|U^\top z^{(i)}\| \leq \sigma \alpha \quad (87)$$

Finally setting $\delta \leftarrow \delta/3/t$ and taking the union bound over three bounds over all the tasks in $[t]$ gets us the desired result. \square

Here we prove the approximate incoherence of the intermediate V and the spectrum of intermediate W .

Lemma A.6 (Incoherence of intermediate $v^{(i)}$). *If $\|H\|_F \leq (1 - \frac{1}{\sqrt{2}})\sqrt{(t/r)\lambda_r((r/t)W^*)}$, $\|H\|_{\infty,2}^2 \leq O(\mu\lambda_r((r/t)W^*))$, $\|Q\| \leq 1$ and $\frac{1}{2} \leq \sigma_{\min}(Q)$, and (61) and (62) are true, then*

$$\|v^{(i)}\| \leq O\left(\mu\lambda_r((r/t)W)\right), \text{ and } \lambda_r((r/t)W^*) \leq 2\lambda_r((r/t)W) \quad (88)$$

Proof of Lemma A.6.

$$\|v^{(i)}\| \leq \|Q^{-1}v^{*(i)}\| + \|v^{(i)} - Q^{-1}v^{*(i)}\| \leq 2\|v^{*(i)}\| + \|h^{(i)}\| \quad (89)$$

$$\implies \|v^{(i)}\|^2 \leq O(\|V^*\|_{\infty,2}^2) + O(\|H\|_{\infty,2}^2) \leq O\left(\mu\lambda_r((r/t)W^*)\right) \quad (90)$$

where the second inequality use the definition $h^{(i)} = v^{(i)} - Q^{-1}v^{*(i)}$ and $\|Q^{-1}\| \leq 2$ (as $\sigma_{\min}(Q) \geq \frac{1}{2}$), the third inequality use the fact that $a + b \leq 2a^2 + 2b^2$ an (62), and the final inequality uses $\|H\|_{\infty,2} \leq \|V\|_{\infty,2}$.

Notice that $W = V^T V$ and $W^* = (V^*)^T V^*$. Thus $\sqrt{\lambda_r((r/t)W)} = \sqrt{(r/t)\sigma_r(V)}$ and $\sqrt{\lambda_r((r/t)W^*)} = \sqrt{(r/t)\sigma_r(W^*)}$, and both W and W^* are positive semi-definite (PSD). Similarly, using $\sigma_{\min}(Q^{-1}) = \sigma_{\min}(((U^*)^\top U)^{-1}) \geq 1$ and Lemma D.1 we can get that

$$\sqrt{\lambda_r((r/t)W^*)} \leq \sqrt{\sigma_{\min}^2(Q^{-1})\lambda_r((r/t)W^*)} \leq \sqrt{(r/t)\lambda_r(Q^{-1}(V^*)^T V^* Q^{-\top})} \leq \sqrt{(r/t)\sigma_r(V^* Q^{-T})} \quad (91)$$

Therefore, instead of analyzing the relation between $\lambda_r(W)$ and $\lambda_r(W^*)$, we can analyze the relation between $\sigma_r(V)$ and $\sigma_r(W^*)$. Notice that $V^* Q^{-T} = V + V^* Q^{-T} - V$. Then by Weyl's inequality (Lemma D.2, by setting $A \leftarrow V^* Q^{-T}$, $B \leftarrow V$, and $C \leftarrow V^* Q^{-T} - V$) we get that

$$\sqrt{\lambda_r((r/t)W^*)} \leq \sqrt{(r/t)\sigma_r(V^* Q^{-T})} \leq \sqrt{(r/t)\sigma_r(V)} + \sqrt{(r/t)\|V - V^* Q^{-T}\|} \quad (92)$$

$$\leq \sqrt{\lambda_r((r/t)W)} + \sqrt{(r/t)\|H\|} \quad (93)$$

$$\leq \sqrt{\lambda_r((r/t)W)} + \sqrt{(r/t)\|H\|_F} \quad (94)$$

$$\leq \sqrt{\lambda_r((r/t)W)} + \left(1 - \frac{1}{\sqrt{2}}\right)\sqrt{\lambda_r((r/t)W^*)} \quad (95)$$

where the last inequality uses $\|H\|_F \leq (1 - \frac{1}{\sqrt{2}})\sqrt{(t/r)\lambda_r((r/t)W^*)}$. Finally we get the desired result by re-arranging the terms. \square

A.3 Analysis of update on U

A.3.1 Proof of Lemma A.2

Proof of Lemma A.2. In this proof for brevity, we will first set that $\mathcal{T}_k \leftarrow [t]$, $|\mathcal{T}_k| = t/K \leftarrow t$, $S_2^{(i)} \leftarrow S^{(i)} = \frac{1}{m} \sum_{j \in [m]} x_j^{(i)} (x_j^{(i)})^\top$. This can be done due to the approximate equivalence of the subset \mathcal{T}_k by Lemma A.4. Finally at the end of the analysis we will reset $\mathcal{T}_k \leftarrow \mathcal{T}_k$, $|\mathcal{T}_k| = t/K \leftarrow t/K$, $S_2^{(i)} \leftarrow S_2^{(i)} = \frac{2}{m} \sum_{j \in [m/2+1, m]} x_j^{(i)} (x_j^{(i)})^\top$.

Recall that

$$\widehat{U} - U^* Q = \mathcal{W}^{-\frac{1}{2}} (\mathcal{I} + \mathcal{E}_1) (-\mathcal{W}^{-\frac{1}{2}} \mathcal{H} + \mathcal{E}_2) (U^* Q) + \mathcal{W}^{-\frac{1}{2}} \left(\sum_{i \in [t]} z^{(i)} (v^{(i)})^\top \right) \quad (96)$$

where $\mathcal{E}_1 = (\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}})^\dagger - \mathcal{I}$ and $\mathcal{E}_2 = \mathcal{W}^{-\frac{1}{2}} \widehat{\mathcal{H}} - \mathcal{W}^{-\frac{1}{2}} \mathcal{H}$, and $F = \widehat{U} - U^* Q + \mathcal{W}^{-1}(\mathcal{H}(U^* Q))$. Therefore

$$\|F\|_F \leq \|\mathcal{W}^{-\frac{1}{2}}\|_F (\|\mathcal{E}_1\|_F \|\mathcal{W}^{-\frac{1}{2}} \mathcal{H}(U^* Q)\|_F + \|\mathcal{I} + \mathcal{E}_1\|_F (\|\mathcal{E}_2(U^* Q)\|_F + \|\mathcal{W}^{-\frac{1}{2}}(\sum_{i \in [t]} z^{(i)}(v^{(i)})^\top)\|_F)) \quad (97)$$

We can trivially bound $\|\mathcal{W}^{-\frac{1}{2}}\|_F$ as follows. For all $\|U\|_F = 1$, the following is true.

$$\|\mathcal{W}^{-\frac{1}{2}}(U)\|_F = \|U \mathcal{W}^{-\frac{1}{2}}\|_F \leq \|U\|_F \|\mathcal{W}^{-\frac{1}{2}}\| \leq \sqrt{\frac{r/t}{\lambda_r}} \quad (98)$$

$\Omega(\mu d r^2 \log(1/\delta)) \leq m t$ and approximate incoherence of intermediate V (19) implies that $\Omega(d r \frac{\|V\|_{\infty,2}^2}{\lambda_r(W)/t} \log(1/\delta)) \leq \Omega(\mu d r^2 \log(1/\delta)) \leq m t$, then by Lemma A.7 we have that, with a probability of at least $1 - \delta/3$

$$\|\mathcal{E}_1\|_F \leq 3c \sqrt{\frac{d r \|V\|_{\infty,2}^2 \log(27/\delta)}{m \lambda_r(W)}} \leq 3c \sqrt{\frac{\mu d r^2 \log(27/\delta)}{m t}} \leq \frac{1}{2} \quad (99)$$

This also implies that

$$\|\mathcal{I} + \mathcal{E}_1\|_F \leq \|\mathcal{I}\| + \|\mathcal{E}_1\|_F \leq 1 + \Delta \leq \frac{3}{2} \quad (100)$$

By Lemma A.8,

$$\|(\mathcal{W}^{-\frac{1}{2}} \mathcal{H})(U^* Q)\|_F \leq \|H\|_F \quad (101)$$

and with a probability of at least $1 - \delta/3$

$$\|\mathcal{E}_2(U^* Q)\|_F \leq c(\min(\|H\|_F \frac{\|V\|_{\infty,2}}{\sqrt{\lambda_r(W)}}, \|H\|_{\infty,2}) \sqrt{\frac{d r \log(15/\delta)}{m}} + \|H\|_{\infty,2} \frac{\|V\|_{\infty,2}}{\sqrt{\lambda_r(W)}} \frac{d r \log(15/\delta)}{m}) \quad (102)$$

Using the approximate incoherence of V (19) in the above inequality, we get that

$$\|\mathcal{E}_2(U^* Q)\|_F \leq c(\min(\|H\|_F \sqrt{\frac{\mu r}{t}}, \|H\|_{\infty,2}) \sqrt{\frac{d r \log(15/\delta)}{m}} + \|H\|_{\infty,2} \sqrt{\frac{\mu r}{t}} \cdot \frac{d r \log(15/\delta)}{m}) \quad (103)$$

By Lemma A.9 with a probability of at least $1 - \delta/3$

$$\|\sum_{i \in [t]} \mathcal{W}^{-\frac{1}{2}}(z^{(i)}(v^{(i)})^\top)\|_F \leq O\left(\sigma \sqrt{\frac{d r}{m} \log\left(\frac{t}{\delta}\right) \log\left(\frac{r}{\delta}\right)}\right) \quad (104)$$

Finally taking union bound over the above results and using Lemma A.1, we can bound each of the terms constituting F . Using (98), (101) and (20) (recall that we set $t \leftarrow t/K$) we get

$$\|\mathcal{W}^{-1} \mathcal{H}(U^* Q)\|_F \leq \|\mathcal{W}^{-\frac{1}{2}}\|_F \|\mathcal{W}^{-\frac{1}{2}} \mathcal{H}(U^* Q)\|_F \quad (105)$$

$$\leq \sqrt{\frac{r}{t}} \frac{\|H\|_F}{\sqrt{\lambda_r}} \leq O\left(\sqrt{\frac{\lambda_1^*}{\lambda_r^*}} \sqrt{\frac{\log(\frac{t}{\delta})}{\log(\frac{1}{\delta})}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{\delta})}{m}}\right) \quad (106)$$

Using (98), (100), (101), and (20) we get

$$\|\mathcal{W}^{-\frac{1}{2}}\|_F \|\mathcal{E}_1\|_F \|\mathcal{W}^{-\frac{1}{2}} \mathcal{H}(U^* Q)\|_F \quad (107)$$

$$\leq O\left(\sqrt{\frac{\mu d r^2 \log(\frac{1}{\delta})}{m t}} \sqrt{\frac{r}{t}} \frac{\|H\|_F}{\sqrt{\lambda_r}}\right) \quad (108)$$

$$\leq O\left(\sqrt{\frac{\lambda_1^*}{\lambda_r^*}} \frac{\mu d r^2 \log(\frac{t}{\delta})}{m t} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \sqrt{\frac{\mu d r^2 \log(\frac{1}{\delta})}{m t}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{\delta})}{m}}\right) \quad (109)$$

