
Multi-task Representation Learning for Pure Exploration in Linear Bandits

Yihan Du¹ Longbo Huang¹ Wen Sun²

Abstract

Despite the recent success of representation learning in sequential decision making, the study of the pure exploration scenario (i.e., identify the best option and minimize the sample complexity) is still limited. In this paper, we study multi-task representation learning for best arm identification in linear bandits (RepBAI-LB) and best policy identification in contextual linear bandits (RepBPI-CLB), two popular pure exploration settings with wide applications, e.g., clinical trials and web content optimization. In these two problems, all tasks share a common low-dimensional linear representation, and our goal is to leverage this feature to accelerate the best arm (policy) identification process for all tasks. For these problems, we design computationally and sample efficient algorithms DouExpDes and C-DouExpDes, which perform double experimental designs to plan optimal sample allocations for learning the global representation. We show that by learning the common representation among tasks, our sample complexity is significantly better than that of the native approach which solves tasks independently. To the best of our knowledge, this is the first work to demonstrate the benefits of representation learning for multi-task pure exploration.

1. Introduction

Multi-task representation learning (Caruana, 1997) is an important problem which aims to learn a common low-dimensional representation from multiple related tasks. Representation learning has received extensive attention in both empirical applications (Ando et al., 2005; Bengio et al., 2013; Li et al., 2014) and theoretical study (Maurer et al., 2016; Du et al., 2021a; Tripuraneni et al., 2021).

¹IIS, Tsinghua University ²Cornell University. Correspondence to: Yihan Du <duyh18@mails.tsinghua.edu.cn>, Longbo Huang <longbohuang@tsinghua.edu.cn>, Wen Sun <ws455@cornell.edu>.

Recently, an emerging number of works (Yang et al., 2021; 2022; Hu et al., 2021; Cella et al., 2022b) investigate representation learning for sequential decision making, and show that if all tasks share a joint low-rank representation, then by leveraging such a joint representation, it is possible to learn faster than treating each task independently. Despite the accomplishments of these works, they mainly focus on the regret minimization setting, where the performance is measured by the cumulative reward gap between the optimal option and the actually chosen options.

However, in real-world applications where obtaining a sample is expensive and time-consuming, e.g., clinical trials (Zhang et al., 2012), it is often desirable to identify the optimal option using as few samples as possible, i.e., we face the *pure exploration* scenario rather than regret minimization. Moreover, in many decision-making applications, we often need to tackle multiple related tasks, e.g., treatment planning for different diseases (Bragman et al., 2018) and content optimization for multiple websites (Agarwal et al., 2009), and there usually exists a common representation among these tasks, e.g., the features of drugs and the representations of website items. Thus, we desire to exploit the shared representation among tasks to expedite learning. For example, in clinical treatment planning, we want to identify the optimal treatment for multiple diseases, and there exists a joint representation of treatments. In this case, since conducting a clinical trial and collecting a sample is time-consuming, we desire to make use of the shared representation and reduce the number of samples required.

Motivated by the above fact, in this paper, we study representation learning for multi-task pure exploration in sequential decision making. Following prior works (Yang et al., 2021; 2022; Hu et al., 2021), we consider the linear bandit setting, which is one of the most popular settings in sequential decision making and has various applications such as clinical trials and recommendation systems. Specifically, we investigate two pure exploration problems, i.e., representation learning for best arm identification in linear bandits (RepBAI-LB) and best policy identification in contextual linear bandits (RepBPI-CLB).

In RepBAI-LB, an agent is given a confidence parameter δ , an arm set $\mathcal{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$ and M tasks. For each task $m \in [M]$, the expected reward of each arm $\mathbf{x} \in \mathcal{X}$

is generated by $\mathbf{x}^\top \boldsymbol{\theta}_m$, where $\boldsymbol{\theta}_m \in \mathbb{R}^d$ is an underlying reward parameter. There exists an unknown global feature extractor $\mathbf{B} \in \mathbb{R}^{d \times k}$ and an underlying prediction parameter \mathbf{w}_m such that $\boldsymbol{\theta}_m = \mathbf{B}\mathbf{w}_m$ for any $m \in [M]$, where $M \gg d \gg k$. We can understand the problem as that all tasks share a joint representation $\mathbf{f}(\mathbf{x}) := \mathbf{B}^\top \mathbf{x}$ for arms, where the dimension of $\mathbf{f}(\mathbf{x})$ is much smaller than that of \mathbf{x} . The agent sequentially selects arms and tasks to sample, and observes noisy rewards. The goal of the agent is to identify the best arm with the maximum expected reward for each task with confidence $1 - \delta$, using as few samples as possible.

The RepBPI-CLB problem is an extension of RepBAI-LB to environments with random and varying contexts. In RepBPI-CLB, there are a context space \mathcal{S} , an action space \mathcal{A} , a known feature mapping $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$ and an *unknown* context distribution \mathcal{D} . For each task $m \in [M]$, the expected reward of each context-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ is generated by $\phi(s, a)^\top \boldsymbol{\theta}_m$, where $\boldsymbol{\theta}_m = \mathbf{B}\mathbf{w}_m$. We can similarly interpret the problem as that all tasks share a low-dimensional context-action representation $\mathbf{B}^\top \phi(s, a) \in \mathbb{R}^k$. At each timestep, the agent first observes a context drawn from \mathcal{D} , and chooses an action and a task to sample, and then observes a random reward. Given a confidence parameter δ and an accuracy parameter ε , the agent aims to identify an ε -optimal policy (i.e., a mapping $\mathcal{S} \mapsto \mathcal{A}$ that gives suboptimality within ε) for each task with confidence $1 - \delta$, while minimizing the number of samples used.

In contrast to existing representation learning works (Yang et al., 2021; 2022; Hu et al., 2021; Cella et al., 2022b), we focus on the pure exploration scenario and face several unique challenges: (i) The sample complexity minimization objective requires us to plan an optimal sample allocation for recovering the low-rank representation, in order to save samples to the highest degree. (ii) Unlike prior works which either assume that the arm set is an ellipsoid/sphere (Yang et al., 2021; 2022) or are computationally inefficient (Hu et al., 2021), we allow an arbitrary arm set that spans \mathbb{R}^d , which poses challenges on how to efficiently schedule samples according to the shapes of arms. (iii) Different from prior works (Huang et al., 2015; Li et al., 2022), we do not assume prior knowledge of the context distribution. This imposes additional difficulties in sample allocation planning and estimator construction. To handle these challenges, we design computationally and sample efficient algorithms, which effectively estimate the context distribution and employ the experimental design approaches to plan samples.

We summarize our contributions in this paper as follows.

- We formulate the problems of multi-task representation learning for best arm identification in linear bandits (RepBAI-LB) and best policy identification in contextual linear bandits (RepBPI-CLB). To the best of our knowledge, this is the first work to study representation

learning in the multi-task pure exploration scenario.

- For RepBAI-LB, we propose an efficient algorithm DouExpDes equipped with *double experimental designs*. The first design optimally schedules samples to learn the joint representation according to arm shapes, and the second design minimizes the estimation error for rewards using low-dimensional representations. Furthermore, we establish a sample complexity guarantee $\tilde{O}(\frac{Mk}{\Delta_{\min}^2})$, which shows superiority over the baseline result $\tilde{O}(\frac{Md}{\Delta_{\min}^2})$ (i.e., solving each task independently). Here Δ_{\min} denotes the minimum reward gap.
- For RepBPI-CLB, we develop C-DouExpDes, an algorithm which efficiently estimates the context distribution and conducts double experimental designs under the estimated context distribution to learn the global representation. A sample complexity result $\tilde{O}(\frac{Mk^2}{\varepsilon^2})$ is also provided for C-DouExpDes, which significantly outperforms the baseline result $\tilde{O}(\frac{Md^2}{\varepsilon^2})$, and demonstrates the power of representation learning.

2. Related Work

In this section, we introduce two lines of related works, and defer a more complete literature review to Appendix A.

Representation Learning. The study of representation learning has been initiated and developed in the supervised learning setting, e.g., (Baxter, 2000; Ando et al., 2005; Maurer et al., 2016; Du et al., 2021a; Tripuraneni et al., 2021).

Recently, representation learning for sequential decision making has attracted extensive attention. Lale et al. (2019); Jun et al. (2019); Lu et al. (2021b); Huang et al. (2021) study linear bandits with a hidden low-rank structure (e.g., bilinear bandits), which is very related to the problem of representation learning. Yang et al. (2021; 2022); Hu et al. (2021); Cella et al. (2022b) consider multi-task representation learning for linear bandits with the regret minimization objective. Yang et al. (2021; 2022) assume that the arm set is an ellipsoid or sphere. Hu et al. (2021) relax this assumption and allow arbitrary arm sets, but their algorithms that build upon a multi-task joint least-square estimator are computationally inefficient. Cella et al. (2022b) design algorithms that do not need to know the dimension of the underlying representation. There are also other works (Lu et al., 2021a; 2022; Pacchiano et al., 2022; Zhang & Wang, 2021; Cheng et al., 2022; Agarwal et al., 2022) which investigate representation learning for reinforcement learning.

Different from the above works which consider regret minimization, we study representation learning for (contextual) linear bandits with the pure exploration objective, which brings unique challenges on how to optimally allocate samples to learn the feature extractor, and motivates us to design

algorithms based on double experimental designs.

Pure Exploration in (Contextual) Linear Bandits. Most existing linear bandit works focus on regret minimization, e.g., (Dani et al., 2008; Chu et al., 2011; Abbasi-Yadkori et al., 2011). Recently, there has been a surge of interests in the pure exploration objective for (contextual) linear bandits. For linear bandits, Soare et al. (2014) firstly apply the experimental design approach to distinguish the optimal arm, and establish sample complexity that heavily depends on the minimum reward gap. Tao et al. (2018) design a novel randomized estimator for the underlying reward parameter, and achieve tighter sample complexity which depends on the reward gaps of the best d arms. Fiez et al. (2019) provide the first near-optimal sample complexity upper and lower bounds for best arm identification in linear bandits. For contextual linear bandits, Zanette et al. (2021) develop a non-adaptive policy to collect data, from which a near-optimal policy can be computed. Li et al. (2022) build instance-optimal sample complexity for best policy identification in contextual linear bandits, with prior knowledge of the context distribution. By contrast, our work studies a multi-task setting where tasks share a common representation, and does not assume any prior knowledge of the context distribution.

3. Problem Formulation

In this section, we present the formal problem formulations of RepBAI-LB and RepBPI-CLB. Before describing the formulations, we first introduce some useful notations.

Notations. We use bold lower-case letters to denote vectors and bold upper-case letters to denote matrices. For any matrix \mathbf{A} , $\|\mathbf{A}\|$ denotes the spectral norm of \mathbf{A} , and $\sigma_{\min}(\mathbf{A})$ denotes the minimum singular value of \mathbf{A} . For any positive semi-definite matrix $\mathbf{A} \in \mathbb{R}^{d' \times d'}$ and vector $\mathbf{x} \in \mathbb{R}^{d'}$, $\|\mathbf{x}\|_{\mathbf{A}} := \sqrt{\mathbf{x}^\top \mathbf{A} \mathbf{x}}$. We use $\text{polylog}(\cdot)$ to denote a polylogarithmic factor in given parameters, and $\tilde{O}(\cdot)$ to denote an expression that hides polylogarithmic factors in all problem parameters except δ and ε .

Representation Learning for Best Arm Identification in Linear Bandits (RepBAI-LB). An agent is given a set of arms $\mathcal{X} := \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \subseteq \mathbb{R}^d$ and M best arm identification tasks. Without loss of generality, we assume that \mathcal{X} spans \mathbb{R}^d , as done in many prior works (Fiez et al., 2019; Katz-Samuels et al., 2020; Degenne et al., 2020). For any $\mathbf{x} \in \mathcal{X}$, $\|\mathbf{x}\| \leq L_x$ for some constant L_x . For each task $m \in [M]$, the expected reward of each arm $\mathbf{x} \in \mathcal{X}$ is $\mathbf{x}^\top \boldsymbol{\theta}_m$, where $\boldsymbol{\theta}_m \in \mathbb{R}^d$ is an unknown reward parameter. Among all tasks, there exists a common underlying feature extractor $\mathbf{B} \in \mathbb{R}^{d \times k}$, which satisfies that for each task $m \in [M]$, $\boldsymbol{\theta}_m = \mathbf{B} \mathbf{w}_m$. Here \mathbf{B} has orthonormal columns, $\mathbf{w}_m \in \mathbb{R}^k$ is an unknown prediction parameter, and $M \gg d \gg k$. For

any $m \in [M]$, $\|\mathbf{w}_m\| \leq L_w$ for some constant L_w .

At each timestep t , the agent chooses an arm $\mathbf{x} \in \mathcal{X}$ and a task $m \in [M]$, to sample arm \mathbf{x} in task m . Then, she observes a random reward $r_t = \mathbf{x}^\top \boldsymbol{\theta}_m + \eta_t = \mathbf{x}^\top \mathbf{B} \mathbf{w}_m + \eta_t$, where η_t is an independent, zero-mean and sub-Gaussian noise. For simplicity of analysis, we assume that $\mathbb{E}[\eta_t^2] = 1$, which can be easily relaxed by using a more carefully-designed estimator in our algorithm. Given a confidence parameter $\delta \in (0, 1)$, the agent aims to identify the best arms $\mathbf{x}_m^* := \arg\max_{\mathbf{x} \in \mathcal{X}} \mathbf{x}^\top \boldsymbol{\theta}_m$ for all tasks $m \in [M]$ with probability at least $1 - \delta$, using as few samples as possible. We define sample complexity as the total number of samples used over all tasks, which is the performance metric considered in our paper.

To efficiently learn the underlying low-dimensional representation, we make the following standard assumptions.

Assumption 3.1 (Diverse Tasks). We assume that $\sigma_{\min}(\frac{1}{M} \sum_{m=1}^M \mathbf{w}_m \mathbf{w}_m^\top) = \Omega(\frac{1}{k})$.

This assumption indicates that the prediction parameters $\mathbf{w}_1, \dots, \mathbf{w}_M$ are uniformly spread out in all directions of \mathbb{R}^k , which was also assumed in (Du et al., 2021a; Tripuraneni et al., 2021; Yang et al., 2021), and is necessary for recovering the feature extractor \mathbf{B} .

For any distribution $\boldsymbol{\lambda} \in \Delta_{\mathcal{X}}$ and $\mathbf{B} \in \mathbb{R}^{d \times k}$, let $\mathbf{A}(\boldsymbol{\lambda}, \mathbf{B}) := \sum_{i=1}^n \lambda(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B}$. For any task $m \in [M]$, let

$$\boldsymbol{\lambda}_m^* := \arg\min_{\boldsymbol{\lambda} \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_m^*\}} \frac{\|\mathbf{B}^\top (\mathbf{x}_m^* - \mathbf{x})\|_{\mathbf{A}(\boldsymbol{\lambda}, \mathbf{B})^{-1}}^2}{((\mathbf{x}_m^* - \mathbf{x})^\top \boldsymbol{\theta}_m)^2}.$$

Here $\boldsymbol{\lambda}_m^*$ denotes the optimal sample allocation that minimizes prediction error of arms (i.e., the solution of G-optimal design (Pukelsheim, 2006)) under the underlying low-dimensional representation.

Assumption 3.2 (Eigenvalue of G-optimal Design Matrix). For any task $m \in [M]$, $\sigma_{\min}(\mathbf{A}(\boldsymbol{\lambda}_m^*, \mathbf{B})) \geq \omega$ for some constant $\omega > 0$.

This assumption implies that the covariance matrix $\mathbf{A}(\boldsymbol{\lambda}_m^*, \mathbf{B})$ under the optimal sample allocation is invertible, which is necessary for estimating \mathbf{w}_m . Note that the quantities introduced in Assumptions 3.1 and 3.2, i.e., $\sigma_{\min}(\frac{1}{M} \sum_{m=1}^M \mathbf{w}_m \mathbf{w}_m^\top)$ and $\sigma_{\min}(\mathbf{A}(\boldsymbol{\lambda}_m^*, \mathbf{B}))$, are both defined on the low-dimensional subspace, which scale as k instead of d .

Representation Learning for Best Policy Identification in Contextual Linear Bandits (RepBPI-CLB). In this problem, there are a context space \mathcal{S} , an action space \mathcal{A} , a feature mapping $\phi(\cdot, \cdot) : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$ and an *unknown* context distribution $\mathcal{D} \in \Delta_{\mathcal{S}}$. For any $(s, a) \in \mathcal{S} \times \mathcal{A}$, $\|\phi(s, a)\| \leq L_\phi$ for some constant L_ϕ . An agent needs to solve M best

policy identification tasks. For each task $m \in [M]$, the expected reward of each context-action pair $(s, a) \in \mathcal{S} \times \mathcal{A}$ is $\phi(s, a)^\top \theta_m$, where $\theta_m \in \mathbb{R}^d$ is an unknown reward parameter. Similar to RepBAI-LB, there exists a global feature extractor $B \in \mathbb{R}^{d \times k}$ with orthonormal columns, such that for each task $m \in [M]$, $\theta_m = Bw_m$. Here $w_m \in \mathbb{R}^k$ is an unknown prediction parameter, $\|w_m\| \leq L_w$ for any $m \in [M]$, and $M \gg d \gg k$.

At each timestep t , the agent first observes a random context s_t , which is i.i.d. drawn from \mathcal{D} . Then, she selects an action $a_t \in \mathcal{A}$ and a task $m \in [M]$, to sample action a_t in context s_t under task m . After sampling, she observes a random reward $r_t = \phi(s_t, a_t)^\top \theta_m + \eta_t = \phi(s_t, a_t)^\top Bw_m + \eta_t$, where η_t is an independent, zero-mean and 1-sub-Gaussian noise.

We define a policy π as a mapping from \mathcal{S} to \mathcal{A} . For each task $m \in [M]$, we say a policy $\hat{\pi}_m$ is ε -optimal if

$$\mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} (\phi(s, a) - \phi(s, \hat{\pi}_m(s))^\top \theta_m) \right] \leq \varepsilon.$$

Given a confidence parameter $\delta \in (0, 1)$ and an accuracy parameter $\varepsilon > 0$, the goal of the agent is to identify an ε -optimal policy $\hat{\pi}_m$ for each task $m \in [M]$ with probability at least $1 - \delta$, and minimize the number of samples used, i.e., sample complexity.

We also make two standard assumptions for RepBPI-CLB: Assumption 3.1 and the following assumption on the context distribution and context-action features.

Assumption 3.3. There exists some $\lambda \in \Delta_{\mathcal{A}}$ such that

$$\sigma_{\min} \left(\sum_{a \in \mathcal{A}} \lambda(a) \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right) \geq \nu$$

for some constant $\nu > 0$.

Assumption 3.3 manifests that there exists at least one sample allocation, under which the expected covariance matrix with respect to random contexts is invertible. This assumption enables one to reveal the feature extractor B , despite stochastic and varying contexts. Note that Assumption 3.3 only assumes the existence of a feasible sample allocation, rather than the knowledge of this sample allocation.

It is worth mentioning that in this work, we do not assume that we can sample arbitrary vectors in an ellipsoid/sphere as in (Yang et al., 2021; 2022), or assume that each arm (action) has zero mean and identity covariance as in (Tripuraneni et al., 2021). In contrast, we allow arbitrary shapes of arms (actions), and efficiently allocate samples according to their different shapes. Moreover, we do not assume prior knowledge of the context distribution as in (Huang et al., 2015; Li et al., 2022). Instead, we design an effective

scheme to estimate the context distribution, and carefully bound the estimation error in our analysis.

Below we will introduce our algorithms and results. We defer all our proofs to Appendix due to space limit.

4. Representation Learning for Best Arm Identification in Linear Bandits

In this section, we design a computationally efficient algorithm DouExpDes for RepBAI-LB, which performs double delicate experimental designs to recover the feature extractor and distinguish the best arms using low-rank representations. Furthermore, we provide sample complexity guarantees that mainly depend on the underlying low dimension.

To better describe our algorithm, we first introduce the notion of *experimental design*. Experimental design is an important problem in statistics (Pukelsheim, 2006). Consider a set of feature vectors and an unknown linear regression parameter. Sampling each feature vector will produce a noisy feedback of the inner-product of this feature vector and the unknown parameter. Experimental design investigates how to schedule samples to maximize the statistical power of estimating the unknown parameter. In our algorithm, we mainly use two popular types of experimental design, i.e., *E-optimal design*, which minimizes the spectral norm of the inverse of sample covariance matrix, and *G-optimal design*, which minimizes the maximum prediction error for feature vectors.

4.1. Algorithm DouExpDes

Now we present our algorithm DouExpDes, whose pseudo-code is provided in Algorithm 1. DouExpDes is a phased elimination algorithm, which first conducts the E-optimal design to optimally schedule samples for learning the feature extractor B , and then performs the G-optimal design with low-dimensional representations to eliminate suboptimal arms.

DouExpDes uses a *rounding procedure* ROUND (Allen-Zhu et al., 2017; Fiez et al., 2019), which transforms a given continuous sample allocation (design) into a discrete sample sequence and maintains important properties (e.g., E-optimality and G-optimality) of the design. ROUND($\{(\mathbf{q}_i, \mathbf{Q}_i)\}_{i=1}^{n'}, \lambda, \zeta, N$) takes n' arm-matrix pairs $(\mathbf{q}_1, \mathbf{Q}_1), \dots, (\mathbf{q}_{n'}, \mathbf{Q}_{n'}) \in \mathcal{X} \times \mathbb{R}^{d' \times d'}$, a distribution $\lambda \in \Delta_{\{\mathbf{q}_1, \dots, \mathbf{q}_{n'}\}}$, a rounding approximation parameter $\zeta > 0$, and the number of samples N such that $N \geq \frac{180d'}{\zeta^2}$ as inputs. It will return a sample sequence $\mathbf{s}_1, \dots, \mathbf{s}_N \in \mathcal{X}$, which correspond to feature matrices $\mathbf{S}_1, \dots, \mathbf{S}_N \in \{\mathbf{Q}_1, \dots, \mathbf{Q}_{n'}\}$, and $\sum_{j=1}^N \mathbf{S}_j$ has similar properties as the covariance matrix of the inputted design $N \sum_{i=1}^{n'} \lambda(\mathbf{q}_i) \mathbf{Q}_i$ (see Appendix B for more details).

Algorithm 1 DouExpDes (Double Experimental Design)

- 1: **Input:** \mathcal{X} , δ , rounding procedure ROUND, rounding approximation parameter $\zeta := \frac{1}{10}$, and the size of sample batch $p := \frac{180d}{\zeta^2}$.
- 2: Let λ^E and ρ^E be the optimal solution and the optimal value of the E-optimal design optimization:

$$\min_{\lambda \in \Delta_{\mathcal{X}}} \left\| \left(\sum_{i=1}^n \lambda(x_i) x_i x_i^\top \right)^{-1} \right\|$$
- 3: $\bar{x}_1, \dots, \bar{x}_p \leftarrow \text{ROUND}(\{(x_i, x_i x_i^\top)\}_{i=1}^n, \lambda^E, \zeta, p)$
- 4: $\hat{\mathcal{X}}_{1,m} \leftarrow \mathcal{X}$ for any $m \in [M]$. $\delta_t \leftarrow \frac{\delta}{2t^2}$ for any $t \geq 1$
- 5: **for** phase $t = 1, 2, \dots$ **do**
- 6: $T_t \leftarrow \lceil \frac{c_1(1+\zeta)^3(\rho^E)^2 k^4 L_x^4 L_w^4}{M} \max\{2^{2t}, \frac{L_x^4}{\omega^2}\} \cdot \text{polylog}(\zeta, \rho^E, p, k, L_x, L_w, \frac{1}{\delta_t}, \frac{1}{\omega}) \rceil$, where c_1 is an absolute constant
- 7: $\hat{B}_t \leftarrow \text{FeatRecover}(T_t, \{\bar{x}_i\}_{i \in [p]})$
- 8: $\{\hat{\mathcal{X}}_{t+1,m}\}_{m \in [M]} \leftarrow \text{EliLowRep}(t, \mathcal{X}, \{\hat{\mathcal{X}}_{t,m}\}_{m \in [M]}, \delta_t, \text{ROUND}, \zeta, \hat{B}_t)$
- 9: **if** $|\hat{\mathcal{X}}_{t+1,m}| = 1, \forall m \in [M]$ **then**
- 10: **return** $\hat{\mathcal{X}}_{t+1,m}$ for all tasks $m \in [M]$
- 11: **end if**
- 12: **end for**

The procedure of DouExpDes is as follows. At the beginning, DouExpDes performs the E-optimal design with raw representations, to plan an optimal sample allocation λ^E for the purpose of recovering the feature extractor B (Line 2). Then, DouExpDes calls ROUND to convert the E-optimal sample allocation λ^E into a discrete sample batch $\bar{x}_1, \dots, \bar{x}_p$, which satisfies that

$$\left\| \left(\sum_{j=1}^p \bar{x}_j \bar{x}_j^\top \right)^{-1} \right\| \leq (1 + \zeta) \left\| \left(p \sum_{i=1}^n \lambda^E(x_i) x_i x_i^\top \right)^{-1} \right\|.$$

Next, DouExpDes enters multiple phases, and maintains a candidate arm set $\hat{\mathcal{X}}_{t,m}$ for each task. The specific value of T_t in Line 6 is presented in Eq. (8) of Appendix C.2.

In each phase t , DouExpDes first calls subroutine FeatRecover to recover the feature extractor B . In FeatRecover (Algorithm 2), we repeatedly sample $\bar{x}_1, \dots, \bar{x}_p$ in all tasks, and construct an estimator Z for $\frac{1}{M} \sum_{i=1}^M \theta_m \theta_m^\top$, which contains the information of underlying reward parameters (Line 9). Then, we perform SVD on Z and obtain the estimated feature extractor \hat{B} (Line 10).

Then, DouExpDes calls subroutine EliLowRep to eliminate suboptimal arms using low-dimensional representations. In EliLowRep (Algorithm 3), we conduct the G-optimal design with the reduced-dimensional representations $\hat{B}^\top x$, and obtain sample allocation λ_m^G for each task (Line 2). We further use ROUND to transform λ_m^G into a sample sequence

Algorithm 2 FeatRecover($T, \{\bar{x}_i\}_{i \in [p]}$)

- 1: **for** task $m \in [M]$ **do**
- 2: **for** round $j \in [T]$ **do**
- 3: **for** arm $i \in [p]$ **do**
- 4: Sample \bar{x}_i , and observe random reward $\alpha_{m,j,i}$
- 5: **end for**
- 6: $\tilde{\theta}_{m,j} \leftarrow (\sum_{i=1}^p \bar{x}_i \bar{x}_i^\top)^{-1} \sum_{i=1}^p \bar{x}_i \alpha_{m,j,i}$
- 7: **end for**
- 8: **end for**
- 9: $Z \leftarrow \frac{1}{MT} \sum_{m=1}^M \sum_{j=1}^T \tilde{\theta}_{m,j} (\tilde{\theta}_{m,j})^\top - (\sum_{i=1}^p \bar{x}_i \bar{x}_i^\top)^{-1}$
- 10: Perform SVD decomposition on Z , and let \hat{B} be the top- k left singular vectors of Z
- 11: **return** \hat{B}

Algorithm 3 EliLowRep($t, \mathcal{X}, \{\hat{\mathcal{X}}_m\}_{m \in [M]}, \delta', \text{ROUND}, \zeta, \hat{B}$)

- 1: **for** task $m \in [M]$ **do**
- 2: Let λ_m^G and ρ_m^G be the optimal solution and the optimal value of the G-optimal design optimization:

$$\text{argmin}_{\lambda \in \Delta_{\mathcal{X}}} \max_{x, x' \in \hat{\mathcal{X}}_m} \left\| \hat{B}^\top (x - x') \right\|_{A(\lambda, \hat{B})}^2$$
- 3: $N_m \leftarrow \lceil \max\{32(1 + \zeta)2^{2t} \rho_m^G \log(\frac{4n^2 M}{\delta'}), \frac{180k}{\zeta^2}\} \rceil$
- 4: $z_{m,1}, \dots, z_{m,N_m} \leftarrow \text{ROUND}(\{(x_i, \hat{B}^\top x_i x_i^\top \hat{B})\}_{i=1}^n, \lambda_m^G, \zeta, N_m)$
- 5: Sample the arms $z_{m,1}, \dots, z_{m,N_m} \in \mathcal{X}$, and observe random rewards $r_{m,1}, \dots, r_{m,N_m}$
- 6: Let $\tilde{z}_{m,j} := \hat{B}^\top z_{m,j}$ for any $j \in [N_m]$
- 7: $\hat{w}_m \leftarrow (\sum_{j=1}^{N_m} \tilde{z}_{m,j} \tilde{z}_{m,j}^\top)^{-1} \sum_{j=1}^{N_m} \tilde{z}_{m,j} r_{m,j}$
- 8: $\hat{\theta}_m \leftarrow \hat{B} \hat{w}_m$
- 9: $\hat{\mathcal{X}}'_m \leftarrow \hat{\mathcal{X}}_m \setminus \{x \in \hat{\mathcal{X}}_m \mid \exists x' \in \hat{\mathcal{X}}_m : (x' - x)^\top \hat{\theta}_m > 2^{-t}\}$
- 10: **end for**
- 11: **return** $\{\hat{\mathcal{X}}'_m\}_{m \in [M]}$

$z_{m,1}, \dots, z_{m,N_m}$, which satisfies that

$$\begin{aligned} & \max_{x, x' \in \hat{\mathcal{X}}_m} \|x - x'\|_{(\sum_{j=1}^{N_m} \hat{B}^\top z_{m,j} z_{m,j}^\top \hat{B})^{-1}}^2 \\ & \leq (1 + \zeta) \max_{x, x' \in \hat{\mathcal{X}}_m} \|x - x'\|_{(N_m \sum_{i=1}^n \lambda_m^G(x_i) \hat{B}^\top x_i x_i^\top \hat{B})^{-1}}^2. \end{aligned}$$

After sampling this sequence, we build estimators $\hat{w}_{t,m}$ and $\hat{\theta}_{t,m}$ for the underlying prediction parameter w_m and reward parameter θ_m , respectively (Lines 7-8). Then, we discard the arms that show large gaps to the estimated optimal arm for each task (Line 9).

4.2. Theoretical Performance of DouExpDes

In this subsection, we provide sample complexity guarantees for DouExpDes. To formally present our sample complexity,

we first revisit existing results for conventional single-task best arm identification in linear bandits (BAI-LB).

For a single-task BAI-LB instance with arm set $\mathcal{X} \in \mathbb{R}^d$ and underlying reward parameter $\theta \in \mathbb{R}^d$, the instance-dependent hardness is defined as (Fiez et al., 2019)

$$\rho^S(\mathcal{X}, \theta) := \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}^*\}} \frac{\|\mathbf{x}^* - \mathbf{x}\|^2_{(\sum_{i=1}^n \lambda(\mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top)^{-1}}}{((\mathbf{x}^* - \mathbf{x})^\top \theta)^2},$$

and the best known sample complexity result is $\tilde{O}(\rho^S(\mathcal{X}, \theta) \log(\frac{1}{\delta})) = \tilde{O}(\frac{d}{(\Delta_{\min}^S)^2} \log(\frac{1}{\delta}))$ (Fiez et al., 2019). Here $\mathbf{x}^* := \arg\max_{\mathbf{x} \in \mathcal{X}} \mathbf{x}^\top \theta$ denotes the best arm, and $\Delta_{\min}^S := \min_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}^*\}} (\mathbf{x}^* - \mathbf{x})^\top \theta$ refers to the minimum reward gap.

It can be seen that a naive algorithm for RepBAI-LB is to run an existing single-task BAI-LB algorithm (Fiez et al., 2019; Katz-Samuels et al., 2020) to solve M tasks independently. Then, the sample complexity of such naive algorithm is

$$\tilde{O}\left(\sum_{m=1}^M \rho^S(\mathcal{X}, \theta_m) \log\left(\frac{1}{\delta}\right)\right) = \tilde{O}\left(\frac{Md}{\Delta_{\min}^2} \log\left(\frac{1}{\delta}\right)\right), \quad (1)$$

where $\Delta_{\min} := \min_{m \in [M], \mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_m^*\}} (\mathbf{x}_m^* - \mathbf{x})^\top \theta_m$ denotes the minimum reward gap among all tasks. In the following, we take Eq. (1) as the baseline to demonstrate the power of representation learning.

Now we state the sample complexity for DouExpDes.

Theorem 4.1. *With probability at least $1 - \delta$, algorithm DouExpDes returns the best arms \mathbf{x}_m^* for all tasks $m \in [M]$, and the number of samples used is bounded by*

$$\begin{aligned} & \tilde{O}\left(\sum_{m=1}^M \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_m^*\}} \frac{\|\mathbf{B}^\top(\mathbf{x}_m^* - \mathbf{x})\|_{\mathbf{A}(\lambda, \mathbf{B})}^2}{((\mathbf{x}_m^* - \mathbf{x})^\top \theta_m)^2} \log\left(\frac{1}{\delta}\right)\right. \\ & \quad \left. + (\rho^E)^2 dk^4 L_x^2 L_w^2 D \log^4\left(\frac{1}{\delta}\right)\right) \\ & = \tilde{O}\left(\frac{Mk}{\Delta_{\min}^2} \log\left(\frac{1}{\delta}\right) + (\rho^E)^2 dk^4 L_x^2 L_w^2 D \log^4\left(\frac{1}{\delta}\right)\right), \end{aligned} \quad (2)$$

where $D := \max\{\frac{1}{\Delta_{\min}^2}, \frac{L_x^4}{\omega^2}\}$.

