## Achieving Limited Adaptivity for Multinomial Logistic Bandits

Sukruta Prakash Midigeshi, Tanmay Goyal, Gaurav Sinha

Keywords: Multinomial Logistic Bandits, Limited Adaptivity, Batched Bandits, Contextual Bandits

## Summary

Multinomial Logistic Bandits have recently attracted much attention due to their ability to model problems with multiple outcomes. In the multinomial model, each decision is associated with many possible outcomes, modeled using a multinomial logit function. Several recent works on multinomial logistic bandits have simultaneously achieved optimal regret and computational efficiency. However, motivated by real-world challenges and practicality, there is a need to develop algorithms with limited adaptivity, wherein we are allowed only M policy updates. To address these challenges, we present two algorithms, B-MNL-CB and RS-MNL, that operate in the batched and rarely-switching paradigms, respectively. The batched setting involves choosing the M policy update rounds at the start of the algorithm, while the rarelyswitching setting can choose these M policy update rounds in an adaptive fashion. Our first algorithm, B-MNL-CB extends the notion of distributional optimal designs to the multinomial setting and achieves  $\tilde{O}(\sqrt{T})$  regret assuming the contexts are generated stochastically when presented with  $\Omega(\log \log T)$  update rounds. Our second algorithm, RS-MNL works with adversarially generated contexts and can achieve  $\tilde{O}(\sqrt{T})$  regret with  $\tilde{O}(\log T)$  policy updates. Further, we conducted experiments that demonstrate that our algorithms (with a fixed number of policy updates) are extremely competitive (and often better) than several state-of-the-art baselines (which update their policy every round), showcasing the applicability of our algorithms in various practical scenarios.

## **Contribution(s)**

- 1. We present an algorithm, B-MNL-CB, that achieves an optimal  $\tilde{O}(\sqrt{T})$  regret with  $\Omega(\log \log T)$  batches in the batched setting. Moreover, the leading term of the regret is independent of  $\kappa$ , an instance-dependent non-linearity parameter.
- **Context:** In the batched setting, the rounds at which the policy is updated are fixed beforehand. Gao et al. (2019) showed that having  $\Omega(\log \log T)$  batches is necessary to achieve the optimal minimax regret. Our algorithm, B-MNL-CB, combines the idea of distributional optimal designs (introduced in Ruan et al. (2021)) with the idea of suitable scalings for arms (introduced in Sawarni et al. (2024)) to the multinomial logistic setting. This requires a natural extension of distributional optimal designs to this setting. Achieving a  $\kappa$ - independent regret is important because Amani & Thrampoulidis (2021) showed that  $\kappa$  scales exponentially in several instance parameters and hence, can increase the regret significantly.
- 2. We present a rarely-switching algorithm RS-MNL that achieves an optimal  $O(\sqrt{T})$  regret (with a  $\kappa$ -free leading term) requiring  $O(\log T)$  switches (policy updates). **Context:** In the rarely-switching setting, the switching rounds (policy-update rounds) are adaptively chosen during the course of the algorithm. The need for the update is decided based on a switching criterion similar to the one in Abbasi-Yadkori et al. (2011). While the algorithm bears similarities to the rarely-switching algorithm presented in Sawarni et al. (2024), an alternate regret decomposition method allows us to get rid of the warm-up criterion, which helps reduce the number of switches from  $O(\log^2 T)$  to  $O(\log T)$ . Further, we also get rid of the *Successive Eliminations* in Sawarni et al. (2024) that determine the arm to be played, and replace it with the simpler UCB-maximization rule of Abbasi-Yadkori et al. (2011), resulting in a more efficient runtime for the algorithm.

# Achieving Limited Adaptivity for Multinomial Logistic Bandits

Sukruta Prakash Midigeshi<sup>1</sup>, Tanmay Goyal<sup>1</sup>, Gaurav Sinha<sup>1</sup>

{t-smidigeshi,t-tangoyal,gauravsinha}@microsoft.com

<sup>1</sup>Microsoft Research India

## Abstract

Multinomial Logistic Bandits have recently attracted much attention due to their ability to model problems with multiple outcomes. In the multinomial model, each decision is associated with many possible outcomes, modeled using a multinomial logit function. Several recent works on multinomial logistic bandits have simultaneously achieved optimal regret and computational efficiency. However, motivated by real-world challenges and practicality, there is a need to develop algorithms with limited adaptivity, wherein we are allowed only M policy updates. To address these challenges, we present two algorithms, B-MNL-CB and RS-MNL, that operate in the batched and rarely-switching paradigms, respectively. The batched setting involves choosing the M policy update rounds at the start of the algorithm, while the rarely-switching setting can choose these M policy update rounds in an adaptive fashion. Our first algorithm, B-MNL-CB extends the notion of distributional optimal designs to the multinomial setting and achieves  $O(\sqrt{T})$  regret assuming the contexts are generated stochastically when presented with  $\Omega(\log \log T)$  update rounds. Our second algorithm, RS-MNL works with adversarially generated contexts and can achieve  $O(\sqrt{T})$  regret with  $O(\log T)$  policy updates. Further, we conducted experiments that demonstrate that our algorithms (with a fixed number of policy updates) are extremely competitive (and often better) than several state-of-the-art baselines (which update their policy every round), showcasing the applicability of our algorithms in various practical scenarios.