Using (98), (100), (103), (20) and (21) we get

$$\|\mathcal{W}^{-\frac{1}{2}}\|_F \|\mathcal{I} + \mathcal{E}_1\|_F (\|\mathcal{E}_2(U^*Q)\|_F) \quad (110)$$

$$\leq O\left(\sqrt{\frac{r}{t}} \min\left(\frac{\|H\|_F}{\sqrt{\lambda_r}} \sqrt{\frac{\mu r}{t}}, \frac{\|H\|_{\infty,2}}{\sqrt{\lambda_r}}\right) \sqrt{\frac{dr \log(\frac{1}{\delta})}{m}} + \sqrt{\frac{r}{t}} \frac{\|H\|_{\infty,2}}{\sqrt{\lambda_r}} \sqrt{\frac{\mu r}{t}} \sqrt{\frac{dr \log(\frac{1}{\delta})}{m}}\right) \quad (111)$$

$$\leq O\left(\min\left(\sqrt{\frac{\lambda_1^* \mu dr^2 \log(\frac{t}{\delta})}{\lambda_r^* mt}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \sqrt{\frac{\mu dr^2 \log(\frac{1}{\delta})}{mt}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{\delta})}{m}}\right), \quad (112)$$

$$\sqrt{\frac{\mu dr^2 \log(\frac{t}{\delta})}{mt}} \|(\mathbf{I} - U^*(U^*)^\top)U\| + \sqrt{\frac{dr \log(\frac{1}{\delta})}{m}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{\delta})}{mt}}) + \quad (113)$$

$$\frac{\mu dr^2 \log(\frac{t}{\delta})}{mt} \|(\mathbf{I} - U^*(U^*)^\top)U\| + \frac{\sqrt{\mu} dr \sqrt{r} \log(\frac{1}{\delta})}{m\sqrt{t}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{\delta})}{mt}}) \quad (114)$$

Using (98), (100), (104), and (19) we get

$$\|\mathcal{W}^{-\frac{1}{2}}\|_F \|\mathcal{I} + \mathcal{E}_1\|_F \left\| \sum_{i \in [t]} \mathcal{W}^{-\frac{1}{2}}(z^{(i)}(v^{(i)})^\top) \right\|_F \leq O\left(\frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{dr^2 \log(\frac{t}{\delta}) \log(\frac{r}{\delta})}{mt}}\right) \quad (115)$$

Substituting (106), (109), (114), and (115) in (97) we get

$$\|F\|_F \leq \|\mathcal{W}^{-\frac{1}{2}}\|_F (\|\mathcal{E}_1\|_F \|\mathcal{W}^{-\frac{1}{2}} \mathcal{H}(U^*Q)\|_F + \|\mathcal{I} + \mathcal{E}_1\|_F (\|\mathcal{E}_2(U^*Q)\|_F + \left\| \sum_{i \in [t]} \mathcal{W}^{-\frac{1}{2}}(z^{(i)}(v^{(i)})^\top) \right\|_F)) \quad (116)$$

$$\leq O\left(\sqrt{\frac{\lambda_1^* \mu dr^2 \log(\frac{t}{\delta})}{\lambda_r^* mt}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \sqrt{\frac{\mu dr^2 \log(\frac{1}{\delta})}{mt}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{\delta})}{m}}\right) + \quad (117)$$

$$O\left(\frac{\mu dr^2 \log(\frac{t}{\delta})}{mt} \|(\mathbf{I} - U^*(U^*)^\top)U\| + \frac{\sqrt{\mu} dr \sqrt{r} \log(\frac{1}{\delta})}{mt} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{\delta})}{m}}\right) + O\left(\frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{dr^2 \log(\frac{t}{\delta}) \log(\frac{r}{\delta})}{mt}}\right) \quad (118)$$

$$\leq O\left(\sqrt{\frac{\lambda_1^* \mu dr^2 \log(\frac{t}{\delta})}{\lambda_r^* mt}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \sqrt{\frac{\mu dr^2 \log(\frac{1}{\delta})}{mt}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{\delta})}{m}}\right) + O\left(\frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{dr^2 \log(\frac{t}{\delta}) \log(\frac{r}{\delta})}{mt}}\right) \quad (119)$$

$$\leq O\left(\sqrt{\frac{\lambda_1^* \mu dr^2 \log(\frac{t}{\delta})}{\lambda_r^* mt}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \sqrt{\frac{\mu dr^2 \log(\frac{t}{\delta}) \log(\frac{r}{\delta})}{mt}} \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{1}{\delta})}{m}}\right) \quad (120)$$

where the second-last inequality used the fact that $mt \geq \Omega(\mu dr^2 \log(\frac{t}{\delta}))$. Finally, by resetting $\mathcal{T}_k \leftarrow \mathcal{T}_k$, $|\mathcal{T}_k| = t/K \leftarrow t/K$, $S_2^{(i)} \leftarrow S_2^{(i)} = \frac{2}{m} \sum_{j \in [m/2+1, m]} x_j^{(i)} (x_j^{(i)})^\top$, we obtain the desired result. \square

A.3.2 Supporting lemmas for the analysis of update on U

Lemma A.7. *If $\max(1, 4c^2) dr \frac{\|V\|_{\infty,2}^2}{\lambda_r(W)t} \log(27/\delta) \leq mt$, then with a probability of at least $1 - \delta/3$,*

$$\|\mathcal{E}_1\|_F \leq 3c \sqrt{\frac{dr \|V\|_{\infty,2}^2 \log(27/\delta)}{m \lambda_r(W)}} \quad (121)$$

Proof of Lemma A.7. Let $\mathcal{S}_F = \{U \in \mathbb{R}^{d \times r} \mid \|U\|_F = 1\}$ be the set of all real matrices of dimensions $d \times r$ with unit Frobenius norm. For $\epsilon \leq 1$, there exists an ϵ -net, $N_\epsilon \subset \mathcal{S}_F$, of size $(1 + 2/\epsilon)^{dr}$ with respect to the Frobenius norm [Ver10, Lemma 5.2]. That is for any $U' \in \mathcal{S}_F$, there exists some $U \in N_\epsilon$ such that $\|U' - U\|_F \leq \epsilon$.

Consider a $U \in N_\epsilon$, such that $\|U\|_F = 1$. Now we will prove with high-probability that $\langle (\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}} - \mathcal{I})(U), U \rangle$ is small. Consider the the following quadratic form

$$\langle (\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}})(U), U \rangle = \left\langle \sum_{i \in [t]} S^{(i)} U W^{-\frac{1}{2}} v^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}}, U \right\rangle \quad (122)$$

$$= \sum_{i \in [t]} \frac{1}{m} \sum_{j \in [m]} (x_j^{(i)})^\top (U W^{-\frac{1}{2}} v^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}} U^\top) x_j^{(i)} \quad (123)$$

where $S^{(i)} = \frac{1}{m} \sum_{j \in [m]} x_j^{(i)} (x_j^{(i)})^\top$ and $x_j^{(i)} \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$ are i.i.d. standard Gaussian random vectors and $W = \sum_{i \in [t]} v^{(i)} (v^{(i)})^\top$ is rank- r matrix. We will use Hanson-Wright inequality (Lemma D.5) to prove that the above quadratic form concentrates around its mean. Notice that the the expectation of $\langle (\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}})(U), U \rangle$ is $\langle \mathcal{I}(U), U \rangle$.

$$\sum_{i \in [t]} \mathbb{E} \left[\left\langle S^{(i)} U W^{-\frac{1}{2}} v^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}}, U \right\rangle \right] = \left\langle U W^{-\frac{1}{2}} \sum_{i \in [t]} v^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}}, U \right\rangle = \langle U, U \rangle = \|U\|_F^2 = 1. \quad (124)$$

We will also need the following bounds to apply the Hanson-Wright inequality. Recall that $\|V\|_{\infty,2} = \max_{i \in [t]} \|v^{(i)}\|$. Then,

$$\max_{i \in [t]} \|U W^{-\frac{1}{2}} v^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}} U^\top\| = \max_{i \in [t]} \|U W^{-\frac{1}{2}} v^{(i)}\|^2 \leq \max_{i \in [t]} \|U\|^2 \|W\|^2 \|v^{(i)}\|^2 \leq \frac{\|V\|_{\infty,2}^2}{\lambda_r(W)} \quad (125)$$

Also note that,

$$\sum_{i \in [t]} \|U W^{-\frac{1}{2}} v^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}} U^\top\|_F^2 = \sum_{i \in [t]} \|U W^{-\frac{1}{2}} v^{(i)}\|^4 \quad (126)$$

$$= \max_{i \in [t]} \|U W^{-\frac{1}{2}} v^{(i)}\|^2 \sum_{i \in [t]} \left\langle U W^{-\frac{1}{2}} v^{(i)}, U W^{-\frac{1}{2}} v^{(i)} \right\rangle \quad (127)$$

$$\leq \frac{\|V\|_{\infty,2}^2}{\lambda_r(W)} \quad (128)$$

where the last inequality used (124) and (125). Then by Hanson-Wright inequality (Lemma D.5), with probability at least $1 - \delta/|N_\epsilon|$

$$\left| \langle (\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}} - \mathcal{I})(U), U \rangle \right| = \left| \left\langle \sum_{i \in [t]} \frac{1}{m} \sum_{j \in [m]} x_j^{(i)} (x_j^{(i)})^\top U W^{-\frac{1}{2}} v^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}}, U \right\rangle - \langle U, U \rangle \right| \leq \Delta_\epsilon \quad (129)$$

where $\Delta_\epsilon = c \max\left(\sqrt{\frac{\|V\|_{\infty,2}^2 \log(|N_\epsilon|/\delta)}{m \lambda_r(W)}}, \frac{\|V\|_{\infty,2}^2 \log(|N_\epsilon|/\delta)}{m \lambda_r(W)}\right)$. Taking union bound over all $U \in N_\epsilon$ implies that with probability at least $1 - \delta$

$$\left| \langle (\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}} - \mathcal{I})(U), U \rangle \right| \leq \Delta_\epsilon, \text{ for all } U \in N_\epsilon. \quad (130)$$

For brevity, let $\mathcal{E}'_1(U) = (\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}} - \mathcal{I})(U)$. Notice that \mathcal{E}'_1 is self-adjoint, therefore it has an eigen decomposition with respect to the Frobenius norm. Then, let $U' \in \mathcal{S}_F \subset \mathbb{R}^{d \times r}$ be the largest ‘‘eigenmatrix’’ of \mathcal{E}'_1 , such that $\langle \mathcal{E}'_1(U), U \rangle = \|\mathcal{E}'_1\|_F = \max_{\|\tilde{U}\|_F=1} \langle \mathcal{E}'_1(\tilde{U}), \tilde{U} \rangle = \max_{\|\tilde{U}\|_F=\|\tilde{U}'\|_F=1} \langle \mathcal{E}'_1(\tilde{U}), \tilde{U}' \rangle$. Then there exists some $U \in N_\epsilon$ such that $\|U' - U\|_F \leq \epsilon$.

$$\|\mathcal{E}'_1\|_F = \langle \mathcal{E}'_1(U'), U' \rangle = \langle \mathcal{E}'_1(U), U \rangle + \langle \mathcal{E}'_1(U' - U), U \rangle + \langle \mathcal{E}'_1(U'), U' - U \rangle \quad (131)$$

$$\leq \langle \mathcal{E}'_1(U), U \rangle + \|\mathcal{E}'_1\|_F \|U' - U\|_F (\|U\|_F + \|U'\|_F) \quad (132)$$

$$\leq \langle \mathcal{E}'_1(U), U \rangle + 2\epsilon \|\mathcal{E}'_1\|_F \quad (133)$$

Re-arranging and setting $\epsilon = 1/4$, and $c \leftarrow 2c$, we get

$$\|\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}} - \mathcal{I}\|_F = \|\mathcal{E}'_1\|_F \leq \Delta_{\frac{1}{4}} = \Delta. \quad (134)$$

where $\Delta = c \max \left(\sqrt{\frac{dr \|V\|_{\infty,2}^2 \log(9/\delta)}{m \lambda_r(W)}}, \frac{dr \|V\|_{\infty,2}^2 \log(9/\delta)}{m \lambda_r(W)} \right)$.