Remark 1. In Theorem 4.1, the factors that have implicit dimensional dependency include $\min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_m^*\}} \frac{\|\mathbf{B}^\top(\mathbf{x}_m^* - \mathbf{x})\|_{\mathbf{A}(\lambda, \mathbf{B})}^2}{((\mathbf{x}_m^* - \mathbf{x})^\top \theta_m)^2}$, ω and ρ^E , which scale as k , $\frac{1}{k}$ and d , respectively.

In our sample complexity bound (Eq. (2)), the first term, $\sum_{m=1}^M \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_m^*\}} \frac{\|\mathbf{B}^\top(\mathbf{x}_m^* - \mathbf{x})\|_{\mathbf{A}(\lambda, \mathbf{B})}^2}{((\mathbf{x}_m^* - \mathbf{x})^\top \theta_m)^2} = O(\frac{Mk}{\Delta_{\min}^2})$, represents the hardness of M k -dimensional linear bandit instances with arm set $\{\mathbf{B}^\top \mathbf{x} : \mathbf{x} \in \mathcal{X}\}$ and underlying reward parameters $\mathbf{w}_1, \dots, \mathbf{w}_M$. This term only

depends on the reduced dimension k , instead of d . In other words, it is an essential price that is needed for solving M low-dimensional tasks, even if one knows the feature extractor \mathbf{B} . The second term $(\rho^E)^2 dk^4 L_x^2 L_w^2 D$, which depends on the raw dimension d , is a cost paid for learning the feature extractor. Note that since this term does not contain M , the cost for learning the underlying features is paid only once, rather than for all tasks.

When $M \gg d \gg k$, the first term dominates the bound, which only depends on the low dimension k . This indicates that algorithm DouExpDes effectively learns the low-dimensional representation, and exploits the intrinsic problem structure to reduce the sample complexity from $\tilde{O}(\frac{Md}{\Delta_{\min}^2} \log(\frac{1}{\delta}))$ (i.e., learning each task independently) to only $\tilde{O}(\frac{Mk}{\Delta_{\min}^2} \log(\frac{1}{\delta}))$. Our result corroborates the benefits of representation learning for multi-task pure exploration.

Technical Novelty. We highlight the novelty in the analysis of Theorem 4.1 as follows. (i) Prior low-rank bandit works (Jun et al., 2019; Lu et al., 2021b) use *arbitrary* sample distributions to recover the low-dimensional subspace, and their results depend on the eigenvalue of an arbitrary sample distribution $\|\mathbf{X}^{-1}\|$, where $\mathbf{X} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(d_1)}]$ is a collection of arbitrary d_1 arms from the arm set. By contrast, we utilize the *E-optimality* of the sample batch $\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_p$ to obtain an optimized dependency $\rho^E \approx \min_{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(d_1)} \in \mathcal{X}} \|\mathbf{X}^{-1}\|$, which is the best one can achieve at the subspace recovery stage. (ii) If one naively applies existing single-task BAI-LB analysis (Fiez et al., 2019; Katz-Samuels et al., 2020) in the estimated subspace $\hat{\mathbf{B}}_t$, one can only obtain a sample complexity $\|\hat{\mathbf{B}}_t^\top(\mathbf{x} - \mathbf{x}')\|_{(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t)^{-1}}^2$ dependent on $\hat{\mathbf{B}}_t$, but this is not a valid upper bound. To tackle this challenge, we connect the low-dimensional sample complexity under the estimated subspace $\|\hat{\mathbf{B}}_t^\top(\mathbf{x} - \mathbf{x}')\|_{(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t)^{-1}}^2$ with that under the true subspace $\|\mathbf{B}^\top(\mathbf{x} - \mathbf{x}')\|_{(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B})^{-1}}^2$, and drive a tight sample complexity.

Lower Bound Conjecture. We conjecture that the lower bound for RepBAI-LB is $\Omega(\sum_{m=1}^M \rho^S(\mathcal{X}, \theta_m) \log(\frac{1}{\delta}))$. We describe the preliminary idea below.

First, the lower bound for single-task BAI-LB with arm set \mathcal{X} and underlying reward parameter θ_m is $\Omega(\rho^S(\mathcal{X}, \theta_m) \log(\frac{1}{\delta}))$ (Fiez et al., 2019). If the global feature extractor \mathbf{B} is known, then the RepBAI-LB problem will reduce to M k -dimensional BAI-LB instances with arm set $\{\mathbf{B}^\top \mathbf{x} : \mathbf{x} \in \mathcal{X}\}$ and underlying reward parameters $\mathbf{w}_1, \dots, \mathbf{w}_M$. Therefore, we conjecture that the lower bound for RepBAI-LB is $\Omega(\sum_{m=1}^M \rho^S(\mathcal{X}, \theta_m) \log(\frac{1}{\delta}))$, which is the cost of solving M k -dimensional BAI-LB instances. However, it is challenging to rigorously analyze

Algorithm 4 C-DouExpDes (Contextual Double Experimental Design)

1: **Input:** $\delta, \varepsilon, \phi(\cdot, \cdot)$, regularization parameter $\gamma \geq 1$, rounding procedure ROUND, rounding approximation parameter $\zeta := \frac{1}{10}$, and the size of sample batch $p := \lceil \frac{c_2(1+\zeta)^2 L_\phi^4}{\nu^2} \text{polylog}(\zeta, M, d, k, L_\phi, L_w, \gamma, \frac{1}{\nu}, \frac{1}{\delta}, \frac{1}{\varepsilon}) \rceil$, where c_2 is an absolute constant.

2: $T_0 \leftarrow \lceil \frac{32^2(1+\zeta)^2 L_\phi^4}{\nu^2} \log^2(\frac{20d|A|}{\delta}) \rceil$. $\hat{\mathcal{D}} \leftarrow \emptyset$

3: **for** $\tau \in [T_0]$ **do**

4: Observe context s_τ , and randomly sample an action

5: $\hat{\mathcal{D}} \leftarrow \hat{\mathcal{D}} \cup \{s_\tau\}$

6: **end for**

7: Let $\lambda_{\hat{\mathcal{D}}}^E$ and $\rho_{\hat{\mathcal{D}}}^E$ be the optimal solution and the optimal value of the E-optimal design optimization:

$$\min_{\lambda \in \Delta_A} \left\| \left(\sum_{a \in A} \lambda(a) \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] \right)^{-1} \right\|$$

8: $\{\bar{a}_i\}_{i \in [p]} \leftarrow \text{ROUND}(\{(a, \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top])\}_{a \in A}, \lambda_{\hat{\mathcal{D}}}^E, \zeta, p)$

9: $T \leftarrow \lceil \frac{c_3(1+\zeta)^2 k^4 L_\phi^4 L_w^4}{M \nu^2 \varepsilon^2} \text{polylog}(\zeta, d, k, L_\phi, L_w, \gamma, \frac{1}{\nu}, \frac{1}{\delta}, \frac{1}{\varepsilon}) \rceil$, where c_3 is an absolute constant

10: $\hat{\mathbf{B}} \leftarrow \text{C-FeatRecover}(T, \{\bar{a}_i\}_{i \in [p]})$

11: $N \leftarrow \lceil \frac{(k^2 + \gamma k L_w^2)}{\varepsilon^2} \log^4(\frac{\gamma k L_w}{\varepsilon \delta}) \rceil$

12: $\{\hat{\theta}_{m,N}\}_{m \in [M]} \leftarrow \text{EstLowRep}(N, \gamma, \hat{\mathbf{B}})$

13: **return** $\hat{\pi}_m(\cdot) := \arg\max_{a \in A} \phi(\cdot, a)^\top \hat{\theta}_{m,N}$ for all tasks $m \in [M]$

the independence of these M k -dimensional instances and drive the summation in our conjectured lower bound. We leave the formal lower bound proof for future work.

When $M \gg d \gg k$, Theorem 4.1 matches our conjectured lower bound, which implies that algorithm DouExpDes performs as well as an oracle that knows the low-rank representation \mathbf{B} in advance.

5. Representation Learning for Best Policy Identification in Contextual Linear Bandits

In this section, we turn to contextual linear bandits. Different from prior contextual linear bandit works, e.g., (Huang et al., 2015; Li et al., 2022), here we do not assume any knowledge of context distribution. As a result, our RepBPI-CLB problem faces several unique challenges: (i) how to plan an efficient sample allocation for recovering the feature extractor in advance under an *unknown* context distribution, and (ii) how to construct an estimator for the feature extractor with a partially observed context space.

We propose algorithm C-DouExpDes, which first (i) efficiently estimates the context distribution and conducts ex-

Algorithm 5 C-FeatRecover($T, \{\bar{a}_i\}_{i \in [p]}$)

1: **for** task $m \in [M]$ **do**

2: **for** round $j \in [T]$ **do**

3: **for** arm $i \in [p]$ **do**

4: Observe context $s_{m,j,i}^{(1)}$, sample action \bar{a}_i in task m , and observe reward $\alpha_{m,j,i}^{(1)}$

5: Observe context $s_{m,j,i}^{(2)}$, sample action \bar{a}_i in task m , and observe reward $\alpha_{m,j,i}^{(2)}$

6: **end for**

7: Let $\phi_{m,j,i}^{(\ell)} := \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)$, $\forall i \in [p], \forall \ell \in \{1, 2\}$

8: $\tilde{\theta}_{m,j}^{(\ell)} \leftarrow (\sum_{i=1}^p \phi_{m,j,i}^{(\ell)} \phi_{m,j,i}^{(\ell)\top})^{-1} \sum_{i=1}^p \phi_{m,j,i}^{(\ell)} \alpha_{m,j,i}^{(\ell)}$, $\forall \ell \in \{1, 2\}$

9: **end for**

10: **end for**

11: $\mathbf{Z} \leftarrow \frac{1}{MT} \sum_{m=1}^M \sum_{j=1}^T \tilde{\theta}_{m,j}^{(1)} (\tilde{\theta}_{m,j}^{(2)})^\top$

12: Perform SVD decomposition on \mathbf{Z} , and let $\hat{\mathbf{B}}$ be the top- k left singular vectors

13: **return** $\hat{\mathbf{B}}$

Algorithm 6 EstLowRep($N, \gamma, \hat{\mathbf{B}}$)

1: $\Sigma_{m,0} \leftarrow \gamma I$ for any $m \in [M]$

2: **for** task $m \in [M]$ **do**

3: **for** timestep $t \in [N]$ **do**

4: Observe context $s_{m,t}$

5: $a_{m,t} \leftarrow \arg\max_{a \in A} \|\hat{\mathbf{B}}^\top \phi(s_{m,t}, a)\|_{\Sigma_{m,t-1}^{-1}}$

6: Sample action $a_{m,t}$, and observe reward $r_{m,t}$

7: $\Sigma_{m,t} \leftarrow \Sigma_{m,t-1} + \hat{\mathbf{B}}^\top \phi(s_{m,t}, a_{m,t}) \phi(s_{m,t}, a_{m,t})^\top \hat{\mathbf{B}}$

8: $\hat{\mathbf{w}}_{m,t} \leftarrow \Sigma_{m,t}^{-1} \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) r_{m,\tau}$

9: $\hat{\theta}_{m,t} \leftarrow \hat{\mathbf{B}} \hat{\mathbf{w}}_{m,t}$

10: **end for**

11: **end for**

12: **return** $\{\hat{\theta}_{m,N}\}_{m \in [M]}$

perimental designs under the estimated context distribution, and then (ii) builds a delicate estimator for the feature extractor using instantaneous contexts. Moreover, we also establish a sample complexity guarantee for C-DouExpDes, which mainly depends on the low dimension of the common representation among tasks.

5.1. Algorithm C-DouExpDes

Algorithm 4 presents the pseudo-code of C-DouExpDes. At the beginning, C-DouExpDes uses T_0 samples to estimate the context distribution \mathcal{D} (Lines 3-6). Then, it performs the E-optimal design under the estimated context distribution $\hat{\mathcal{D}}$, and obtains an efficient sample allocation $\lambda_{\hat{\mathcal{D}}}^E$ for the purpose of recovering the feature extractor \mathbf{B} (Line 7). Further, C-DouExpDes calls the rounding procedure ROUND

to transform $\lambda_{\mathcal{D}}^E$ into a sample batch $\bar{a}_1, \dots, \bar{a}_p$, such that

$$\begin{aligned} & \left\| \left(\sum_{j=1}^p \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_j) \phi(s, \bar{a}_j)^\top] \right)^{-1} \right\| \\ & \leq (1 + \zeta) \left\| \left(p \sum_{a \in \mathcal{A}} \lambda_{\mathcal{D}}^E(a) \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] \right)^{-1} \right\|. \end{aligned}$$

The specific values of p and T in Lines 1, 9 are provided in Eq. (19) of Appendix D.1 and Eq. (29) of Appendix D.2, respectively.

Next, C-DouExpDes runs subroutine C-FeatRecover to estimate the feature extractor \mathbf{B} using the sample batch $\bar{a}_1, \dots, \bar{a}_p$. In C-FeatRecover (Algorithm 5), we repeatedly sample $\bar{a}_1, \dots, \bar{a}_p$ in all tasks with random contexts. In Lines 4-5, we sample this batch twice, and the superscripts (1) and (2) denotes the first and second samples, respectively. After sampling, we carefully establish an estimator \mathbf{Z} for the reward parameter related matrix $\frac{1}{M} \sum_{m=1}^M \boldsymbol{\theta}_m \boldsymbol{\theta}_m^\top$, using instantaneous context-action features $\phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top$. We then perform SVD decomposition on \mathbf{Z} to obtain the estimated feature extractor $\hat{\mathbf{B}}$ (Lines 11-12).

Then, C-DouExpDes calls subroutine EstLowRep, which adapts existing reward-free-exploration algorithm in (Zanette et al., 2021) with low-rank representations to estimate $\boldsymbol{\theta}_m$. In EstLowRep (Algorithm 6), we employ the estimated representation $\hat{\mathbf{B}}^\top \phi(s, a)$ to sample the actions with the maximum uncertainty under the observed contexts. After that, we construct estimators $\hat{\mathbf{w}}_{m,t}$ and $\hat{\boldsymbol{\theta}}_{m,t}$ for the prediction parameter $\hat{\mathbf{w}}_m$ and reward parameter $\hat{\boldsymbol{\theta}}_m$ (Lines 8-9). At last, C-DouExpDes returns the greedy policy with respect to the estimated reward parameter $\hat{\boldsymbol{\theta}}_{m,N}$ for each task.

5.2. Theoretical Performance of C-DouExpDes

Next, we establish sample complexity guarantees for algorithm C-DouExpDes. In order to illustrate the advantages of representation learning, we first review existing results for traditional single-task best policy identification in contextual linear bandits (BPI-CLB). For a single BPI-CLB instance with context-action features $\phi(s, a) \in \mathbb{R}^d$ and reward parameter $\boldsymbol{\theta} \in \mathbb{R}^d$, the best known sample complexity is $\tilde{O}(\frac{d^2}{\varepsilon^2} \log(\frac{1}{\delta}))$ (Zanette et al., 2021; Li et al., 2022).

Apparently, if one naively solves the RepBPI-CLB problem by running single-task BPI-CLB algorithms to tackle M tasks independently, one will have a sample complexity

$$\tilde{O}\left(\frac{Md^2}{\varepsilon^2} \log\left(\frac{1}{\delta}\right)\right),$$

which heavily depends on the raw dimension d of context-action features. The goal of representation learning is to leverage the common representation among tasks to alleviate the dependency of dimension and save samples.

Now we present the sample complexity for C-DouExpDes.

Theorem 5.1. *With probability at least $1 - \delta$, C-DouExpDes returns an ε -optimal policy $\hat{\pi}_m$ such that $\mathbb{E}_{s \sim \mathcal{D}}[\max_{a \in \mathcal{A}} (\phi(s, a) - \phi(s, \hat{\pi}_m(s))^\top \boldsymbol{\theta}_m)] \leq \varepsilon$ for each task $m \in [M]$, and the number of samples used is*

$$\tilde{O}\left(\frac{M(k^2 + \gamma k L_w^2)}{\varepsilon^2} + \frac{k^4 L_\phi^8 L_w^4}{\nu^4 \varepsilon^2}\right).$$

Remark 2. In this result, only factor ν has implicit dimensional dependency, which scales as $\frac{1}{d}$. The first term $\frac{M(k^2 + \gamma k L_w^2)}{\varepsilon^2}$ is a cost of identifying optimal policies for M tasks with k -dimensional features $\mathbf{B}^\top \phi(s, a)$. The second term $\frac{k^4 L_\phi^8 L_w^4}{\nu^4 \varepsilon^2}$ is a price paid for learning global feature extractor \mathbf{B} and does not depend on M . This indicates that we only need to pay this price once, and then enjoy the benefits of dimension reduction for all M tasks.

When $M \gg \frac{1}{\nu} \gg k$, this result becomes $\tilde{O}(\frac{Mk^2}{\varepsilon^2})$ and only depends on the low dimension k , which implies that C-DouExpDes performs as well as an oracle that knows the underlying low-rank subspace \mathbf{B} . This sample complexity significantly outperforms the baseline result $\tilde{O}(\frac{Md^2}{\varepsilon^2})$ (i.e., solving M tasks independently), and demonstrates the power of representation learning.

Analytical Novelty. Below we elaborate the novelty in the proof of Theorem 5.1. (i) We carefully bound the deviation between the context-action features under the estimated context distribution $\mathbb{E}_{s \sim \hat{\mathcal{D}}}[\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]$ and those under the true context distribution $\mathbb{E}_{s \sim \mathcal{D}}[\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]$. We further bound the distance between $\mathbb{E}_{s \sim \mathcal{D}}[\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]$ and the context-action features under actual instantaneous contexts $\phi(s_{m,j,i}^{(\ell)}, \bar{a}_i) \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top$. (ii) We leverage the E-optimality of the sample batch $\bar{a}_1, \dots, \bar{a}_p$ to bound $\|(\sum_{i=1}^p \phi_{m,j,i}^{(\ell)} \phi_{m,j,i}^{(\ell)\top})^{-1}\|$. Then, we establish a concentration inequality for $\|\mathbf{Z} - \frac{1}{M} \sum_{m=1}^M \boldsymbol{\theta}_m \boldsymbol{\theta}_m^\top\|$ using the bounded $\|(\sum_{i=1}^p \phi_{m,j,i}^{(\ell)} \phi_{m,j,i}^{(\ell)\top})^{-1}\|$ and matrix Bernstein inequality with truncated noises. (iii) Furthermore, we decompose the prediction error $\phi(s, a)^\top (\hat{\boldsymbol{\theta}}_{m,t} - \boldsymbol{\theta}_m)$ into three components, including the sample variance and bias of $\hat{\mathbf{w}}_{m,t}$, and the estimation error of $\hat{\mathbf{B}}$. This prediction error is bounded via self-normalized concentration inequalities with the reduced dimension k .

6. Experiments

In this section, we present experiments to evaluate the empirical performance of our algorithms.

In our experiments, we set $\delta = 0.005$, $d = 5$, $k = 2$ and $M \in [50, 230]$, where k divides M . In RepBAI-LB, \mathcal{X} is

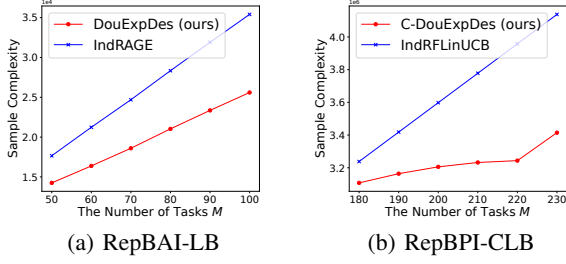


Figure 1. Experimental results for RepBAI-LB and RepBPI-CLB. The two figures compare the sample complexities of our algorithms with the naive algorithms which treat each task independently.

the canonical basis of \mathbb{R}^d . In RepBPI-CLB, we set $\varepsilon = 0.1$, $|\mathcal{S}| = 5$ and $|\mathcal{A}| = 5$. \mathcal{D} is the uniform distribution on \mathcal{S} . For any $s \in \mathcal{S}$, $\{\phi(s, a)\}_{a \in \mathcal{A}}$ is the canonical basis of \mathbb{R}^d . In both problems, $\mathbf{B} = [I_k; \mathbf{0}]$, where I_k denotes the $k \times k$ identity matrix. $\mathbf{w}_1, \dots, \mathbf{w}_M$ are divided into k groups, with $\frac{M}{k}$ same members in each group. The members in the i -th group ($i \in [k]$), i.e., $\mathbf{w}_{(M/k) \times (i-1) + 1}, \dots, \mathbf{w}_{(M/k) \times i}$, have 1 in the i -th coordinate and 0 in all other coordinates. For any $m \in [M]$, $\theta_m = \mathbf{B}\mathbf{w}_m$. We vary M and perform 50 independent runs to report the average sample complexity across runs.

For RepBAI-LB, we compare algorithm DouExpDes with the baseline IndRAGE which runs the state-of-the-art single-task BAI-LB algorithm RAGE (Fiez et al., 2019) to solve M tasks independently. Figure 1(a) shows the empirical results for RepBAI-LB. From Figure 1(a), we can see that DouExpDes has a better sample complexity than IndRAGE, and as the number of tasks M increases, the sample complexity of DouExpDes increases at a lower rate than that of IndRAGE. This demonstrates that DouExpDes effectively utilize the shared representation among tasks to reduce the number of samples needed for multi-task learning.

For RepBPI-CLB, our algorithm C-DouExpDes is compared with the baseline IndRFLinUCB, which tackles M tasks independently by calling the state-of-the-art single-task BPI-CLB algorithm Reward-free LinUCB (Zanette et al., 2021). As presented in Figure 1(b), C-DouExpDes achieves a significantly lower sample complexity than IndRFLinUCB. In addition, the slope of the sample complexity curve of C-DouExpDes with respect to M is much smaller than that of IndRFLinUCB, which validates that C-DouExpDes enjoys a lighter dependency on dimension in multi-task learning. These empirical results match our theoretical bounds, and corroborate the power of representation learning.

7. Conclusion and Future Work

In this paper, we investigate representation learning for pure exploration in multi-task (contextual) linear bandits. We

propose two efficient algorithms which conduct double experimental designs to optimally allocate samples for learning the low-rank representation. The sample complexities of our algorithms mainly depend on the low dimension of the underlying joint representation among tasks, instead of the raw high dimension. Our theoretical and experimental results demonstrate the benefit of representation learning for pure exploration in multi-task bandits. There are many interesting directions for further exploration. One direction is to establish lower bounds to validate the optimality of our algorithms. Another direction is to extend this work to more complex (nonlinear) representation settings.

Acknowledgements

The work of Yihan Du and Longbo Huang is supported by the Technology and Innovation Major Project of the Ministry of Science and Technology of China under Grant 2020AAA0108400 and 2020AAA0108403 and the Tsinghua Precision Medicine Foundation 10001020109. Wen Sun acknowledges funding support from NSF IIS-2154711.

References

- Abbasi-Yadkori, Y., Pál, D., and Szepesvári, C. Improved algorithms for linear stochastic bandits. *Advances in Neural Information Processing Systems*, 24, 2011.
- Agarwal, A., Song, Y., Sun, W., Wang, K., Wang, M., and Zhang, X. Provable benefits of representational transfer in reinforcement learning. *arXiv preprint arXiv:2205.14571*, 2022.
- Agarwal, D., Chen, B.-C., and Elango, P. Explore/exploit schemes for web content optimization. In *International Conference on Data Mining*, pp. 1–10. IEEE, 2009.
- Allen-Zhu, Z., Li, Y., Singh, A., and Wang, Y. Near-optimal design of experiments via regret minimization. In *International Conference on Machine Learning*, pp. 126–135. PMLR, 2017.
- Ando, R. K., Zhang, T., and Bartlett, P. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6(11), 2005.
- Baxter, J. A model of inductive bias learning. *Journal of Artificial Intelligence Research*, 12:149–198, 2000.
- Ben-David, S. and Schuller, R. Exploiting task relatedness for multiple task learning. In *Learning Theory and Kernel Machines*, pp. 567–580. Springer, 2003.
- Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives. *IEEE Transac-*

- tions on Pattern Analysis and Machine Intelligence, 35 (8):1798–1828, 2013.
- Bhatia, R. *Matrix analysis*, volume 169. Springer Science & Business Media, 2013.
- Bragman, F. J., Tanno, R., Eaton-Rosen, Z., Li, W., Hawkes, D. J., Ourselin, S., Alexander, D. C., McClelland, J. R., and Cardoso, M. J. Uncertainty in multitask learning: joint representations for probabilistic MR-only radiotherapy planning. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 3–11. Springer, 2018.
- Caruana, R. Multitask learning. *Machine Learning*, 28(1): 41–75, 1997.
- Cavallanti, G., Cesa-Bianchi, N., and Gentile, C. Linear algorithms for online multitask classification. *Journal of Machine Learning Research*, 11:2901–2934, 2010.
- Cella, L., Lounici, K., and Pontil, M. Meta representation learning with contextual linear bandits. *arXiv preprint arXiv:2205.15100*, 2022a.
- Cella, L., Lounici, K., and Pontil, M. Multi-task representation learning with stochastic linear bandits. *arXiv preprint arXiv:2202.10066*, 2022b.
- Cheng, Y., Feng, S., Yang, J., Zhang, H., and Liang, Y. Provable benefit of multitask representation learning in reinforcement learning. In *Advances in Neural Information Processing Systems*, 2022.
- Chu, W., Li, L., Reyzin, L., and Schapire, R. Contextual bandits with linear payoff functions. In *International Conference on Artificial Intelligence and Statistics*, pp. 208–214, 2011.
- Dani, V., Hayes, T. P., and Kakade, S. M. Stochastic linear optimization under bandit feedback. In *Conference on Learning Theory*, 2008.
- Degenne, R., Ménard, P., Shang, X., and Valko, M. Gamification of pure exploration for linear bandits. In *International Conference on Machine Learning*, pp. 2432–2442. PMLR, 2020.
- Du, S. S., Hu, W., Kakade, S. M., Lee, J. D., and Lei, Q. Few-shot learning via learning the representation, provably. In *International Conference on Learning Representations*, 2021a.
- Du, Y., Kuroki, Y., and Chen, W. Combinatorial pure exploration with full-bandit or partial linear feedback. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 7262–7270, 2021b.
- Fiez, T., Jain, L., Jamieson, K. G., and Ratliff, L. Sequential experimental design for transductive linear bandits. *Advances in Neural Information Processing Systems*, 32, 2019.
- Hu, J., Chen, X., Jin, C., Li, L., and Wang, L. Near-optimal representation learning for linear bandits and linear rl. In *International Conference on Machine Learning*, pp. 4349–4358. PMLR, 2021.
- Huang, B., Huang, K., Kakade, S., Lee, J. D., Lei, Q., Wang, R., and Yang, J. Optimal gradient-based algorithms for non-concave bandit optimization. *Advances in Neural Information Processing Systems*, 34:29101–29115, 2021.
- Huang, T.-K., Agarwal, A., Hsu, D. J., Langford, J., and Schapire, R. E. Efficient and parsimonious agnostic active learning. *Advances in Neural Information Processing Systems*, 28, 2015.
- Jedra, Y. and Proutiere, A. Optimal best-arm identification in linear bandits. *Advances in Neural Information Processing Systems*, 33:10007–10017, 2020.
- Jun, K.-S., Willett, R., Wright, S., and Nowak, R. Bilinear bandits with low-rank structure. In *International Conference on Machine Learning*, pp. 3163–3172. PMLR, 2019.
- Katz-Samuels, J., Jain, L., Jamieson, K. G., et al. An empirical process approach to the union bound: Practical algorithms for combinatorial and linear bandits. *Advances in Neural Information Processing Systems*, 33:10371–10382, 2020.
- Kiefer, J. and Wolfowitz, J. The equivalence of two extremum problems. *Canadian Journal of Mathematics*, 12: 363–366, 1960.
- Lale, S., Azizzadenesheli, K., Anandkumar, A., and Hassibi, B. Stochastic linear bandits with hidden low rank structure. *arXiv preprint arXiv:1901.09490*, 2019.
- Lattimore, T. and Hao, B. Bandit phase retrieval. *Advances in Neural Information Processing Systems*, 34:18801–18811, 2021.
- Li, J., Zhang, H., Zhang, L., Huang, X., and Zhang, L. Joint collaborative representation with multitask learning for hyperspectral image classification. *IEEE Transactions on Geoscience and Remote Sensing*, 52(9):5923–5936, 2014.
- Li, Z., Ratliff, L., Nassif, H., Jamieson, K., and Jain, L. Instance-optimal PAC algorithms for contextual bandits. *Advances in Neural Information Processing Systems*, 2022.

- Lu, R., Huang, G., and Du, S. S. On the power of multitask representation learning in linear MDP. *arXiv preprint arXiv:2106.08053*, 2021a.
- Lu, R., Zhao, A., Du, S. S., and Huang, G. Provable general function class representation learning in multitask bandits and MDPs. *Advances in Neural Information Processing Systems*, 2022.
- Lu, Y., Meisami, A., and Tewari, A. Low-rank generalized linear bandit problems. In *International Conference on Artificial Intelligence and Statistics*, pp. 460–468. PMLR, 2021b.
- Maurer, A. Bounds for linear multi-task learning. *Journal of Machine Learning Research*, 7:117–139, 2006.
- Maurer, A., Pontil, M., and Romera-Paredes, B. The benefit of multitask representation learning. *Journal of Machine Learning Research*, 17(81):1–32, 2016.
- Pacchiano, A., Nachum, O., Tripuraneni, N., and Bartlett, P. Joint representation training in sequential tasks with shared structure. *arXiv preprint arXiv:2206.12441*, 2022.
- Pukelsheim, F. *Optimal design of experiments*. SIAM, 2006.
- Qin, Y., Menara, T., Oymak, S., Ching, S., and Pasqualetti, F. Non-stationary representation learning in sequential linear bandits. *IEEE Open Journal of Control Systems*, 2022.
- Rivasplata, O. Subgaussian random variables: An expository note. *Internet Publication, PDF*, 5, 2012.
- Rusmevichientong, P. and Tsitsiklis, J. N. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.
- Soare, M., Lazaric, A., and Munos, R. Best-arm identification in linear bandits. *Advances in Neural Information Processing Systems*, 27, 2014.
- Tao, C., Blanco, S., and Zhou, Y. Best arm identification in linear bandits with linear dimension dependency. In *International Conference on Machine Learning*, pp. 4877–4886. PMLR, 2018.
- Tripuraneni, N., Jin, C., and Jordan, M. Provable meta-learning of linear representations. In *International Conference on Machine Learning*, pp. 10434–10443. PMLR, 2021.
- Tropp, J. A. et al. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.
- Xu, L., Honda, J., and Sugiyama, M. A fully adaptive algorithm for pure exploration in linear bandits. In *International Conference on Artificial Intelligence and Statistics*, pp. 843–851. PMLR, 2018.
- Yang, J., Hu, W., Lee, J. D., and Du, S. S. Impact of representation learning in linear bandits. In *ICLR*, 2021.
- Yang, J., Lei, Q., Lee, J. D., and Du, S. S. Nearly minimax algorithms for linear bandits with shared representation. *arXiv preprint arXiv:2203.15664*, 2022.
- Zanette, A., Dong, K., Lee, J. N., and Brunskill, E. Design of experiments for stochastic contextual linear bandits. *Advances in Neural Information Processing Systems*, 34: 22720–22731, 2021.
- Zhang, C. and Wang, Z. Provably efficient multi-task reinforcement learning with model transfer. *Advances in Neural Information Processing Systems*, 34:19771–19783, 2021.
- Zhang, D., Shen, D., Initiative, A. D. N., et al. Multi-modal multi-task learning for joint prediction of multiple regression and classification variables in Alzheimer’s disease. *NeuroImage*, 59(2):895–907, 2012.

Appendix

A. Related Work

In this section, we present a full literature review for two lines of related works, i.e., representation learning and pure exploration in (contextual) linear bandits.

Representation Learning. The study of representation learning has been initiated and developed in the supervised learning setting, e.g., (Baxter, 2000; Ben-David & Schuller, 2003; Ando et al., 2005; Maurer, 2006; Cavallanti et al., 2010; Maurer et al., 2016; Du et al., 2021a; Tripuraneni et al., 2021). A most related work is (Tripuraneni et al., 2021), which proposes a method-of-moments estimator for recovering the feature extractor, and establishes error guarantees for transferring the learned representation from past tasks to a new task.

Recently, representation learning for sequential decision making (bandits and reinforcement learning) has attracted extensive attention. We first introduce several works on low-rank bandits, which is a very similar topic to representation learning for bandits. Lale et al. (2019) study linear bandits with a hidden low-rank structure, and provide a regret bound dependent on the eigenvalue of the action distribution covariance. Jun et al. (2019); Lu et al. (2021b) also investigate low-rank linear bandits (bilinear bandits), and design algorithms which run traditional linear bandit algorithm LinUCB (Abbasi-Yadkori et al., 2011) in the estimated low-dimensional subspace. Lattimore & Hao (2021) consider an instantiation of low-rank bandits, called bandit phase retrieval. Huang et al. (2021) study a large family of bandit problems with non-concave reward functions, including low-rank linear bandits. They design a stochastic gradient-based algorithm that achieves an improved regret bound over those in (Jun et al., 2019; Lu et al., 2021b).

Now we introduce related works on representation learning for bandits. Yang et al. (2021; 2022) study multi-task representation learning for linear bandits with the regret minimization objective, and assume that the action set at each timestep is an ellipsoid or sphere. Hu et al. (2021) further relax this assumption and allow arbitrary action sets, but their algorithms equipped with a multi-task joint least-square estimator are computationally inefficient. Cella et al. (2022a;b) also investigate the problem in (Yang et al., 2021) and propose algorithms which do not need to know the dimension of the underlying representation. Qin et al. (2022) study multi-task representation learning for linear bandits in a non-stationary environment, and develop algorithms that learn and transfer non-stationary representations adaptively.

There are also other works studying multi-task representation learning for reinforcement learning (RL). Lu et al. (2021a; 2022) consider multi-task representation learning for linear MDPs, where the agent learns a shared representation function from a given function class. Pacchiano et al. (2022) investigate multi-task RL with a joint low-dimensional linear representation, and design a computationally efficient algorithm using a bilinear optimization oracle. Zhang & Wang (2021) consider multi-task (multi-player) RL in tabular MDPs, where the relatedness of MDPs are measured by the similarity of reward functions and transition distributions. Cheng et al. (2022); Agarwal et al. (2022) study multi-task representation learning and representational transfer for low-rank MDPs, where multiple low-rank MDPs share a common state-action feature mapping.

Different from the above works which consider regret minimization, we study representation learning for (contextual) linear bandits with the pure exploration objective, which imposes unique challenges on how to optimally allocate samples to learn the feature extractor, and motivates us to design algorithms based on double experimental designs.

Pure Exploration in (Contextual) Linear Bandits. Most linear bandit studies consider regret minimization, e.g., (Dani et al., 2008; Rusmevichientong & Tsitsiklis, 2010; Chu et al., 2011; Abbasi-Yadkori et al., 2011). Recently, there is a surge of interests in pure exploration for (contextual) linear bandits, e.g., (Soare et al., 2014; Tao et al., 2018; Xu et al., 2018; Fiez et al., 2019; Katz-Samuels et al., 2020; Degenne et al., 2020; Jedra & Proutiere, 2020; Du et al., 2021b; Zanette et al., 2021; Li et al., 2022). For linear bandits, Soare et al. (2014) firstly apply the G-optimal design to identify the best arm, and provide a sample complexity result that heavily depends on the minimum reward gap. Tao et al. (2018) design a novel randomized estimator for the underlying reward parameter, and achieve tighter sample complexity which depends on the reward gaps of the best d arms. Du et al. (2021b) further extend the algorithm in (Tao et al., 2018) to develop a polynomial-time algorithm for combinatorially large arm sets. Xu et al. (2018) propose a fully-adaptive algorithm which changes the arm selection strategy at each timestep. Fiez et al. (2019) establish the first near-optimal sample complexity upper and lower bounds for best arm identification in linear bandits. Katz-Samuels et al. (2020) further extend the algorithm in (Fiez et al., 2019) and use empirical processes to avoid an explicit union bound over the number of arms. Degenne et al. (2020); Jedra & Proutiere (2020) develop asymptotically optimal algorithms using the track-and-stop approaches. For contextual linear

bandits, Zanette et al. (2021) design a single non-adaptive policy to collect a dataset, from which a near-optimal policy can be computed. Li et al. (2022) build the first instance-dependent upper and lower bounds for best policy identification in contextual linear bandits, with the prior knowledge of the context distribution. By contrast, our work studies multi-task best arm/policy identification in (contextual) linear bandits with a shared representation among tasks, and does not assume any prior knowledge of the context distribution.

B. Rounding Procedure

In this section, we introduce the rounding procedure ROUND in detail.

Let $\mathcal{X}^+ := \mathcal{X} \cup \mathcal{A}$ denote the union space of arm set \mathcal{X} and action space \mathcal{A} . There are n arms or actions $p_1, \dots, p_n \in \mathcal{X}^+$ and n positive semi-definite matrices $\mathbf{Q}_1, \dots, \mathbf{Q}_n \in \mathbb{S}_+^d$, where \mathbf{Q}_i represents the feature of arm or action p_i for any $i \in [n]$. Denote $\mathcal{P} := \{p_1, \dots, p_n\}$ and $\mathcal{Q} := \{\mathbf{Q}_1, \dots, \mathbf{Q}_n\}$.

The rounding procedure $\text{ROUND}(\{(p_i, \mathbf{Q}_i)\}_{i=1}^n, \lambda, \zeta, N)$ (Allen-Zhu et al., 2017; Fiez et al., 2019) takes n arm-matrix or action-matrix pairs $(p_1, \mathbf{Q}_1), \dots, (p_n, \mathbf{Q}_n) \in \mathcal{X}^+ \times \mathbb{S}_+^d$, a distribution $\lambda \in \Delta_{\mathcal{P}}$ (or equivalently, $\lambda \in \Delta_{\mathcal{Q}}$), an approximation parameter $\zeta > 0$, and the number of samples N which satisfies that $N \geq \frac{180d}{\zeta^2}$ as inputs. Roughly speaking, it will find a N -length discrete arm or action sequence whose associated feature matrices maintain the similar property (e.g., G-optimality and E-optimality) as the continuous sample allocation λ .

Formally, $\text{ROUND}(\{(p_i, \mathbf{Q}_i)\}_{i=1}^n, \lambda, \zeta, N)$ returns a discrete sample sequence $s_1, \dots, s_N \in \mathcal{P}^N$ associated with feature matrices $\mathbf{S}_1, \dots, \mathbf{S}_N \in \mathbb{S}_+^d$, which satisfy the following properties:

(i) If λ is an E-optimal design, i.e., λ is the optimal solution of the optimization

$$\min_{\lambda \in \Delta_{\mathcal{Q}}} \left\| \left(\sum_{i=1}^n \lambda(\mathbf{Q}_i) \mathbf{Q}_i \right)^{-1} \right\|,$$

then $\mathbf{S}_1, \dots, \mathbf{S}_N$ satisfy that

$$\left\| \left(\sum_{j=1}^N \mathbf{S}_j \right)^{-1} \right\| \leq (1 + \zeta) \left\| \left(N \sum_{i=1}^n \lambda(\mathbf{Q}_i) \mathbf{Q}_i \right)^{-1} \right\|.$$

(ii) If λ is a G-optimal design, i.e., for a given prediction set $\mathcal{Y} \subseteq \mathbb{R}^d$, λ is the optimal solution of the optimization

$$\min_{\lambda \in \Delta_{\mathcal{Q}}} \max_{\mathbf{y} \in \mathcal{Y}} \|\mathbf{y}\|^2_{\left(\sum_{i=1}^n \lambda(\mathbf{Q}_i) \mathbf{Q}_i \right)^{-1}},$$

then $\mathbf{S}_1, \dots, \mathbf{S}_N$ satisfy that

$$\max_{\mathbf{y} \in \mathcal{Y}} \|\mathbf{y}\|^2_{\left(\sum_{j=1}^N \mathbf{S}_j \right)^{-1}} \leq (1 + \zeta) \max_{\mathbf{y} \in \mathcal{Y}} \|\mathbf{y}\|^2_{\left(N \sum_{i=1}^n \lambda(\mathbf{Q}_i) \mathbf{Q}_i \right)^{-1}}.$$

We implement ROUND by setting $\pi^* = N\lambda$, $k = r = N$ and $\mathbf{x}_i \mathbf{x}_i^\top = \left(\sum_{i=1}^n \pi^*(\mathbf{Q}_i) \mathbf{Q}_i \right)^{-\frac{1}{2}} \mathbf{Q}_i \left(\sum_{i=1}^n \pi^*(\mathbf{Q}_i) \mathbf{Q}_i \right)^{-\frac{1}{2}}$ for any $i \in [n]$ in Algorithm 1 of (Allen-Zhu et al., 2017). Note that Algorithm 1 in (Allen-Zhu et al., 2017) only needs to access the feature matrix $\mathbf{x}_i \mathbf{x}_i^\top$ rather than the separate feature vector \mathbf{x}_i , which allows us to apply it to our problem. We refer interested readers to (Allen-Zhu et al., 2017) and Appendix B in (Fiez et al., 2019) for more implementation details of this rounding procedure.

C. Proofs for Algorithm DouExpDes

In this section, we provide the proofs for Algorithm DouExpDes.

Throughout our proofs, we use L_θ to denote the upper bound of $\|\theta_m\|$ for any $m \in [M]$. Since $\theta_m = \mathbf{B} \mathbf{w}_m$ for any $m \in [M]$, we have that $\|\theta_m\| \leq \|\mathbf{B}\| \|\mathbf{w}_m\| \leq \|\mathbf{w}_m\| \leq L_w$, and thus $L_\theta \leq L_w$.

C.1. Sample Batch Planning

Recall that

$$\lambda^E := \operatorname{argmin}_{\lambda \in \Delta_{\mathcal{X}}} \left\| \left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right\|$$

and

$$\rho^E := \min_{\lambda \in \Delta_{\mathcal{X}}} \left\| \left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right\|$$

are the optimal solution and the optimal value of the E-optimal design optimization, respectively (Line 2 in Algorithm 1). $\bar{\mathbf{x}}_1, \dots, \bar{\mathbf{x}}_p$ is an arm sequence generated according to sample allocation λ^E via rounding procedure ROUND (Line 3 in Algorithm 1).

Let

$$\mathbf{X}_{\text{batch}} := \begin{bmatrix} \bar{\mathbf{x}}_1^\top \\ \vdots \\ \bar{\mathbf{x}}_p^\top \end{bmatrix},$$

and

$$\mathbf{X}_{\text{batch}}^+ := (\mathbf{X}_{\text{batch}}^\top \mathbf{X}_{\text{batch}})^{-1} \mathbf{X}_{\text{batch}}^\top.$$

According to the fact that \mathcal{X} spans \mathbb{R}^d , the definition of E-optimal design and the guarantee of ROUND, we have that $\mathbf{X}_{\text{batch}}^\top \mathbf{X}_{\text{batch}}$ is invertible.

Now, we first give an upper bound of $\|\mathbf{X}_{\text{batch}}^+\|$.

Lemma C.1. *It holds that*

$$\|\mathbf{X}_{\text{batch}}^+\| \leq \sqrt{\frac{(1 + \zeta) \rho^E}{p}}.$$

Proof of Lemma C.1. We have

$$\begin{aligned} \|\mathbf{X}_{\text{batch}}^+\| &= \left\| (\mathbf{X}_{\text{batch}}^\top \mathbf{X}_{\text{batch}})^{-1} \mathbf{X}_{\text{batch}}^\top \right\| \\ &= \sqrt{\left\| (\mathbf{X}_{\text{batch}}^\top \mathbf{X}_{\text{batch}})^{-1} \mathbf{X}_{\text{batch}}^\top \mathbf{X}_{\text{batch}} (\mathbf{X}_{\text{batch}}^\top \mathbf{X}_{\text{batch}})^{-1} \right\|} \\ &= \sqrt{\left\| (\mathbf{X}_{\text{batch}}^\top \mathbf{X}_{\text{batch}})^{-1} \right\|} \\ &= \sqrt{\left\| \left(\sum_{i=1}^p \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \right)^{-1} \right\|} \\ &\leq \sqrt{(1 + \zeta) \left\| \left(p \sum_{i=1}^n \lambda^E(\mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right\|} \\ &= \sqrt{\frac{(1 + \zeta) \rho^E}{p}}. \end{aligned}$$

□

C.2. Global Feature Extractor Recovery

For clarity of notation, we add subscript t to the notations in subroutine `FeatRecover` to denote the quantities generated in phase t . Specifically, we use $\alpha_{t,m,j,i}$, $\tilde{\theta}_{t,m,j}$, \mathbf{Z}_t and $\hat{\mathbf{B}}_t$ to denote the random reward, estimator of reward parameter, estimator of $\frac{1}{M} \sum_{i=1}^M \theta_m \theta_m^\top$ and estimator of feature extractor in phase t , respectively.

For any phase $t > 0$, task $m \in [M]$, round $j \in [T_t]$ and arm $i \in [p]$, let $\eta_{t,m,j,i}$ denote the noise of the sample on arm \bar{x}_i in the j -th round for task m , during the execution of `FeatRecover` in phase t (Line 4 in Algorithm 2). The noise $\eta_{t,m,j,i}$ is zero-mean and sub-Gaussian, and has variance 1. $\eta_{t,m,j,i}$ is independent for different t, m, j, i .

For any phase $t > 0$, task $m \in [M]$, round $j \in [T_t]$, let $\alpha_{t,m,j} := [\alpha_{t,m,j,1}, \dots, \alpha_{t,m,j,p}]^\top$. Then, we have that

$$\tilde{\theta}_{t,m,j} = \mathbf{X}_{\text{batch}}^+ \alpha_{t,m,j},$$

and

$$\mathbf{Z}_t = \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \tilde{\theta}_{t,m,j} (\tilde{\theta}_{t,m,j})^\top - \mathbf{X}_{\text{batch}}^+ (\mathbf{X}_{\text{batch}}^+)^{\top}.$$

Lemma C.2 (Expectation of \mathbf{Z}_t). *It holds that*

$$\mathbb{E}[\mathbf{Z}_t] = \frac{1}{M} \sum_{m=1}^M \theta_m \theta_m^\top.$$

Proof of Lemma C.2. \mathbf{Z}_t can be written as

$$\begin{aligned} \mathbf{Z}_t &= \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \tilde{\theta}_{t,m,j} (\tilde{\theta}_{t,m,j})^\top - \mathbf{X}_{\text{batch}}^+ (\mathbf{X}_{\text{batch}}^+)^{\top} \\ &= \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{X}_{\text{batch}}^+ \begin{bmatrix} \alpha_{t,m,j,1} \\ \vdots \\ \alpha_{t,m,j,p} \end{bmatrix} [\alpha_{t,m,j,1}, \dots, \alpha_{t,m,j,p}]^\top (\mathbf{X}_{\text{batch}}^+)^{\top} - \mathbf{X}_{\text{batch}}^+ (\mathbf{X}_{\text{batch}}^+)^{\top} \\ &= \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{X}_{\text{batch}}^+ \begin{bmatrix} \bar{x}_1^\top \theta_m + \eta_{t,m,j,1} \\ \vdots \\ \bar{x}_p^\top \theta_m + \eta_{t,m,j,p} \end{bmatrix} [\bar{x}_1^\top \theta_m + \eta_{t,m,j,1}, \dots, \bar{x}_p^\top \theta_m + \eta_{t,m,j,p}]^\top (\mathbf{X}_{\text{batch}}^+)^{\top} - \mathbf{X}_{\text{batch}}^+ (\mathbf{X}_{\text{batch}}^+)^{\top} \\ &= \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{X}_{\text{batch}}^+ \begin{bmatrix} (\bar{x}_1^\top \theta_m + \eta_{t,m,j,1})^2 & \cdots & (\bar{x}_1^\top \theta_m + \eta_{t,m,j,1})(\bar{x}_p^\top \theta_m + \eta_{t,m,j,p}) \\ \vdots & \ddots & \vdots \\ (\bar{x}_p^\top \theta_m + \eta_{t,m,j,p})(\bar{x}_1^\top \theta_m + \eta_{t,m,j,1}) & \cdots & (\bar{x}_p^\top \theta_m + \eta_{t,m,j,p})^2 \end{bmatrix} (\mathbf{X}_{\text{batch}}^+)^{\top} \\ &\quad - \mathbf{X}_{\text{batch}}^+ (\mathbf{X}_{\text{batch}}^+)^{\top} \\ &= \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{X}_{\text{batch}}^+ \left(\begin{bmatrix} (\bar{x}_1^\top \theta_m)^2 & \cdots & \bar{x}_1^\top \theta_m \bar{x}_p^\top \theta_m \\ \vdots & \ddots & \vdots \\ \bar{x}_1^\top \theta_m \bar{x}_p^\top \theta_m & \cdots & (\bar{x}_p^\top \theta_m)^2 \end{bmatrix} \right. \\ &\quad + \begin{bmatrix} 2\bar{x}_1^\top \theta_m \eta_{t,m,j,1} & \cdots & \bar{x}_1^\top \theta_m \eta_{t,m,j,p} + \bar{x}_p^\top \theta_m \eta_{t,m,j,1} \\ \vdots & \ddots & \vdots \\ \bar{x}_1^\top \theta_m \eta_{t,m,j,p} + \bar{x}_p^\top \theta_m \eta_{t,m,j,1} & \cdots & 2\bar{x}_p^\top \theta_m \eta_{t,m,j,p} \end{bmatrix} \\ &\quad \left. + \begin{bmatrix} (\eta_{t,m,j,1})^2 & \cdots & \eta_{t,m,j,1} \eta_{t,m,j,p} \\ \vdots & \ddots & \vdots \\ \eta_{t,m,j,1} \eta_{t,m,j,p} & \cdots & (\eta_{t,m,j,p})^2 \end{bmatrix} \right) (\mathbf{X}_{\text{batch}}^+)^{\top} - \mathbf{X}_{\text{batch}}^+ (\mathbf{X}_{\text{batch}}^+)^{\top}. \end{aligned} \tag{3}$$

Then, taking the expectation on \mathbf{Z}_t , we have

$$\mathbb{E}[\mathbf{Z}_t] = \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{X}_{\text{batch}}^+ \left(\begin{bmatrix} (\bar{x}_1^\top \theta_m)^2 & \cdots & \bar{x}_1^\top \theta_m \bar{x}_p^\top \theta_m \\ \vdots & \ddots & \vdots \\ \bar{x}_1^\top \theta_m \bar{x}_p^\top \theta_m & \cdots & (\bar{x}_p^\top \theta_m)^2 \end{bmatrix} + \mathbf{I}_d \right) (\mathbf{X}_{\text{batch}}^+)^{\top} - \mathbf{X}_{\text{batch}}^+ (\mathbf{X}_{\text{batch}}^+)^{\top}$$

$$\begin{aligned}
 &= \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{X}_{\text{batch}}^+ \begin{bmatrix} (\bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m)^2 & \cdots & \bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \\ \vdots & \cdots & \vdots \\ \bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m & \cdots & (\bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m)^2 \end{bmatrix} (\mathbf{X}_{\text{batch}}^+)^{\top} \\
 &= \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{X}_{\text{batch}}^+ \begin{bmatrix} \bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \\ \vdots \\ \bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \end{bmatrix} [\bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m, \dots, \bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m]^\top (\mathbf{X}_{\text{batch}}^+)^{\top} \\
 &= \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{X}_{\text{batch}}^+ \mathbf{X}_{\text{batch}} \boldsymbol{\theta}_m \boldsymbol{\theta}_m^\top \mathbf{X}_{\text{batch}}^\top (\mathbf{X}_{\text{batch}}^+)^{\top} \\
 &= \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \boldsymbol{\theta}_m \boldsymbol{\theta}_m^\top \\
 &= \frac{1}{M} \sum_{m=1}^M \boldsymbol{\theta}_m \boldsymbol{\theta}_m^\top.
 \end{aligned}$$

□

Recall that for any $t > 0$, $\delta_t := \frac{\delta}{2t^2}$.

For any phase $t > 0$, define events

$$\mathcal{E}_t := \left\{ \|\mathbf{Z}_t - \mathbb{E}[\mathbf{Z}_t]\| \leq \frac{96 \|\mathbf{X}_{\text{batch}}^+\|^2 p L_x L_\theta \log\left(\frac{16p}{\delta_t}\right) \log\left(\frac{16pMT_t}{\delta_t}\right)}{\sqrt{MT_t}} \right\},$$

and

$$\mathcal{E} := \cap_{t=1}^{\infty} \mathcal{E}_t.$$

Lemma C.3 (Concentration of \mathbf{Z}_t). *It holds that*

$$\Pr[\mathcal{E}] \geq \frac{\delta}{2}.$$

Proof of Lemma C.3. According to Eq. (3), we have

$$\begin{aligned}
 \mathbf{Z}_t - \mathbb{E}[\mathbf{Z}_t] &= \frac{1}{MT_t} \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{X}_{\text{batch}}^+ \left(\begin{bmatrix} 2\bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \eta_{t,m,j,1} & \cdots & \bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \eta_{t,m,j,p} + \bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \eta_{t,m,j,1} \\ \vdots & \cdots & \vdots \\ \bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \eta_{t,m,j,p} + \bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \eta_{t,m,j,1} & \cdots & 2\bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \eta_{t,m,j,p} \end{bmatrix} \right. \\
 &\quad - \mathbb{E} \begin{bmatrix} 2\bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \eta_{t,m,j,1} & \cdots & \bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \eta_{t,m,j,p} + \bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \eta_{t,m,j,1} \\ \vdots & \cdots & \vdots \\ \bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \eta_{t,m,j,p} + \bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \eta_{t,m,j,1} & \cdots & 2\bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \eta_{t,m,j,p} \end{bmatrix} \\
 &\quad \left. + \begin{bmatrix} (\eta_{t,m,j,1})^2 & \cdots & \eta_{t,m,j,1} \eta_{t,m,j,p} \\ \vdots & \cdots & \vdots \\ \eta_{t,m,j,1} \eta_{t,m,j,p} & \cdots & (\eta_{t,m,j,p})^2 \end{bmatrix} - \mathbb{E} \begin{bmatrix} (\eta_{t,m,j,1})^2 & \cdots & \eta_{t,m,j,1} \eta_{t,m,j,p} \\ \vdots & \cdots & \vdots \\ \eta_{t,m,j,1} \eta_{t,m,j,p} & \cdots & (\eta_{t,m,j,p})^2 \end{bmatrix} \right) (\mathbf{X}_{\text{batch}}^+)^{\top}.
 \end{aligned}$$

Define the following matrices:

$$\begin{aligned}
 \mathbf{A}_{t,m,j} &:= \frac{1}{MT_t} \begin{bmatrix} 2\bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \eta_{t,m,j,1} & \cdots & \bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \eta_{t,m,j,p} + \bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \eta_{t,m,j,1} \\ \vdots & \cdots & \vdots \\ \bar{\mathbf{x}}_1^\top \boldsymbol{\theta}_m \eta_{t,m,j,p} + \bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \eta_{t,m,j,1} & \cdots & 2\bar{\mathbf{x}}_p^\top \boldsymbol{\theta}_m \eta_{t,m,j,p} \end{bmatrix}, \\
 \mathbf{A}_t &:= \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{A}_{t,m,j},
 \end{aligned}$$

$$\begin{aligned} \mathbf{C}_{t,m,j} &:= \frac{1}{MT_t} \begin{bmatrix} (\eta_{t,m,j,1})^2 & \cdots & \eta_{t,m,j,1}\eta_{t,m,j,p} \\ \cdots & \cdots & \cdots \\ \eta_{t,m,j,1}\eta_{t,m,j,p} & \cdots & (\eta_{t,m,j,p})^2 \end{bmatrix}, \\ \mathbf{C}_t &:= \sum_{m=1}^M \sum_{j=1}^{T_t} \mathbf{C}_{t,m,j}. \end{aligned}$$

Then, we can write $\mathbf{Z}_t - \mathbb{E}[\mathbf{Z}_t]$ as

$$\mathbf{Z}_t - \mathbb{E}[\mathbf{Z}_t] = \mathbf{X}_{\text{batch}}^+ (\mathbf{A}_t - \mathbb{E}[\mathbf{A}_t] + \mathbf{C}_t - \mathbb{E}[\mathbf{C}_t]) (\mathbf{X}_{\text{batch}}^+)^{\top},$$

and thus,

$$\|\mathbf{Z}_t - \mathbb{E}[\mathbf{Z}_t]\| \leq \|\mathbf{X}_{\text{batch}}^+\|^2 (\|\mathbf{A}_t - \mathbb{E}[\mathbf{A}_t]\| + \|\mathbf{C}_t - \mathbb{E}[\mathbf{C}_t]\|). \quad (4)$$

Next, we analyze $\|\mathbf{A}_t - \mathbb{E}[\mathbf{A}_t]\|$ and $\|\mathbf{C}_t - \mathbb{E}[\mathbf{C}_t]\|$. In order to use the truncated matrix Bernstein inequality (Lemma E.2), we define the truncated noise and truncated matrices as follows.

Let $R > 0$ be a truncation level of noises, which will be chosen later. For any $t > 0$, $m \in [M]$, $j \in [T_t]$ and $i \in [p]$, let $\tilde{\eta}_{t,m,j,i} = \eta_{t,m,j,i} \mathbb{1}\{|\eta_{t,m,j,i}| \leq R\}$ denote the truncated noise. Then, we define the following truncated matrices:

$$\begin{aligned} \tilde{\mathbf{A}}_{t,m,j} &:= \frac{1}{MT_t} \begin{bmatrix} 2\bar{\mathbf{x}}_1^{\top} \boldsymbol{\theta}_m \tilde{\eta}_{t,m,j,1} & \cdots & \bar{\mathbf{x}}_1^{\top} \boldsymbol{\theta}_m \tilde{\eta}_{t,m,j,p} + \bar{\mathbf{x}}_p^{\top} \boldsymbol{\theta}_m \tilde{\eta}_{t,m,j,1} \\ \cdots & \cdots & \cdots \\ \bar{\mathbf{x}}_1^{\top} \boldsymbol{\theta}_m \tilde{\eta}_{t,m,j,p} + \bar{\mathbf{x}}_p^{\top} \boldsymbol{\theta}_m \tilde{\eta}_{t,m,j,1} & \cdots & 2\bar{\mathbf{x}}_p^{\top} \boldsymbol{\theta}_m \tilde{\eta}_{t,m,j,p} \end{bmatrix} \\ \tilde{\mathbf{A}}_t &:= \sum_{m=1}^M \sum_{j=1}^{T_t} \tilde{\mathbf{A}}_{t,m,j}, \\ \tilde{\mathbf{C}}_{t,m,j} &:= \frac{1}{MT_t} \begin{bmatrix} (\tilde{\eta}_{t,m,j,1})^2 & \cdots & \tilde{\eta}_{t,m,j,1} \tilde{\eta}_{t,m,j,p} \\ \cdots & \cdots & \cdots \\ \tilde{\eta}_{t,m,j,1} \tilde{\eta}_{t,m,j,p} & \cdots & (\tilde{\eta}_{t,m,j,p})^2 \end{bmatrix} \\ \tilde{\mathbf{C}}_t &:= \sum_{m=1}^M \sum_{j=1}^{T_t} \tilde{\mathbf{C}}_{t,m,j} \end{aligned} \quad (5)$$

First, we bound $\|\mathbf{A}_t - \mathbb{E}[\mathbf{A}_t]\|$. Since for any $t > 0$, $m \in [M]$, $j \in [T_t]$ and $i \in [p]$, $|\tilde{\eta}_{t,m,j,i}| \leq R$ and $|\bar{\mathbf{x}}_i^{\top} \boldsymbol{\theta}_m| \leq L_x L_{\theta}$, we have $\|\tilde{\mathbf{A}}_{t,m,j}\| \leq \frac{1}{MT_t} \cdot 2pL_x L_{\theta} R$.

Recall that for any $t > 0$, $m \in [M]$, $j \in [T_t]$ and $i \in [p]$, $\eta_{t,m,j,i}$ is 1-sub-Gaussian. Using a union bound over $i \in [p]$, we have that for any $t > 0$, $m \in [M]$, $j \in [T_t]$, with probability at least $1 - 2p \exp(-\frac{R^2}{2})$, $|\eta_{t,m,j,i}| \leq R$ for all $i \in [p]$. Thus, with probability at least $1 - 2p \exp(-\frac{R^2}{2})$, $\|\mathbf{A}_{t,m,j}\| \leq \frac{1}{MT_t} \cdot 2pL_x L_{\theta} R$.

Then, we have

$$\begin{aligned} \left\| \mathbb{E}[\mathbf{A}_{t,m,j}] - \mathbb{E}[\tilde{\mathbf{A}}_{t,m,j}] \right\| &\leq \left\| \mathbb{E} \left[\mathbf{A}_{t,m,j} \cdot \mathbb{1} \left\{ \|\mathbf{A}_{t,m,j}\| \geq \frac{2pL_x L_{\theta} R}{MT_t} \right\} \right] \right\| \\ &\leq \mathbb{E} \left[\|\mathbf{A}_{t,m,j}\| \cdot \mathbb{1} \left\{ \|\mathbf{A}_{t,m,j}\| \geq \frac{2pL_x L_{\theta} R}{MT_t} \right\} \right] \\ &= \mathbb{E} \left[\frac{2pL_x L_{\theta} R}{MT_t} \cdot \mathbb{1} \left\{ \|\mathbf{A}_{t,m,j}\| \geq \frac{2pL_x L_{\theta} R}{MT_t} \right\} \right] \\ &\quad + \mathbb{E} \left[\left(\|\mathbf{A}_{t,m,j}\| - \frac{2pL_x L_{\theta} R}{MT_t} \right) \cdot \mathbb{1} \left\{ \|\mathbf{A}_{t,m,j}\| \geq \frac{2pL_x L_{\theta} R}{MT_t} \right\} \right] \\ &= \frac{2pL_x L_{\theta} R}{MT_t} \cdot \Pr \left[\|\mathbf{A}_{t,m,j}\| \geq \frac{2pL_x L_{\theta} R}{MT_t} \right] + \int_0^{\infty} \Pr \left[\|\mathbf{A}_{t,m,j}\| - \frac{2pL_x L_{\theta} R}{MT_t} > x \right] dx \end{aligned}$$

$$\begin{aligned}
 &\leq \frac{2pL_xL_\theta R}{MT_t} \cdot 2p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2pL_xL_\theta}{MT_t} \int_R^\infty \Pr\left[\|\mathbf{A}_{t,m,j}\| > \frac{2pL_xL_\theta y}{MT_t}\right] dy \\
 &\leq \frac{2pL_xL_\theta R}{MT_t} \cdot 2p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2pL_xL_\theta}{MT_t} \int_R^\infty 2p \exp\left(-\frac{y^2}{2}\right) dy \\
 &\leq \frac{2pL_xL_\theta R}{MT_t} \cdot 2p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2pL_xL_\theta}{MT_t} \cdot 2p \cdot \frac{1}{R} \cdot \exp\left(-\frac{R^2}{2}\right) \\
 &= \frac{2pL_xL_\theta}{MT_t} \cdot 2p \cdot \left(R + \frac{1}{R}\right) \exp\left(-\frac{R^2}{2}\right).
 \end{aligned}$$

Let $\delta' \in (0, 1)$ be a confidence parameter which will be chosen later. Using the truncated matrix Bernstein inequality (Lemma E.2) with $n = MT_t$, $R = \sqrt{2 \log\left(\frac{2pMT_t}{\delta'}\right)}$, $n \Pr[\|\mathbf{A}_{t,m,j}\| \geq \frac{1}{MT_t} \cdot 2pL_xL_\theta R] \leq \delta'$, $U = \frac{2pL_xL_\theta \sqrt{2 \log\left(\frac{2pMT_t}{\delta'}\right)}}{MT_t}$, $\sigma^2 = MT_t U^2$, $\tau = \frac{4 \cdot 2pL_xL_\theta \sqrt{2 \log\left(\frac{2pMT_t}{\delta'}\right)} \log\left(\frac{2p}{\delta'}\right)}{\sqrt{MT_t}} + \frac{4 \cdot 2pL_xL_\theta \sqrt{2 \log\left(\frac{2pMT_t}{\delta'}\right)} \log\left(\frac{2p}{\delta'}\right)}{MT_t}$ and $\Delta = \frac{2pL_xL_\theta \cdot 2 \sqrt{2 \log\left(\frac{2pMT_t}{\delta'}\right)}}{MT_t} \cdot \frac{\delta'}{MT_t}$, we have that with probability at least $1 - 2\delta'$,

$$\begin{aligned}
 \|\mathbf{A}_t - \mathbb{E}[\mathbf{A}_t]\| &\leq \frac{4 \cdot 2pL_xL_\theta \sqrt{2 \log\left(\frac{2pMT_t}{\delta'}\right)} \log\left(\frac{2p}{\delta'}\right)}{\sqrt{MT_t}} + \frac{4 \cdot 2pL_xL_\theta \sqrt{2 \log\left(\frac{2pMT_t}{\delta'}\right)} \log\left(\frac{2p}{\delta'}\right)}{MT_t} \\
 &\leq \frac{8 \cdot 2pL_xL_\theta \sqrt{2 \log\left(\frac{2pMT_t}{\delta'}\right)} \log\left(\frac{2p}{\delta'}\right)}{\sqrt{MT_t}}.
 \end{aligned} \tag{6}$$

Now we investigate $\|\mathbf{C}_t - \mathbb{E}[\mathbf{C}_t]\|$. Recall that in Eq. (5), for any $t > 0$, $m \in [M]$, $j \in [T_t]$ and $i \in [p]$, $|\tilde{\eta}_{t,m,j,i}| \leq R$. Then, we have $\|\tilde{\mathbf{C}}_{t,m,j}\| \leq \frac{1}{MT_t} \cdot pR^2$.