For brevity, let $\widehat{\mathcal{A}}(U) = (\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}})(U)$. Notice that $\widehat{\mathcal{A}}$ is self-adjoint, therefore it has an eigen decomposition with respect to the Frobenius norm. Then, let $U' \in \mathcal{S}_F \subset \mathbb{R}^{d \times r}$ be the smallest ‘‘eigenmatrix’’ of $\widehat{\mathcal{A}}$, such that $\langle \widehat{\mathcal{A}}(U), U \rangle = \lambda_{\min}(\widehat{\mathcal{A}}) = \min_{\|\tilde{U}\|_F=1} \langle \widehat{\mathcal{A}}(\tilde{U}), \tilde{U} \rangle = \min_{\|\tilde{U}\|_F=\|U'\|_F=1} \langle \widehat{\mathcal{A}}(\tilde{U}), \tilde{U} \rangle$. Then there exists some $U \in N_\epsilon$ such that $\|U' - U\|_F \leq \epsilon$.

$$\lambda_{\min}(\widehat{\mathcal{A}}) = \langle \widehat{\mathcal{A}}(U'), U' \rangle = \langle \mathcal{I}(U), U \rangle + \langle (\widehat{\mathcal{A}} - \mathcal{I})(U), U \rangle + \langle \widehat{\mathcal{A}}(U' - U), U \rangle + \langle \widehat{\mathcal{A}}(U'), U' - U \rangle \quad (135)$$

$$\geq 1 - |\langle (\widehat{\mathcal{A}} - \mathcal{I})(U), U \rangle| - \lambda_{\min}(\widehat{\mathcal{A}}) \|U' - U\|_F (\|U\|_F + \|U'\|_F) \quad (136)$$

$$\geq 1 - \Delta_\epsilon - 2\epsilon \lambda_{\min}(\widehat{\mathcal{A}}) \quad (137)$$

Re-arranging and setting $\epsilon = 1/4$, and $c \leftarrow 2c$, we get that $\lambda_{\min}(\widehat{\mathcal{A}}) \geq \frac{2}{3}(1 - \Delta)$. Therefore,

$$\|(\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}})^\dagger\|_F = \frac{1}{\lambda_{\min}(\widehat{\mathcal{A}})} \leq \frac{3}{2(1 - \Delta)}. \quad (138)$$

where $\Delta = c \max \left(\sqrt{\frac{dr \|V\|_{\infty,2}^2 \log(9/\delta)}{m \lambda_r(W)}}, \frac{dr \|V\|_{\infty,2}^2 \log(9/\delta)}{m \lambda_r(W)} \right)$. If $\max(1, 4c^2) dr \frac{\|V\|_{\infty,2}^2}{\lambda_r(W)/t} \log(27/\delta) \leq mt$, we get that $\Delta \leq c \sqrt{\frac{dr \|V\|_{\infty,2}^2 \log(9/\delta)}{m \lambda_r(W)}} \leq \frac{1}{2}$.

By setting $A + B = \mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}}$ and $A = \mathcal{I}$ such that $\mathcal{E}_1 = (A + B)^{-1} - B^{-1}$, in the Woodbury matrix inverse identity (235) (Lemma D.3) we get that, with a probability of at least $1 - \delta$

$$\|(A + B)^{-1} - A^{-1}\|_F \leq \|A^{-1}\|_F \|B\|_F \|(A + B)^{-1}\|_F \quad (139)$$

$$\implies \|\mathcal{E}_1\|_F \leq \|(\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}})^\dagger - \mathcal{I}\|_F \leq \|\mathcal{I}^\dagger\|_F \|\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}} - \mathcal{I}\|_F \|(\mathcal{W}^{-\frac{1}{2}} \mathcal{A} \mathcal{W}^{-\frac{1}{2}})^\dagger\|_F \quad (140)$$

$$\leq 1 \cdot \Delta \cdot \frac{3}{2(1 - \Delta)} \leq 3\Delta \leq 3c \sqrt{\frac{dr \|V\|_{\infty,2}^2 \log(9/\delta)}{m \lambda_r(W)}} \quad (141)$$

Finally, setting $\delta \leftarrow \delta/3$ get us the desired result. \square

Lemma A.8. $\|(\mathcal{W}^{-\frac{1}{2}} \mathcal{H})(U^* Q)\|_F \leq \|H\|_F$ and with a probability of at least $1 - \delta/3$

$$\|\mathcal{E}_2(U^* Q)\|_F \leq c(\min(\|H\|_F \frac{\|V\|_{\infty,2}}{\sqrt{\lambda_r(W)}}, \|H\|_{\infty,2}) \sqrt{\frac{dr \log(15/\delta)}{m}} + \|H\|_{\infty,2} \frac{\|V\|_{\infty,2}}{\sqrt{\lambda_r(W)}} \frac{dr \log(15/\delta)}{m}) \quad (142)$$

Proof of Lemma A.8. First we prove that the expected value $\mathbb{E}[(\mathcal{W}^{-\frac{1}{2}} \widehat{\mathcal{H}})(U^* Q)] = (\mathcal{W}^{-\frac{1}{2}} \mathcal{H})(U^* Q)$ is bounded.

$$\|(\mathcal{W}^{-\frac{1}{2}} \mathcal{H})(U^* Q)\|_F = \max_{\|U\|_F=1} \langle (\mathcal{W}^{-\frac{1}{2}} \mathcal{H})(U^* Q), U \rangle \quad (143)$$

$$= \max_{\|U\|_F=1} \sum_{i \in [t]} \langle U^* Q h^{(i)}(v^{(i)})^\top \mathcal{W}^{-\frac{1}{2}}, U \rangle \quad (144)$$

$$= \max_{\|U\|_F=1} \sum_{i \in [t]} \langle U^* Q h^{(i)}, U \mathcal{W}^{-\frac{1}{2}} v^{(i)} \rangle \quad (145)$$

$$\leq \max_{\|U\|_F=1} \sqrt{\sum_{i \in [t]} \|U^* Q h^{(i)}\|^2} \sqrt{\sum_{i \in [t]} \langle U \mathcal{W}^{-\frac{1}{2}} v^{(i)}, U \mathcal{W}^{-\frac{1}{2}} v^{(i)} \rangle} \quad (146)$$

$$\leq \max_{\|U\|_F=1} \|Q\| \sqrt{\sum_{i \in [t]} \|h^{(i)}\|^2} \sqrt{\left\langle U \sum_{i \in [t]} \mathcal{W}^{-\frac{1}{2}} v^{(i)} (v^{(i)})^\top \mathcal{W}^{-\frac{1}{2}}, U \right\rangle} \quad (147)$$

$$\leq \max_{\|U\|_F=1} \|H\|_F \|U\|_F = \|H\|_F \quad (148)$$

where used the fact that $\langle AB, C \rangle = \langle A, CB^\top \rangle$ and $(U^*)^\top U^* = \mathbf{I}$.

Let $\mathcal{S}_F = \{U \in \mathbb{R}^{d \times r} \mid \|U\|_F = 1\}$ be the set of all real matrices of dimensions $d \times r$ with unit Frobenius norm. For $\epsilon \leq 1$, there exists an ϵ -net, $N_\epsilon \subset \mathcal{S}_F$, of size $(1+2/\epsilon)^{dr}$ with respect to the Frobenius norm [Ver10, Lemma 5.2]. That is for any $U' \in \mathcal{S}_F$, there exists some $U \in N_\epsilon$ such that $\|U' - U\|_F \leq \epsilon$.

Consider a $U \in N_\epsilon$, such that $\|U\|_F = 1$. Now we will prove with high-probability that $\langle (\mathcal{W}^{-\frac{1}{2}} \mathcal{H})(U^* Q)(U) - \mathcal{W}^{-\frac{1}{2}} (\sum_{i \in [t]} S^{(i)} U^* Q h^{(i)}(v^{(i)})^\top), U \rangle$ is small. Consider the the following quadratic form

$$\langle \mathcal{W}^{-\frac{1}{2}} (\sum_{i \in [t]} S^{(i)} U^* Q h^{(i)}(v^{(i)})^\top), U \rangle = \left\langle \sum_{i \in [t]} S^{(i)} U^* Q h^{(i)}(v^{(i)})^\top \mathcal{W}^{-\frac{1}{2}}, U \right\rangle \quad (149)$$

$$= \sum_{i \in [t]} \frac{1}{m} \sum_{j \in [m]} (x_j^{(i)})^\top (U^* Q h^{(i)}(v^{(i)})^\top \mathcal{W}^{-\frac{1}{2}} U^\top) x_j^{(i)} \quad (150)$$

where $S^{(i)} = \frac{1}{m} \sum_{j \in [m]} x_j^{(i)} (x_j^{(i)})^\top$ and $x_j^{(i)} \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$ are i.i.d. standard Gaussian random vectors and $W = \sum_{i \in [t]} v^{(i)} (v^{(i)})^\top$ is rank- r matrix. We will use Hanson-Wright inequality (Lemma D.5) to prove that the above quadratic form concentrates around its mean. Notice that the the expectation of $\langle \mathcal{W}^{-\frac{1}{2}} (\sum_{i \in [t]} S^{(i)} U^* Q h^{(i)}(v^{(i)})^\top), U \rangle$ is $\langle \mathcal{W}^{-\frac{1}{2}} \mathcal{H}(U), U \rangle$.

$$\mathbb{E}[\mathcal{W}^{-\frac{1}{2}} (\sum_{i \in [t]} S^{(i)} U^* Q h^{(i)}(v^{(i)})^\top)] = \mathcal{W}^{-\frac{1}{2}} (\sum_{i \in [t]} U^* Q h^{(i)}(v^{(i)})^\top) = (\mathcal{W}^{-\frac{1}{2}} \mathcal{H})(U^* Q). \quad (151)$$

We will also need the following bounds to apply the Hanson-Wright inequality. Recall that $\|H\|_{\infty,2} = \max_{i \in [t]} \|h^{(i)}\|$ and $\|V\|_{\infty,2} = \max_{i \in [t]} \|v^{(i)}\|$. Then,

$$\max_{i \in [t]} \|U^* Q h^{(i)}(v^{(i)})^\top \mathcal{W}^{-\frac{1}{2}} U^\top\| \leq \max_{i \in [t]} \|U^*\| \|Q\| \|h^{(i)}\| \max_{i \in [t]} \frac{\|v^{(i)}\|}{\sqrt{\lambda_r(W)}} \|U\| \leq \|H\|_{\infty,2} \frac{\|V\|_{\infty,2}}{\sqrt{\lambda_r(W)}} \quad (152)$$