Recall that for any $t > 0$, $m \in [M]$ and $j \in [T_t]$, with probability at least $1 - 2p \exp(-\frac{R^2}{2})$, $|\eta_{t,m,j,i}| \leq R$ for all $i \in [p]$. Thus, with probability at least $1 - 2p \exp(-\frac{R^2}{2})$, $\|\mathbf{C}_{t,m,j}\| \leq \frac{1}{MT_t} \cdot pR^2$. Then, we have

$$\begin{aligned}
 \|\mathbb{E}[\mathbf{C}_{t,m,j}] - \mathbb{E}[\tilde{\mathbf{C}}_{t,m,j}]\| &\leq \left\| \mathbb{E} \left[\mathbf{C}_{t,m,j} \cdot \mathbb{1} \left\{ \|\mathbf{C}_{t,m,j}\| \geq \frac{pR^2}{MT_t} \right\} \right] \right\| \\
 &\leq \mathbb{E} \left[\|\mathbf{C}_{t,m,j}\| \cdot \mathbb{1} \left\{ \|\mathbf{C}_{t,m,j}\| \geq \frac{pR^2}{MT_t} \right\} \right] \\
 &= \mathbb{E} \left[\frac{pR^2}{MT_t} \cdot \mathbb{1} \left\{ \|\mathbf{C}_{t,m,j}\| \geq \frac{pR^2}{MT_t} \right\} \right] + \mathbb{E} \left[\left(\|\mathbf{C}_{t,m,j}\| - \frac{pR^2}{MT_t} \right) \cdot \mathbb{1} \left\{ \|\mathbf{C}_{t,m,j}\| \geq \frac{pR^2}{MT_t} \right\} \right] \\
 &= \frac{pR^2}{MT_t} \cdot \Pr \left[\|\mathbf{C}_{t,m,j}\| \geq \frac{pR^2}{MT_t} \right] + \int_0^\infty \Pr \left[\|\mathbf{C}_{t,m,j}\| - \frac{pR^2}{MT_t} > x \right] dx \\
 &\leq \frac{pR^2}{MT_t} \cdot 2p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2p}{MT_t} \int_R^\infty \mathbf{y} \cdot \Pr \left[\|\mathbf{C}_{t,m,j}\| > \frac{dy^2}{MT_t} \right] dy \\
 &\leq \frac{pR^2}{MT_t} \cdot 2p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2p}{MT_t} \int_R^\infty \mathbf{y} \cdot 2p \exp\left(-\frac{y^2}{2}\right) dy \\
 &\leq \frac{pR^2}{MT_t} \cdot 2p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2p}{MT_t} \cdot 2p \cdot \exp\left(-\frac{R^2}{2}\right) \\
 &= \frac{p}{MT_t} \cdot 2p \cdot (R^2 + 2) \exp\left(-\frac{R^2}{2}\right).
 \end{aligned}$$

Using the truncated matrix Bernstein inequality (Lemma E.2) with $n = MT_t$, $R = \sqrt{2 \log\left(\frac{2pMT_t}{\delta'}\right)}$, $n \Pr[\|\mathbf{C}_{t,m,j}\| \geq \frac{1}{MT_t} \cdot pR^2] \leq \delta'$, $U = \frac{p \cdot 2 \log\left(\frac{2pMT_t}{\delta'}\right)}{MT_t}$, $\sigma^2 = \frac{32p}{MT_t}$, $\tau = \frac{4 \cdot p \cdot 2 \log\left(\frac{2pMT_t}{\delta'}\right) \log\left(\frac{2p}{\delta'}\right)}{\sqrt{MT_t}} + \frac{4 \cdot p \cdot 2 \log\left(\frac{2pMT_t}{\delta'}\right) \log\left(\frac{2p}{\delta'}\right)}{MT_t}$ and $\Delta =$

$\frac{p \cdot 2 \cdot 2 \log\left(\frac{2pMT_t}{\delta'}\right)}{MT_t} \cdot \frac{\delta'}{MT_t}$, we have that with probability at least $1 - 2\delta'$,

$$\begin{aligned} \|C_t - \mathbb{E}[C_t]\| &\leq \frac{4 \cdot 2p \log\left(\frac{2pMT_t}{\delta'}\right) \log\left(\frac{2p}{\delta'}\right)}{\sqrt{MT_t}} + \frac{4 \cdot 2p \log\left(\frac{2pMT_t}{\delta'}\right) \log\left(\frac{2p}{\delta'}\right)}{MT_t} \\ &\leq \frac{8 \cdot 2p \log\left(\frac{2pMT_t}{\delta'}\right) \log\left(\frac{2p}{\delta'}\right)}{\sqrt{MT_t}} \end{aligned} \quad (7)$$

Plugging Eqs. (6) and (7) into Eq. (4), we have that with probability at least $1 - 4\delta'$,

$$\begin{aligned} \|Z_t - \mathbb{E}[Z_t]\| &\leq \|X_{\text{batch}}^+\|^2 (\|A_t - \mathbb{E}[A_t]\| + \|C_t - \mathbb{E}[C_t]\|) \\ &\leq \|X_{\text{batch}}^+\|^2 \left(\frac{8 \cdot 2pL_xL_\theta \sqrt{2 \log\left(\frac{2pMT_t}{\delta'}\right) \log\left(\frac{2p}{\delta'}\right)}}{\sqrt{MT_t}} + \frac{8 \cdot 2p \log\left(\frac{2pMT_t}{\delta'}\right) \log\left(\frac{2p}{\delta'}\right)}{\sqrt{MT_t}} \right) \\ &\leq \frac{96 \|X_{\text{batch}}^+\|^2 pL_xL_\theta \log\left(\frac{2p}{\delta'}\right) \log\left(\frac{2pMT_t}{\delta'}\right)}{\sqrt{MT_t}}. \end{aligned}$$

Let $\delta' = \frac{\delta_t}{8}$. Then, we obtain that with probability at least $1 - \frac{\delta_t}{2}$,

$$\|Z_t - \mathbb{E}[Z_t]\| \leq \frac{96 \|X_{\text{batch}}^+\|^2 pL_xL_\theta \log\left(\frac{16p}{\delta_t}\right) \log\left(\frac{16pMT_t}{\delta_t}\right)}{\sqrt{MT_t}},$$

which implies that $\Pr[\mathcal{E}_t] \geq 1 - \frac{\delta_t}{2}$.

Taking a union bound over all phases $t \geq 1$ and recalling $\delta_t := \frac{\delta}{2t^2}$, we obtain

$$\begin{aligned} \Pr[\mathcal{E}] &\geq 1 - \sum_{t=1}^{\infty} \Pr[\bar{\mathcal{E}}_t] \\ &\geq 1 - \sum_{t=1}^{\infty} \frac{\delta_t}{2} \\ &= 1 - \sum_{t=1}^{\infty} \frac{\delta}{4t^2} \\ &\geq 1 - \frac{\delta}{2}. \end{aligned}$$

□

For any matrix $A \in \mathbb{R}^{m \times n}$ with $m \geq n$, let $\sigma_{\max}(A)$ and $\sigma_{\min}(A)$ denote the maximum and minimum singular values of A , respectively. For any $i \in [m]$, let $\sigma_i(A)$ denote the i -th singular value of A .

For any matrix $A \in \mathbb{R}^{m \times n}$ with $m \geq n$, let A_\perp denote the orthogonal complement matrix of A , where the columns of A_\perp are the orthogonal complement of those of A . Then, it holds that $AA^\top + A_\perp A_\perp^\top = I_m$, where I_m is the $m \times m$ identity matrix.

According to Assumption 3.1, there exists an absolute constant c_0 which satisfies that $\sigma_{\min}\left(\frac{1}{M} \sum_{m=1}^M \mathbf{w}_m \mathbf{w}_m^\top\right) = \sigma_{\min}\left(\frac{1}{M} \sum_{m=1}^M \boldsymbol{\theta}_m \boldsymbol{\theta}_m^\top\right) \geq \frac{c_0}{k}$.

Lemma C.4 (Concentration of \hat{B}_t). *Suppose that event \mathcal{E} holds. Then, for any phase $t > 0$,*

$$\left\| \hat{B}_{t,\perp}^\top B \right\| \leq \frac{192 \|X_{\text{batch}}^+\|^2 kpL_xL_\theta \log\left(\frac{16p}{\delta_t}\right) \log\left(\frac{16pMT_t}{\delta_t}\right)}{\sqrt{MT_t}}.$$

Furthermore, for any phase $t > 0$, if

$$T_t = \left\lceil \frac{68 \cdot 192^2 \cdot 8^2 (1 + \zeta)^3 (\rho^E)^2 k^4 L_x^4 L_\theta^2 L_w^2}{c_0^2 M} \cdot \max \left\{ 2^{2t}, \frac{L_x^4}{\omega^2} \right\} \cdot \log^2 \left(\frac{16p}{\delta_t} \right) \right. \\ \left. \log^2 \left(\frac{192 \cdot 16 \cdot 8 (1 + \zeta)^{\frac{3}{2}} \rho^E k^2 p L_x^2 L_\theta L_w}{c_0} \cdot \max \left\{ 2^t, \frac{L_x^2}{\omega} \right\} \cdot \frac{1}{\delta_t} \cdot \log \left(\frac{16p}{\delta_t} \right) \right) \right\rceil, \quad (8)$$

then

$$\|\hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B}\| \leq \min \left\{ \frac{1}{8kL_x L_w \cdot 2^t \sqrt{1 + \zeta}}, \frac{\omega}{6L_x^2} \right\}.$$

Proof of Lemma C.4. From Assumption 3.1, $\sigma_k(\mathbb{E}[\mathbf{Z}_t]) - \sigma_{k+1}(\mathbb{E}[\mathbf{Z}_t]) = \sigma_{\min}(\frac{1}{M} \sum_{m=1}^M \boldsymbol{\theta}_m \boldsymbol{\theta}_m^\top) \geq \frac{c_0}{k}$. Using the Davis-Kahan sin θ Theorem (Bhatia, 2013) and letting T_t be large enough to satisfy $\|\mathbf{Z}_t - \mathbb{E}[\mathbf{Z}_t]\| \leq \frac{c_0}{2k}$, we have

$$\|\hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B}\| \leq \frac{\|\mathbf{Z}_t - \mathbb{E}[\mathbf{Z}_t]\|}{\sigma_k(\mathbb{E}[\mathbf{Z}_t]) - \sigma_{k+1}(\mathbb{E}[\mathbf{Z}_t]) - \|\mathbf{Z}_t - \mathbb{E}[\mathbf{Z}_t]\|} \\ \leq \frac{2k}{c_0} \|\mathbf{Z}_t - \mathbb{E}[\mathbf{Z}_t]\| \\ \stackrel{(a)}{\leq} \frac{192 \|\mathbf{X}_{\text{batch}}^+\|^2 k p L_x L_\theta \log \left(\frac{16p}{\delta_t} \right)}{c_0 \sqrt{MT_t}} \log \left(\frac{16pMT_t}{\delta_t} \right),$$

where inequality (a) uses the definition of event \mathcal{E} .

Using Lemma E.3 with $A = \frac{192 \|\mathbf{X}_{\text{batch}}^+\|^2 k p L_x L_\theta}{c_0} \log \left(\frac{16p}{\delta_t} \right)$, $B = \frac{16p}{\delta_t}$ and $\kappa = \min \left\{ \frac{1}{8kL_x L_w \cdot 2^t \sqrt{1 + \zeta}}, \frac{\omega}{6L_x^2} \right\}$, we have that if

$$MT_t \geq 68 \left(\frac{192 \|\mathbf{X}_{\text{batch}}^+\|^2 k p L_x L_\theta}{c_0} \log \left(\frac{16p}{\delta_t} \right) \right)^2 \cdot \max \left\{ \left(8kL_x L_w \cdot 2^t \sqrt{1 + \zeta} \right)^2, \frac{6^2 L_x^4}{\omega^2} \right\} \\ \log^2 \left(\frac{192 \|\mathbf{X}_{\text{batch}}^+\|^2 k p L_x L_\theta}{c_0} \log \left(\frac{16p}{\delta_t} \right) \cdot \frac{16p}{\delta_t} \cdot \max \left\{ 8kL_x L_w \cdot 2^t \sqrt{1 + \zeta}, \frac{6L_x^2}{\omega} \right\} \right),$$

$$\text{then } \|\hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B}\| \leq \min \left\{ \frac{1}{8kL_x L_w \cdot 2^t \sqrt{1 + \zeta}}, \frac{\omega}{6L_x^2} \right\}.$$

According to Lemma C.1, we have $\|\mathbf{X}_{\text{batch}}^+\| \leq \sqrt{\frac{(1+\zeta)\rho^E}{p}}$.

Then, further enlarging MT_t , we have that if

$$MT_t \geq \frac{68 \cdot 192^2 \cdot 8^2 (1 + \zeta)^3 (\rho^E)^2 k^4 L_x^4 L_\theta^2 L_w^2}{c_0^2} \cdot \max \left\{ 2^{2t}, \frac{L_x^4}{\omega^2} \right\} \cdot \log^2 \left(\frac{16p}{\delta_t} \right) \\ \log^2 \left(\frac{192 \cdot 16 \cdot 8 (1 + \zeta)^{\frac{3}{2}} \rho^E k^2 p L_x^2 L_\theta L_w}{c_0} \cdot \max \left\{ 2^t, \frac{L_x^2}{\omega} \right\} \cdot \frac{1}{\delta_t} \cdot \log \left(\frac{16p}{\delta_t} \right) \right),$$

then

$$\|\hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B}\| \leq \min \left\{ \frac{1}{8kL_x L_w \cdot 2^t \sqrt{1 + \zeta}}, \frac{\omega}{6L_x^2} \right\}.$$

□

C.3. Elimination with Low-dimensional Representations

For clarity of notation, we also add subscript t to the notations in subroutine `EliLowRep` to denote the quantities generated in phase t . Specifically, we use the notations $\hat{\mathbf{B}}_t$, $\hat{\mathcal{X}}_{t,m}$, $\lambda_{t,m}^G$, $\rho_{t,m}^G$, $N_{t,m}$, $\{\mathbf{z}_{t,m,i}\}_{i \in [N_{t,m}]}$, $\{r_{t,m,i}\}_{i \in [N_{t,m}]}$, $\hat{\mathbf{w}}_{t,m}$ and $\hat{\boldsymbol{\theta}}_{t,m}$ to denote the corresponding quantities used in `EliLowRep` in phase t .

Before analyzing the sample complexity of `EliLowRep`, we first prove that there exists a sample allocation $\boldsymbol{\lambda} \in \Delta_{\mathcal{X}}$ such that $\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t$ is invertible, i.e., the G-optimal design optimization with $\hat{\mathbf{B}}_t$ is non-vacuous (Line 2 in Algorithm 3).

For any task $m \in [M]$, let

$$\boldsymbol{\lambda}_m^* := \operatorname{argmin}_{\boldsymbol{\lambda} \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_m^*\}} \frac{\|\mathbf{B}^\top \mathbf{x}_m^* - \mathbf{B}^\top \mathbf{x}\|^2_{(\sum_{i=1}^n \lambda(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B})^{-1}}}{((\mathbf{x}_m^* - \mathbf{x})^\top \boldsymbol{\theta}_m)^2}.$$

$\boldsymbol{\lambda}_m^*$ is the optimal solution of the G-optimal design optimization with true feature extractor \mathbf{B} .

Lemma C.5. *For any phase $t > 0$ and task $m \in [M]$, if $\|\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp\| \leq \frac{\omega}{6L_x^2}$, we have*

$$\sigma_{\min} \left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right) > 0.$$

Proof of Lemma C.5. For any task $m \in [M]$, let $\mathbf{A}_m := \sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{x}_i \mathbf{x}_i^\top$. Then, for any phase $t > 0$ and task $m \in [M]$, we have

$$\begin{aligned} \sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t &= \hat{\mathbf{B}}_t^\top \mathbf{A}_m \hat{\mathbf{B}}_t \\ &= \hat{\mathbf{B}}_t^\top (\mathbf{B} \mathbf{B}^\top + \mathbf{B}_\perp \mathbf{B}_\perp^\top) \mathbf{A}_m (\mathbf{B} \mathbf{B}^\top + \mathbf{B}_\perp \mathbf{B}_\perp^\top) \hat{\mathbf{B}}_t \\ &= \hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{A}_m \mathbf{B} \mathbf{B}^\top \hat{\mathbf{B}}_t + \hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{A}_m \mathbf{B}_\perp \mathbf{B}_\perp^\top \hat{\mathbf{B}}_t \\ &\quad + \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{A}_m \mathbf{B} \mathbf{B}^\top \hat{\mathbf{B}}_t + \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{A}_m \mathbf{B}_\perp \mathbf{B}_\perp^\top \hat{\mathbf{B}}_t. \end{aligned}$$

Hence, we have

$$\begin{aligned} \sigma_{\min} \left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right) &\geq \sigma_{\min} \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{A}_m \mathbf{B} \mathbf{B}^\top \hat{\mathbf{B}}_t \right) - \sigma_{\max} \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{A}_m \mathbf{B}_\perp \mathbf{B}_\perp^\top \hat{\mathbf{B}}_t \right) \\ &\quad - \sigma_{\max} \left(\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{A}_m \mathbf{B} \mathbf{B}^\top \hat{\mathbf{B}}_t \right) - \sigma_{\max} \left(\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{A}_m \mathbf{B}_\perp \mathbf{B}_\perp^\top \hat{\mathbf{B}}_t \right) \\ &\geq \sigma_{\min} \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \right) \sigma_{\min} \left(\mathbf{B}^\top \mathbf{A}_m \mathbf{B} \right) \sigma_{\min} \left(\mathbf{B}^\top \hat{\mathbf{B}}_t \right) - \left\| \mathbf{B}_\perp^\top \hat{\mathbf{B}}_t \right\| \left\| \mathbf{A}_m \right\| \\ &\quad - \left\| \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \right\| \left\| \mathbf{A}_m \right\| - \left\| \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \right\| \left\| \mathbf{A}_m \right\| \\ &\geq \sigma_{\min}^2 \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \right) \sigma_{\min} \left(\mathbf{B}^\top \mathbf{A}_m \mathbf{B} \right) - 3 \left\| \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \right\| L_x^2 \\ &\stackrel{(a)}{\geq} \left(1 - \left\| \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \right\|^2 \right) \omega - 3 \left\| \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \right\| L_x^2, \end{aligned}$$

where inequality (a) uses the fact that $\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \hat{\mathbf{B}}_t + \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \hat{\mathbf{B}}_t = \hat{\mathbf{B}}_t^\top (\mathbf{B} \mathbf{B}^\top + \mathbf{B}_\perp \mathbf{B}_\perp^\top) \hat{\mathbf{B}}_t = \hat{\mathbf{B}}_t^\top \hat{\mathbf{B}}_t = \mathbf{I}_k$, and thus, $\sigma_{\min}^2(\hat{\mathbf{B}}_t^\top \mathbf{B}) = 1 - \left\| \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \right\|^2$.

Let $\left\| \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \right\| \leq \frac{\omega}{6L_x^2}$. Then, we have

$$\begin{aligned} \sigma_{\min} \left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right) &\geq \left(1 - \frac{\omega^2}{36L_x^4} \right) \omega - \frac{\omega}{2} \\ &= \frac{\omega}{2} - \frac{\omega^3}{36L_x^4} \end{aligned}$$

> 0 ,

where the last inequality is due to $\omega \leq L_x^2 < \sqrt{18}L_x^2$. \square

Next, we bound the optimal value $\rho_{t,m}^G$ of the G-optimal design optimization with the estimated feature extractor \hat{B}_t .

For any $\mathcal{Z} \subseteq \mathcal{X}$, let $\mathcal{Y}(\mathcal{Z}) := \{\mathbf{x} - \mathbf{x}' : \forall \mathbf{x}, \mathbf{x}' \in \mathcal{Z}, \mathbf{x} \neq \mathbf{x}'\}$. Recall that in Line 2 of Algorithm 3, for any phase $t > 0$ and task $m \in [M]$,

$$\rho_{t,m}^G := \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{y} \in \mathcal{Y}(\hat{\mathcal{X}}_{t,m})} \|\hat{B}_t^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2.$$

Lemma C.6. For any phase $t > 0$ and task $m \in [M]$,

$$\rho_{t,m}^G \leq 4k.$$

Proof of Lemma C.6. For any phase $t > 0$ and task $m \in [M]$, we have that $\hat{\mathcal{X}}_{t,m} \subseteq \mathcal{X}$ and $\mathcal{Y}(\hat{\mathcal{X}}_{t,m}) \subseteq \mathcal{Y}(\mathcal{X})$.

For any fixed $\lambda \in \Delta_{\mathcal{X}}$,

$$\begin{aligned} \max_{\mathbf{y} \in \mathcal{Y}(\hat{\mathcal{X}}_{t,m})} \|\hat{B}_t^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2 &\leq \max_{\mathbf{y} \in \mathcal{Y}(\mathcal{X})} \|\hat{B}_t^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2 \\ &= \|\hat{B}_t^\top (\mathbf{x}'_1 - \mathbf{x}'_2)\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2 \\ &\leq \left(\|\hat{B}_t^\top \mathbf{x}'_1\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}} + \|\hat{B}_t^\top \mathbf{x}'_2\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}} \right)^2 \\ &\leq 2\|\hat{B}_t^\top \mathbf{x}'_1\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2 + 2\|\hat{B}_t^\top \mathbf{x}'_2\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2 \\ &\leq 4 \max_{\mathbf{x} \in \mathcal{X}} \|\hat{B}_t^\top \mathbf{x}\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2, \end{aligned}$$

where \mathbf{x}'_1 and \mathbf{x}'_2 are the arms which satisfy that $\mathbf{y} = \mathbf{x}'_1 - \mathbf{x}'_2$ achieves the maximum value $\max_{\mathbf{y} \in \mathcal{Y}(\mathcal{X})} \|\hat{B}_t^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2$.

Since $\hat{B}_t^\top \mathbf{x} \in \mathbb{R}^k$, according to the Equivalence Theorem in (Kiefer & Wolfowitz, 1960), we have

$$\min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X}} \|\hat{B}_t^\top \mathbf{x}\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2 = k.$$

Therefore, we have

$$\begin{aligned} 4k &= 4 \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X}} \|\hat{B}_t^\top \mathbf{x}\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2 \\ &= 4 \max_{\mathbf{x} \in \mathcal{X}} \|\hat{B}_t^\top \mathbf{x}\|_{\left(\sum_{i=1}^n \lambda'(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2 \\ &\geq \max_{\mathbf{y} \in \mathcal{Y}(\hat{\mathcal{X}}_{t,m})} \|\hat{B}_t^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda'(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2 \\ &\geq \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{y} \in \mathcal{Y}(\hat{\mathcal{X}}_{t,m})} \|\hat{B}_t^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2 \\ &= \rho_{t,m}^G, \end{aligned}$$

where $\lambda' := \operatorname{argmin}_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X}} \|\hat{B}_t^\top \mathbf{x}\|_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{B}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{B}_t\right)^{-1}}^2$. \square

Now we analyze the estimation error of the estimated reward parameter $\hat{\theta}_{t,m} = \hat{B}_t \hat{w}_{t,m}$ in `EliLowRep`.

For any phase $t > 0$, task $m \in [M]$ and arm $j \in [N_{t,m}]$, let $\xi_{t,m,j}$ denote the noise of the sample on arm $\mathbf{z}_{t,m,j}$ for task m , during the execution of `EliLowRep` in phase t (Line 5 in Algorithm 3).

For any phase $t > 0$, define events

$$\begin{aligned} \mathcal{F}_t &:= \left\{ \mathbf{y}^\top \hat{\mathbf{B}}_t \left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1} \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \xi_{t,m,j} \right. \\ &\quad \left. \leq \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|_{\left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1}} \sqrt{2 \log \left(\frac{4n^2 M}{\delta_t} \right)}, \forall m \in [M], \forall \mathbf{y} \in \mathcal{Y}(\hat{\mathcal{X}}_{t,m}) \right\}, \end{aligned} \quad (9)$$

and

$$\mathcal{F} := \cap_{t=1}^{\infty} \mathcal{F}_t.$$

Lemma C.7 (Concentration of the Variance Term). *It holds that*

$$\Pr[\mathcal{F}] \geq 1 - \frac{\delta}{2}.$$

Proof of Lemma C.7. Let $\Sigma_{t,m} := \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t$. Then, we can write

$$\mathbf{y}^\top \hat{\mathbf{B}}_t \left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1} \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \xi_{t,m,j} = \sum_{j=1}^{N_{t,m}} \mathbf{y}^\top \hat{\mathbf{B}}_t \Sigma_{t,m}^{-1} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \xi_{t,m,j}.$$

For any phase $t > 0$, task $m \in [M]$ and arm $j \in [N_{t,m}]$, $\hat{\mathbf{B}}_t$, $\Sigma_{t,m}$ and $\{\mathbf{z}_{t,m,j}\}_{j=1}^{N_{t,m}}$ are fixed before the sampling in **EliLowRep**, and the noise $\xi_{t,m,j}$ is 1-sub-Gaussian (Line 5 in Algorithm 3). Thus, we have that for any $t > 0$, $m \in [M]$ and $j \in [N_{t,m}]$, $\mathbf{y}^\top \hat{\mathbf{B}}_t \Sigma_{t,m}^{-1} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \xi_{t,m,j}$ is $(\mathbf{y}^\top \hat{\mathbf{B}}_t \Sigma_{t,m}^{-1} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j})$ -sub-Gaussian.

Using Hoeffding's inequality and taking a union bound over all $m \in [M]$ and $\mathbf{y} \in \mathcal{Y}(\hat{\mathcal{X}}_{t,m})$, we have that with probability at least $1 - \frac{\delta_t}{2}$,

$$\begin{aligned} &\sum_{j=1}^{N_{t,m}} \mathbf{y}^\top \hat{\mathbf{B}}_t \Sigma_{t,m}^{-1} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \xi_{t,m,j} \\ &\leq \sqrt{2 \sum_{j=1}^{N_{t,m}} \left(\mathbf{y}^\top \hat{\mathbf{B}}_t \Sigma_{t,m}^{-1} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \right)^2 \cdot \log \left(\frac{4n^2 M}{\delta_t} \right)} \\ &= \sqrt{2 \sum_{j=1}^{N_{t,m}} \mathbf{y}^\top \hat{\mathbf{B}}_t \Sigma_{t,m}^{-1} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \Sigma_{t,m}^{-1} \hat{\mathbf{B}}_t^\top \mathbf{y} \cdot \log \left(\frac{4n^2 M}{\delta_t} \right)} \\ &= \sqrt{2 \mathbf{y}^\top \hat{\mathbf{B}}_t \Sigma_{t,m}^{-1} \left(\hat{\mathbf{B}}_t^\top \sum_{j=1}^{N_{t,m}} \mathbf{z}_{t,m,j} \cdot \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right) \Sigma_{t,m}^{-1} \hat{\mathbf{B}}_t^\top \mathbf{y} \cdot \log \left(\frac{4n^2 M}{\delta_t} \right)} \\ &= \sqrt{2 \mathbf{y}^\top \hat{\mathbf{B}}_t \Sigma_{t,m}^{-1} \hat{\mathbf{B}}_t^\top \mathbf{y} \cdot \log \left(\frac{4n^2 M}{\delta_t} \right)} \\ &= \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|_{\Sigma_{t,m}^{-1}} \sqrt{2 \log \left(\frac{4n^2 M}{\delta_t} \right)}, \end{aligned}$$

which implies that

$$\Pr[\mathcal{F}_t] \geq 1 - \frac{\delta_t}{2}.$$

Taking a union bound over all phases $t \geq 1$ and recalling $\delta_t := \frac{\delta}{2t^2}$, we obtain

$$\begin{aligned} \Pr[\mathcal{F}] &\geq 1 - \sum_{t=1}^{\infty} \Pr[\bar{\mathcal{F}}_t] \\ &\geq 1 - \sum_{t=1}^{\infty} \frac{\delta_t}{2} \\ &= 1 - \sum_{t=1}^{\infty} \frac{\delta}{4t^2} \\ &\geq 1 - \frac{\delta}{2}. \end{aligned}$$

□

Lemma C.8 (Concentration of $\hat{\theta}_{t,m}$). *Suppose that event $\mathcal{E} \cap \mathcal{F}$ holds. Then, for any phase $t > 0$, task $m \in [M]$ and $\mathbf{y} \in \mathcal{Y}(\mathcal{X}_{t,m})$,*

$$\left| \mathbf{y}^\top (\hat{\theta}_{t,m} - \theta_m) \right| \leq \frac{1}{2t}.$$

Proof of Lemma C.8. For any phase $t > 0$, task $m \in [M]$ and $\mathbf{y} \in \mathcal{Y}(\mathcal{X}_{t,m})$,

$$\begin{aligned} \mathbf{y}^\top (\hat{\theta}_{t,m} - \theta_m) &= \mathbf{y}^\top \hat{\mathbf{B}}_t \hat{\mathbf{w}}_{t,m} - \mathbf{y}^\top (\hat{\mathbf{B}}_t \hat{\mathbf{B}}_t^\top + \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top) \theta_m \\ &= \mathbf{y}^\top \hat{\mathbf{B}}_t (\hat{\mathbf{w}}_{t,m} - \hat{\mathbf{B}}_t^\top \theta_m) - \mathbf{y}^\top \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top \theta_m. \end{aligned} \quad (10)$$

Here, $\hat{\mathbf{w}}_{t,m}$ can be written as

$$\begin{aligned} \hat{\mathbf{w}}_{t,m} &= \left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1} \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot r_{t,m,j} \\ &= \left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1} \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot (\mathbf{z}_{t,m,j}^\top \theta_m + \xi_{t,m,j}) \\ &= \left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1} \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot (\mathbf{z}_{t,m,j}^\top (\hat{\mathbf{B}}_t \hat{\mathbf{B}}_t^\top + \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top) \theta_m + \xi_{t,m,j}) \\ &= \hat{\mathbf{B}}_t^\top \theta_m + \left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1} \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top \theta_m \\ &\quad + \left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1} \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \xi_{t,m,j}. \end{aligned} \quad (11)$$

Plugging Eq. (11) into Eq. (10), we can decompose the estimation error of $\hat{\theta}_{t,m}$ in EliLowRep into three parts as

$$\mathbf{y}^\top (\hat{\theta}_{t,m} - \theta_m) = \underbrace{\mathbf{y}^\top \hat{\mathbf{B}}_t \left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1} \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top \theta_m}_{\text{Bias}} + \mathbf{y}^\top \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top \theta_m + \mathbf{y}^\top \hat{\mathbf{B}}_t \left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1} \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \xi_{t,m,j}$$

$$\begin{aligned}
 & + \underbrace{\mathbf{y}^\top \hat{\mathbf{B}}_t \left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1} \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \xi_{t,m,j}}_{\text{Variance}} - \underbrace{\mathbf{y}^\top \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \mathbf{w}_m}_{\text{Estimation error of } \hat{\mathbf{B}}_t}.
 \end{aligned}$$

Taking the absolute value on both sides, and using the Cauchy–Schwarz inequality and definition of event \mathcal{F} (Eq. (9)), we have

$$\begin{aligned}
 & \left| \mathbf{y}^\top \hat{\boldsymbol{\theta}}_{t,m} - \mathbf{y}^\top \boldsymbol{\theta}_m \right| \\
 & \leq \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|_{\left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1}} \cdot \left\| \sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \cdot \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \mathbf{w}_m \right\|_{\left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1}} \\
 & \quad + \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|_{\left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1}} \sqrt{2 \log \left(\frac{4n^2 M}{\delta_t} \right)} + \left| \mathbf{y}^\top \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \mathbf{w}_m \right| \\
 & \stackrel{(a)}{\leq} \frac{\sqrt{1+\zeta} \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|_{\left(\sum_{i=1}^n \lambda_{t,m}^G(\mathbf{x}_i) \cdot \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right)^{-1}}}{\sqrt{N_{t,m}}} \cdot L_x L_w \left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| \cdot \sum_{j=1}^{N_{t,m}} \left\| \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \right\|_{\left(\sum_{j=1}^{N_{t,m}} \hat{\mathbf{B}}_t^\top \mathbf{z}_{t,m,j} \mathbf{z}_{t,m,j}^\top \hat{\mathbf{B}}_t \right)^{-1}} \\
 & \quad + \frac{\sqrt{1+\zeta} \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|_{\left(\sum_{i=1}^n \lambda_{t,m}^G(\mathbf{x}_i) \cdot \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right)^{-1}}}{\sqrt{N_{t,m}}} \cdot \sqrt{2 \log \left(\frac{4n^2 M}{\delta_t} \right)} + 2L_x L_w \left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| \\
 & \stackrel{(b)}{\leq} \frac{\sqrt{1+\zeta} \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|_{\left(\sum_{i=1}^n \lambda_{t,m}^G(\mathbf{x}_i) \cdot \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right)^{-1}}}{\sqrt{N_{t,m}}} \cdot L_x L_w \left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| \cdot \sqrt{k N_{t,m}} \\
 & \quad + \frac{\sqrt{1+\zeta} \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|_{\left(\sum_{i=1}^n \lambda_{t,m}^G(\mathbf{x}_i) \cdot \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right)^{-1}}}{\sqrt{N_{t,m}}} \cdot \sqrt{2 \log \left(\frac{4n^2 M}{\delta_t} \right)} + 2L_x L_w \left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| \\
 & \leq \sqrt{1+\zeta} \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|_{\left(\sum_{i=1}^n \lambda_{t,m}^G(\mathbf{x}_i) \cdot \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right)^{-1}} \cdot L_x L_w \left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| \cdot \sqrt{k} \\
 & \quad + \frac{\sqrt{1+\zeta} \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|_{\left(\sum_{i=1}^n \lambda_{t,m}^G(\mathbf{x}_i) \cdot \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right)^{-1}}}{\sqrt{N_{t,m}}} \cdot \sqrt{2 \log \left(\frac{4n^2 M}{\delta_t} \right)} + 2L_x L_w \left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| \\
 & \leq \sqrt{(1+\zeta) \cdot k \cdot \rho_{t,m}^G} \cdot L_x L_w \left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| + \frac{\sqrt{(1+\zeta) \cdot \rho_{t,m}^G \cdot 2 \log \left(\frac{4n^2 M}{\delta_t} \right)}}{\sqrt{N_{t,m}}} + 2L_x L_w \left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| \\
 & \stackrel{(c)}{\leq} \sqrt{(1+\zeta) \cdot 4k^2} \cdot L_x L_w \left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| + \frac{\sqrt{(1+\zeta) \cdot \rho_{t,m}^G \cdot 2 \log \left(\frac{4n^2 M}{\delta_t} \right)}}{\sqrt{N_{t,m}}} + 2L_x L_w \left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| \\
 & \stackrel{(d)}{\leq} \sqrt{(1+\zeta) \cdot 4k^2} \cdot L_x L_w \cdot \frac{1}{8k L_x L_w \cdot 2^t \sqrt{1+\zeta}} + \frac{1}{4 \cdot 2^t} + 2L_x L_w \cdot \frac{1}{8k L_x L_w \cdot 2^t \sqrt{1+\zeta}} \\
 & \leq \frac{1}{4 \cdot 2^t} + \frac{1}{4 \cdot 2^t} + \frac{1}{4 \cdot 2^t} \\
 & \leq \frac{1}{2^t}.
 \end{aligned}$$

Here inequality (a) is due to the guarantee of rounding procedure ROUND and the triangle inequality. Inequality (b) uses Lemma E.5, and inequality (c) follows from Lemma C.6. Inequality (d) comes from Lemma C.4 and $N_{t,m} := \max\{\lceil 32 \cdot 2^{2t} (1+\zeta) \rho_{t,m}^G \log(\frac{4n^2 M}{\delta_t}) \rceil, \frac{180k}{\zeta^2}\}$. \square

For any task $m \in [M]$ and arm $\mathbf{x} \in \mathcal{X}$, let $\Delta_m(\mathbf{x}) := (\mathbf{x}_m^* - \mathbf{x})^\top \boldsymbol{\theta}_m$ denote the reward gap between the optimal arm \mathbf{x}_m^*

and arm \mathbf{x} in task m . For any phase $t > 0$ and task $m \in [M]$, let $\mathcal{Z}_{t,m} := \{\mathbf{x} \in \mathcal{X} : \Delta_m(\mathbf{x}) \leq 4 \cdot 2^{-t}\}$.