Also note that

$$\sum_{i \in [t]} \|U^* Q h^{(i)}(v^{(i)})^\top \mathcal{W}^{-\frac{1}{2}} U^\top\|_F^2 = \sum_{i \in [t]} \|U^* Q h^{(i)}\|^2 \|U \mathcal{W}^{-\frac{1}{2}} v^{(i)}\|^2 \quad (153)$$

$$\leq (\sum_{i \in [t]} \|U^* Q h^{(i)}\|^2) (\max_{i \in [t]} \|U \mathcal{W}^{-\frac{1}{2}} v^{(i)}\|^2) \quad (154)$$

$$\leq (\|Q\|^2 \sum_{i \in [t]} \|h^{(i)}\|^2) (\max_{i \in [t]} \|U\|^2 \|W^{-\frac{1}{2}}\|^2 \|v^{(i)}\|^2) \quad (155)$$

$$\leq \|H\|_F^2 \frac{\|V\|_{\infty,2}^2}{\lambda_r(W)} \quad (156)$$

and

$$\sum_{i \in [t]} \|U^* Q h^{(i)}(v^{(i)})^\top \mathcal{W}^{-\frac{1}{2}} U^\top\|_F^2 = \sum_{i \in [t]} \|U^* Q h^{(i)}\|^2 \|U \mathcal{W}^{-\frac{1}{2}} v^{(i)}\|^2 \quad (157)$$

$$\leq (\max_{i \in [t]} \|U^* Q h^{(i)}\|^2) \text{tr}(U \mathcal{W}^{-\frac{1}{2}} \sum_{i \in [t]} v^{(i)} (v^{(i)})^\top \mathcal{W}^{-\frac{1}{2}} U^\top) \quad (158)$$

$$\leq \|Q\| \max_{i \in [t]} \|h^{(i)}\|^2 \|U\|_F^2 \quad (159)$$

$$= \|H\|_{\infty,2}^2. \quad (160)$$

Therefore, $\sum_{i \in [t]} \|U^* Q h^{(i)}(v^{(i)})^\top \mathcal{W}^{-\frac{1}{2}} U^\top\|_F^2 \leq \min\{\|H\|_F^2 \frac{\|V\|_{\infty,2}^2}{\lambda_r(W)}, \|H\|_{\infty,2}^2\}$. For brevity, let $\mathcal{E}_2(U) = \mathcal{W}^{-\frac{1}{2}} (\sum_{i \in [t]} S^{(i)} U^* Q h^{(i)}(v^{(i)})^\top) - (\mathcal{W}^{-\frac{1}{2}} \mathcal{H})(U)$. Then by Hanson-Wright inequality (Lemma D.5), with probability at least $1 - \delta/|N_\epsilon|$

$$|\langle \mathcal{E}_2(U^* Q), U \rangle| = \left| \left\langle \sum_{i \in [t]} \frac{1}{m} \sum_{j \in [m]} x_j^{(i)} (x_j^{(i)})^\top U^* Q h^{(i)}(v^{(i)})^\top \mathcal{W}^{-\frac{1}{2}}, U \right\rangle - \left\langle (\mathcal{W}^{-\frac{1}{2}} \mathcal{H})(U^* Q), U \right\rangle \right| \leq \Delta_\epsilon \quad (161)$$

where $\Delta_\epsilon = c(\min(\|H\|_F \frac{\|V\|_{\infty,2}}{\sqrt{\lambda_r(W)}}, \|H\|_{\infty,2}) \sqrt{\frac{\log(|N_\epsilon|/\delta)}{m}} + \|H\|_{\infty,2} \frac{\|V\|_{\infty,2}}{\sqrt{\lambda_r(W)}} \frac{\log(|N_\epsilon|/\delta)}{m})$. Taking union bound over all $U \in N_\epsilon$ implies that with probability at least $1 - \delta$

$$|\langle \mathcal{E}_2(U), U \rangle| \leq \Delta_\epsilon, \text{ for all } U \in N_\epsilon. \quad (162)$$

Let $U' \in \mathcal{S}_F \subset \mathbb{R}^{d \times r}$ be the matrix “parallel” to \mathcal{E}_1 , that is $\|\mathcal{E}_2(U^*Q)\|_F = \max_{\|\tilde{U}\|_F=1} \langle \mathcal{E}_1(U^*Q), \tilde{U} \rangle = \langle \mathcal{E}_2(U^*Q), U' \rangle$. Then there exists some $U \in N_\epsilon$ such that $\|U' - U\|_F \leq \epsilon$.

$$\|\mathcal{E}_2(U^*Q)\|_F = \langle \mathcal{E}_2(U^*Q), U' \rangle = \langle \mathcal{E}_2(U^*Q), U \rangle + \langle \mathcal{E}_2(U^*Q), U' - U \rangle \quad (163)$$

$$\leq \langle \mathcal{E}_1(U), U \rangle + \|\mathcal{E}_2(U^*Q)\|_F \|U' - U\|_F \quad (164)$$

$$\leq \langle \mathcal{E}_1(U), U \rangle + \epsilon \|\mathcal{E}_2(U^*Q)\|_F \quad (165)$$

Re-arranging and setting $\epsilon = 1/2$, and $c \leftarrow 2c$, we get

$$\|\mathcal{E}_2(U^*Q)\|_F \leq \Delta_{\frac{1}{2}} = c(\min(\|H\|_F \frac{\|V\|_{\infty,2}}{\sqrt{\lambda_r(W)}}, \|H\|_{\infty,2}) \sqrt{\frac{dr \log(5/\delta)}{m}} + \|H\|_{\infty,2} \frac{\|V\|_{\infty,2}}{\sqrt{\lambda_r(W)}} \frac{dr \log(5/\delta)}{m}) \quad (166)$$

Finally setting $\delta \leftarrow \delta/3$ get us the desired result. \square

Lemma A.9. *With a probability of at least $1 - \delta/3$*

$$\left\| \sum_{i \in [t]} \mathcal{W}^{-\frac{1}{2}}(z^{(i)}(v^{(i)})^\top) \right\|_F \leq O\left(\sigma \sqrt{\frac{dr}{m} \log\left(\frac{t}{\delta}\right) \log\left(\frac{r}{\delta}\right)}\right) \quad (167)$$

Proof of Lemma A.9. Notice that $z^{(i)}$ (defined in Appendix A) is a Gaussian random vector of the following form

$$z^{(i)} = \frac{1}{m} \sum_{j \in [m]} \varepsilon_j^{(i)} x_j^{(i)} = \frac{1}{m} \|\varepsilon^{(i)}\| g^{(i)}, g^{(i)} \sim \mathcal{N}(0, \mathbf{I}_{d \times d}) \quad (168)$$

Using Hanson-Wright inequality (Lemma D.5, by setting $m \leftarrow 1$, $x_1 \leftarrow \varepsilon^{(i)}$, and $A_1 \leftarrow \mathbf{I}_{m \times m}$) and taking union bound over all tasks, we get that, with probability of at least $1 - \frac{\delta}{2}$

$$\|\varepsilon^{(i)}\|^2 \leq \sigma^2 m \left(1 + c \sqrt{\frac{\log(\frac{2t}{\delta})}{m}} + c \frac{\log(\frac{2t}{\delta})}{m}\right) \leq 2c\sigma^2 m \log\left(\frac{2t}{\delta}\right), \text{ for all } i \in [t] \quad (169)$$

where used the fact that $m \geq 1$ and $\log\left(\frac{2t}{\delta}\right) \geq 1$.

Let $\hat{v}^{(i)} = W^{-\frac{1}{2}} v^{(i)}$, then

$$\sum_{i \in [t]} \|\hat{v}^{(i)}\|^2 = \sum_{i \in [t]} \text{tr}((v^{(i)})^\top W^{-1} v^{(i)}) = \sum_{i \in [t]} \text{tr}(W^{-1} v^{(i)} (v^{(i)})^\top) = r \quad (170)$$

Notice that $\sum_{i \in [t]} \frac{1}{m} \|\varepsilon^{(i)}\| g^{(i)} \hat{v}_j^{(i)}$ is a Gaussian random vector of the following form

$$\sum_{i \in [t]} \frac{1}{m} \|\varepsilon^{(i)}\| g^{(i)} \hat{v}_j^{(i)} = \frac{1}{m} \sqrt{\sum_{i \in [t]} \|\varepsilon^{(i)}\|^2 (\hat{v}_j^{(i)})^2} \hat{g}_j, \hat{g}_j \sim \mathcal{N}(0, \mathbf{I}_{d \times d}) \quad (171)$$

Using Hanson-Wright inequality (Lemma D.5, by setting $m \leftarrow 1$, $x_1 \leftarrow \hat{g}_j$, and $A_1 \leftarrow \mathbf{I}_{d \times d}$) and taking union bound over all $j \in [r]$, we get that, with probability of at least $1 - \frac{\delta}{2}$

$$\|\hat{g}_j\|^2 \leq d \left(1 + c \sqrt{\frac{\log(\frac{2r}{\delta})}{d}} + c \frac{\log(\frac{2r}{\delta})}{d}\right) \leq 2cd \log\left(\frac{2r}{\delta}\right), \text{ for all } j \in [r] \quad (172)$$

where used the fact that $d \geq 1$ and $\log\left(\frac{2r}{\delta}\right) \geq 1$.

Combining the above results and using union bound, we get that, with a probability of at least $1 - \delta$,

$$\left\| \sum_{i \in [t]} \mathcal{W}^{-\frac{1}{2}} (z^{(i)} (v^{(i)})^\top) \right\|_F^2 = \left\| \sum_{i \in [t]} z^{(i)} (v^{(i)})^\top W^{-\frac{1}{2}} \right\|_F^2 = \left\| \sum_{i \in [t]} \frac{1}{m} \|\varepsilon^{(i)}\| g^{(i)} (\widehat{v}^{(i)})^\top \right\|_F^2 \quad (173)$$

$$= \sum_{j \in [r]} \left\| \sum_{i \in [t]} \frac{1}{m} \|\varepsilon^{(i)}\| g_j^{(i)} \widehat{v}_j^{(i)} \right\|^2 \quad (174)$$

$$\leq \sum_{j \in [r]} \sum_{i \in [t]} \frac{\|\varepsilon^{(i)}\|^2}{m^2} (\widehat{v}_j^{(i)})^2 \|\widehat{g}_j\|^2 \quad (175)$$

$$\leq \sum_{j \in [r]} \sum_{i \in [t]} O\left(\frac{m\sigma^2}{m^2} \log\left(\frac{t}{\delta}\right)\right) (\widehat{v}_j^{(i)})^2 O\left(d \log\left(\frac{r}{\delta}\right)\right) \quad (176)$$

$$\leq O\left(\frac{d\sigma^2}{m} \log\left(\frac{t}{\delta}\right) \log\left(\frac{r}{\delta}\right)\right) \sum_{i \in [t]} \|\widehat{v}^{(i)}\|^2 \quad (177)$$

$$\leq O\left(\frac{\sigma^2 dr}{m} \log\left(\frac{t}{\delta}\right) \log\left(\frac{r}{\delta}\right)\right). \quad (178)$$