Lemma C.9. Suppose that event $\mathcal{E} \cap \mathcal{F}$ holds. For any phase $t > 0$ and task $m \in [M]$,

$$\mathbf{x}_m^* \in \hat{\mathcal{X}}_{t,m},$$

and for any phase $t \geq 2$ and task $m \in [M]$,

$$\hat{\mathcal{X}}_{t,m} \subseteq \mathcal{Z}_{t,m}.$$

Proof of Lemma C.9. This proof follows a similar analytical procedure as that of Lemma 2 in (Fiez et al., 2019).

First, we prove $\mathbf{x}_m^* \in \hat{\mathcal{X}}_{t,m}$ for any phase $t > 0$ and task $m \in [M]$ by contradiction.

Suppose that for some $t > 0$ and some $m \in [M]$, \mathbf{x}_m^* is eliminated from $\hat{\mathcal{X}}_{t,m}$ in phase t . Then, we have that there exists some $\mathbf{x}' \in \hat{\mathcal{X}}_{t,m}$ such that

$$(\mathbf{x}' - \mathbf{x}_m^*)^\top \hat{\boldsymbol{\theta}}_{t,m} > 2^{-t}.$$

Then, we have

$$\begin{aligned} (\mathbf{x}' - \mathbf{x}_m^*)^\top \boldsymbol{\theta}_m &= (\mathbf{x}' - \mathbf{x}_m^*)^\top \hat{\boldsymbol{\theta}}_{t,m} - (\mathbf{x}' - \mathbf{x}_m^*)^\top (\hat{\boldsymbol{\theta}}_{t,m} - \boldsymbol{\theta}_m) \\ &\geq (\mathbf{x}' - \mathbf{x}_m^*)^\top \hat{\boldsymbol{\theta}}_{t,m} - 2^{-t} \\ &> 2^{-t} - 2^{-t} \\ &= 0, \end{aligned}$$

which contradicts the definition of \mathbf{x}_m^* . Thus, we obtain that $\mathbf{x}_m^* \in \hat{\mathcal{X}}_{t,m}$ for any phase $t > 0$ and task $m \in [M]$.

Next, we prove $\hat{\mathcal{X}}_{t,m} \subseteq \mathcal{Z}_{t,m}$ for any phase $t \geq 2$ and task $m \in [M]$, i.e., each $\mathbf{x} \in \hat{\mathcal{X}}_{t,m}$ satisfies that $\Delta_m(\mathbf{x}) \leq 4 \cdot 2^{-t}$.

Suppose that there exists some phase t , some task m and some $\mathbf{x} \in \hat{\mathcal{X}}_{t,m}$ such that $\Delta_m(\mathbf{x}) > 4 \cdot 2^{-t}$. Then, in phase $t - 1 \geq 1$, we have

$$\begin{aligned} (\mathbf{x}_m^* - \mathbf{x})^\top \hat{\boldsymbol{\theta}}_{t-1,m} &= (\mathbf{x}_m^* - \mathbf{x})^\top \boldsymbol{\theta}_m - (\mathbf{x}_m^* - \mathbf{x})^\top (\boldsymbol{\theta}_m - \hat{\boldsymbol{\theta}}_{t-1,m}) \\ &\geq (\mathbf{x}_m^* - \mathbf{x})^\top \boldsymbol{\theta}_m - 2^{-(t-1)} \\ &> 4 \cdot 2^{-t} - 2^{-(t-1)} \\ &= 2^{-(t-1)}, \end{aligned}$$

which implies that \mathbf{x} should have been eliminated from $\hat{\mathcal{X}}_{t,m}$ in phase $t - 1$, and contradicts our supposition. Thus, we complete the proof. \square

C.4. Proof of Theorem 4.1

Before proving Theorem 4.1, we first introduce a useful lemma.

For any task $m \in [M]$, let

$$\lambda_m^* := \operatorname{argmin}_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_m^*\}} \frac{\|B^\top \mathbf{x}_m^* - B^\top \mathbf{x}\|^2_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) B^\top \mathbf{x}_i \mathbf{x}_i^\top B\right)^{-1}}}{((\mathbf{x}_m^* - \mathbf{x})^\top \boldsymbol{\theta}_m)^2},$$

and

$$\rho_m^* := \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_m^*\}} \frac{\|B^\top \mathbf{x}_m^* - B^\top \mathbf{x}\|^2_{\left(\sum_{i=1}^n \lambda(\mathbf{x}_i) B^\top \mathbf{x}_i \mathbf{x}_i^\top B\right)^{-1}}}{((\mathbf{x}_m^* - \mathbf{x})^\top \boldsymbol{\theta}_m)^2}.$$

λ_m^* and ρ_m^* are the optimal solution and the optimal value of the G-optimal design optimization with true feature extractor B , respectively.

Lemma C.10. Suppose that event $\mathcal{E} \cap \mathcal{F}$ holds. For any task $m \in [M]$ and $\mathbf{y} \in \mathbb{R}^d$,

$$\|\hat{\mathbf{B}}_t^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t\right)^{-1}}^2 \leq \|\mathbf{B}^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B}\right)^{-1}}^2 + \frac{11L_x^4}{k\omega^2 \cdot 2^t}.$$

Proof of Lemma C.10. We first handle the term $\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t\right)^{-1}$.

For any task $m \in [M]$, we have

$$\begin{aligned} & \sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \\ &= \sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i + \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i \right) \cdot \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i + \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i \right)^\top \\ &= \sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \left(\left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i \right) \cdot \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i \right)^\top + \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i \right) \cdot \left(\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i \right)^\top \right. \\ & \quad \left. + \left(\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i \right) \cdot \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i \right)^\top + \left(\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i \right) \cdot \left(\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i \right)^\top \right) \\ &= \sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i \right) \cdot \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i \right)^\top + \sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \left(\left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i \right) \cdot \left(\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i \right)^\top \right. \\ & \quad \left. + \left(\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i \right) \cdot \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i \right)^\top + \left(\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i \right) \cdot \left(\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i \right)^\top \right). \end{aligned}$$

Let $\mathbf{P}_t := \sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) (\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i) \cdot (\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i)^\top$. Let $\mathbf{Q}_t := \sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) ((\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i) \cdot (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i)^\top + (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i) \cdot (\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{x}_i)^\top + (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i) \cdot (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{x}_i)^\top)$. Then, we have $\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t = \mathbf{P}_t + \mathbf{Q}_t$.

From Assumption 3.2, we have that for any task $m \in [M]$, $\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B}$ is invertible. Since $\hat{\mathbf{B}}_t^\top \mathbf{B}$ is also invertible, we have that \mathbf{P}_t is invertible. According to Lemmas C.4 and C.5, we have that $\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t$ is also invertible. Thus, we can write $\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t\right)^{-1}$ as follows.

$$\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right)^{-1} = \mathbf{P}_t^{-1} - (\mathbf{P}_t + \mathbf{Q}_t)^{-1} \mathbf{Q}_t \mathbf{P}_t^{-1}$$

Hence, for any task $m \in [M]$ and $\mathbf{y} \in \mathbb{R}^d$, we have

$$\begin{aligned} \|\hat{\mathbf{B}}_t^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t\right)^{-1}}^2 &= \left(\hat{\mathbf{B}}_t^\top \mathbf{y} \right)^\top \left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t \right)^{-1} \hat{\mathbf{B}}_t^\top \mathbf{y} \\ &= \underbrace{\left(\hat{\mathbf{B}}_t^\top \mathbf{y} \right)^\top \mathbf{P}_t^{-1} \hat{\mathbf{B}}_t^\top \mathbf{y}}_{\text{Term 1}} - \underbrace{\left(\hat{\mathbf{B}}_t^\top \mathbf{y} \right)^\top (\mathbf{P}_t + \mathbf{Q}_t)^{-1} \mathbf{Q}_t \mathbf{P}_t^{-1} \hat{\mathbf{B}}_t^\top \mathbf{y}}_{\text{Term 2}}. \end{aligned} \quad (12)$$

From Lemma C.4, we have

$$\|\hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B}\| \leq \min \left\{ \frac{1}{8k \cdot 2^t \sqrt{1+\zeta}}, \frac{\omega}{6L_x^2} \right\} \leq \min \left\{ \frac{1}{8k \cdot 2^t}, \frac{\omega}{6L_x^2} \right\}.$$

Since $\mathbf{B}^\top \hat{\mathbf{B}}_t \hat{\mathbf{B}}_t^\top \mathbf{B} + \mathbf{B}^\top \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} = \mathbf{B}^\top (\hat{\mathbf{B}}_t \hat{\mathbf{B}}_t^\top + \hat{\mathbf{B}}_{t,\perp} \hat{\mathbf{B}}_{t,\perp}^\top) \mathbf{B} = \mathbf{B}^\top \mathbf{B} = \mathbf{I}_k$, we have $\sigma_{\min}^2(\hat{\mathbf{B}}_t^\top \mathbf{B}) = 1 - \|\hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B}\|^2$.

Thus, we have

$$\sigma_{\min}(\hat{B}_t^\top B) = \sqrt{1 - \|\hat{B}_{t,\perp}^\top B\|^2} \geq \sqrt{1 - \min\left\{\frac{1}{64k^2 \cdot 2^{2t}}, \frac{\omega^2}{36L_x^4}\right\}} > 0,$$

which implies that $\hat{B}_t^\top B$ is invertible.

Now, we first analyze Term 1 in Eq. (12).

$$\begin{aligned} \text{Term 1} &= \left(\hat{B}_t^\top y\right)^\top P_t^{-1} \hat{B}_t^\top y \\ &= \left(\hat{B}_t^\top BB^\top y + \hat{B}_t^\top B_\perp B_\perp^\top y\right)^\top P_t^{-1} \left(\hat{B}_t^\top BB^\top y + \hat{B}_t^\top B_\perp B_\perp^\top y\right) \\ &= \underbrace{\left(\hat{B}_t^\top BB^\top y\right)^\top P_t^{-1} \left(\hat{B}_t^\top BB^\top y\right)}_{\text{Term 1-1}} + \underbrace{\left(\hat{B}_t^\top BB^\top y\right)^\top P_t^{-1} \left(\hat{B}_t^\top B_\perp B_\perp^\top y\right)}_{\text{Term 1-2}} \\ &\quad + \underbrace{\left(\hat{B}_t^\top B_\perp B_\perp^\top y\right)^\top P_t^{-1} \left(\hat{B}_t^\top BB^\top y\right)}_{\text{Term 1-3}} + \underbrace{\left(\hat{B}_t^\top B_\perp B_\perp^\top y\right)^\top P_t^{-1} \left(\hat{B}_t^\top B_\perp B_\perp^\top y\right)}_{\text{Term 1-4}}. \end{aligned}$$

In the following, we bound Terms 1-1, 1-2, 1-3 and 1-4, respectively.

First, we have

$$\begin{aligned} \text{Term 1-1} &= \left(\hat{B}_t^\top BB^\top y\right)^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) \left(\hat{B}_t^\top BB^\top x_i\right) \cdot \left(\hat{B}_t^\top BB^\top x_i\right)^\top\right)^{-1} \hat{B}_t^\top BB^\top y \\ &= \left(\hat{B}_t^\top BB^\top y\right)^\top \left(\hat{B}_t^\top B \left(\sum_{i=1}^n \lambda_m^*(x_i) B^\top x_i x_i^\top B\right) \left(\hat{B}_t^\top B\right)^\top\right)^{-1} \hat{B}_t^\top BB^\top y \\ &= \left(\hat{B}_t^\top BB^\top y\right)^\top \left(\left(\hat{B}_t^\top B\right)^{-1}\right)^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) B^\top x_i x_i^\top B\right)^{-1} \left(\hat{B}_t^\top B\right)^{-1} \hat{B}_t^\top BB^\top y \\ &= \left(B^\top y\right)^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) B^\top x_i x_i^\top B\right)^{-1} B^\top y \\ &= \|B^\top y\|_{\left(\sum_{i=1}^n \lambda_m^*(x_i) B^\top x_i x_i^\top B\right)^{-1}}^2. \end{aligned}$$

We note that since $\hat{B}_t^\top BB^\top \hat{B}_t + \hat{B}_t^\top B_\perp B_\perp^\top \hat{B}_t = \hat{B}_t^\top (BB^\top + B_\perp B_\perp^\top) \hat{B}_t = \hat{B}_t^\top \hat{B}_t = I_k$, $\sigma_{\min}^2(\hat{B}_t^\top B) = 1 - \|\hat{B}_t^\top B_\perp\|^2$. In addition, $\left\|\left(\hat{B}_t^\top B\right)^{-1}\right\| = \frac{1}{\sigma_{\min}(\hat{B}_t^\top B)} = \frac{1}{\sqrt{1 - \|\hat{B}_t^\top B_\perp\|^2}}$.

Then, second, we have

$$\begin{aligned} \text{Term 1-2} &= \left(\hat{B}_t^\top BB^\top y\right)^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) \left(\hat{B}_t^\top BB^\top x_i\right) \cdot \left(\hat{B}_t^\top BB^\top x_i\right)^\top\right)^{-1} \hat{B}_t^\top B_\perp B_\perp^\top y \\ &= \left(\hat{B}_t^\top BB^\top y\right)^\top \left(\hat{B}_t^\top B \left(\sum_{i=1}^n \lambda_m^*(x_i) B^\top x_i x_i^\top B\right) \left(\hat{B}_t^\top B\right)^\top\right)^{-1} \hat{B}_t^\top B_\perp B_\perp^\top y \\ &= \left(\hat{B}_t^\top BB^\top y\right)^\top \left(\left(\hat{B}_t^\top B\right)^{-1}\right)^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) B^\top x_i x_i^\top B\right)^{-1} \left(\hat{B}_t^\top B\right)^{-1} \hat{B}_t^\top B_\perp B_\perp^\top y \end{aligned}$$

$$\begin{aligned}
 &= (\mathbf{B}^\top \mathbf{y})^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) \mathbf{B}^\top x_i x_i^\top \mathbf{B} \right)^{-1} (\hat{\mathbf{B}}_t^\top \mathbf{B})^{-1} \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y} \\
 &\leq 2L_x \cdot \frac{1}{\omega} \cdot \frac{1}{\sqrt{1 - \|\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp\|^2}} \cdot \|\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp\| \cdot 2L_x \\
 &\leq 4L_x^2 \cdot \frac{1}{\omega} \cdot \frac{1}{\sqrt{1 - \left(\frac{1}{8k \cdot 2^t}\right)^2}} \cdot \frac{1}{8k \cdot 2^t} \\
 &\leq 4L_x^2 \cdot \frac{1}{\omega} \cdot \frac{1}{\sqrt{1 - \frac{3}{4}}} \cdot \frac{1}{8k \cdot 2^t} \\
 &= \frac{L_x^2}{k\omega \cdot 2^t}.
 \end{aligned}$$

Third, we have

$$\begin{aligned}
 \text{Term 1-3} &= (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y})^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) (\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top x_i) \cdot (\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top x_i)^\top \right)^{-1} \hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{y} \\
 &= (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y})^\top \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \left(\sum_{i=1}^n \lambda_m^*(x_i) \mathbf{B}^\top x_i x_i^\top \mathbf{B} \right) (\hat{\mathbf{B}}_t^\top \mathbf{B})^\top \right)^{-1} \hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{y} \\
 &= (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y})^\top \left((\hat{\mathbf{B}}_t^\top \mathbf{B})^{-1} \right)^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) \mathbf{B}^\top x_i x_i^\top \mathbf{B} \right)^{-1} (\hat{\mathbf{B}}_t^\top \mathbf{B})^{-1} \hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top \mathbf{y} \\
 &= \left((\hat{\mathbf{B}}_t^\top \mathbf{B})^{-1} \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y} \right)^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) \mathbf{B}^\top x_i x_i^\top \mathbf{B} \right)^{-1} \mathbf{B}^\top \mathbf{y}, \\
 &\leq \frac{1}{\sqrt{1 - \|\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp\|^2}} \cdot \|\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp\| \cdot 2L_x \cdot \frac{1}{\omega} \cdot 2L_x \\
 &\leq 4L_x^2 \cdot \frac{1}{\omega} \cdot \frac{1}{\sqrt{1 - \left(\frac{1}{8k \cdot 2^t}\right)^2}} \cdot \frac{1}{8k \cdot 2^t} \\
 &\leq 4L_x^2 \cdot \frac{1}{\omega} \cdot \frac{1}{\sqrt{1 - \frac{3}{4}}} \cdot \frac{1}{8k \cdot 2^t} \\
 &= \frac{L_x^2}{k\omega \cdot 2^t}.
 \end{aligned}$$

Finally, we have

$$\begin{aligned}
 \text{Term 1-4} &= (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y})^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) (\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top x_i) \cdot (\hat{\mathbf{B}}_t^\top \mathbf{B} \mathbf{B}^\top x_i)^\top \right)^{-1} \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y} \\
 &= (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y})^\top \left(\hat{\mathbf{B}}_t^\top \mathbf{B} \left(\sum_{i=1}^n \lambda_m^*(x_i) \mathbf{B}^\top x_i x_i^\top \mathbf{B} \right) (\hat{\mathbf{B}}_t^\top \mathbf{B})^\top \right)^{-1} \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y} \\
 &= (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y})^\top \left((\hat{\mathbf{B}}_t^\top \mathbf{B})^{-1} \right)^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) \mathbf{B}^\top x_i x_i^\top \mathbf{B} \right)^{-1} (\hat{\mathbf{B}}_t^\top \mathbf{B})^{-1} \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y} \\
 &= \left((\hat{\mathbf{B}}_t^\top \mathbf{B})^{-1} \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y} \right)^\top \left(\sum_{i=1}^n \lambda_m^*(x_i) \mathbf{B}^\top x_i x_i^\top \mathbf{B} \right)^{-1} (\hat{\mathbf{B}}_t^\top \mathbf{B})^{-1} \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \mathbf{B}_\perp^\top \mathbf{y},
 \end{aligned}$$

$$\begin{aligned}
 &\leq \left(\frac{1}{\sqrt{1 - \|\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp\|^2}} \cdot \|\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp\| \cdot 2L_x \right)^2 \cdot \frac{1}{\omega} \\
 &\leq \left(2L_x \cdot \frac{1}{\sqrt{1 - \left(\frac{1}{8k \cdot 2^t}\right)^2}} \cdot \frac{1}{8k \cdot 2^t} \right)^2 \cdot \frac{1}{\omega} \\
 &\leq \left(2L_x \cdot \frac{1}{\sqrt{1 - \frac{3}{4}}} \cdot \frac{1}{8k \cdot 2^t} \right)^2 \cdot \frac{1}{\omega} \\
 &= \frac{L_x^2}{4k^2 \omega \cdot 2^{2t}}.
 \end{aligned}$$

Thus, we have

$$\begin{aligned}
 \text{Term 1} &= \left(\hat{\mathbf{B}}_t^\top \mathbf{y} \right)^\top \mathbf{P}_t^{-1} \hat{\mathbf{B}}_t^\top \mathbf{y} \\
 &\leq \|\mathbf{B}^\top \mathbf{y}\|^2_{\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B} \right)^{-1}} + \frac{2L_x^2}{k\omega \cdot 2^t} + \frac{L_x^2}{4k^2 \omega \cdot 2^{2t}} \\
 &\leq \|\mathbf{B}^\top \mathbf{y}\|^2_{\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B} \right)^{-1}} + \frac{3L_x^2}{k\omega \cdot 2^t}.
 \end{aligned} \tag{13}$$

Next, we investigate Term 2. In order to bound Term 2, we first bound the minimum singular value of \mathbf{P}_t and the maximum singular value of \mathbf{Q}_t .

Since $\mathbf{P}_t = \hat{\mathbf{B}}_t^\top \mathbf{B} \left(\sum_{i=1}^n \lambda^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B} \right) (\hat{\mathbf{B}}_t^\top \mathbf{B})^\top$, we have

$$\begin{aligned}
 \sigma_{\min}(\mathbf{P}_t) &\geq \sigma_{\min}^2(\hat{\mathbf{B}}_t^\top \mathbf{B}) \cdot \omega \\
 &= \left(1 - \|\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp\|^2 \right) \omega \\
 &\geq \left(1 - \frac{1}{8^2 k^2 \cdot 2^{2t}} \right) \omega \\
 &\geq \frac{3}{4} \omega.
 \end{aligned}$$

Since $\mathbf{Q}_t = \hat{\mathbf{B}}_t^\top \mathbf{B} \left(\sum_{i=1}^n \lambda^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B}_\perp \right) (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp)^\top + \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \left(\sum_{i=1}^n \lambda^*(\mathbf{x}_i) \mathbf{B}_\perp^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B} \right) (\hat{\mathbf{B}}_t^\top \mathbf{B})^\top + \hat{\mathbf{B}}_t^\top \mathbf{B}_\perp \left(\sum_{i=1}^n \lambda^*(\mathbf{x}_i) \mathbf{B}_\perp^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B}_\perp \right) (\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp)^\top$, we have

$$\begin{aligned}
 \sigma_{\max}(\mathbf{Q}_t) &\leq 3L_x^2 \|\hat{\mathbf{B}}_t^\top \mathbf{B}_\perp\| \\
 &\leq \min \left\{ \frac{3L_x^2}{8k \cdot 2^t}, \frac{\omega}{2} \right\}.
 \end{aligned}$$

Then, we can bound Term 2 as

$$\begin{aligned}
 \text{Term 2} &= \left(\hat{\mathbf{B}}_t^\top \mathbf{y} \right)^\top (\mathbf{P}_t + \mathbf{Q}_t)^{-1} \mathbf{Q}_t \mathbf{P}_t^{-1} \hat{\mathbf{B}}_t^\top \mathbf{y} \\
 &\leq \left\| \hat{\mathbf{B}}_t^\top \mathbf{y} \right\|^2 \cdot \left\| (\mathbf{P}_t + \mathbf{Q}_t)^{-1} \right\| \cdot \|\mathbf{Q}_t\| \cdot \|\mathbf{P}_t^{-1}\| \\
 &\leq \frac{4L_x^2 \cdot \sigma_{\max}(\mathbf{Q}_t)}{\sigma_{\min}(\mathbf{P}_t + \mathbf{Q}_t) \cdot \sigma_{\min}(\mathbf{P}_t)} \\
 &\leq \frac{4L_x^2 \cdot \sigma_{\max}(\mathbf{Q}_t)}{(\sigma_{\min}(\mathbf{P}_t) - \sigma_{\max}(\mathbf{Q}_t)) \cdot \sigma_{\min}(\mathbf{P}_t)}
 \end{aligned}$$

$$\begin{aligned}
 &\leq \frac{4L_x^2 \cdot \frac{3L_x^2}{8k \cdot 2^t}}{\left(\frac{3}{4}\omega - \frac{1}{2}\omega\right) \cdot \frac{3}{4}\omega} \\
 &= \frac{8L_x^4}{k\omega^2 \cdot 2^t}.
 \end{aligned} \tag{14}$$

Plugging Eqs. (13) and (14) into Eq. (12), we have

$$\begin{aligned}
 \|\hat{\mathbf{B}}_t^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t\right)^{-1}}^2 &\leq \|\mathbf{B}^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B}\right)^{-1}}^2 + \frac{3L_x^2}{k\omega \cdot 2^t} + \frac{8L_x^4}{k\omega^2 \cdot 2^t} \\
 &\leq \|\mathbf{B}^\top \mathbf{y}\|_{\left(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B}\right)^{-1}}^2 + \frac{11L_x^4}{k\omega^2 \cdot 2^t}.
 \end{aligned}$$

□

Below we prove the sample complexity for algorithm DouExpDes (Theorem 4.1).

Proof of Theorem 4.1. According to Lemmas C.3 and C.7, we have $\Pr[\mathcal{E} \cap \mathcal{F}] \geq 1 - \delta$. Below, supposing that event $\mathcal{E} \cap \mathcal{F}$ holds, we prove the correctness and sample complexity.

We first prove the correctness.

For any task $m \in [M]$, let t_m^* denote the first phase which satisfies $|\hat{\mathcal{X}}_{t,m}| = 1$. Let $t_* = \max_{m \in [M]} t_m^*$ denote the total number of phases used. For any task $m \in [M]$, let $\Delta_{m,\min} := \min_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_m^*\}} (\mathbf{x}_m^* - \mathbf{x})^\top \boldsymbol{\theta}_m$ denote the minimum reward gap for task m . Let $\Delta_{\min} := \min_{m \in [M]} \Delta_{m,\min}$ denote the minimum reward gap among all tasks.

From Lemma C.9, we can obtain the following facts: (i) For any task $m \in [M]$, the optimal arm \mathbf{x}_m^* will never be eliminated. (ii) $t_m^* \leq \lceil \log(\frac{4}{\Delta_{m,\min}}) \rceil + 1$, and thus, $t_* \leq \lceil \log(\frac{4}{\Delta_{\min}}) \rceil + 1$. Therefore, after at most $\lceil \log(\frac{4}{\Delta_{\min}}) \rceil + 1$ phases, algorithm DouExpDes will return the optimal arms \mathbf{x}_m^* for all tasks $m \in [M]$.

Now we prove the sample complexity. In the following, we first prove that the sample complexity of algorithm DouExpDes is bounded by $\tilde{O}(\frac{Mk}{\Delta_{\min}^2} \log(\delta^{-1}) + (\rho^E)^2 dk^4 L_x^2 L_w^2 D \log^4(\delta^{-1}))$.

Recall that $p = \frac{180d}{\zeta^2}$ and $\zeta = \frac{1}{10}$. Then, summing the number of samples used in subroutines `FeatRecover` and `ElLowRep` in all phases (Line 4 in Algorithm 2, Line 5 in Algorithm 3), we have that the total number of samples is

$$\begin{aligned}
 &\sum_{t=1}^{t_*} pMT_t + \sum_{m=1}^M \sum_{t=1}^{t_m^*} N_{t,m} \\
 &= \sum_{t=1}^{t_*} p \cdot O\left((1+\zeta)^3 (\rho^E)^2 k^4 L_x^2 L_\theta^2 \max\left\{2^{2t}, \frac{L_x^4}{\omega^2}\right\} \log^2\left(\frac{p}{\delta_t}\right) \cdot \right. \\
 &\quad \left. \log^2\left((1+\zeta) \rho^E k p L_x L_\theta \max\left\{2^t, \frac{L_x}{\omega}\right\} \frac{1}{\delta_t} \log\left(\frac{p}{\delta_t}\right)\right)\right) \\
 &\quad + \sum_{m=1}^M \sum_{t=1}^{t_m^*} O\left(2^{2t} (1+\zeta) \rho_{t,m}^G \log\left(\frac{n^2 M}{\delta_t}\right) + \frac{k}{\zeta^2}\right) \\
 &= \sum_{t=1}^{O(\log(\Delta_{\min}^{-1}))} O\left((\rho^E)^2 k^4 d L_x^2 L_\theta^2 \max\left\{2^{2t}, \frac{L_x^4}{\omega^2}\right\} \log^2\left(\frac{d \log(\Delta_{\min}^{-1})}{\delta}\right) \cdot \right. \\
 &\quad \left. \log^2\left(\rho^E k d L_x L_\theta \max\left\{\Delta_{\min}^{-1}, \frac{L_x}{\omega}\right\} \frac{\log(\Delta_{\min}^{-1})}{\delta} \log\left(\frac{d \log(\Delta_{\min}^{-1})}{\delta}\right)\right)\right) \\
 &\quad + \sum_{m=1}^M \sum_{t=1}^{O(\log(\Delta_{m,\min}^{-1}))} O\left(2^{2t} \rho_{t,m}^G \log\left(\frac{n^2 M \log(\Delta_{m,\min}^{-1})}{\delta}\right) + k\right)
 \end{aligned} \tag{15}$$

$$\begin{aligned}
 &= O\left((\rho^E)^2 k^4 d L_x^2 L_\theta^2 \max\left\{\Delta_{\min}^{-2}, \frac{L_x^4 \log(\Delta_{\min}^{-1})}{\omega^2}\right\} \log^2\left(\frac{d \log(\Delta_{\min}^{-1})}{\delta}\right) \right. \\
 &\quad \left. \log^2\left(\rho^E k d L_x L_\theta \max\left\{\Delta_{\min}^{-1}, \frac{L_x}{\omega}\right\} \frac{\log(\Delta_{\min}^{-1})}{\delta} \log\left(\frac{d \log(\Delta_{\min}^{-1})}{\delta}\right)\right)\right) \\
 &\quad + O\left(M k \Delta_{\min}^{-2} \log\left(\frac{n^2 M \log(\Delta_{\min}^{-1})}{\delta}\right)\right).
 \end{aligned}$$

□

Next, we prove that the sample complexity of algorithm DouExpDes is bounded by $\tilde{O}(\sum_{m=1}^M \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}_m^*\}} \frac{\|\mathbf{B}^\top(\mathbf{x}_m^* - \mathbf{x})\|_{\mathbf{A}(\lambda)^{-1}}^2}{((\mathbf{x}_m^* - \mathbf{x})^\top \boldsymbol{\theta}_m)^2} \log(\delta^{-1}) + (\rho^E)^2 d k^4 L_x^2 L_w^2 D \log^4(\delta^{-1}))$.

From Eq. (15), we have that with probability $1 - \delta$, the number of samples used by algorithm DouExpDes is bounded by

$$\tilde{O}\left(\sum_{t=1}^{\log(\Delta_{\min}^{-1})} (\rho^E)^2 k^4 d L_x^2 L_\theta^2 \max\left\{2^{2t}, \frac{L_x^4}{\omega^2}\right\} + \sum_{m=1}^M \sum_{t=1}^{\log(\Delta_{m,\min}^{-1})} 2^{2t} \rho_{t,m}^G + M k\right). \quad (16)$$

For any $\mathcal{Z} \subseteq \mathcal{X}$, $\mathcal{Y}(\mathcal{Z}) := \{\mathbf{x} - \mathbf{x}' : \forall \mathbf{x}, \mathbf{x}' \in \mathcal{Z}, \mathbf{x} \neq \mathbf{x}'\}$ and $\mathcal{Y}_m^*(\mathcal{Z}) := \{\mathbf{x}_m^* - \mathbf{x} : \forall \mathbf{x} \in \mathcal{Z}, \mathbf{x} \neq \mathbf{x}_m^*\}$. Then, we have that for any task $m \in [M]$ and phase $t \geq 2$,

$$\begin{aligned}
 (2^t)^2 \rho_{t,m}^G &= (2^t)^2 \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{y} \in \mathcal{Y}(\hat{\mathcal{X}}_{t,m})} \|\hat{\mathbf{B}}_t^\top \mathbf{y}\|_{(\sum_{i=1}^n \lambda(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t)^{-1}}^2 \\
 &\leq (2^t)^2 \max_{\mathbf{y} \in \mathcal{Y}(\hat{\mathcal{X}}_{t,m})} \|\hat{\mathbf{B}}_t^\top \mathbf{y}\|_{(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t)^{-1}}^2 \\
 &\stackrel{(a)}{\leq} (2^t)^2 \max_{\mathbf{y} \in \mathcal{Y}(\mathcal{Z}_{t,m})} \|\hat{\mathbf{B}}_t^\top \mathbf{y}\|_{(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t)^{-1}}^2 \\
 &\stackrel{(b)}{\leq} 4 (2^t)^2 \max_{\mathbf{y} \in \mathcal{Y}_m^*(\mathcal{Z}_{t,m})} \|\hat{\mathbf{B}}_t^\top \mathbf{y}\|_{(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \hat{\mathbf{B}}_t^\top \mathbf{x}_i \mathbf{x}_i^\top \hat{\mathbf{B}}_t)^{-1}}^2 \\
 &\stackrel{(c)}{\leq} 4 (2^t)^2 \left(\max_{\mathbf{y} \in \mathcal{Y}_m^*(\mathcal{Z}_{t,m})} \|\mathbf{B}^\top \mathbf{y}\|_{(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B})^{-1}}^2 + \frac{11 L_x^4}{\omega^2 k \cdot 2^t} \right) \\
 &= 4 \left(\frac{16 \max_{\mathbf{y} \in \mathcal{Y}_m^*(\mathcal{Z}_{t,m})} \|\mathbf{B}^\top \mathbf{y}\|_{(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B})^{-1}}^2}{(4 \cdot 2^{-t})^2} + \frac{11 L_x^4 \cdot 2^t}{\omega^2 k} \right) \\
 &\stackrel{(d)}{\leq} 4 \left(16 \max_{\mathbf{y} \in \mathcal{Y}_m^*(\mathcal{Z}_{t,m})} \frac{\|\mathbf{B}^\top \mathbf{y}\|_{(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B})^{-1}}^2}{(\mathbf{y}^\top \boldsymbol{\theta}_m)^2} + \frac{11 L_x^4 \cdot 2^t}{\omega^2 k} \right) \\
 &\leq 4 \left(16 \max_{\mathbf{y} \in \mathcal{Y}_m^*(\mathcal{X})} \frac{\|\mathbf{B}^\top \mathbf{y}\|_{(\sum_{i=1}^n \lambda_m^*(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B})^{-1}}^2}{(\mathbf{y}^\top \boldsymbol{\theta}_m)^2} + \frac{11 L_x^4 \cdot 2^t}{\omega^2 k} \right) \\
 &\stackrel{(e)}{=} 4 \left(16 \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{y} \in \mathcal{Y}_m^*(\mathcal{X})} \frac{\|\mathbf{B}^\top \mathbf{y}\|_{(\sum_{i=1}^n \lambda(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B})^{-1}}^2}{(\mathbf{y}^\top \boldsymbol{\theta}_m)^2} + \frac{11 L_x^4 \cdot 2^t}{\omega^2 k} \right). \quad (17)
 \end{aligned}$$

Here inequality (a) is due to $\hat{\mathcal{X}}_{t,m} \subseteq \mathcal{Z}_{t,m}$ (from Lemma C.9). Inequality (b) uses the fact that for any $\mathbf{y} = \mathbf{x}_i - \mathbf{x}_j \in \mathcal{Y}(\mathcal{Z}_{t,m})$, we can write $\mathbf{y} = (\mathbf{x}_m^* - \mathbf{x}_j) - (\mathbf{x}_m^* - \mathbf{x}_i)$, and the triangle inequality. Inequality (c) follows from Lemma C.10, and inequality (d) is due to that for any $\mathbf{y} \in \mathcal{Y}_m^*(\mathcal{Z}_{t,m})$, $\mathbf{y}^\top \boldsymbol{\theta}_m \leq 4 \cdot 2^{-t}$ (from the definition of $\mathcal{Z}_{t,m}$). Equality (e) comes from the definition of λ_m^* .

Let $L := \log^2\left(\frac{d \log(\Delta_{\min}^{-1})}{\delta}\right) \cdot \log^2(\rho^E k d L_x L_\theta \max\{\Delta_{\min}^{-1}, \frac{L_x}{\omega}\} \frac{\log(\Delta_{\min}^{-1})}{\delta} \log(\frac{d \log(\Delta_{\min}^{-1})}{\delta}))$. Plugging Eq. (17) into Eq. (16), we have that with probability $1 - \delta$, the number of samples used by algorithm DouExpDes is bounded by

$$\begin{aligned}
 & O\left(\sum_{m=1}^M \sum_{t=1}^{\log(\Delta_{m,\min}^{-1})} 2^{2t} \rho_{t,m}^G \log\left(\frac{n^2 M \log(\Delta_{m,\min}^{-1})}{\delta}\right) + Mk \log(\Delta_{\min}^{-1}) + \sum_{t=1}^{\log(\Delta_{\min}^{-1})} (\rho^E)^2 k^4 d L_x^2 L_\theta^2 \max\left\{2^{2t}, \frac{L_x^4}{\omega^2}\right\} L\right) \\
 &= O\left(\sum_{m=1}^M \sum_{t=2}^{\log(\Delta_{m,\min}^{-1})} \left(\min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{y} \in \mathcal{Y}_m^*(\mathcal{X})} \frac{\|\mathbf{B}^\top \mathbf{y}\|^2 (\sum_{i=1}^n \lambda(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B})^{-1}}{(\mathbf{y}^\top \boldsymbol{\theta}_m)^2} + \frac{L_x^4 \cdot 2^t}{\omega^2 k}\right) \cdot \log\left(\frac{n^2 M \log(\Delta_{m,\min}^{-1})}{\delta}\right)\right. \\
 &\quad \left.+ \sum_{m=1}^M \rho_{1,m}^* \log\left(\frac{n^2 M \log(\Delta_{m,\min}^{-1})}{\delta}\right) + Mk \log(\Delta_{\min}^{-1}) + \sum_{t=1}^{\log(\Delta_{\min}^{-1})} (\rho^E)^2 k^4 d L_x^2 L_\theta^2 \max\left\{2^{2t}, \frac{L_x^4}{\omega^2}\right\} \cdot L\right) \\
 &\stackrel{(a)}{=} O\left(\sum_{m=1}^M \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{y} \in \mathcal{Y}_m^*(\mathcal{X})} \frac{\|\mathbf{B}^\top \mathbf{y}\|^2 (\sum_{i=1}^n \lambda(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B})^{-1}}{(\mathbf{y}^\top \boldsymbol{\theta}_m)^2} \cdot \log\left(\frac{n^2 M \log(\Delta_{m,\min}^{-1})}{\delta}\right) \cdot \log(\Delta_{m,\min}^{-1})\right. \\
 &\quad \left.+ \frac{M L_x^4}{\omega^2 k \cdot \Delta_{\min}} \cdot \log\left(\frac{n^2 M \log(\Delta_{\min}^{-1})}{\delta}\right) + Mk \cdot \log\left(\frac{n^2 M \log(\Delta_{\min}^{-1})}{\delta}\right) \cdot \log(\Delta_{\min}^{-1})\right. \\
 &\quad \left.+ (\rho^E)^2 k^4 d L_x^2 L_\theta^2 \max\left\{\Delta_{\min}^{-2}, \frac{L_x^4 \cdot \log(\Delta_{\min}^{-1})}{\omega^2}\right\} \cdot L\right),
 \end{aligned}$$

where equality (a) uses Lemma C.6.

When $L_x = \omega = \Theta(1)$, we have that with probability $1 - \delta$, the sample complexity of algorithm DouExpDes is bounded by

$$\tilde{O}\left(\sum_{m=1}^M \min_{\lambda \in \Delta_{\mathcal{X}}} \max_{\mathbf{y} \in \mathcal{Y}_m^*(\mathcal{X})} \frac{\|\mathbf{B}^\top \mathbf{y}\|^2 (\sum_{i=1}^n \lambda(\mathbf{x}_i) \mathbf{B}^\top \mathbf{x}_i \mathbf{x}_i^\top \mathbf{B})^{-1}}{(\mathbf{y}^\top \boldsymbol{\theta}_m)^2} \log\left(\frac{1}{\delta}\right) + (\rho^E)^2 k^4 d L_x^2 L_\theta^2 \max\left\{\Delta_{\min}^{-2}, \frac{L_x^4}{\omega^2}\right\} \log^4\left(\frac{1}{\delta}\right)\right).$$

D. Proofs for Algorithm C-DouExpDes

In this section, we present the proofs for Algorithm C-DouExpDes.

D.1. Context Distribution Estimation and Sample Batch Planning

Define λ_D^E and ρ_D^E as the optimal solution and the optimal value of the following E-optimal design optimization:

$$\min_{\lambda \in \Delta_{\mathcal{A}}} \left\| \left(\sum_{a \in \mathcal{A}} \lambda(a) \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right)^{-1} \right\|. \quad (18)$$

Lemma D.1. *It holds that*

$$\rho_D^E \leq \frac{1}{\nu}.$$

Proof of Lemma D.1. The optimization in Eq. (18) is equivalent to maximize the minimum singular value of the matrix $\sum_{a \in \mathcal{A}} \lambda(a) \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top]$.

Thus, λ_D^E is the optimal solution of the following optimization:

$$\max_{\lambda \in \Delta_{\mathcal{A}}} \sigma_{\min} \left(\sum_{a \in \mathcal{A}} \lambda(a) \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right).$$

Using Assumption 3.3, we have

$$\sigma_{\min} \left(\sum_{a \in \mathcal{A}} \lambda_D^E \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right) \geq \nu.$$

Then, we have

$$\begin{aligned}\rho_{\mathcal{D}}^E &= \left\| \left(\sum_{a \in \mathcal{A}} \lambda_{\mathcal{D}}^E \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right)^{-1} \right\| \\ &= \frac{1}{\sigma_{\min} \left(\sum_{a \in \mathcal{A}} \lambda_{\mathcal{D}}^E \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right)} \\ &\leq \frac{1}{\nu}.\end{aligned}$$

□

Define event

$$\mathcal{K} := \left\{ \left\| \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] - \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right\| \leq \frac{8L_\phi^2 \log \left(\frac{20d|\mathcal{A}|}{\delta} \right)}{\sqrt{T_0}}, \forall a \in \mathcal{A} \right\}.$$

Lemma D.2. *It holds that*

$$\Pr[\mathcal{K}] \geq 1 - \frac{\delta}{5}.$$

Furthermore, if event \mathcal{K} holds and

$$T_0 = \left\lceil \frac{32^2(1+\zeta)^2 L_\phi^4 \log^2 \left(\frac{20d|\mathcal{A}|}{\delta} \right)}{\nu^2} \right\rceil,$$

we have that for any $a \in \mathcal{A}$,

$$\left\| \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] - \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right\| \leq \frac{\nu}{4(1+\zeta)}.$$

Proof of Lemma D.2. For any $(s, a) \in \mathcal{S} \times \mathcal{A}$, $\|\phi(s, a) \phi(s, a)^\top\| \leq L_\phi^2$. Then, using the matrix Bernstein inequality (Lemma E.2) and a union bound over $a \in \mathcal{A}$, we have that with probability $1 - \frac{\delta}{5}$, for any $a \in \mathcal{A}$,

$$\begin{aligned}\left\| \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] - \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right\| &\leq 4L_\phi^2 \sqrt{\frac{\log \left(\frac{10 \cdot 2d|\mathcal{A}|}{\delta} \right)}{T_0}} + \frac{4L_\phi^2 \log \left(\frac{10 \cdot 2d|\mathcal{A}|}{\delta} \right)}{T_0} \\ &\leq \frac{8L_\phi^2 \log \left(\frac{20d|\mathcal{A}|}{\delta} \right)}{\sqrt{T_0}}.\end{aligned}$$

If $T_0 \geq 32^2(1+\zeta)^2 \nu^{-2} L_\phi^4 \log^2 \left(\frac{20d|\mathcal{A}|}{\delta} \right)$, we have

$$\left\| \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] - \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right\| \leq \frac{\nu}{4(1+\zeta)},$$

which completes the proof. □

Define event

$$\begin{aligned}\mathcal{L} &:= \left\{ \left\| \sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)} \bar{a}_i) \phi(s_{m,j,i}^{(\ell)} \bar{a}_i)^\top - \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] \right\| \leq 8L_\phi^2 \sqrt{p} \log \left(\frac{40dMT}{\delta} \right), \right. \\ &\quad \left. \forall m \in [M], \forall j \in [T], \forall \ell \in \{1, 2\} \right\}.\end{aligned}$$

Lemma D.3. *It holds that*

$$\Pr[\mathcal{L}] \geq 1 - \frac{\delta}{5}.$$

Furthermore, if event \mathcal{L} holds and

$$p = \left\lceil \frac{32^2(1+\zeta)^2 L_\phi^4}{\nu^2} \log^2 \left(\frac{40dMT}{\delta} \right) \right\rceil, \quad (19)$$

we have that for any $m \in [M]$, $j \in [T]$ and $\ell \in \{1, 2\}$,

$$\left\| \sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i) \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top - \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] \right\| \leq \frac{p\nu}{4(1+\zeta)}.$$

Here, the value of T is specified in Eq. (29).

Proof of Lemma D.3. For any $(s, a) \in \mathcal{S} \times \mathcal{A}$, $\|\phi(s, a) \phi(s, a)^\top\| \leq L_\phi^2$. Then, using the matrix Bernstein inequality (Lemma E.2) and a union bound over $m \in [M]$, $j \in [T]$ and $\ell \in \{1, 2\}$, we have that with probability $1 - \frac{\delta}{5}$, for any $m \in [M]$, $j \in [T]$ and $\ell \in \{1, 2\}$,

$$\begin{aligned} \left\| \sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i) \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top - \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] \right\| &\leq 4L_\phi^2 \sqrt{p \log \left(\frac{10 \cdot 4dMT}{\delta} \right)} + 4L_\phi^2 \log \left(\frac{10 \cdot 4dMT}{\delta} \right) \\ &\leq 8L_\phi^2 \sqrt{p} \log \left(\frac{40dMT}{\delta} \right). \end{aligned}$$

In addition, if $p \geq 32^2(1+\zeta)^2 \nu^{-2} L_\phi^4 \log^2 \left(\frac{40dMT}{\delta} \right)$, we have that

$$8L_\phi^2 \sqrt{p} \log \left(\frac{40dMT}{\delta} \right) \leq \frac{p\nu}{4(1+\zeta)}$$

and thus,

$$\left\| \sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i) \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top - \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] \right\| \leq \frac{p\nu}{4(1+\zeta)},$$

which completes the proof. \square

For any task $m \in [M]$, round $j \in [T]$ and $\ell \in \{1, 2\}$, let

$$\Phi_{m,j}^{(\ell)} = \begin{bmatrix} \phi(s_{m,j,1}^{(\ell)}, \bar{a}_1)^\top \\ \vdots \\ \phi(s_{m,j,p}^{(\ell)}, \bar{a}_p)^\top \end{bmatrix},$$

and

$$(\Phi_{m,j}^{(\ell)})^+ = ((\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)})^{-1} (\Phi_{m,j}^{(\ell)})^\top.$$

Lemma D.4. *Suppose that event $\mathcal{K} \cap \mathcal{L}$ holds. Then, for any $m \in [M]$, $j \in [T]$ and $\ell \in \{1, 2\}$,*

$$\|(\Phi_{m,j}^{(\ell)})^+\| \leq 2 \sqrt{\frac{(1+\zeta)}{p\nu}}.$$

Proof of Lemma D.4. We first assume that $(\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)}$ is invertible. In our later analysis, we will prove that as long as T_0 and p are large enough, $(\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)}$ is invertible.

For any $m \in [M]$, $j \in [T]$ and $\ell \in \{1, 2\}$, we have

$$\begin{aligned}
 \|(\Phi_{m,j}^{(\ell)})^+\| &= \|((\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)})^{-1} (\Phi_{m,j}^{(\ell)})^\top\| \\
 &= \sqrt{\|((\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)})^{-1} (\Phi_{m,j}^{(\ell)})^\top ((\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)})^{-1}\|} \\
 &= \sqrt{\|((\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)})^{-1}\|} \\
 &= \frac{1}{\sqrt{\sigma_{\min}((\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)})}}.
 \end{aligned} \tag{20}$$

In addition, we have

$$\begin{aligned}
 &\sigma_{\min}((\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)}) \\
 &= \sigma_{\min}\left(\sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i) \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top\right) \\
 &= \sigma_{\min}\left(\sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] + \sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i) \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top - \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right) \\
 &\geq \sigma_{\min}\left(\sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right) - \left\|\sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i) \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top - \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right\| \\
 &= \sigma_{\min}\left(\sum_{i=1}^p \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] + \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] - \sum_{i=1}^p \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right) \\
 &\quad - \left\|\sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i) \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top - \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right\| \\
 &\geq \sigma_{\min}\left(\sum_{i=1}^p \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right) - \left\|\sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] - \sum_{i=1}^p \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right\| \\
 &\quad - \left\|\sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i) \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top - \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right\| \\
 &\geq \sigma_{\min}\left(\sum_{i=1}^p \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right) - \sum_{i=1}^p \|\mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] - \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\| \\
 &\quad - \left\|\sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i) \phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top - \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right\| \\
 &\geq \sigma_{\min}\left(\sum_{i=1}^p \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right) - \frac{p\nu}{4(1+\zeta)} - \frac{p\nu}{4(1+\zeta)},
 \end{aligned} \tag{21}$$

where the last inequality uses Lemmas D.2 and D.3.

In the following, we analyze $\sigma_{\min}(\sum_{i=1}^p \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top])$. According to the guarantee of the rounding procedure ROUND, we have

$$\left\|\left(\sum_{i=1}^p \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top]\right)^{-1}\right\| \leq (1+\zeta) \left\|\left(p \sum_{a \in \mathcal{A}} \lambda_{\mathcal{D}}^E(a) \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top]\right)^{-1}\right\|$$

$$\leq (1 + \zeta) \left\| \left(p \sum_{a \in \mathcal{A}} \lambda^E(a) \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] \right)^{-1} \right\|,$$

which implies that

$$\begin{aligned} & \sigma_{\min} \left(\sum_{i=1}^p \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] \right) \\ & \geq \frac{p}{1 + \zeta} \sigma_{\min} \left(\sum_{a \in \mathcal{A}} \lambda^E(a) \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] \right) \\ & \geq \frac{p}{1 + \zeta} \sigma_{\min} \left(\sum_{a \in \mathcal{A}} \lambda^E(a) \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right. \\ & \quad \left. + \sum_{a \in \mathcal{A}} \lambda^E(a) \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] - \sum_{a \in \mathcal{A}} \lambda^E(a) \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right) \\ & \geq \frac{p}{1 + \zeta} \left(\sigma_{\min} \left(\sum_{a \in \mathcal{A}} \lambda^E(a) \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right) \right. \\ & \quad \left. - \left\| \sum_{a \in \mathcal{A}} \lambda^E(a) \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] - \sum_{a \in \mathcal{A}} \lambda^E(a) \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right\| \right) \\ & \geq \frac{p}{1 + \zeta} \left(\frac{1}{\rho_{\mathcal{D}}^E} - \sum_{a \in \mathcal{A}} \lambda^E(a) \left\| \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] - \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right\| \right) \\ & \stackrel{(a)}{\geq} \frac{p}{1 + \zeta} \left(\nu - \frac{\nu}{4(1 + \zeta)} \right) \\ & \geq \frac{3p\nu}{4(1 + \zeta)}, \end{aligned} \tag{22}$$

where inequality (a) uses Lemmas D.1 and D.2.

Plugging Eq. (22) into Eq. (21), we have

$$\begin{aligned} \sigma_{\min} \left((\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)} \right) & \geq \frac{3p\nu}{4(1 + \zeta)} - \frac{p\nu}{4(1 + \zeta)} - \frac{p\nu}{4(1 + \zeta)} \\ & = \frac{p\nu}{4(1 + \zeta)}. \end{aligned} \tag{23}$$

Equations (21) and (23) show that if T_0 and p are large enough to satisfy that $\left\| \mathbb{E}_{s \sim \hat{\mathcal{D}}} [\phi(s, a) \phi(s, a)^\top] - \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, a) \phi(s, a)^\top] \right\| \leq \frac{\nu}{4(1 + \zeta)}$ for any $a \in \mathcal{A}$ and $\left\| \sum_{i=1}^p \phi(s_{m,j,i}^{(\ell)} \bar{a}_i) \phi(s_{m,j,i}^{(\ell)} \bar{a}_i)^\top - \sum_{i=1}^p \mathbb{E}_{s \sim \mathcal{D}} [\phi(s, \bar{a}_i) \phi(s, \bar{a}_i)^\top] \right\| \leq \frac{p\nu}{4(1 + \zeta)}$ for any $m \in [M]$, $j \in [T]$ and $\ell \in \{1, 2\}$, respectively, then we have that $(\Phi_{m,j}^{(\ell)})^\top \Phi_{m,j}^{(\ell)}$ is invertible.

Continuing with Eq. (20), we have

$$\left\| (\Phi_{m,j}^{(\ell)})^+ \right\| \leq 2 \sqrt{\frac{(1 + \zeta)}{p\nu}}.$$

□

D.2. Global Feature Extractor Recovery with Stochastic Contexts

In subroutine C-FeatRecover, for any $m \in [M]$, $j \in [T]$, $i \in [p]$ and $\ell \in \{1, 2\}$, let $s_{m,j,i}^{(\ell)}$ and $\eta_{m,j,i}^{(\ell)}$ denote the random context and noise of the ℓ -th sample on action \bar{a}_i in the j -th round for task m , respectively. Here, the superscript $\ell \in \{1, 2\}$ refers to the first sample (Line 4 in Algorithm 5) or the second sample (Line 5 in Algorithm 5) on an action \bar{a}_i .

In C-FeatRecover, for any $m \in [M]$, $j \in [T]$, $i \in [p]$ and $\ell \in \{1, 2\}$, let $\alpha_{m,j}^{(\ell)} \leftarrow [\alpha_{m,j,1}^{(\ell)}, \dots, \alpha_{m,j,p}^{(\ell)}]^\top$, and then, $\tilde{\theta}_{m,j}^{(\ell)} = (\Phi_{m,j}^{(\ell)})^+ \alpha_{m,j}^{(\ell)}$. Recall that $Z = \frac{1}{MT} \sum_{m=1}^M \sum_{j=1}^T \tilde{\theta}_{m,j}^{(1)} (\tilde{\theta}_{m,j}^{(2)})^\top$.

Lemma D.5 (Expectation of Z). *It holds that*

$$\mathbb{E}[Z] = \frac{1}{M} \sum_{m=1}^M \theta_m \theta_m^\top.$$

Proof of Lemma D.5. Z can be written as

$$\begin{aligned} Z &= \frac{1}{MT} \sum_{m=1}^M \sum_{j=1}^T \tilde{\theta}_{m,j}^{(1)} (\tilde{\theta}_{m,j}^{(2)})^\top \\ &= \frac{1}{MT} \sum_{m=1}^M \sum_{j=1}^T (\Phi_{m,j}^{(1)})^+ \begin{bmatrix} \alpha_{m,j,1}^{(1)} \\ \vdots \\ \alpha_{m,j,p}^{(1)} \end{bmatrix} [\alpha_{m,j,1}^{(2)}, \dots, \alpha_{m,j,p}^{(2)}] ((\Phi_{m,j}^{(2)})^+)^{\top} \\ &= \frac{1}{MT} \sum_{m=1}^M \sum_{j=1}^T (\Phi_{m,j}^{(1)})^+ \left(\begin{bmatrix} (\phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \theta_m) (\phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \theta_m) \dots (\phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \theta_m) (\phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \theta_m) \\ \vdots \\ (\phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \theta_m) (\phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \theta_m) \dots (\phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \theta_m) (\phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \theta_m) \end{bmatrix} \right. \\ &\quad + \begin{bmatrix} \phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \theta_m \cdot \eta_{m,j,1}^{(2)} + \eta_{m,j,1}^{(1)} \cdot \phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \theta_m \dots \phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \theta_m \cdot \eta_{m,j,p}^{(2)} + \eta_{m,j,1}^{(1)} \cdot \phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \theta_m \\ \vdots \\ \phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \theta_m \cdot \eta_{m,j,1}^{(2)} + \eta_{m,j,p}^{(1)} \cdot \phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \theta_m \dots \phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \theta_m \cdot \eta_{m,j,p}^{(2)} + \eta_{m,j,p}^{(1)} \cdot \phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \theta_m \end{bmatrix} \\ &\quad \left. + \begin{bmatrix} \eta_{m,j,1}^{(1)} \cdot \eta_{m,j,1}^{(2)} & \dots & \eta_{m,j,1}^{(1)} \cdot \eta_{m,j,p}^{(2)} \\ \vdots & \ddots & \vdots \\ \eta_{m,j,p}^{(1)} \cdot \eta_{m,j,1}^{(2)} & \dots & \eta_{m,j,p}^{(1)} \cdot \eta_{m,j,p}^{(2)} \end{bmatrix} \right) ((\Phi_{m,j}^{(2)})^+)^{\top}. \end{aligned} \quad (24)$$

For any task $m \in [M]$, $j \in [T]$, $i \in [p]$, the sample on action a_i in the first round (i.e., $s_{m,j,i}^{(1)}$ and $\eta_{m,j,i}^{(1)}$) is independent of that in the second round (i.e., $s_{m,j,i}^{(2)}$ and $\eta_{m,j,i}^{(2)}$). Hence, taking the expectation on Z , we obtain

$$\begin{aligned} \mathbb{E}[Z] &= \frac{1}{MT} \sum_{m=1}^M \sum_{j=1}^T \mathbb{E} \left[(\Phi_{m,j}^{(1)})^+ \cdot \begin{bmatrix} (\phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \theta_m) (\phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \theta_m) \dots (\phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \theta_m) (\phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \theta_m) \\ \vdots \\ (\phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \theta_m) (\phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \theta_m) \dots (\phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \theta_m) (\phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \theta_m) \end{bmatrix} ((\Phi_{m,j}^{(2)})^+)^{\top} \right] \\ &= \frac{1}{MT} \sum_{m=1}^M \sum_{j=1}^T \mathbb{E} \left[((\Phi_{m,j}^{(1)})^\top \Phi_{m,j}^{(1)})^{-1} (\Phi_{m,j}^{(1)})^\top \begin{bmatrix} \phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \theta_m \\ \vdots \\ \phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \theta_m \end{bmatrix} \begin{bmatrix} \phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \theta_m, \dots, \phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \theta_m \end{bmatrix} \cdot \Phi_{m,j}^{(2)} ((\Phi_{m,j}^{(2)})^\top \Phi_{m,j}^{(2)})^{-1} \right] \\ &= \frac{1}{MT} \sum_{m=1}^M \sum_{j=1}^T \mathbb{E} \left[((\Phi_{m,j}^{(1)})^\top \Phi_{m,j}^{(1)})^{-1} (\Phi_{m,j}^{(1)})^\top \cdot \Phi_{m,j}^{(1)} \theta_m (\theta_m)^\top (\Phi_{m,j}^{(2)})^\top \cdot \Phi_{m,j}^{(2)} ((\Phi_{m,j}^{(2)})^\top \Phi_{m,j}^{(2)})^{-1} \right] \\ &= \frac{1}{MT} \sum_{m=1}^M \sum_{j=1}^T \theta_m (\theta_m)^\top \end{aligned}$$

$$= \frac{1}{M} \sum_{m=1}^M \boldsymbol{\theta}_m \boldsymbol{\theta}_m^\top.$$

□

Define event

$$\mathcal{G} := \left\{ \|\mathbf{Z} - \mathbb{E}[\mathbf{Z}]\| \leq \frac{256(1+\zeta)L_\phi L_\theta \log\left(\frac{50d}{\delta}\right)}{\nu\sqrt{MT}} \log\left(\frac{100pMT}{\delta}\right) \right\}.$$

Lemma D.6 (Concentration of \mathbf{Z}). *Suppose that $\mathcal{K} \cap \mathcal{L}$ holds. Then, it holds that*

$$\Pr[\mathcal{G}] \geq 1 - \frac{\delta}{5}.$$

Proof of Lemma D.6. Define the following matrices:

$$\begin{aligned} \mathbf{D}_{m,j} &:= \frac{1}{MT} (\boldsymbol{\Phi}_{m,j}^{(1)})^+. \\ &\begin{bmatrix} \left(\phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \boldsymbol{\theta}_m \right) \left(\phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \boldsymbol{\theta}_m \right) & \dots & \left(\phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \boldsymbol{\theta}_m \right) \left(\phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \right) \\ \vdots & \ddots & \vdots \\ \left(\phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \right) \left(\phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \boldsymbol{\theta}_m \right) & \dots & \left(\phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \right) \left(\phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \right) \end{bmatrix} ((\boldsymbol{\Phi}_{m,j}^{(2)})^+)^{\top} \\ \mathbf{D} &:= \sum_{m=1}^M \sum_{j=1}^T \mathbf{D}_{m,j} \\ \mathbf{E}_{m,j} &:= \frac{1}{MT} (\boldsymbol{\Phi}_{m,j}^{(1)})^+. \\ &\begin{bmatrix} \phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \boldsymbol{\theta}_m \cdot \eta_{m,j,1}^{(2)} + \eta_{m,j,1}^{(1)} \cdot \phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \boldsymbol{\theta}_m & \dots & \phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \boldsymbol{\theta}_m \cdot \eta_{m,j,p}^{(2)} + \eta_{m,j,1}^{(1)} \cdot \phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \\ \vdots & \ddots & \vdots \\ \phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \cdot \eta_{m,j,1}^{(2)} + \eta_{m,j,p}^{(1)} \cdot \phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \boldsymbol{\theta}_m & \dots & \phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \cdot \eta_{m,j,p}^{(2)} + \eta_{m,j,p}^{(1)} \cdot \phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \end{bmatrix} \\ &((\boldsymbol{\Phi}_{m,j}^{(2)})^+)^{\top} \\ \mathbf{E} &:= \sum_{m=1}^M \sum_{j=1}^T \mathbf{E}_{m,j} \\ \mathbf{F}_{m,j} &:= \frac{1}{MT} (\boldsymbol{\Phi}_{m,j}^{(1)})^+ \begin{bmatrix} \eta_{m,j,1}^{(1)} \cdot \eta_{m,j,1}^{(2)} & \dots & \eta_{m,j,1}^{(1)} \cdot \eta_{m,j,p}^{(2)} \\ \vdots & \ddots & \vdots \\ \eta_{m,j,p}^{(1)} \cdot \eta_{m,j,1}^{(2)} & \dots & \eta_{m,j,p}^{(1)} \cdot \eta_{m,j,p}^{(2)} \end{bmatrix} ((\boldsymbol{\Phi}_{m,j}^{(2)})^+)^{\top} \\ \mathbf{F} &:= \sum_{m=1}^M \sum_{j=1}^T \mathbf{F}_{m,j} \end{aligned}$$

From Eq. (24), we can bound $\|\mathbf{Z} - \mathbb{E}[\mathbf{Z}]\|$ as

$$\|\mathbf{Z} - \mathbb{E}[\mathbf{Z}]\| \leq \|\mathbf{D} - \mathbb{E}[\mathbf{D}]\| + \|\mathbf{E} - \mathbb{E}[\mathbf{E}]\| + \|\mathbf{F} - \mathbb{E}[\mathbf{F}]\|. \quad (25)$$

Similar to the proof of Lemma C.3, in order to use the truncated matrix Bernstein inequality (Lemma E.2), we define the truncated noise and some truncated matrices as follows.

Let $R > 0$ be a truncation parameter of noises which will be chosen later. For any $m \in [M]$, $j \in [T]$, $i \in [p]$ and $\ell \in \{1, 2\}$, let $\tilde{\eta}_{m,j,i}^{(\ell)} = \eta_{m,j,i}^{(\ell)} \mathbb{1}_{\{|\eta_{m,j,i}^{(\ell)}| \leq R\}}$ denote the truncated noise. Furthermore, we define the following matrices with truncated noises:

$$\tilde{\mathbf{E}}_{m,j} := \frac{1}{MT} (\boldsymbol{\Phi}_{m,j}^{(1)})^+.$$

$$\begin{aligned}
 & \begin{bmatrix} \phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \boldsymbol{\theta}_m \cdot \tilde{\eta}_{m,j,1}^{(2)} + \tilde{\eta}_{m,j,1}^{(1)} \cdot \phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \boldsymbol{\theta}_m \dots \phi(s_{m,j,1}^{(1)}, \bar{a}_1)^\top \boldsymbol{\theta}_m \cdot \tilde{\eta}_{m,j,p}^{(2)} + \tilde{\eta}_{m,j,1}^{(1)} \cdot \phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \\ \vdots \\ \phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \cdot \tilde{\eta}_{m,j,1}^{(2)} + \tilde{\eta}_{m,j,p}^{(1)} \cdot \phi(s_{m,j,1}^{(2)}, \bar{a}_1)^\top \boldsymbol{\theta}_m \dots \phi(s_{m,j,p}^{(1)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \cdot \tilde{\eta}_{m,j,p}^{(2)} + \tilde{\eta}_{m,j,p}^{(1)} \cdot \phi(s_{m,j,p}^{(2)}, \bar{a}_p)^\top \boldsymbol{\theta}_m \end{bmatrix} \\
 & ((\boldsymbol{\Phi}_{m,j}^{(2)})^+)^{\top} \\
 \tilde{\mathbf{E}} &:= \sum_{m=1}^M \sum_{j=1}^T \tilde{\mathbf{E}}_{m,j} \\
 \tilde{\mathbf{F}}_{m,j} &:= \frac{1}{MT} (\boldsymbol{\Phi}_{m,j}^{(1)})^+ \begin{bmatrix} \tilde{\eta}_{m,j,1}^{(1)} \cdot \tilde{\eta}_{m,j,1}^{(2)} & \dots & \tilde{\eta}_{m,j,1}^{(1)} \cdot \tilde{\eta}_{m,j,p}^{(2)} \\ \vdots & \ddots & \vdots \\ \tilde{\eta}_{m,j,p}^{(1)} \cdot \tilde{\eta}_{m,j,1}^{(2)} & \dots & \tilde{\eta}_{m,j,p}^{(1)} \cdot \tilde{\eta}_{m,j,p}^{(2)} \end{bmatrix} ((\boldsymbol{\Phi}_{m,j}^{(2)})^+)^{\top} \\
 \tilde{\mathbf{F}} &:= \sum_{m=1}^M \sum_{j=1}^T \tilde{\mathbf{F}}_{m,j}
 \end{aligned}$$

Recall that from Lemma D.4, we have that for any $m \in [M]$, $j \in [T]$ and $\ell \in \{1, 2\}$, $\|(\boldsymbol{\Phi}_{m,j}^{(\ell)})^+\| \leq 2\sqrt{\frac{(1+\zeta)}{p\nu}}$. Let $B_\Phi := 2\sqrt{\frac{(1+\zeta)}{p\nu}}$.