Finally, we get the desired result by setting $\delta \leftarrow \delta/3$. \square

A.4 Analysis of QR decomposition

Proof of Lemma A.3.

$$\sigma_{\min}(R) \geq \min_{\|z\|=1} \|Rz\| = \min_{\|z\|=1} \|U^+ Rz\| = \min_{\|z\|=1} \|\widehat{U}z\| \quad (179)$$

$$\geq \min_{\|z\|=1} \|(U^*Q + \mathcal{W}^\dagger \mathcal{H}(U^*Q) + F)z\| \quad (180)$$

$$\geq \min_{\|z\|=1} \sqrt{z^\top Q^\top Q z} - \|\mathcal{W}^\dagger \mathcal{H}(U^*Q)\| - \|F\| \quad (181)$$

$$\geq \min_{\|z\|=1} \sigma_{\min}(Q) - \|\mathcal{W}^\dagger \mathcal{H}(U^*Q)\| - \|F\| \quad (182)$$

$$\geq \frac{1}{2} - \frac{1}{8} - \frac{1}{8} \geq \frac{1}{4} \quad (183)$$

There fore R is invertible and $\|R^{-1}\| = (\sigma_{\min}(R))^{-1} \leq 4$ \square

A.5 Analysis of shuffling and partitioning

Proof of Lemma A.4. We will assume that the set of tasks $[t]$ is shuffled. We will prove that incoherence holds for the all subset $\mathcal{T}_k = [1 + \frac{t(k-1)}{K}, \frac{tk}{K}]$ of size t/K . Shuffling and K -way partitioning to get \mathcal{T}_k is equivalent to uniformly sampling without replacement t/K elements from $[t]$. We prove that incoherence holds for the first subset \mathcal{T}_1 , then this is equivalent to proving that incoherence holds for the k -th partition \mathcal{T}_k by symmetry. Let the tasks sampled for \mathcal{T}_1 without replacement be $\{i_l\}_{l=1}^{t/k}$, where i_l is the l -th sample.

Let $\mathcal{S}_F = \{z \in \mathbb{R}^r \mid \|z\| = 1\}$ be the set of all real vectors of dimensions r with unit Euclidean norm. For $\epsilon \leq 1$, there exists an ϵ -net, $N_\epsilon \subset \mathcal{S}_F$, of size $(1 + 2/\epsilon)^r$ with respect to the Euclidean norm [Ver10, Lemma 5.2]. That is for any $z' \in \mathcal{S}_F$, there exists some $z \in N_\epsilon$ such that $\|z' - z\| \leq \epsilon$.

Consider a $z \in N_\epsilon$, such that $\|z\| = 1$. Now we will prove with high-probability that $z^\top (\sum_{l=1}^{t/K} v^{*(i_l)} (v^{*(i_l)})^\top) z$ is approximately equal to $z^\top \mathbb{E}[\sum_{l=1}^{t/K} v^{*(i_l)} (v^{*(i_l)})^\top] z$. Now consider the martingale X_l , such that $X_0 = 0$ and $X_l = X_{l-1} + z^\top (v^{*(i_l)} (v^{*(i_l)})^\top - \mathbb{E}[v^{*(i_l)} (v^{*(i_l)})^\top \mid X_0, \dots, X_{l-1}]) z$, for all $l \in [t/K]$. Clearly this is a martingale as $\mathbb{E}[X_l \mid X_0, \dots, X_{l-1}] = 0$, for all $l \in [t/K]$. The maximum difference two consecutive steps is $\max_l |X_l - X_{l-1}| \leq 2\|v^{*(i_l)}\|^2 \leq 2\|V^*\|_{\infty,2}^2$. Therefore by Azuma-Hoeffding martingale inequality,

$$\left| \sum_{l=1}^{t/K} z^\top v^{*(i_l)} (v^{*(i_l)})^\top z - z^\top \mathbb{E}\left[\sum_{l=1}^{t/K} v^{*(i_l)} (v^{*(i_l)})^\top\right] z \right| = |X_{t/K}| \leq \sqrt{\frac{2t}{K} \|V\|_{\infty,2}^4 \log\left(\frac{2|N_\epsilon|}{\delta}\right)} \quad (184)$$

with a probability of at least $1 - \delta/|N_\epsilon|$.

For brevity, let $E = \sum_{l=1}^{t/K} v^{*(i_l)}(v^{*(i_l)})^\top - \mathbb{E}[\sum_{l=1}^{t/K} v^{*(i_l)}(v^{*(i_l)})^\top]$. Notice that E is a real symmetric matrix, therefore it has an eigen decomposition. Then, let $v' \in \mathcal{S} \subset \mathbb{R}^r$ be the largest ‘‘eigenvector’’ of E , such that $(v')^\top E v' = \|E\| = \max_{\|\tilde{v}\|=1} \tilde{v}^\top E \tilde{v} = \max_{\|\tilde{v}\|_F=1} \tilde{v}^\top E \tilde{v}'$. Then there exists some $v \in N_\epsilon$ such that $\|v' - v\| \leq \epsilon$.

$$\|E\|_F = (v')^\top E v = v^\top E v + (v' - v)^\top E v + (v')^\top E (v' - v) \quad (185)$$

$$\leq v^\top E v + \|v' - v\| \|E\| \|v\| + \|v'\| \|E\| \|v' - v\| \quad (186)$$

$$\leq v^\top E v + 2\epsilon \|E\| \quad (187)$$

Re-arranging and setting $\epsilon = 1/4$, and $c \leftarrow 2c$, we get

$$\left\| \sum_{l=1}^{t/K} v^{*(i_l)}(v^{*(i_l)})^\top - \mathbb{E}\left[\sum_{l=1}^{t/K} v^{*(i_l)}(v^{*(i_l)})^\top\right] \right\| = \|E\| \leq \sqrt{\frac{2tr}{K} \|V\|_{\infty,2}^4 \log\left(\frac{18}{\delta}\right)} \leq \frac{1}{2} \lambda_r(\mathbb{E}[\sum_{l=1}^{t/K} v^{*(i_l)}(v^{*(i_l)})^\top]). \quad (188)$$

with probability at least $1 - \delta/k$, where the last inequality used the fact that $t \geq \Omega(\mu^2 r^3 K \log(1/\delta))$. Additionally note that $\mathbb{E}[\sum_{l=1}^{t/k} v^{*(i_l)}(v^{*(i_l)})^\top] = \frac{1}{K} \sum_{i=1}^t v^{*(i)}(v^{*(i)})^\top = \frac{1}{K} (V^*)^\top V^*$, Therefore

$$\lambda_{r'}\left(\sum_{i \in \mathcal{T}_k} v^{*(i)}(v^{*(i)})^\top\right) = \frac{1}{K} \Theta(\lambda_{r'}((V^*)^\top V^*)) \text{ for all } r' \in [r] \quad (189)$$

where $\lambda_i(\cdot)$ is the r' -th largest eigenvalue matrix operator. □

B Analysis of MLLAMS (Algorithm 2) with subset selection

Initialized at U , the k -th step of alternating minimization-based MLLAMS (Algorithm 2) is:

$$\mathcal{T}_k = \{i \in [1 + (k-1)t/K, tk/K] \mid \sigma_{\min}(U^\top S^{(i)} U) \geq 1/2 \text{ and } \sigma_{\max}(U^\top S^{(i)} U) \leq 2\} \quad (190)$$

$$v^{(i)} \leftarrow (U^\top S^{(i)} U)^\dagger ((U^\top S^{(i)} U^*) v^{*(i)} + U^\top z^{(i)}), \quad \text{for } i \in \mathcal{T}_k \quad (191)$$

$$\hat{U} \leftarrow \mathcal{A}^\dagger \left(\sum_{i \in \mathcal{T}} S^{(i)} U^* v^{*(i)} (v^{(i)})^\top + z^{(i)} (v^{(i)})^\top \right), \quad (192)$$

$$U^+ \leftarrow \text{QR}(\hat{U}), \quad (193)$$

where U^+ is the next iterate, $S_1^{(i)} = \frac{2}{m} \sum_{j \in [1, m/2]} x_j^{(i)} (x_j^{(i)})^\top$, $S_2^{(i)} = \frac{2}{m} \sum_{j \in [1+m/2, m]} x_j^{(i)} (x_j^{(i)})^\top$, $z^{(i)} \triangleq (1/m) \sum_{j \in [m]} \varepsilon_j^{(i)} x_j^{(i)}$ and $\mathcal{A} : \mathbb{R}^{d \times r} \rightarrow \mathbb{R}^{d \times r}$ is a self-adjoint linear operator such that $\mathcal{A}(U) = \sum_{i \in \mathcal{T}} S^{(i)} U v^{*(i)} (v^{(i)})^\top$.

Remark: Note the subset \mathcal{T}_k , which we analyze, is slightly different from that of Algorithm 2. This is done to save some polylog factors in the final error-bound and sample complexity. However, the analysis will remain almost the same even if eliminate the subset selection criterion, $\sigma_{\max}(U^\top S^{(i)} U) \geq 2$ for all $i \in \mathcal{T}_k$.

Theorem 6. *Let there be t linear regression tasks, each with m samples satisfying Assumptions 1 and 2, and $K = \lceil \log_2(\frac{\lambda_r^* \lambda_1^* mt}{\mu dr^2}) \rceil$, $\|(\mathbf{I} - U^*(U^*)^\top) U_{\text{init}}\|_F \leq \min\left(\frac{3}{4}, O\left(\sqrt{\frac{\lambda_r^*}{\lambda_1^*} \frac{1}{\log(t/K)}}\right)\right)$, $m \geq \Omega\left(\left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 r^2 \log\left(\frac{t}{\delta}\right) + r^2 \log\left(\frac{K}{\delta}\right) + \log(\mu r)\right)$, $t \geq \Omega(\mu^2 r^3 K \log\left(\frac{K}{\delta}\right))$ and $mt \geq \Omega\left(\mu dr^2 K \frac{\lambda_r^*}{\lambda_1^*} \left(\log\left(\frac{t}{\delta}\right) + \left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 \log^2\left(\frac{t}{\delta}\right) \log\left(\frac{rK}{\delta}\right)\right)\right)$. Then, for any $0 < \delta < 1$, after K iterations, MLLAMS (Algorithm 2) returns an orthonormal matrix $U \in \mathbb{R}^{d \times r}$, such that with a probability of at least $1 - \delta$*

$$\frac{1}{\sqrt{r}} \|(\mathbf{I} - U^*(U^*)^\top) U\|_F \leq O\left(\frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{\mu dr K \log\left(\frac{t}{\delta}\right) \log\left(\frac{rK}{\delta}\right)}{mt}}\right) \quad (194)$$

and the algorithm uses an additional memory of size $O(d^2 r^2)$.

A proof is in Section B.1.

B.1 Analysis

First, in the following lemma, we prove that the task subset \mathcal{T}_k has similar properties as the full task partition $[1 + t(k-1)/K, tk/K]$.

Lemma B.1 (Subset selection). *If $m \geq \Omega(r + \log(\mu r))$ and $t \geq \Omega(\mu^2 r^2 K \log(\frac{1}{\delta}))$, then with a probability of at least $1 - \delta/3$,*

$$|\mathcal{T}_k| = \Theta\left(\frac{t}{K}\right), \quad \text{and} \quad \|V^*\|_{\infty,2}^2 \leq O\left(\frac{\mu r}{|\mathcal{T}_k|} \lambda_r \left(\sum_{i \in \mathcal{T}} v^{*(i)} (v^{*(i)})^\top\right)\right) \quad (195)$$

$$\lambda_r \left(\sum_{i \in \mathcal{T}} v^{*(i)} (v^{*(i)})^\top\right) = \Theta\left(\lambda_r \left(\sum_{i \in \mathcal{P}_k} v^{*(i)} (v^{*(i)})^\top\right)\right), \quad \text{and} \quad \lambda_1 \left(\sum_{i \in \mathcal{T}} v^{*(i)} (v^{*(i)})^\top\right) = \Theta\left(\lambda_1 \left(\sum_{i \in \mathcal{P}_k} v^{*(i)} (v^{*(i)})^\top\right)\right) \quad (196)$$

where $\mathcal{P}_k = [1 + t(k-1)/K, tk/K]$ is the k -th K -way partition of $[t]$ after shuffling.

A proof is in Section B.2. Therefore, assuming that the above high-probability event holds, in the rest of the proof we can consider that \mathcal{T}_k is equivalent to \mathcal{P}_k .

In the rest of the proof, when compared to the proof of Theorem 5, only the following Lemma (corresponding to Lemma A.1) analyzing the V -update changes in its necessary condition.

Lemma B.2. *If $\|(\mathbf{I} - U^*(U^*)^\top)U\|_F \leq \min\left(\frac{3}{4}, O\left(\sqrt{\frac{\lambda_r^*}{\lambda_1^*} \frac{1}{\log(t/K)}}\right)\right)$ and $m \geq \Omega\left(\left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 r^2 \log\left(\frac{t}{K\delta}\right) + r \log\left(\frac{1}{\delta}\right)\right)$, then with a probability of at least $1 - \delta/3$,*

$$\|v^{(i)}\| \leq O\left(\mu \lambda_r\right), \quad \text{and} \quad \lambda_r^* \leq 2\lambda_r \quad (197)$$

and

$$\sqrt{\frac{rK}{t}} \frac{\|H\|_F}{\sqrt{\lambda_r}} \leq O\left(\sqrt{\frac{\log(\frac{t}{K\delta})}{\log(\frac{1}{\delta})}} \sqrt{\frac{\lambda_1^*}{\lambda_r^*}} \|(\mathbf{I} - U^*(U^*)^\top)U\|_F + \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 \log(\frac{t}{K\delta})}{m}}\right) \quad (198)$$

$$\sqrt{\frac{rK}{t}} \frac{\|H\|_{\infty,2}}{\sqrt{\lambda_r}} \leq O\left(\sqrt{\frac{\log(\frac{t}{K\delta})}{\log(\frac{1}{\delta})}} \|(\mathbf{I} - U^*(U^*)^\top)U\| \sqrt{\frac{\mu r K}{t}} + \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{r^2 K \log(\frac{t}{K\delta})}{mt}}\right) \quad (199)$$

A proof is in Section B.3.1. We omit the rest of the proof, as it is same as that of Theorem 5.

B.2 Analysis of task subset selection

Proof of Lemma B.1 (Subset selection). Let $\mathcal{P}_k = [1 + (k-1)t/K, tk/K]$ and

$$\mathcal{T}_k = \left\{ i \in [1 + (k-1)t/K, tk/K] \mid \sigma_{\min}(U^\top S^{(i)} U) \geq 1/2 \text{ and } \sigma_{\max}(U^\top S^{(i)} U) \leq 2 \right\}. \quad (200)$$

For all $i \in \mathcal{P}_k$, $X_i = \mathbb{I}(\sigma_{\min}(U^\top S^{(i)} U) \geq 1/2 \text{ and } \sigma_{\max}(U^\top S^{(i)} U) \leq 2)$ be the indicator variable denoting whether index i was select into the subset $\widehat{\mathcal{T}}$.

By Lemma D.7 (by setting $a_j \leftarrow 1$, $x_j \leftarrow U^\top x_j^{(i)}$ for all $j \in [m]$, and $\delta \leftarrow 1/4\mu r$) X_i are i.i.d. Bernoulli random variables with mean $p \geq 1 - \frac{1}{4\mu r}$, if $c \max\left(\sqrt{\frac{r \log(9) + \log(4\mu r)}{m}}, \frac{r \log(9) + \log(4\mu r)}{m}\right) \leq 1/2$, which is satisfied by $m \geq \Omega(r + \log(\mu r))$, for all $i \in \mathcal{P}_k$.

By Hoeffding inequality for Bernoulli random variables, with a probability of at least $1 - \delta/3$

$$|\mathcal{T}_k| - pt/K = \left| \sum_{i \in \mathcal{P}_k} X_i - \left(1 - \frac{1}{4\mu r}\right) \frac{t}{K} \right| \leq \frac{t}{K} \sqrt{\frac{K \log(\frac{3}{\delta})}{2t}} \leq \frac{t}{K} O\left(\frac{1}{4\mu r}\right) \quad (201)$$

where we used the fact that $t \geq \Omega(8K\mu^2 r^2 \log(\frac{3}{\delta}))$. Therefore

$$\frac{t}{K} - |\mathcal{T}_k| \leq \frac{t}{K} O\left(\frac{1}{2\mu r}\right), \quad \text{and} \quad |\mathcal{T}_k| \leq \Theta\left(\frac{t}{K}\right) \quad (202)$$

where we used the fact that $\mu \geq 1$ and $r \geq 1$.

$$\frac{r}{t} \left| z^\top \left(\sum_{i \in \mathcal{T}_k} v^{*(i)} (v^{*(i)})^\top \right) z - z^\top \left(\sum_{i \in \mathcal{P}_k} v^{*(i)} (v^{*(i)})^\top \right) z \right| \leq \frac{r}{t} (t - \hat{t}) \|V^*\|_{\infty,2}^2 \leq \frac{r}{t} O\left(\frac{t}{2\mu r}\right) \cdot \|V^*\|_{\infty,2}^2 \leq \frac{\lambda_r}{2}, \quad (203)$$

for all $z \in \mathbb{R}^r$, where $\lambda_r = \lambda_r(\sum_{i \in \mathcal{P}_k} v^{*(i)} (v^{*(i)})^\top)$. Therefore

$$\lambda_r \left(\sum_{i \in \mathcal{T}} v^{*(i)} (v^{*(i)})^\top \right) = \Theta \left(\lambda_r \left(\sum_{i \in \mathcal{P}_k} v^{*(i)} (v^{*(i)})^\top \right) \right), \text{ and } \lambda_1 \left(\sum_{i \in \mathcal{T}} v^{*(i)} (v^{*(i)})^\top \right) = \Theta \left(\lambda_1 \left(\sum_{i \in \mathcal{P}_k} v^{*(i)} (v^{*(i)})^\top \right) \right) \quad (204)$$

Using approximate incoherence of the partition \mathcal{P}_k (Lemma A.4) we get

$$\|V^*\|_{\infty,2}^2 \leq O\left(\frac{\mu r K}{t}\right) \lambda_r \left(\sum_{i \in \mathcal{P}_k} v^{*(i)} (v^{*(i)})^\top \right) = O\left(\frac{\mu r K}{t}\right) \min_{\|z\|=1} z^\top \left(\sum_{i \in \mathcal{P}_k} v^{*(i)} (v^{*(i)})^\top \right) z \quad (205)$$

$$\leq O\left(\frac{\mu r K}{t}\right) \min_{\|z\|=1} z^\top \left(\sum_{i \in \mathcal{T}_k} v^{*(i)} (v^{*(i)})^\top \right) z + O\left(\frac{\mu r K}{t}\right) \left(\frac{t}{K} - |\mathcal{T}_k|\right) \|V^*\|_{\infty,2}^2 \quad (206)$$

$$\leq O\left(\frac{\mu r K}{t}\right) \lambda_r \left(\sum_{i \in \mathcal{T}_k} v^{*(i)} (v^{*(i)})^\top \right) + \frac{1}{2} \|V^*\|_{\infty,2}^2 \quad (207)$$

$$(208)$$

This implies that approximate incoherence holds for \mathcal{T}_k , $\|V^*\|_{\infty,2}^2 \leq O\left(\frac{\mu r K}{t}\right) \lambda_r(\sum_{i \in \mathcal{T}_k} v^{*(i)} (v^{*(i)})^\top) \leq O\left(\frac{\mu r}{|\mathcal{T}_k|} \lambda_r(\sum_{i \in \mathcal{T}_k} v^{*(i)} (v^{*(i)})^\top)\right)$. \square

B.3 Analysis of update on V

B.3.1 Proof of Lemma B.2

Proof of Lemma B.2. The proof is similar to that of Lemma A.1, but instead of using Lemma A.5 to bound some linear operators, we use the definition of selected task subset \mathcal{T}_k and Lemma B.3 to get that $\|(U^\top S^{(i)} U)^\dagger\| \leq 2$ for all $i \in \mathcal{T}_k$ and with a probability of at least $1 - \delta$,

$$\left. \begin{aligned} \|U^\top S^{(i)} U_\perp U_\perp^\top U^* v^{*(i)}\| &\leq \alpha \|U_\perp^\top U^* v^{*(i)}\|, \text{ and } \\ \|U^\top z^{(i)}\| &\leq \sigma \alpha, \end{aligned} \right\} \text{ for all } i \in \mathcal{T}_k \quad (209)$$

where $\alpha = c \sqrt{\frac{r \log(10t/\delta)}{m}}$. We omit the rest of the proof, as it is same as that of Lemma A.1. \square

Here we bound the linear operators in the $v^{(i)}$ update.

Lemma B.3. *With a probability of at least $1 - \delta$, the following are true for all $i \in [t]$*

$$\|U^\top S^{(i)} U_\perp (U_\perp)^\top U^* v^{*(i)}\| \leq \sqrt{\frac{2cr \log(10t/\delta)}{m}} \|U_\perp U^* v^{*(i)}\|, \text{ and} \quad (210)$$

$$\|U^\top z^{(i)}\| \leq \sigma \sqrt{\frac{2cr \log(10t/\delta)}{m}} \quad (211)$$

Proof. Let $i \in [t]$. Let $b = (U_\perp)^\top U^* v^{*(i)} \in \mathbb{R}^r$

Let $\mathcal{S} = \{v \in \mathbb{R}^r \mid \|v\| = 1\}$ be the set of all real vectors of dimension r with unit Euclidean norm. For $\epsilon \leq 1$, there exists an ϵ -net, $N_\epsilon \subset \mathcal{S}$, of size $(1 + 2/\epsilon)^r$ with respect to the Euclidean norm [Ver10, Lemma 5.2]. That is for any $v' \in \mathcal{S}$, there exists some $v \in N_\epsilon$ such that $\|v' - v\|_F \leq \epsilon$.

Consider a $v \in N_\epsilon$, such that $\|v\|_F = 1$. Now we will prove with high-probability that $\langle (U^\top S^{(i)} U_\perp) v, b \rangle$ is small. Consider the the following quadratic form

$$v^\top (U^\top S^{(i)} U_\perp) b = \frac{1}{m} \sum_{j \in [m]} v^\top (U^\top x_j^{(i)} (x_j^{(i)})^\top U_\perp) b \stackrel{d}{=} \|b\| \frac{1}{m} \sum_{j \in [m]} \tilde{x}_j g_j \quad (212)$$

where $g_j \sim \mathcal{N}(0, 1)$ are i.i.d. standard Gaussian random variables and $\tilde{x}_j = v^\top U^\top x_j^{(i)} \in \mathbb{R}^d$. This follows from the fact that sets of columns of U and U_\perp forms an orthonormal basis.