We first analyze $\|\mathbf{D} - \mathbb{E}[\mathbf{D}]\|$. Since $|\phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top \boldsymbol{\theta}_m| \leq L_\phi L_\theta$ for any $m \in [M]$, $j \in [T]$, $i \in [p]$ and $\ell \in \{1, 2\}$, we have that $\|\mathbf{D}_{m,j}\| \leq \frac{1}{MT} \cdot p L_\phi L_\theta B_\Phi^2$ and $\|\sum_{m=1}^M \sum_{j=1}^T \mathbb{E}[\mathbf{D}_{m,j}^2]\| \leq MT \cdot \frac{1}{M^2 T^2} \cdot p^2 L_\phi^2 L_\theta^2 B_\Phi^4 = \frac{1}{MT} \cdot p^2 L_\phi^2 L_\theta^2 B_\Phi^4$ for any $m \in [M]$ and $j \in [T]$.

Let $\delta' \in (0, 1)$ be a confidence parameter which will be chosen later. Using the matrix Bernstein inequality (Lemma E.2), we have that with probability at least $1 - \delta'$,

$$\begin{aligned}
 \|\mathbf{D} - \mathbb{E}[\mathbf{D}]\| &\leq 4\sqrt{\frac{p^2 L_\phi^2 L_\theta^2 B_\Phi^4 \log\left(\frac{2d}{\delta'}\right)}{MT}} + \frac{4p L_\phi L_\theta B_\Phi^2 \log\left(\frac{2d}{\delta'}\right)}{MT} \\
 &\leq \frac{8 \cdot 4p L_\phi L_\theta B_\Phi^2 \log\left(\frac{2d}{\delta'}\right)}{\sqrt{MT}}.
 \end{aligned} \tag{26}$$

Next, we bound $\|\mathbf{E} - \mathbb{E}[\mathbf{E}]\|$. Since $|\phi(s_{m,j,i}^{(\ell)}, \bar{a}_i)^\top \boldsymbol{\theta}_m| \leq L_\phi L_\theta$ and $|\tilde{\eta}_{m,j,i}^{(\ell)}| \leq R$ for any $m \in [M]$, $j \in [T]$, $i \in [p]$ and $\ell \in \{1, 2\}$, we have that $\|\tilde{\mathbf{E}}_{m,j}\| \leq \frac{1}{MT} \cdot 2p R L_\phi L_\theta B_\Phi^2$ and $\left\|\sum_{m=1}^M \sum_{j=1}^T \mathbb{E}[\tilde{\mathbf{E}}_{m,j}^2]\right\| \leq \frac{1}{MT} \cdot 4p^2 R^2 L_\phi^2 L_\theta^2 B_\Phi^4$ for any $m \in [M]$ and $j \in [T]$.

Since $\eta_{m,j,i}^{(\ell)}$ is 1-sub-Gaussian for any $m \in [M]$, $j \in [T]$, $i \in [p]$ and $\ell \in \{1, 2\}$, using a union bound over $i \in [p]$ and $\ell \in \{1, 2\}$, we have that for any $m \in [M]$ and $j \in [T]$, with probability at least $1 - 4p \exp(-\frac{R^2}{2})$, $|\eta_{m,j,i}^{(\ell)}| \leq R$ for all $i \in [p]$ and $\ell \in \{1, 2\}$, and thus, $\|\mathbf{E}_{m,j}\| \leq \frac{1}{MT} \cdot 2p R L_\phi L_\theta B_\Phi^2$. Then, we have

$$\begin{aligned}
 & \left\| \mathbb{E}[\mathbf{E}_{m,j}] - \mathbb{E}[\tilde{\mathbf{E}}_{m,j}] \right\| \\
 & \leq \left\| \mathbb{E} \left[\mathbf{E}_{m,j} \cdot \mathbb{1} \left\{ \|\mathbf{E}_{m,j}\| \geq \frac{2p R L_\phi L_\theta B_\Phi^2}{MT} \right\} \right] \right\| \\
 & \leq \mathbb{E} \left[\|\mathbf{E}_{m,j}\| \cdot \mathbb{1} \left\{ \|\mathbf{E}_{m,j}\| \geq \frac{2p R L_\phi L_\theta B_\Phi^2}{MT} \right\} \right] \\
 & = \mathbb{E} \left[\frac{2p R L_\phi L_\theta B_\Phi^2}{MT} \cdot \mathbb{1} \left\{ \|\mathbf{E}_{m,j}\| \geq \frac{2p R L_\phi L_\theta B_\Phi^2}{MT} \right\} \right] + \mathbb{E} \left[\left(\|\mathbf{E}_{m,j}\| - \frac{2p R L_\phi L_\theta B_\Phi^2}{MT} \right) \cdot \mathbb{1} \left\{ \|\mathbf{E}_{m,j}\| \geq \frac{2p R L_\phi L_\theta B_\Phi^2}{MT} \right\} \right] \\
 & = \frac{2p R L_\phi L_\theta B_\Phi^2}{MT} \cdot \Pr \left[\|\mathbf{E}_{m,j}\| \geq \frac{2p R L_\phi L_\theta B_\Phi^2}{MT} \right] + \int_0^\infty \Pr \left[\|\mathbf{E}_{m,j}\| - \frac{2p R L_\phi L_\theta B_\Phi^2}{MT} > x \right] dx
 \end{aligned}$$

$$\begin{aligned}
 &\leq \frac{2pRL_\phi L_\theta B_\Phi^2}{MT} \cdot 4p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2pL_\phi L_\theta B_\Phi^2}{MT} \int_R^\infty \Pr\left[\|\mathbf{E}_{m,j}\| > \frac{2pL_\phi L_\theta B_\Phi^2 y}{MT}\right] dy \\
 &\leq \frac{2pRL_\phi L_\theta B_\Phi^2}{MT} \cdot 4p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2pL_\phi L_\theta B_\Phi^2}{MT} \int_R^\infty 4p \exp\left(-\frac{y^2}{2}\right) dy \\
 &\leq \frac{2pRL_\phi L_\theta B_\Phi^2}{MT} \cdot 4p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2pL_\phi L_\theta B_\Phi^2}{MT} \cdot 4p \cdot \frac{1}{R} \cdot \exp\left(-\frac{R^2}{2}\right) \\
 &= \frac{2pL_\phi L_\theta B_\Phi^2}{MT} \cdot 4p \cdot \left(R + \frac{1}{R}\right) \exp\left(-\frac{R^2}{2}\right).
 \end{aligned}$$

Using the truncated matrix Bernstein inequality (Lemma E.2) with $n = MT$, $R = \sqrt{2 \log\left(\frac{4pMT}{\delta'}\right)}$, $U = \frac{2pL_\phi L_\theta B_\Phi^2 \sqrt{2 \log\left(\frac{4pMT}{\delta'}\right)}}{MT}$, $\sigma^2 = \frac{(2pL_\phi L_\theta B_\Phi^2 \sqrt{2 \log\left(\frac{4pMT}{\delta'}\right)})^2}{MT}$, $\tau = 4\sqrt{\frac{(2pL_\phi L_\theta B_\Phi^2 \sqrt{2 \log\left(\frac{4pMT}{\delta'}\right)})^2 \cdot \log\left(\frac{2d}{\delta'}\right)}{MT}} + \frac{4 \cdot 2pL_\phi L_\theta B_\Phi^2 \sqrt{2 \log\left(\frac{4pMT}{\delta'}\right)} \cdot \log\left(\frac{2d}{\delta'}\right)}{MT}$ and $\Delta = \frac{2pL_\phi L_\theta B_\Phi^2 \cdot 2 \sqrt{2 \log\left(\frac{4pMT}{\delta'}\right)}}{MT} \cdot \frac{\delta'}{MT}$, we have that with probability at least $1 - 2\delta'$,

$$\|\mathbf{E} - \mathbb{E}[\mathbf{E}]\| \leq \frac{8 \cdot 2pL_\phi L_\theta B_\Phi^2 \sqrt{2 \log\left(\frac{4pMT}{\delta'}\right)} \cdot \log\left(\frac{2d}{\delta'}\right)}{\sqrt{MT}}. \quad (27)$$

Now we investigate $\|\mathbf{F} - \mathbb{E}[\mathbf{F}]\|$. Since $|\tilde{\eta}_{m,j,i}^{(\ell)}| \leq R$ for any $m \in [M]$, $j \in [T]$, $i \in [p]$ and $\ell \in \{1, 2\}$, we have that $\|\tilde{\mathbf{F}}_{m,j}\| \leq \frac{1}{MT} \cdot pR^2 B_\Phi^2$ and $\left\|\sum_{m=1}^M \sum_{j=1}^T \mathbb{E}[\tilde{\mathbf{F}}_{m,j}^2]\right\| \leq \frac{1}{MT} \cdot p^2 R^4 B_\Phi^4$.

Recall that for any $m \in [M]$ and $j \in [T]$, with probability at least $1 - 4p \exp(-\frac{R^2}{2})$, $|\eta_{m,j,i}^{(\ell)}| \leq R$ for all $i \in [p]$ and $\ell \in \{1, 2\}$, and thus, $\|\mathbf{F}_{m,j}\| \leq \frac{1}{MT} \cdot pB_\Phi^2 R^2$. Then, we have

$$\begin{aligned}
 \left\|\mathbb{E}[\mathbf{F}_{m,j}] - \mathbb{E}[\tilde{\mathbf{F}}_{m,j}]\right\| &\leq \left\|\mathbb{E}\left[\mathbf{F}_{m,j} \cdot \mathbb{1}\left\{\|\mathbf{F}_{m,j}\| \geq \frac{pB_\Phi^2 R^2}{MT}\right\}\right]\right\| \\
 &\leq \mathbb{E}\left[\|\mathbf{F}_{m,j}\| \cdot \mathbb{1}\left\{\|\mathbf{F}_{m,j}\| \geq \frac{pB_\Phi^2 R^2}{MT}\right\}\right] \\
 &= \mathbb{E}\left[\frac{pB_\Phi^2 R^2}{MT} \cdot \mathbb{1}\left\{\|\mathbf{F}_{m,j}\| \geq \frac{pB_\Phi^2 R^2}{MT}\right\}\right] + \left[\left(\|\mathbf{F}_{m,j}\| - \frac{pB_\Phi^2 R^2}{MT}\right) \cdot \mathbb{1}\left\{\|\mathbf{F}_{m,j}\| \geq \frac{pB_\Phi^2 R^2}{MT}\right\}\right] \\
 &= \frac{pB_\Phi^2 R^2}{MT} \cdot \Pr\left[\|\mathbf{F}_{m,j}\| \geq \frac{pB_\Phi^2 R^2}{MT}\right] + \int_0^\infty \Pr\left[\|\mathbf{F}_{m,j}\| - \frac{pB_\Phi^2 R^2}{MT} > x\right] dx \\
 &\leq \frac{pB_\Phi^2 R^2}{MT} \cdot 4p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2pB_\Phi^2}{MT} \int_R^\infty y \cdot \Pr\left[\|\mathbf{F}_{m,j}\| > \frac{pB_\Phi^2 y^2}{MT}\right] dy \\
 &\leq \frac{pB_\Phi^2 R^2}{MT} \cdot 4p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2pB_\Phi^2}{MT} \int_R^\infty y \cdot 4p \exp\left(-\frac{y^2}{2}\right) dy \\
 &\leq \frac{pB_\Phi^2 R^2}{MT} \cdot 4p \cdot \exp\left(-\frac{R^2}{2}\right) + \frac{2pB_\Phi^2}{MT} \cdot 4p \cdot \exp\left(-\frac{R^2}{2}\right) \\
 &= \frac{pB_\Phi^2}{MT} \cdot 4p \cdot (R^2 + 2) \exp\left(-\frac{R^2}{2}\right).
 \end{aligned}$$

Using the truncated matrix Bernstein inequality (Lemma E.2) with $n = MT$, $R = \sqrt{2 \log\left(\frac{4pMT}{\delta'}\right)}$, $U = \frac{pB_\Phi^2 \cdot 2 \log\left(\frac{4pMT}{\delta'}\right)}{MT}$, $\sigma^2 = \frac{(pB_\Phi^2 \cdot 2 \log\left(\frac{4pMT}{\delta'}\right))^2}{MT}$, $\tau = 4\sqrt{\frac{(pB_\Phi^2 \cdot 2 \log\left(\frac{4pMT}{\delta'}\right))^2 \cdot \log\left(\frac{2d}{\delta'}\right)}{MT}} + \frac{4 \cdot pB_\Phi^2 \cdot 2 \log\left(\frac{4pMT}{\delta'}\right) \cdot \log\left(\frac{2d}{\delta'}\right)}{MT}$ and $\Delta = \frac{pB_\Phi^2 \cdot 2 \cdot 2 \log\left(\frac{4pMT}{\delta'}\right)}{MT} \cdot \frac{\delta'}{MT}$, we have that with probability at least $1 - 2\delta'$,

$$\|\mathbf{F} - \mathbb{E}[\mathbf{F}]\| \leq \frac{8 \cdot pB_\Phi^2 \cdot 2 \log\left(\frac{4pMT}{\delta'}\right) \cdot \log\left(\frac{2d}{\delta'}\right)}{\sqrt{MT}}. \quad (28)$$

Plugging Eqs. (26)-(28) into Eq. (25), we have that with probability at least $1 - 5\delta'$,

$$\begin{aligned} \|\mathbf{Z} - \mathbb{E}[\mathbf{Z}]\| &\leq \|\mathbf{D} - \mathbb{E}[\mathbf{D}]\| + \|\mathbf{E} - \mathbb{E}[\mathbf{E}]\| + \|\mathbf{F} - \mathbb{E}[\mathbf{F}]\| \\ &\leq \frac{64pL_\phi L_\theta B_\Phi^2 \log\left(\frac{4pMT}{\delta'}\right) \log\left(\frac{2d}{\delta'}\right)}{\sqrt{MT}}. \end{aligned}$$

Let $\delta' = \frac{\delta}{25}$. Recall that $B_\Phi := 2\sqrt{\frac{(1+\zeta)}{p\nu}}$. Then, we obtain that with probability at least $1 - \frac{\delta}{5}$,

$$\|\mathbf{Z} - \mathbb{E}[\mathbf{Z}]\| \leq \frac{256(1+\zeta)L_\phi L_\theta \log\left(\frac{50d}{\delta}\right) \log\left(\frac{100pMT}{\delta}\right)}{\nu\sqrt{MT}},$$

which implies that $\Pr[\mathcal{G}] \geq 1 - \frac{\delta}{5}$. \square

According to Assumption 3.1, there exists an absolute constant c_0 which satisfies that $\sigma_{\min}(\frac{1}{M} \sum_{m=1}^M \mathbf{w}_m \mathbf{w}_m^\top) = \sigma_{\min}(\frac{1}{M} \sum_{m=1}^M \boldsymbol{\theta}_m \boldsymbol{\theta}_m^\top) \geq \frac{c_0}{k}$.

Lemma D.7 (Concentration of $\hat{\mathbf{B}}$). *Suppose that event \mathcal{G} holds. Then,*

$$\|\hat{\mathbf{B}}_\perp^\top \mathbf{B}\| \leq \frac{2048(1+\zeta)kL_\phi L_\theta \log\left(\frac{50d}{\delta}\right) \log\left(\frac{135(1+\zeta)dL_\phi MT}{\nu\delta}\right)}{c_0\nu\sqrt{MT}}.$$

Furthermore, if

$$T = \left\lceil \frac{68 \cdot 2048^2 \cdot 96^2 (1+\zeta)^2 k^4 L_\phi^4 L_\theta^2 L_w^2}{c_0^2 \nu^2 \varepsilon^2 M} \log^6 \left(\frac{2048 \cdot 135 \cdot 96 \cdot 50 \cdot 5(1+\zeta)^2 k^2 d^2 L_\phi^3 L_\theta L_w N}{c_0 \nu^2 \delta^3 \varepsilon} \right) \right\rceil, \quad (29)$$

we have

$$\|\hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B}\| \leq \frac{\varepsilon}{96k \log\left(\frac{5N}{\delta}\right) L_\phi L_w}.$$

Proof of Lemma D.7. First, we have that $\sigma_k(\mathbb{E}[\mathbf{Z}]) - \sigma_{k+1}(\mathbb{E}[\mathbf{Z}]) = \sigma_{\min}(\frac{1}{M} \sum_{m=1}^M \boldsymbol{\theta}_m \boldsymbol{\theta}_m^\top) \geq \frac{c_0}{k}$. Let $p := \lceil 32^2(1+\zeta)^2 \nu^{-2} L_\phi^4 \log^2\left(\frac{40dMT}{\delta}\right) \rceil$. Then, using the Davis-Kahan sin θ Theorem (Bhatia, 2013) and letting T_t be large enough to satisfy that $\|\mathbf{Z} - \mathbb{E}[\mathbf{Z}]\| \leq \frac{c_0}{2k}$, we have

$$\begin{aligned} \|\hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B}\| &\leq \frac{\|\mathbf{Z} - \mathbb{E}[\mathbf{Z}]\|}{\sigma_k(\mathbb{E}[\mathbf{Z}]) - \sigma_{k+1}(\mathbb{E}[\mathbf{Z}]) - \|\mathbf{Z} - \mathbb{E}[\mathbf{Z}]\|} \\ &\leq \frac{2k}{c_0} \|\mathbf{Z} - \mathbb{E}[\mathbf{Z}]\| \\ &\leq \frac{512(1+\zeta)kL_\phi L_\theta \log\left(\frac{50d}{\delta}\right) \log\left(\frac{100pMT}{\delta}\right)}{c_0\nu\sqrt{MT}} \\ &\leq \frac{512(1+\zeta)kL_\phi L_\theta \log\left(\frac{50d}{\delta}\right) \log\left(\frac{100MT}{\delta} \cdot \frac{2 \cdot 32^2(1+\zeta)^2 L_\phi^4}{\nu^2} \log^2\left(\frac{40dMT}{\delta}\right)\right)}{c_0\nu\sqrt{MT}} \\ &\leq \frac{512(1+\zeta)kL_\phi L_\theta \log\left(\frac{50d}{\delta}\right) \log\left(\frac{2 \cdot 100 \cdot 32^2 \cdot 40^2(1+\zeta)^2 d^2 L_\phi^4 M^3 T^3}{\nu^2 \delta^3}\right)}{c_0\nu\sqrt{MT}} \\ &\leq \frac{2048(1+\zeta)kL_\phi L_\theta \log\left(\frac{50d}{\delta}\right) \log\left(\frac{135(1+\zeta)dL_\phi MT}{\nu\delta}\right)}{c_0\nu\sqrt{MT}}. \end{aligned}$$

Using Lemma E.3 with $A = 2048(1+\zeta)kc_0^{-1}\nu^{-1}L_\phi L_\theta \log\left(\frac{50d}{\delta}\right)$, $B = \frac{135(1+\zeta)dL_\phi}{\nu\delta}$ and $\kappa = \frac{\varepsilon}{96k \log\left(\frac{5N}{\delta}\right) L_\phi L_w}$, we have that if

$$MT \geq \frac{68 \cdot 2048^2 \cdot 96^2 (1+\zeta)^2 k^4 L_\phi^4 L_\theta^2 L_w^2}{c_0^2 \nu^2 \varepsilon^2}.$$

$$\log^2 \left(\frac{50d}{\delta} \right) \log^2 \left(\frac{5N}{\delta} \right) \log^2 \left(\frac{2048 \cdot 135 \cdot 96(1+\zeta)^2 k^2 d L_\phi^3 L_\theta L_w}{c_0 \nu^2 \delta \varepsilon} \log \left(\frac{50d}{\delta} \right) \log \left(\frac{5N}{\delta} \right) \right),$$

then $\left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| \leq \frac{\varepsilon}{96k \log\left(\frac{5N}{\delta}\right) L_\phi L_w}.$

Further enlarging MT , if

$$MT \geq \frac{68 \cdot 2048^2 \cdot 96^2 (1+\zeta)^2 k^4 L_\phi^4 L_\theta^2 L_w^2}{c_0^2 \nu^2 \varepsilon^2} \log^6 \left(\frac{2048 \cdot 135 \cdot 96 \cdot 50 \cdot 5(1+\zeta)^2 k^2 d^2 L_\phi^3 L_\theta L_w N}{c_0 \nu^2 \delta^3 \varepsilon} \right),$$

then

$$\left\| \hat{\mathbf{B}}_{t,\perp}^\top \mathbf{B} \right\| \leq \frac{\varepsilon}{96k \log\left(\frac{5N}{\delta}\right) L_\phi L_w}.$$

□

D.3. Estimation with Low-dimensional Representations

Lemma D.8. *In subroutine EstLowRep (Algorithm 6), for any $m \in [M]$ and $t > 0$, we have*

$$\log \left(\frac{\det \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)}{\det(\gamma I)} \right) \leq k \log \left(1 + \frac{t}{\gamma k} \right).$$

Proof of Lemma D.8. This proof uses a similar idea as Lemma 11 in (Abbasi-Yadkori et al., 2011).

It holds that

$$\begin{aligned} & \log \left(\frac{\det \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)}{\det(\gamma I)} \right) \\ & \leq \log \left(\frac{\left(\frac{\text{Trace}(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}})}{k} \right)^k}{\gamma^k} \right) \\ & = k \log \left(\frac{\text{Trace}(\gamma I) + \sum_{\tau=1}^t \text{Trace} \left(\hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)}{\gamma k} \right) \\ & = k \log \left(\frac{\gamma k + \sum_{\tau=1}^t \left\| \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \right\|^2}{\gamma k} \right) \\ & \leq k \log \left(1 + \frac{t}{\gamma k} \right). \end{aligned}$$

□

Lemma D.9. *In subroutine EstLowRep (Algorithm 6), for any $m \in [M]$ and $t \geq 0$, we have*

$$\mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_{m,t}^{-1}} \right] \geq \mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_{m,t+1}^{-1}} \right].$$

Proof of Lemma D.9. This proof is similar to that of Lemma 6 in (Zanette et al., 2021).

For any $m \in [M]$ and $t \geq 0$, since $\Sigma_{m,t+1} \succeq \Sigma_{m,t}$, we have $\Sigma_{m,t}^{-1} \succeq \Sigma_{m,t+1}^{-1}$. Hence, for any $m \in [M]$, $t \geq 0$, $s \in \mathcal{S}$ and $a \in \mathcal{A}$, we have

$$\phi(s, a)^\top \hat{B} \Sigma_{m,t}^{-1} \hat{B}^\top \phi(s, a) \geq \phi(s, a)^\top \hat{B} \Sigma_{m,t+1}^{-1} \hat{B}^\top \phi(s, a),$$

which implies that

$$\left\| \hat{B}^\top \phi(s, a) \right\|_{\Sigma_{m,t}^{-1}} \geq \left\| \hat{B}^\top \phi(s, a) \right\|_{\Sigma_{m,t+1}^{-1}}.$$

Therefore, for any $m \in [M]$ and $t \geq 0$, we have

$$\mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{B}^\top \phi(s, a) \right\|_{\Sigma_{m,t}^{-1}} \right] \geq \mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{B}^\top \phi(s, a) \right\|_{\Sigma_{m,t+1}^{-1}} \right].$$

□

In subroutine `EstLowRep`, for any $m \in [M]$ and $t > 0$, let $\xi_{m,t}$ denote the noise of the sample at timestep t for task m (Line 6 in Algorithm 6).

Define event

$$\mathcal{H} := \left\{ \left\| \sum_{\tau=1}^t \hat{B}^\top \phi(s_{m,\tau}, a_{m,\tau}) \xi_{m,\tau} \right\|_{(\gamma I + \sum_{\tau=1}^t \hat{B}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{B})^{-1}} \leq k \log \left(1 + \frac{t}{\gamma k} \right) + 2 \log \left(\frac{5}{\delta} \right), \forall m \in [M], \forall t > 0 \right\}.$$

Lemma D.10 (Martingale Concentration of the Variance Term). *It holds that*

$$\Pr[\mathcal{H}] \geq 1 - \frac{\delta}{5}.$$

Proof of Lemma D.10. Let δ' be a confidence parameter which will be chosen later. Since \hat{B} is fixed before sampling $(s_{m,\tau}, a_{m,\tau})$ for all $m \in [M]$ and $\tau > 0$, using Lemma E.7, we have that with probability at least $1 - \delta'$, for any task $m \in [M]$ and $t > 0$,

$$\begin{aligned} & \left\| \sum_{\tau=1}^t \hat{B}^\top \phi(s_{m,\tau}, a_{m,\tau}) \xi_{m,\tau} \right\|_{(\gamma I + \sum_{\tau=1}^t \hat{B}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{B})^{-1}} \\ & \leq 2 \log \left(\frac{\det \left(\gamma I + \sum_{\tau=1}^t \hat{B}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{B} \right)^{\frac{1}{2}}}{\det(\gamma I)^{\frac{1}{2}} \cdot \delta'} \right) \\ & \leq \log \left(\frac{\det \left(\gamma I + \sum_{\tau=1}^t \hat{B}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{B} \right)}{\det(\gamma I)} \right) + 2 \log \left(\frac{1}{\delta'} \right) \\ & \stackrel{(a)}{\leq} k \log \left(1 + \frac{t}{\gamma k} \right) + 2 \log \left(\frac{1}{\delta'} \right), \end{aligned}$$

where inequality (a) uses Lemma D.8.

Letting $\delta' = \frac{\delta}{5}$, we obtain this lemma. □

Define event

$$\mathcal{J} := \left\{ \sum_{t=1}^N \mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{B}^\top \phi(s, a) \right\|_{\Sigma_{t-1}^{-1}} \right] \leq \right.$$

$$\frac{1}{4} \left(2\sqrt{\log\left(\frac{5}{\delta}\right)} + \sqrt{4\log\left(\frac{5}{\delta}\right) + 4\left(\sum_{t=1}^N \max_{a \in \mathcal{A}} \|\hat{\mathbf{B}}^\top \phi(s_t, a)\|_{\Sigma_{t-1}^{-1}} + 2\log\left(\frac{5}{\delta}\right)\right)} \right)^2 \Bigg\}.$$

Lemma D.11. *It holds that*

$$\Pr[\mathcal{J}] \geq 1 - \frac{\delta}{5}.$$

Proof of Lemma D.11. Using Lemma E.8, we can obtain this lemma. \square

Lemma D.12. *Suppose that event $\mathcal{K} \cap \mathcal{L} \cap \mathcal{G} \cap \mathcal{H} \cap \mathcal{J}$ holds. For any task $m \in [M]$, we have*

$$\begin{aligned} \mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left| \phi(s, a)^\top (\hat{\boldsymbol{\theta}}_{m,N} - \boldsymbol{\theta}_m) \right| \right] &\leq \left(2\sqrt{\frac{2k \log\left(1 + \frac{N}{\gamma k}\right)}{N}} + \frac{8 \log\left(\frac{5}{\delta}\right)}{N} \right) \\ &\quad \left(\|\hat{\mathbf{B}}_\perp^\top \mathbf{B}\| \sqrt{Nk} + \sqrt{k \log\left(1 + \frac{N}{\gamma k}\right) + 2\log\left(\frac{5}{\delta}\right) + \sqrt{\gamma}} \right) + \|\hat{\mathbf{B}}_\perp^\top \mathbf{B}\|. \end{aligned}$$

Furthermore, if

$$N = \left\lceil \frac{4^2 \cdot 26^4 \cdot 24^2 \cdot 2(k^2 + k\gamma L_\theta^2) \log^4\left(\frac{240(k + \sqrt{k\gamma} L_\theta)}{\varepsilon \delta}\right)}{\varepsilon^2} \right\rceil,$$

then

$$\mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left| \phi(s, a)^\top (\hat{\boldsymbol{\theta}}_{m,N} - \boldsymbol{\theta}_m) \right| \right] \leq \frac{\varepsilon}{2}$$

Proof of Lemma D.12. For any task $m \in [M]$ and $t \in [N]$,

$$\begin{aligned} \hat{\mathbf{w}}_{m,t} &= \Sigma_{m,t}^{-1} \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) r_{m,\tau} \\ &= \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)^{-1} \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) (\phi(s_{m,\tau}, a_{m,\tau})^\top \boldsymbol{\theta}_m + \xi_{m,j}) \\ &= \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)^{-1} \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \cdot \\ &\quad (\phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \hat{\mathbf{B}}^\top \boldsymbol{\theta}_m + \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}}_\perp \hat{\mathbf{B}}_\perp^\top \boldsymbol{\theta}_m + \xi_{m,j}) \\ &\quad + \gamma \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)^{-1} \hat{\mathbf{B}}^\top \boldsymbol{\theta}_m \\ &\quad - \gamma \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)^{-1} \hat{\mathbf{B}}^\top \boldsymbol{\theta}_m \\ &= \hat{\mathbf{B}}^\top \boldsymbol{\theta}_m + \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)^{-1} \cdot \\ &\quad \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}}_\perp \hat{\mathbf{B}}_\perp^\top \mathbf{B} \mathbf{w}_m \end{aligned}$$

$$\begin{aligned}
 & + \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)^{-1} \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \xi_{m,j} \\
 & - \gamma \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)^{-1} \hat{\mathbf{B}}^\top \boldsymbol{\theta}_m.
 \end{aligned}$$

Hence, for any task $m \in [M]$, $t \in [N]$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned}
 \phi(s, a)^\top (\hat{\boldsymbol{\theta}}_{m,t} - \boldsymbol{\theta}_m) & = \phi(s, a)^\top \hat{\mathbf{B}} \hat{\mathbf{w}}_{m,t} - \phi(s, a)^\top (\hat{\mathbf{B}} \hat{\mathbf{B}}^\top + \hat{\mathbf{B}}_\perp \hat{\mathbf{B}}_\perp^\top) \boldsymbol{\theta}_m \\
 & = \phi(s, a)^\top \hat{\mathbf{B}} (\hat{\mathbf{w}}_{m,t} - \hat{\mathbf{B}}^\top \boldsymbol{\theta}_m) - \phi(s, a)^\top \hat{\mathbf{B}}_\perp \hat{\mathbf{B}}_\perp^\top \boldsymbol{\theta}_m \\
 & = \phi(s, a)^\top \hat{\mathbf{B}} \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)^{-1} \\
 & \quad \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}}_\perp \hat{\mathbf{B}}_\perp^\top \mathbf{B} \mathbf{w}_m \\
 & \quad + \phi(s, a)^\top \hat{\mathbf{B}} \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)^{-1} \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \xi_{m,j} \\
 & \quad - \gamma \phi(s, a)^\top \hat{\mathbf{B}} \left(\gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)^{-1} \hat{\mathbf{B}}^\top \boldsymbol{\theta}_m - \phi(s, a)^\top \hat{\mathbf{B}}_\perp \hat{\mathbf{B}}_\perp^\top \mathbf{B} \mathbf{w}_m.
 \end{aligned}$$

For any $m \in [M]$, let $\boldsymbol{\Sigma}_{m,0} := \gamma I$. For any $m \in [M]$ and $t \geq 1$, let $\boldsymbol{\Sigma}_{m,t} := \gamma I + \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}}$.

Taking the absolute value on both sides and using the Cauchy–Schwarz inequality, we obtain that for any $m \in [M]$, $t \in [N]$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\begin{aligned}
 & \left| \phi(s, a)^\top (\hat{\boldsymbol{\theta}}_{m,t} - \boldsymbol{\theta}_m) \right| \\
 & \leq \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \left\| \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}}_\perp \hat{\mathbf{B}}_\perp^\top \mathbf{B} \mathbf{w}_m \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \\
 & \quad + \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \left\| \sum_{\tau=1}^t \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \xi_{m,j} \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \\
 & \quad + \gamma \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \left\| \hat{\mathbf{B}}^\top \boldsymbol{\theta}_m \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \\
 & \quad + \left| \phi(s, a)^\top \hat{\mathbf{B}}_\perp \hat{\mathbf{B}}_\perp^\top \mathbf{B} \mathbf{w}_m \right| \\
 & \stackrel{(a)}{\leq} \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \sum_{\tau=1}^t \left| \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}}_\perp \hat{\mathbf{B}}_\perp^\top \mathbf{B} \mathbf{w}_m \right| \cdot \left\| \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \\
 & \quad + \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \sqrt{k \log \left(1 + \frac{t}{\gamma k} \right) + 2 \log \left(\frac{5}{\delta} \right)} \\
 & \quad + \gamma \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \cdot \frac{1}{\sqrt{\gamma}} \cdot \left\| \hat{\mathbf{B}}^\top \boldsymbol{\theta}_m \right\| + \left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\| L_\phi L_w \\
 & \leq \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}} \cdot \left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\| L_\phi L_w \cdot \sum_{\tau=1}^t \left\| \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \right\|_{\boldsymbol{\Sigma}_{m,t}^{-1}}
 \end{aligned}$$

$$\begin{aligned}
 & + \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_{m,t}^{-1}} \sqrt{k \log \left(1 + \frac{t}{\gamma k} \right) + 2 \log \left(\frac{5}{\delta} \right)} \\
 & + \sqrt{\gamma} L_\theta \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_{m,t}^{-1}} + \left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\| L_\phi L_w \\
 & \stackrel{(b)}{\leq} \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_{m,t}^{-1}} \cdot \left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\| L_\phi L_w \cdot \sqrt{tk} \\
 & + \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_{m,t}^{-1}} \sqrt{k \log \left(1 + \frac{t}{\gamma k} \right) + 2 \log \left(\frac{5}{\delta} \right)} \\
 & + \sqrt{\gamma} L_\theta \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_{m,t}^{-1}} + \left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\| L_\phi L_w \\
 & = \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_{m,t}^{-1}} \left(\left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\| L_\phi L_w \cdot \sqrt{tk} + \sqrt{k \log \left(1 + \frac{t}{\gamma k} \right) + 2 \log \left(\frac{5}{\delta} \right) + \gamma L_\theta} \right) \\
 & + \left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\| L_\phi L_w,
 \end{aligned}$$

where inequality (a) uses the triangle inequality and the definition of event \mathcal{H} , and inequality (b) is due to Lemma E.6.

Taking the maximum over $a \in \mathcal{A}$ and taking the expectation on $s \sim \mathcal{D}$, we have that for any task $m \in [M]$,

$$\begin{aligned}
 \mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left| \phi(s, a)^\top (\hat{\boldsymbol{\theta}}_{m,N} - \boldsymbol{\theta}_m) \right| \right] & \leq \mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_N^{-1}} \right] \cdot \\
 & \quad \left(\left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\| L_\phi L_w \cdot \sqrt{Nk} + \sqrt{k \log \left(1 + \frac{N}{\gamma k} \right) + 2 \log \left(\frac{5}{\delta} \right) + \gamma L_\theta} \right) \\
 & + \left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\| L_\phi L_w.
 \end{aligned} \tag{30}$$

According to Lemma D.9, $\mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_t^{-1}} \right]$ is non-increasing with respect to t . Hence, we have

$$\begin{aligned}
 \mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_N^{-1}} \right] & \leq \frac{1}{N} \sum_{t=1}^N \mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_t^{-1}} \right] \\
 & \leq \frac{1}{N} \sum_{t=1}^N \mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_{t-1}^{-1}} \right] \\
 & \stackrel{(a)}{\leq} \frac{1}{4N} \left(2 \sqrt{\log \left(\frac{5}{\delta} \right)} \right. \\
 & \quad \left. + \sqrt{4 \log \left(\frac{5}{\delta} \right) + 4 \left(\sum_{t=1}^N \max_{a \in \mathcal{A}} \left\| \hat{\mathbf{B}}^\top \phi(s_t, a) \right\|_{\Sigma_{t-1}^{-1}} + 2 \log \left(\frac{5}{\delta} \right) \right)} \right)^2 \\
 & = \frac{1}{4N} \left(2 \sqrt{\log \left(\frac{5}{\delta} \right)} + \sqrt{4 \log \left(\frac{5}{\delta} \right) + 4 \left(\sum_{t=1}^N \left\| \hat{\mathbf{B}}^\top \phi(s_t, a_t) \right\|_{\Sigma_{t-1}^{-1}} + 2 \log \left(\frac{5}{\delta} \right) \right)} \right)^2,
 \end{aligned} \tag{31}$$

where inequality (a) is due to the definition of event \mathcal{J} .

In addition, we have

$$\sum_{t=1}^N \left\| \hat{\mathbf{B}}^\top \phi(s_t, a_t) \right\|_{\Sigma_{t-1}^{-1}} \leq \sqrt{N \cdot \sum_{t=1}^N \left\| \hat{\mathbf{B}}^\top \phi(s_t, a_t) \right\|_{\Sigma_{t-1}^{-1}}^2}$$

$$\begin{aligned}
 &\stackrel{(a)}{\leq} \sqrt{2N \log \left(\frac{\det \left(\gamma I + \sum_{\tau=1}^N \hat{\mathbf{B}}^\top \phi(s_{m,\tau}, a_{m,\tau}) \phi(s_{m,\tau}, a_{m,\tau})^\top \hat{\mathbf{B}} \right)}{\det(\gamma I)} \right)} \\
 &\stackrel{(b)}{\leq} \sqrt{2Nk \log \left(1 + \frac{N}{\gamma k} \right)},
 \end{aligned} \tag{32}$$

where inequality (a) uses Lemma E.9, and inequality (b) is due to Lemma D.8.

Combining Eqs. (31) and (32), we have

$$\begin{aligned}
 \mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left\| \hat{\mathbf{B}}^\top \phi(s, a) \right\|_{\Sigma_N^{-1}} \right] &\leq \frac{1}{4N} \left(2\sqrt{\log \left(\frac{5}{\delta} \right)} + \sqrt{4 \log \left(\frac{5}{\delta} \right) + 4 \left(\sqrt{2Nk \log \left(1 + \frac{N}{\gamma k} \right)} + 2 \log \left(\frac{5}{\delta} \right) \right)} \right)^2 \\
 &\stackrel{(a)}{\leq} \frac{1}{2N} \left(4 \log \left(\frac{5}{\delta} \right) + 4 \log \left(\frac{5}{\delta} \right) + 4 \left(\sqrt{2Nk \log \left(1 + \frac{N}{\gamma k} \right)} + 2 \log \left(\frac{5}{\delta} \right) \right) \right) \\
 &= \frac{1}{N} \left(2\sqrt{2Nk \log \left(1 + \frac{N}{\gamma k} \right)} + 8 \log \left(\frac{5}{\delta} \right) \right) \\
 &= 2\sqrt{\frac{2k \log \left(1 + \frac{N}{\gamma k} \right)}{N}} + \frac{8 \log \left(\frac{5}{\delta} \right)}{N},
 \end{aligned} \tag{33}$$

where inequality (a) uses the Cauchy–Schwarz inequality.

Furthermore, plugging Eq. (33) into Eq. (30) and using $\gamma \geq 1$, we have that for $N \geq 1$ and $\sqrt{k} \log(2N) \geq 1$,

$$\begin{aligned}
 &\mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left| \phi(s, a)^\top (\hat{\boldsymbol{\theta}}_{m,N} - \boldsymbol{\theta}_m) \right| \right] \\
 &\leq \left(2\sqrt{\frac{2k \log \left(1 + \frac{N}{\gamma k} \right)}{N}} + \frac{8 \log \left(\frac{5}{\delta} \right)}{N} \right) \\
 &\quad \left(\left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\|_{L_\phi L_w \sqrt{Nk}} + \sqrt{k \log \left(1 + \frac{N}{\gamma k} \right) + 2 \log \left(\frac{5}{\delta} \right) + \sqrt{\gamma} L_\theta} + \left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\|_{L_\phi L_w} \right) \\
 &\leq \frac{12\sqrt{k} \log \left(\frac{5N}{\delta} \right)}{\sqrt{N}} \left(\left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\|_{L_\phi L_w \sqrt{Nk}} + 2\sqrt{k} \log \left(\frac{5N}{\delta} \right) + \sqrt{\gamma} L_\theta \right) + \left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\|_{L_\phi L_w} \\
 &\leq \frac{(24k + 12\sqrt{k\gamma} L_\theta) \log^2 \left(\frac{5N}{\delta} \right)}{\sqrt{N}} + 24k \log \left(\frac{5N}{\delta} \right) \left\| \hat{\mathbf{B}}_\perp^\top \mathbf{B} \right\|_{L_\phi L_w}.
 \end{aligned} \tag{34}$$

Using Lemma E.4 with $A = 24k + 12\sqrt{k\gamma} L_\theta$, $B = \frac{5}{\delta}$ and $\kappa = \frac{\varepsilon}{4}$, we have that if

$$N \geq \frac{26^4 (24k + 12\sqrt{k\gamma} L_\theta)^2 \log^4 \left(\frac{2 \cdot 5(24k + 12\sqrt{k\gamma} L_\theta)}{\varepsilon \delta} \right)}{\left(\frac{\varepsilon}{4} \right)^2},$$

then $\frac{(24k + 12\sqrt{k\gamma} L_\theta) \log^2 \left(\frac{5N}{\delta} \right)}{\sqrt{N}} \leq \frac{\varepsilon}{4}$.

Further enlarging N , if

$$N \geq \frac{4^2 \cdot 26^4 \cdot 24^2 \cdot 2 (k^2 + k\gamma L_\theta^2) \log^4 \left(\frac{240(k + \sqrt{k\gamma} L_\theta)}{\varepsilon \delta} \right)}{\varepsilon^2}, \tag{36}$$

then

$$\frac{(24k + 12\sqrt{k\gamma}L_\theta) \log^2\left(\frac{5N}{\delta}\right)}{\sqrt{N}} \leq \frac{\varepsilon}{4}.$$

According to Lemma D.7, we have $\|\hat{B}_{t,\perp}^\top \mathbf{B}\| \leq \frac{\varepsilon}{96k \log\left(\frac{5N}{\delta}\right) L_\phi L_w}$.

Thus, setting N as the value in Eq. (36), and continuing with Eq. (35), we have

$$\mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left| \phi(s, a)^\top (\hat{\boldsymbol{\theta}}_{m,N} - \boldsymbol{\theta}_m) \right| \right] \leq \frac{\varepsilon}{4} + \frac{\varepsilon}{4} = \frac{\varepsilon}{2}.$$

□

D.4. Proof of Theorem 5.1

Proof of Theorem 5.1. Combining Lemmas D.2, D.3, D.6, D.10 and D.11, we have that $\Pr[\mathcal{K} \cap \mathcal{L} \cap \mathcal{G} \cap \mathcal{H} \cap \mathcal{J}] \geq 1 - \delta$. Suppose that event $\mathcal{K} \cap \mathcal{L} \cap \mathcal{G} \cap \mathcal{H} \cap \mathcal{J}$ holds.

First, we use a similar analytical procedure as that in (Zanette et al., 2021) to prove the correctness.

Using Lemma D.12, we have that for any task $m \in [M]$,

$$\mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} \left| \phi(s, a)^\top (\hat{\boldsymbol{\theta}}_{m,N} - \boldsymbol{\theta}_m) \right| \right] \leq \frac{\varepsilon}{2}.$$

For any $m \in [M]$ and $s \in \mathcal{S}$, let $\beta_m(s) := \max_{a \in \mathcal{A}} |\phi(s, a)^\top (\hat{\boldsymbol{\theta}}_{m,N} - \boldsymbol{\theta}_m)|$ and $\pi_m^*(s) := \operatorname{argmax}_{a \in \mathcal{A}} \phi(s, a)^\top \boldsymbol{\theta}_m$.

For any $m \in [M]$ and $s \in \mathcal{S}$, we have

$$\begin{aligned} \phi(s, \hat{\pi}_m(s))^\top \boldsymbol{\theta}_m &\geq \phi(s, \hat{\pi}_m(s))^\top \hat{\boldsymbol{\theta}}_{m,N} - \beta_m(s) \\ &\stackrel{(a)}{\geq} \phi(s, \pi_m^*(s))^\top \hat{\boldsymbol{\theta}}_{m,N} - \beta_m(s) \\ &\geq \phi(s, \pi_m^*(s))^\top \boldsymbol{\theta}_m - 2\beta_m(s), \end{aligned}$$

where inequality (a) is due to that $\hat{\pi}_m(s)$ is greedy with respect to $\hat{\boldsymbol{\theta}}_{m,N}$.

Rearranging the above equation and taking the expectation of s on both sides, we have

$$\mathbb{E}_{s \sim \mathcal{D}} \left[\max_{a \in \mathcal{A}} (\phi(s, a) - \phi(s, \hat{\pi}_m(s)))^\top \boldsymbol{\theta}_m \right] \leq 2\mathbb{E}_{s \sim \mathcal{D}} [\beta_m(s)] \leq \varepsilon.$$

Now we prove the sample complexity. Summing the number of samples used in the main algorithm of C-DouExpDes and subroutines C-FeatRecover and EstLowRep (Line ?? in Algorithm 4, Lines 4-5 in Algorithm 5 and Line 6 in Algorithm 6), we have that the total number of samples is bounded by

$$\begin{aligned} &T_0 + 2MTp + MN \\ &= O\left(\frac{L_\phi^4}{\nu^2} \log^2\left(\frac{d|\mathcal{A}|}{\delta}\right) + \frac{k^4 L_\phi^4 L_\theta^2 L_w^2}{\nu^2 \varepsilon^2} \log^6\left(\frac{kdL_\phi L_\theta L_w N}{\nu \delta \varepsilon}\right) \cdot \frac{L_\phi^4}{\nu^2} \log^2\left(\frac{dMT}{\delta}\right)\right. \\ &\quad \left.+ M \cdot \frac{(k^2 + k\gamma L_\theta^2) \log^4\left(\frac{k + \sqrt{k\gamma} L_\theta}{\varepsilon \delta}\right)}{\varepsilon^2}\right) \\ &= O\left(\frac{k^4 L_\phi^4 L_\theta^2 L_w^2}{\nu^2 \varepsilon^2} \log^6\left(\frac{|\mathcal{A}| kdL_\phi L_\theta L_w N}{\nu \delta \varepsilon}\right) \cdot \frac{L_\phi^4}{\nu^2} \log^2\left(\frac{dMT}{\delta}\right)\right) \end{aligned}$$

$$\begin{aligned}
 &+ M \cdot \frac{(k^2 + k\gamma L_\theta^2) \log^4 \left(\frac{k + \sqrt{k\gamma} L_\theta}{\varepsilon \delta} \right)}{\varepsilon^2} \\
 &= \tilde{O} \left(\frac{k^4 L_\phi^8 L_\theta^2 L_w^2}{\nu^4 \varepsilon^2} + \frac{M (k^2 + k\gamma L_\theta^2)}{\varepsilon^2} \right).
 \end{aligned}$$

□

E. Technical Tools

In this section, we provide some useful technical tools.

Lemma E.1 (Matrix Bernstein Inequality - Average, Lemma 31 in (Tripuraneni et al., 2021)). *Consider a truncation level $U > 0$. If $\{\mathbf{Z}_1, \dots, \mathbf{Z}_n\}$ is a sequence of $d_1 \times d_2$ independent random matrices and $\mathbf{Z}'_i = \mathbf{Z}_i \cdot \mathbb{1}\{\|\mathbf{Z}_i\| \leq U\}$ for any $i \in [n]$, then*

$$\Pr \left[\left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{Z}_i - \mathbb{E}[\mathbf{Z}_i]) \right\| \geq t \right] \leq \Pr \left[\left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{Z}'_i - \mathbb{E}[\mathbf{Z}'_i]) \right\| \geq t - \Delta \right] + n \Pr [\|\mathbf{Z}_i\| \geq U],$$

where $\Delta \geq \|\mathbb{E}[\mathbf{Z}_i] - \mathbb{E}[\mathbf{Z}'_i]\|$ for any $i \in [n]$.

In addition, for $t \geq \Delta$, we have

$$\Pr \left[\left\| \frac{1}{n} \sum_{i=1}^n (\mathbf{Z}'_i - \mathbb{E}[\mathbf{Z}'_i]) \right\| \geq t - \Delta \right] \leq (d_1 + d_2) \exp \left(-\frac{n^2(t - \Delta)^2}{2\sigma^2 + \frac{2Un(t - \Delta)}{3}} \right),$$

where

$$\begin{aligned}
 \sigma^2 &= \max \left\{ \left\| \sum_{i=1}^n \mathbb{E}[(\mathbf{Z}'_i - \mathbb{E}[\mathbf{Z}'_i])^\top (\mathbf{Z}'_i - \mathbb{E}[\mathbf{Z}'_i])] \right\|, \left\| \sum_{i=1}^n \mathbb{E}[(\mathbf{Z}'_i - \mathbb{E}[\mathbf{Z}'_i])(\mathbf{Z}'_i - \mathbb{E}[\mathbf{Z}'_i])^\top] \right\| \right\} \\
 &\leq \max \left\{ \left\| \sum_{i=1}^n \mathbb{E}[\mathbf{Z}'_i{}^\top \mathbf{Z}'_i] \right\|, \left\| \sum_{i=1}^n \mathbb{E}[\mathbf{Z}'_i \mathbf{Z}'_i{}^\top] \right\| \right\}.
 \end{aligned}$$

Lemma 31 in (Tripuraneni et al., 2021) gives a truncated matrix Bernstein inequality for symmetric random matrices. Here we extend it to general random matrices.

Lemma E.1 can be obtained by combining the truncation argument in the proof of Lemma 31 in (Tripuraneni et al., 2021) and Theorem 6.1.1 in (Tropp et al., 2015) (classic matrix Bernstein inequality for general random matrices).

Lemma E.2 (Matrix Bernstein Inequality - Summation). *Consider a truncation level $U > 0$. If $\{\mathbf{Z}_1, \dots, \mathbf{Z}_n\}$ is a sequence of $d_1 \times d_2$ independent random matrices, and $\mathbf{Z}'_i = \mathbf{Z}_i \cdot \mathbb{1}\{\|\mathbf{Z}_i\| \leq U\}$ and $\Delta \geq \|\mathbb{E}[\mathbf{Z}_i] - \mathbb{E}[\mathbf{Z}'_i]\|$ for any $i \in [n]$, then for $\tau \geq 2n\Delta$,*

$$\Pr \left[\left\| \sum_{i=1}^n (\mathbf{Z}_i - \mathbb{E}[\mathbf{Z}_i]) \right\| \geq \tau \right] \leq (d_1 + d_2) \exp \left(-\frac{1}{4} \cdot \frac{\tau^2}{2\sigma^2 + \frac{U\tau}{3}} \right) + n \Pr [\|\mathbf{Z}_i\| \geq U],$$

where

$$\begin{aligned}
 \sigma^2 &= \max \left\{ \left\| \sum_{i=1}^n \mathbb{E}[(\mathbf{Z}'_i - \mathbb{E}[\mathbf{Z}'_i])^\top (\mathbf{Z}'_i - \mathbb{E}[\mathbf{Z}'_i])] \right\|, \left\| \sum_{i=1}^n \mathbb{E}[(\mathbf{Z}'_i - \mathbb{E}[\mathbf{Z}'_i])(\mathbf{Z}'_i - \mathbb{E}[\mathbf{Z}'_i])^\top] \right\| \right\} \\
 &\leq \max \left\{ \left\| \sum_{i=1}^n \mathbb{E}[\mathbf{Z}'_i{}^\top \mathbf{Z}'_i] \right\|, \left\| \sum_{i=1}^n \mathbb{E}[\mathbf{Z}'_i \mathbf{Z}'_i{}^\top] \right\| \right\}.
 \end{aligned}$$

Furthermore, we have

$$\Pr \left[\left\| \sum_{i=1}^n (\mathbf{Z}_i - \mathbb{E}[\mathbf{Z}_i]) \right\| \geq 4\sqrt{\sigma^2 \log \left(\frac{d_1 + d_2}{\delta} \right)} + 4U \log \left(\frac{d_1 + d_2}{\delta} \right) \right] \leq \delta + n \Pr [\|\mathbf{Z}_i\| \geq U].$$

Proof of Lemma E.2. Using Lemma E.1 and defining $\tau := nt$, we have that for $\tau > n\Delta$,

$$\Pr \left[\left\| \sum_{i=1}^n (\mathbf{Z}_i - \mathbb{E}[\mathbf{Z}_i]) \right\| \geq \tau \right] \leq (d_1 + d_2) \exp \left(-\frac{(\tau - n\Delta)^2}{2\sigma^2 + \frac{2U(\tau - n\Delta)}{3}} \right) + n \Pr [\|\mathbf{Z}_i\| \geq U].$$

If $\tau > 2n\Delta$, then $\tau - n\Delta > \frac{1}{2}\tau$ and we have

$$\begin{aligned} \Pr \left[\left\| \sum_{i=1}^n (\mathbf{Z}_i - \mathbb{E}[\mathbf{Z}_i]) \right\| \geq \tau \right] &\leq (d_1 + d_2) \exp \left(-\frac{(\frac{1}{2}\tau)^2}{2\sigma^2 + \frac{2U(\frac{1}{2}\tau)}{3}} \right) + n \Pr [\|\mathbf{Z}_i\| \geq U] \\ &\leq (d_1 + d_2) \exp \left(-\frac{1}{4} \cdot \frac{\tau^2}{2\sigma^2 + \frac{U\tau}{3}} \right) + n \Pr [\|\mathbf{Z}_i\| \geq U]. \end{aligned} \quad (37)$$

Plugging $\tau = 4\sqrt{\sigma^2 \log \left(\frac{d_1 + d_2}{\delta} \right)} + 4U \log \left(\frac{d_1 + d_2}{\delta} \right)$ into Eq. (37), we have

$$\begin{aligned} &\Pr \left[\left\| \sum_{i=1}^n (\mathbf{Z}_i - \mathbb{E}[\mathbf{Z}_i]) \right\| \geq 4\sqrt{\sigma^2 \log \left(\frac{d_1 + d_2}{\delta} \right)} + 4U \log \left(\frac{d_1 + d_2}{\delta} \right) \right] \\ &\leq (d_1 + d_2) \exp \left(-\frac{1}{4} \cdot \frac{16\sigma^2 \log \left(\frac{d_1 + d_2}{\delta} \right) + 16U^2 \log^2 \left(\frac{d_1 + d_2}{\delta} \right) + 32U \log \left(\frac{d_1 + d_2}{\delta} \right) \sqrt{\sigma^2 \log \left(\frac{d_1 + d_2}{\delta} \right)}}{2\sigma^2 + \frac{1}{3} \left(4U \sqrt{\sigma^2 \log \left(\frac{d_1 + d_2}{\delta} \right)} + 4U^2 \log \left(\frac{d_1 + d_2}{\delta} \right) \right)} \right) \\ &\quad + n \Pr [\|\mathbf{Z}_i\| \geq U] \\ &\leq (d_1 + d_2) \exp \left(-\frac{1}{4} \cdot 4 \log \left(\frac{d_1 + d_2}{\delta} \right) \right) + n \Pr [\|\mathbf{Z}_i\| \geq U] \\ &= \delta + n \Pr [\|\mathbf{Z}_i\| \geq U]. \end{aligned}$$

□

Lemma E.3. For any $A, B > 1$, $\kappa \in (0, 1)$ and $T > 0$ such that $\log \left(\frac{AB}{\kappa} \right) > 1$ and $\log(BT) > 2$, if

$$T \geq \frac{68A^2 \log^2 \left(\frac{AB}{\kappa} \right)}{\kappa^2},$$

then

$$\frac{A}{\sqrt{T}} \log(BT) \leq \kappa.$$

Proof of Lemma E.3. If $T = \frac{68A^2 \log^2 \left(\frac{AB}{\kappa} \right)}{\kappa^2}$, we have

$$\begin{aligned} \frac{A}{\sqrt{T}} \log(BT) &= \frac{A\kappa}{\sqrt{68A \log \left(\frac{AB}{\kappa} \right)}} \log \left(\frac{68A^2 B \log^2 \left(\frac{AB}{\kappa} \right)}{\kappa^2} \right) \\ &= \frac{\kappa}{\sqrt{68 \log \left(\frac{AB}{\kappa} \right)}} \left(\log(68) + \log \left(\frac{A^2 B}{\kappa^2} \right) + \log \left(\log^2 \left(\frac{AB}{\kappa} \right) \right) \right) \\ &\leq \frac{\kappa}{\sqrt{68 \log \left(\frac{AB}{\kappa} \right)}} \left(\log(68) + 2 \log \left(\frac{AB}{\kappa} \right) + 2 \log \left(\frac{AB}{\kappa} \right) \right) \\ &\leq \frac{\kappa}{\sqrt{68 \log \left(\frac{AB}{\kappa} \right)}} \left(\log(68) \log \left(\frac{AB}{\kappa} \right) + 4 \log \left(\frac{AB}{\kappa} \right) \right) \\ &\leq \kappa. \end{aligned}$$

Let $f(T) = \frac{A}{\sqrt{T}} \log(BT)$. Then, the derivative of $f(T)$ is

$$f'(T) = \frac{2A - A \log(BT)}{2T\sqrt{T}}.$$

If $\log(BT) > 2$, then $f'(T) < 0$, and thus $f(T)$ is decreasing with respect to T .

Therefore, if $T \geq \frac{68A^2 \log^2(\frac{AB}{\kappa})}{\kappa^2}$, we have

$$\frac{A}{\sqrt{T}} \log(BT) \leq \kappa.$$

□

Lemma E.4. For any $A, B > 1$ and $\kappa \in (0, 1)$ such that $\log(\frac{AB}{\kappa}) > 1$ and $\log(BN) > 4$, if

$$N \geq \frac{26^4 A^2 \log^4(\frac{AB}{\kappa})}{\kappa^2},$$

then

$$\frac{A \log^2(BN)}{\sqrt{N}} \leq \kappa.$$

Proof of Lemma E.4. If $N = \frac{26^4 A^2 \log^4(\frac{AB}{\kappa})}{\kappa^2}$, we have $\kappa\sqrt{N} = 26^2 A \log^2(\frac{AB}{\kappa})$, and

$$\begin{aligned} A \log^2(BN) &= A \log^2\left(\frac{26^4 A^2 B \log^4(\frac{AB}{\kappa})}{\kappa^2}\right) \\ &\leq A \log^2\left(\frac{26^4 A^2 B}{\kappa^2} \cdot \frac{A^4 B^4}{\kappa^4}\right) \\ &\leq 36A \log^2\left(\frac{26AB}{\kappa}\right) \\ &= 36A \left(\log(26) + \log\left(\frac{AB}{\kappa}\right)\right)^2 \\ &\leq 36A \left(\log(26) \log\left(\frac{AB}{\kappa}\right) + \log\left(\frac{AB}{\kappa}\right)\right)^2 \\ &= 36(\log(26) + 1)^2 A \log^2\left(\frac{AB}{\kappa}\right) \\ &\leq 26^2 A \log^2\left(\frac{AB}{\kappa}\right) \\ &= \kappa\sqrt{N}, \end{aligned}$$

and thus $\frac{A \log^2(BN)}{\sqrt{N}} \leq \kappa$.

Let $f(N) = \frac{A \log^2(BN)}{\sqrt{N}}$. Then, the derivative function of $f(N)$ is

$$f'(N) = \frac{4A \log(BN) - A \log^2(BN)}{2N\sqrt{N}} = \frac{A \log(BN) \cdot (4 - \log(BN))}{2N\sqrt{N}}.$$

If $\log(BN) > 4$, then $f'(N) < 0$, and thus $f(N)$ is decreasing with respect to N .

Therefore, if $N \geq \frac{26^4 A^2 \log^4(\frac{AB}{\kappa})}{\kappa^2}$, we have $\frac{A \log^2(BN)}{\sqrt{N}} \leq \kappa$.

□

Lemma E.5. For any $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^k$, we have

$$\sum_{j=1}^n \|\mathbf{x}_j\|_{(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top)^{-1}} \leq \sqrt{nk}.$$

Proof of Lemma E.5. It holds that

$$\begin{aligned} \sum_{j=1}^n \|\mathbf{x}_j\|_{(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top)^{-1}} &= \sum_{j=1}^n \sqrt{\mathbf{x}_j^\top \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \mathbf{x}_j} \\ &\leq \sqrt{n \cdot \sum_{j=1}^n \mathbf{x}_j^\top \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \mathbf{x}_j} \\ &\leq \sqrt{n \cdot \sum_{j=1}^n \text{Trace} \left(\mathbf{x}_j^\top \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \mathbf{x}_j \right)} \\ &= \sqrt{n \cdot \sum_{j=1}^n \text{Trace} \left(\mathbf{x}_j \mathbf{x}_j^\top \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right)} \\ &= \sqrt{n \cdot \text{Trace} \left(\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right)} \\ &= \sqrt{n \cdot \text{Trace}(\mathbf{I}_k)} \\ &= \sqrt{nk} \end{aligned}$$

□

Lemma E.6. For any $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^k$ and $\gamma > 0$, we have

$$\sum_{j=1}^n \|\mathbf{x}_j\|_{(\gamma \mathbf{I} + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top)^{-1}} \leq \sqrt{nk}.$$

Proof of Lemma E.6. It holds that

$$\begin{aligned} \sum_{j=1}^n \|\mathbf{x}_j\|_{(\gamma \mathbf{I} + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top)^{-1}} &= \sum_{j=1}^n \sqrt{\mathbf{x}_j^\top \left(\gamma \mathbf{I} + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \mathbf{x}_j} \\ &\leq \sqrt{n \cdot \sum_{j=1}^n \mathbf{x}_j^\top \left(\gamma \mathbf{I} + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \mathbf{x}_j} \\ &= \sqrt{n \cdot \sum_{j=1}^n \text{Trace} \left(\mathbf{x}_j^\top \left(\gamma \mathbf{I} + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \mathbf{x}_j \right)} \\ &= \sqrt{n \cdot \sum_{j=1}^n \text{Trace} \left(\mathbf{x}_j \mathbf{x}_j^\top \left(\gamma \mathbf{I} + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right)} \\ &= \sqrt{n \cdot \text{Trace} \left(\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \left(\gamma \mathbf{I} + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right)} \end{aligned}$$

$$\begin{aligned}
 &\stackrel{(a)}{\leq} \sqrt{n \cdot \left(\text{Trace} \left(\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \left(\gamma I + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right) + \text{Trace} \left(\gamma \left(\gamma I + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right) \right)} \\
 &= \sqrt{n \cdot \text{Trace} \left(\sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \left(\gamma I + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} + \gamma \left(\gamma I + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right)} \\
 &= \sqrt{n \cdot \text{Trace} \left(\left(\gamma I + \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \right) \left(\gamma I + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \right)^{-1} \right)} \\
 &= \sqrt{n \cdot \text{Trace}(\mathbf{I}_k)} \\
 &= \sqrt{nk},
 \end{aligned}$$

where inequality (a) is due to that $(\gamma I + \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top)$ is a positive definite matrix. \square

Lemma E.7 (Self-normalized Concentration for Martingales, Theorem 1 in (Abbasi-Yadkori et al., 2011)). *Let $\{\mathcal{F}_t\}_{t=0}^\infty$ be a filtration such that for any $t \geq 1$, the selected action $\mathbf{X}_t \in \mathbb{R}^k$ is \mathcal{F}_{t-1} -measurable, the noise $\eta_t \in \mathbb{R}$ is \mathcal{F}_t -measurable, and conditioning on \mathcal{F}_{t-1} , η_t is zero-mean and R -sub-Gaussian. Let $\mathbf{V}_0 \in \mathbb{R}^{k \times k}$ be a positive definite matrix and let $\mathbf{V}_t = \sum_{i=1}^t \mathbf{X}_i \mathbf{X}_i^\top$ for any $t \geq 1$. Then, for any $\delta > 0$, with probability at least $1 - \delta$, for all $t \geq 1$,*

$$\left\| \sum_{i=1}^t \mathbf{X}_i \cdot \eta_i \right\|_{(\mathbf{V}_0 + \mathbf{V}_t)^{-1}}^2 \leq 2R^2 \log \left(\frac{\det(\mathbf{V}_t)^{\frac{1}{2}}}{\det(\mathbf{V}_0)^{\frac{1}{2}} \cdot \delta} \right).$$

Lemma E.8 (Reverse Bernstein Inequality for Martingales, Theorem 3 in (Zanette et al., 2021)). *Let $(\Sigma, \mathcal{F}, \Pr[\cdot])$ be a probability space and consider the stochastic process $\{\mathbf{X}_t\}$ adapted to the filtration $\{\mathcal{F}_t\}$. Let $\mathbb{E}_t[\mathbf{X}_t] := \mathbb{E}[\mathbf{X}_t | \mathcal{F}_{t-1}]$ be the conditional expectation of \mathbf{X}_t given \mathcal{F}_{t-1} . If $0 \leq \mathbf{X}_t \leq 1$ then it holds that*

$$\Pr \left[\sum_{t=1}^T \mathbb{E}_t[\mathbf{X}_t] \geq \frac{1}{4} \left(2\sqrt{\log \left(\frac{1}{\delta} \right)} + \sqrt{4 \log \left(\frac{1}{\delta} \right) + 4 \left(\sum_{t=1}^T \mathbf{X}_t + 2 \log \left(\frac{1}{\delta} \right) \right)} \right)^2 \right] \leq \delta.$$

Lemma E.9 (Elliptical Potential Lemma, Lemma 11 in (Abbasi-Yadkori et al., 2011)). *Let $\{\mathbf{X}_t\}_{t=1}^\infty$ be a sequence in \mathbb{R}^k . Let \mathbf{V}_0 be a $k \times k$ positive definite matrix and let $\mathbf{V}_t = \mathbf{V}_0 + \sum_{i=1}^t \mathbf{X}_i \mathbf{X}_i^\top$ such that for any $t \geq 1$, $\|\mathbf{X}_t\|_{\mathbf{V}_{t-1}^{-1}}^2 \leq 1$. Then, we have that*

$$\sum_{t=1}^n \|\mathbf{X}_t\|_{\mathbf{V}_{t-1}^{-1}}^2 \leq 2 \log \frac{\det(\mathbf{V}_n)}{\det(\mathbf{V}_0)}.$$

Lemma E.10 (Moments of Sub-Gaussian Random Variables, Proposition 3.2 in (Rivasplata, 2012)). *For a σ^2 -sub-Gaussian random variable \mathbf{X} which satisfies*

$$\mathbb{E}[\exp(\mu \mathbf{X})] \leq \exp\left(\frac{\sigma^2 \mu^2}{2}\right), \quad \forall \mu \in \mathbb{R},$$

we have that for any integer $n \geq 1$,

$$\mathbb{E}[|\mathbf{X}|^n] \leq (2\sigma^2)^{\frac{n}{2}} n \cdot \Gamma\left(\frac{n}{2}\right),$$

where $\Gamma(n) := (n-1)!$ for any integer $n \geq 1$.