Note that g_j and \tilde{x}_j are independent, as U and U_\perp are orthogonal and $U^\top S^{(i)} U$, does not depend on $U_\perp x_j^{(i)}$. We will use the properties of Gaussian random variables to prove that $\|\frac{1}{m} \sum_{j \in [m]} \tilde{x}_j g_j\|$ concentrates around zero. Note that

$$\frac{1}{m} \sum_{j \in [m]} \tilde{x}_j g_j \stackrel{d}{=} \frac{1}{m} \|\tilde{x}\| g, \quad \text{where } g \sim \mathcal{N}(0, 1) \quad (213)$$

Then with probability at least $1 - \delta/2t/|N_\epsilon|$, $|g|^2 \leq c \log(2t|N_\epsilon|/\delta)$. Additionally, by definition of \mathcal{T}_k we have

$$\frac{1}{m} \|\tilde{x}\|^2 = \frac{1}{m} \sum_{j \in [m]} \tilde{x}_j^2 = v^\top U^\top \left(\frac{1}{m} \sum_{j \in [m]} x_j^{(i)} (x_j^{(i)})^\top \right) U v \leq \sigma_{\max}(U^\top S^{(i)} U) \leq 2 \quad (214)$$

Therefore

$$v^\top (U^\top S^{(i)} U_\perp) b \leq \frac{1}{\sqrt{m}} \|b\| \sqrt{2c \log(2t|N_\epsilon|/\delta)} \quad (215)$$

For brevity, let $e = (U^\top S^{(i)} U_\perp) b$. Let $v' \in \mathcal{S} \subset \mathbb{R}^r$ be the unit vector parallel to e , such that $(v')^\top e = \|e\| = \max_{\|\tilde{v}\|=1} \tilde{v}^\top e$. Then there exists some $v \in N_\epsilon$ such that $\|v' - v\| \leq \epsilon$.

$$\|e\| = (v')^\top e = v^\top e + (v' - v)^\top e \leq v^\top e + \|v' - v\| \|e\| \leq v^\top e + \epsilon \|e\| \quad (216)$$

Re-arranging and setting $\epsilon = 1/2$, and $c \leftarrow 2c$, we get

$$\|(U^\top S^{(i)} U_\perp) b\| \leq \|b\| \sqrt{\frac{2cr \log(10t/\delta)}{m}}, \quad \text{with a probability of at least } 1 - \delta/2t \quad (217)$$

Using similar arguments we can also prove that with a probability of at least $1 - \delta$

$$\|U^\top z^{(i)}\| = \left\| \frac{1}{m} U^\top x_j^{(i)} \varepsilon_j^{(i)} \right\| \leq \sigma \sqrt{\frac{2cr \log(10t/\delta)}{m}}, \quad \text{with a probability of at least } 1 - \delta/2t \quad (218)$$

Finally taking the union bound over the two bounds over all the tasks in \mathcal{T} gets us the desired result. \square

C Corollaries of known results

Theorem 7 (Theorem 3, Tripuraneni et al. 2020). *Let there be t linear regression tasks, each with m samples satisfying Assumptions 1 and 2, and*

$$mt \geq \tilde{\Omega} \left(\frac{\lambda_1^*}{\lambda_r^*} \mu dr + \left(\frac{\sigma}{\sqrt{\lambda_r^*}} \right)^4 dr^2 \right) \quad (219)$$

then with a high probability of at least $1 - O((mt)^{-100})$, Method-of-Moments [TJJ20, Algorithm 1] outputs an orthonormal matrix $U \in \mathbb{R}^{d \times r}$ such that

$$\|(\mathbf{I} - U^*(U^*)^\top) U\|_2 \leq \tilde{O} \left(\sqrt{\frac{\lambda_1^* \mu dr}{\lambda_r^* mt}} + \left(\frac{\sigma}{\sqrt{\lambda_r^*}} \right)^2 \sqrt{\frac{dr^2}{mt}} \right) \quad (220)$$

and

$$\|(\mathbf{I} - U^*(U^*)^\top) U\|_F \leq \tilde{O} \left(\sqrt{\frac{\lambda_1^* \mu dr^2}{\lambda_r^* mt}} + \left(\frac{\sigma}{\sqrt{\lambda_r^*}} \right)^2 \sqrt{\frac{dr^3}{mt}} \right). \quad (221)$$

Proof. From the details of the proof of Theorem 3 in [TJJ20] we can derive that, with a high probability of at least $1 - O((mt)^{-100})$,

$$\|(\mathbf{I} - U^*(U^*)^\top)U\|_2 \quad (222)$$

$$\leq \tilde{O}\left(\sqrt{\frac{dr^2 \text{tr}(W^*) \|V^*\|_{\infty,2}^2}{\lambda_r^{*2} mt^2}} + \frac{dr \|V^*\|_{\infty,2}^2}{\lambda_r^* mt} + \sigma \left(\sqrt{\frac{dr^2 \text{tr}(W^*)}{\lambda_r^{*2} mt^2}} + \frac{dr \|V^*\|_{\infty,2}}{\lambda_r^* mt}\right) + \sigma^2 \left(\sqrt{\frac{dr^2}{\lambda_r^{*2} mt}} + \frac{dr}{\lambda_r^* mt}\right)\right) \quad (223)$$

$$\leq \tilde{O}\left(\sqrt{\frac{\lambda_1^* \mu dr}{\lambda_r^* mt}} + \frac{\mu dr}{mt} + \frac{\sigma}{\sqrt{\lambda_r^*}} \left(\sqrt{\frac{\lambda_1^* dr}{\lambda_r^* mt}} + \frac{\sqrt{\mu} dr}{mt}\right) + \left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 \left(\sqrt{\frac{dr^2}{mt}} + \frac{dr}{mt}\right)\right) \quad (224)$$

$$\leq \tilde{O}\left(\sqrt{\frac{\lambda_1^* \mu dr}{\lambda_r^* mt}} + \frac{\sigma}{\sqrt{\lambda_r^*}} \sqrt{\frac{\lambda_1^* dr}{\lambda_r^* mt}} + \left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 \sqrt{\frac{dr^2}{mt}}\right) \quad (225)$$

$$\leq \tilde{O}\left(\sqrt{\frac{\lambda_1^* \mu dr}{\lambda_r^* mt}} + \left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^2 \sqrt{\frac{dr^2}{mt}}\right) \quad (226)$$

where $\|V\|_{\infty,2} = \max_{i \in [t]} \|v^{(i)}\|$, and the second-last inequality uses the fact that $mt \geq \tilde{\Omega}(\mu dr)$ and last inequality uses the fact that $\frac{\lambda_1^*}{\lambda_r^*} \leq \mu r$. Additionally we require that

$$mt \geq \tilde{\Omega}\left(\frac{\lambda_1^*}{\lambda_r^*} \mu dr + \left(\frac{\sigma}{\sqrt{\lambda_r^*}}\right)^4 dr^2\right) \quad (227)$$

□

Theorem 8. [TJJ20, Theorem 5] Let $r \leq d/2$ and $mt \geq r(d-r)$, then for all V^* , w.p. $\geq 1/2$

$$\inf_{\hat{U}} \sup_{U \in G_{r,d}} \frac{\|(\mathbf{I} - U^*(U^*)^\top)\hat{U}\|_F}{\sqrt{r}} \geq \Omega\left(\left(\frac{\lambda_r^*}{\lambda_1^*} \frac{\sigma}{\sqrt{\lambda_r^*}}\right) \sqrt{\frac{dr}{mt}}\right),$$

where $G_{r,d}$ is the Grassmannian manifold of r -dimensional subspaces in \mathbb{R}^d , the infimum for \hat{U} is taken over the set of all measurable functions that takes mt samples in total from the model in Section 2 satisfying Assumption 1 and 2.

Proof. The proof is very similar to that of Theorem 5 of Tripuraneni, Jin, and Jordan [TJJ20]. The main difference is that instead of lower bounding error in spectral norm we have to bound the distance in the Frobenius norm. However, the rest of the details are almost the same, hence we omit a full proof. □

D Technical Lemmas

This section contains some technical lemmas used in this paper.

Lemma D.1. For a real matrix $A \in \mathbb{R}^{m \times n}$ and a real symmetric positive semi-definite (PSD) matrix $B \in \mathbb{R}^{n \times n}$, the following holds true: $\sigma_{\min}^2(A) \lambda_{\min}(B) \leq \lambda_{\min}(ABA^\top)$, where $\sigma_{\min}(\cdot)$ and $\lambda_{\min}(\cdot)$ represents the minimum singular value and minimum eigenvalue operators respectively.

Proof. The proof directly follows from the definitions of σ_{\min} and λ_{\min} . Since B is a PSD matrix, therefore ABA^\top is also PSD, i.e. $\lambda_{\min}(ABA^\top) \geq 0$. This is because since B is PSD, it has a PSD matrix square root $B^{1/2}$ such that $B = (B^{1/2})^\top B^{1/2}$ and $B^{1/2}$ is PSD. Then

$$z^\top ABA^\top z = z^\top A(B^{1/2})^\top B^{1/2} A^\top z = \|B^{1/2} A^\top z\|^2 \geq 0 \quad (228)$$

First assume that $\sigma_{\min}(A) > 0$, then

$$\lambda_{\min}(ABA^\top) = \min_{\|z\|=1} z^\top ABA^\top z \quad (229)$$

$$= \sigma_{\min}^2(A) \min_{\|z\|=1} \left(\frac{A^\top z}{\sigma_{\min}(A)} \right)^\top B \left(\frac{A^\top z}{\sigma_{\min}(A)} \right) \quad (230)$$

$$\geq \sigma_{\min}^2(A) \min_{1 \leq \|z\| \leq \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}} z^\top B z \quad (231)$$

$$\geq \sigma_{\min}^2(A) \min_{\|z\|=1} z^\top B z \quad (232)$$

$$= \sigma_{\min}^2(A) \lambda_{\min}(B) \quad (233)$$

The second last inequality above follows from the fact that B is a PSD matrix, i.e. $\min_{\|z\|=1} z^\top B z = \lambda_{\min}(B) \geq 0$. Secondly if $\sigma_{\min}(A) = 0$, then A is rank deficient and hence ABA^\top is also rank deficient, i.e. $\lambda_{\min}(ABA^\top) = 0$. Therefore $\lambda_{\min}(ABA^\top) = 0 = \sigma_{\min}^2(A) \lambda_{\min}(B)$. \square

Lemma D.2 (Weyl's inequality [AM56]). *For three real r -rank matrices, satisfying $A - B = C$, Weyl's inequality [AM56, Theorem 3.6], tells that*

$$\sigma_k(A) - \sigma_k(B) \leq \|C\|, \text{ for all } k \in [r] \quad (234)$$

where $\sigma_k(\cdot)$ is the k -th largest singular value operator.

Lemma D.3 (a variant of Woodburry matrix identity [HS81]). *For linear operators A and B such that A and $A + B$ are invertible, then*

$$(A + B)^{-1} - A^{-1} = -A^{-1}B(A + B)^{-1} \quad (235)$$

Lemma D.4. *Let $U \in \mathbb{R}^{d \times r}$ and $U^* \in \mathbb{R}^{d \times r}$ be two orthonormal matrices. Let $\{\sin \theta_j(U, U^*)\}_{j=1}^r$ be the singular values of $(U^*)^\top U$. Then following are true.*

$$\|U - U^*(U^*)^\top U\|_F \geq \|\mathbf{I} - (U^*)^\top U\|_F, \quad (236)$$

$$\|U - U^*(U^*)^\top U\|_F \geq r - \|(U^*)^\top U\|_F^2 \geq \sum_{k \in [r]} \sin^2 \theta_k(U, U^*), \quad (237)$$

$$\|(\mathbf{I} - U^*(U^*)^\top)U\| = \|(U_\perp^\top)^\top U\| = \|U_\perp^\top U^*\| = \|(\mathbf{I} - U(U)^\top)U^*\|, \quad (238)$$

$$\|(\mathbf{I} - U^*(U^*)^\top)U\|_F = \|(U_\perp^\top)^\top U\|_F = \|U_\perp^\top U^*\|_F = \|(\mathbf{I} - U(U)^\top)U^*\|_F, \text{ and} \quad (239)$$

$$\sigma_r((U^*)^\top U) \geq \sqrt{1 - \|(\mathbf{I} - U^*(U^*)^\top)U\|} \quad (240)$$

Proof.

$$\|U - U^*(U^*)^\top U\|_F^2 = \langle U - U^*(U^*)^\top U, U - U^*(U^*)^\top U \rangle \quad (241)$$

$$= \langle U, U \rangle - 2 \langle U^*(U^*)^\top U, U \rangle + \langle U^*(U^*)^\top U, U^*(U^*)^\top U \rangle \quad (242)$$

$$= r - 2 \text{tr}(((U^*)^\top U)^\top ((U^*)^\top U)) + \text{tr}(((U^*)^\top U)^\top ((U^*)^\top U)) \quad (243)$$

$$= r - \text{tr}(((U^*)^\top U)^\top ((U^*)^\top U)) \quad (244)$$

$$= r - \sum_{k \in [r]} \cos^2 \theta_k(U, U^*) = \sum_{k \in [r]} \sin^2 \theta_k(U, U^*) \geq \sin^2 \theta_1(U, U^*) \quad (245)$$

$$\geq \sum_{k \in [r]} (1 - \cos^2 \theta_k(U, U^*)) \quad (246)$$

$$\geq \sum_{k \in [r]} (1 - \cos \theta_k(U, U^*))^2 \quad (247)$$

$$= \|\mathbf{I} - (U^*)^\top U\|_F^2 \quad (248)$$

$$\|U_\perp^\top U^*\| = \sigma_{\max}(U_\perp^\top U^*) = \sqrt{\lambda_{\max}((U^*)^\top U_\perp U_\perp^\top U^*)} \quad (249)$$

$$= \sqrt{\lambda_{\max}((U^*)^\top U_\perp U_\perp^\top U_\perp U_\perp^\top U^*)} = \|U_\perp U_\perp^\top U^*\| = \|(\mathbf{I} - UU^\top)U^*\| \quad (250)$$

Note that for $\|z\| = 1$

$$1 = z^\top U^\top U z = z^\top U^\top U^* (U^*)^\top U z + z^\top U^\top U_\perp^* (U_\perp^*)^\top U z \quad (251)$$

$$\implies 1 - z^\top U^\top U^* (U^*)^\top U z = z^\top U^\top U_\perp^* (U_\perp^*)^\top U z \quad (252)$$

$$\implies 1 - \min_{\|z\|=1} z^\top U^\top U^* (U^*)^\top U z = \max_{\|z\|=1} z^\top U^\top U_\perp^* (U_\perp^*)^\top U z \quad (253)$$

$$\implies 1 - \sigma_{\min}^2((U^*)^\top U) = \|(U_\perp^*)^\top U\|^2 \quad (254)$$

Therefore

$$\sigma_{\min}^2(U^\top U^*) + \|U_\perp^\top U^*\|^2 = 1 = \sigma_{\min}^2((U^*)^\top U) + \|(U_\perp^*)^\top U\|^2 \implies \|U_\perp^\top U^*\| = \|(U_\perp^*)^\top U\| \quad (255)$$

Rest of the equality can be obtained in a similar fashion using the above two relations.

$$\|U_\perp^\top U^*\|_F^2 = \text{tr}((U^*)^\top U_\perp U_\perp^\top U^*) = \text{tr}((U^*)^\top (\mathbf{I} - U U^\top) U^*) \quad (256)$$

$$= \text{tr}((U^*)^\top (\mathbf{I} - U U^\top)^2 U^*) \quad (257)$$

$$= \|(\mathbf{I} - U U^\top) U^*\|_F^2 \quad (258)$$

$$= \|(\mathbf{I} - U^* (U^*)^\top) U\|_F^2 = \|(U_\perp^*)^\top U\|_F^2 \quad (259)$$

Let $E = (\mathbf{I} - U^* (U^*)^\top) U$ and $Q = (U^*)^\top U$. Then $U^\top E = \mathbf{I} - Q^\top Q$. Then by Weyl's inequality (Lemma D.2, by setting $A \leftarrow \mathbf{I}$, $B \leftarrow Q^\top Q$, and $C \leftarrow U^\top E$) we get that

$$1 - \sigma_r(Q)^2 = \sigma_r(\mathbf{I}) - \sigma_r(Q^\top Q) \leq \|U^\top E\| \leq \|U\| \|E\| \leq \|(\mathbf{I} - U^* (U^*)^\top) U\| \quad (260)$$

This implies that $\sigma_r((U^*)^\top U) \geq \sqrt{1 - \|(\mathbf{I} - U^* (U^*)^\top) U\|}$ \square

Lemma D.5 (Hanson-Wright inequality, Theorem 6.2.1 [Ver18]). *Let $x_1, \dots, x_m \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$ be m i.i.d. standard isotropic Gaussian random vectors of dimension d . Then, for some universal constant $c \geq 0$, the following holds true with a probability of at least $1 - \delta$.*

$$\left| \frac{1}{m} \sum_{j=1}^m x_j^\top A_j x_j - \frac{1}{m} \sum_{j=1}^m \text{tr} A_j \right| \leq c \max \left(\sqrt{\sum_{j=1}^m \|A_j\|_F^2 \frac{\log(1/\delta)}{m^2}}, \max_{j=1, \dots, m} \|A_j\|_2 \frac{\log(1/\delta)}{m} \right) \quad (261)$$

Lemma D.6. *Let $x_1, \dots, x_m \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$ be m i.i.d. standard isotropic Gaussian random vectors of dimension d . Then, for some universal constant $c \geq 0$, the following holds true with a probability of at least $1 - \delta$.*

$$\left| \frac{1}{m} \sum_{j=1}^m a^\top (x_j x_j^\top) b - a^\top b \right| \leq c \|a\| \|b\| \max \left(\sqrt{\frac{\log(1/\delta)}{m}}, \frac{\log(1/\delta)}{m} \right) \quad (262)$$

Proof. First notice that $a^\top (x_j x_j^\top) b = \text{tr}(a^\top (x_j x_j^\top) b) = \text{tr}(x_j^\top b a^\top x_j) = x_j^\top b a^\top x_j$ and $a^\top b = \text{tr}(b a^\top)$. Then desired result follows from Lemma D.5, by setting $A_j = b a^\top$. \square

Lemma D.7. *Let $x_1, \dots, x_m \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$ be m i.i.d. standard isotropic Gaussian random vectors of dimension d . Then, for some universal constant $c \geq 0$, the following holds true with a probability of at least $1 - \delta$.*

$$\left\| \frac{1}{m} \sum_{j=1}^m a_j x_j x_j^\top - \frac{1}{m} \sum_{j=1}^m a_j \mathbf{I} \right\| \leq c \max \left(\frac{\|a\|_2}{\sqrt{m}} \sqrt{\frac{d \log(9) + \log(1/\delta)}{m}}, \|a\|_\infty \frac{d \log(9) + \log(1/\delta)}{m} \right) \quad (263)$$

Proof. For $\epsilon \leq 1$, consider a unit vector $u \in N_\epsilon$ from the ϵ -net of size $|N_\epsilon| = (1 + 2/\epsilon)^d$, of the sphere \mathbb{S}^{d-1} [Ver10, Lemma 5.2]. That is for any $u' \in \mathbb{S}^{d-1}$, there exists some $u \in N_\epsilon$ such that $\|u' - u\| \leq \epsilon$.

Now we will prove a concentration for $\frac{1}{m} \sum_{j=1}^m a_j u^\top x_j x_j^\top u - \frac{1}{m} \sum_{j=1}^m a_j$. Notice that, $a_j u^\top (x_j x_j^\top) u = a_j \text{tr}(u^\top x_j x_j^\top u) = a_j \text{tr}(x_j^\top u u^\top x_j) = x_j^\top (a_j u u^\top) x_j$ and $\text{tr}(a_j u u^\top) = a_j$. Then, by Hanson-Wright inequality

(Lemma D.5), for some universal constant $c \geq 0$, the following holds true with a probability of at least $1 - \delta'$.

$$\left| \frac{1}{m} \sum_{j=1}^m a_j u^\top x_j x_j^\top u - \frac{1}{m} \sum_{j=1}^m a_j \right| \leq c \max \left(\frac{\|a\|_2}{\sqrt{m}} \sqrt{\frac{\log(1/\delta')}{m}}, \|a\|_\infty \frac{\log(1/\delta')}{m} \right) \quad (264)$$

This implies that, through union bound, for the matrix $A' = \frac{1}{m} \sum_{j=1}^m a_j x_j x_j^\top - \frac{1}{m} \sum_{j=1}^m a_j \mathbf{I}$ the following holds true with probability at least $1 - \delta$

$$u^\top A' u \leq c \max \left(\frac{\|a\|_2}{\sqrt{m}} \sqrt{\frac{\log(|N_\epsilon|/\delta)}{m}}, \|a\|_\infty \frac{\log(|N_\epsilon|/\delta)}{m} \right), \quad \text{any } u \in N_\epsilon \quad (265)$$

Let $u' \in \mathbb{S}^{d-1}$ be the top singular-value of A' , then there exists some $u \in N_\epsilon$ such that $\|u' - u\| \leq \epsilon$.

$$\sigma_{\max}(A') = (u')^\top A' u' = (u' - u)^\top A' u' + u^\top A' u' + u^\top A' u \quad (266)$$

$$\leq \|u' - u\| \sigma_{\max}(A') \|u'\| + \|u\| \sigma_{\max}(A') \|u' - u\| + u^\top A' u \quad (267)$$

$$(268)$$

Re-arranging and setting $\epsilon = 1/4$ and setting $c \leftarrow 2c$, we get

$$\sigma_{\max}(A') \leq \frac{u^\top A' u}{1 - 2\epsilon} \leq 2c \max \left(\frac{\|a\|_2}{\sqrt{m}} \sqrt{\frac{d \log(9) + \log(1/\delta)}{m}}, \|a\|_\infty \frac{d \log(9) + \log(1/\delta)}{m} \right) \quad (269)$$

□