### **1** Introduction and Prior Works

Contextual Bandits help incorporate additional information that a learner may have with the standard Multi-Armed Bandit (MAB) setting. In this setting, at each round, the learner is presented with a set of arms and is expected to choose an arm. She is also presented with a context vector that helps guide the decisions she makes. For each decision, the learner receives a reward, which is generated using a hidden optimal parameter. The goal of the learner is to minimize her cumulative regret (or equivalently, maximize her cumulative reward), over a specified number of rounds T. Contextual Bandits have long been studied under various notions of reward models and settings. For instance, one of the simplest models is to assume that the expected reward is a linear function of the arms and the hidden parameter (Abbasi-Yadkori et al., 2011; Auer, 2003; Chu et al., 2011). This was later extended to non-linear settings such as the logistic setting (Faury et al., 2020; Abeille et al., 2021; Faury et al., 2022), generalized linear setting (Filippi et al., 2010; Li et al., 2017), and the multinomial setting (Amani & Thrampoulidis, 2021; Zhang & Sugiyama, 2023). In this work, we specifically focus on the multinomial setting that can model problems with multiple outcomes, which makes this setting incredibly useful in the fields of machine and reinforcement learning, as well as, in real life.

Though significant progress has been made in designing algorithms for the contextual setting, the algorithms do not demonstrate a lot of applicability. There has been growing interest in constraining the budget available for algorithmic updates. This limited adaptivity setting is crucial in real-world applications, where frequent updates can hinder parallelism and large-scale deployment. Additionally, practical and computational constraints may make it infeasible to make policy updates at every time step. For example, in clinical trials (Group et al., 1997), the treatments made available to the patients cannot be changed with every patient. Thus, the updates are made after administering the treatment to a group of patients, observing the effects and outcomes, and then updating the treatment. We observe a similar tendency in online advertising and recommendations, where it is difficult to update the policy at each round due to resource constraints. A recent line of work (Ruan et al., 2021; Sawarni et al., 2024) has introduced algorithms for contextual bandits in the linear and generalized linear settings, respectively. They introduce algorithms for two different settings: the *batched* setting, wherein the policy update rounds are fixed at the start of the algorithm, and the *rarely-switching* algorithm, wherein the policy update rounds are decided in an adaptive fashion. Since multinomial logistic bandits are not generalized linear models, it is not clear if the algorithms developed in past works would apply in this setting. Hence, the major focus of this work is to develop algorithms with limited adaptivity for the multinomial setting. We now list our contributions:

#### 1.1 Contributions

- We propose a new algorithm B-MNL-CB, which operates in the batched setting where the contexts are generated stochastically. The algorithm achieves  $\tilde{O}(\sqrt{T})$  regret with high probability, with  $\Omega(\log \log T)$  policy updates. In order to accommodate time-varying contexts, we adapt the recently introduced concept of distributional optimal designs (Ruan et al., 2021) to the multinomial logistic setting. This is done by introducing a new scaling technique to counter the non-linearity associated with the reward function. Note that the leading term of the regret bound is free of the instance-dependent non-linearity parameter  $\kappa$ , which can scale exponentially with the instance parameters (refer to Section 2 for more details).
- Our second algorithm, RS-MNL operates in the rarely-switching setting, where the contexts are generated adversarially. The algorithm achieves  $\tilde{O}(\sqrt{T})$  regret while performing  $\tilde{O}(\log T)$  policy updates, each determined by a simple switching criterion. Further, our algorithm does not require a warmup switching criterion, unlike the rarely-switching algorithm in Sawarni et al. (2024), which helps in reducing the number of switches from  $\tilde{O}(\log^2 T)$  to  $\tilde{O}(\log T)$ .
- We empirically demonstrate the performance of our rarely-switching algorithm RS-MNL. Across a range of randomly selected instances, our algorithm achieves regret comparable to, and often better than, several logistic and multinomial logistic state-of-the-art baseline algorithms. Our algorithm manages to do so with a limited number of policy updates as compared to the baselines, which perform an update at each time round. We also empirically show that the number of switches made by our algorithm is  $\tilde{O}(\log T)$ , which is in agreement with our theoretical results.

#### 1.2 Related works

The multinomial logistic setting was first studied by Amani & Thrampoulidis (2021). They proposed an algorithm that achieved a regret bound of  $\tilde{O}(\sqrt{\kappa T})$ , where  $\kappa$  is the instance-dependent nonlinearity parameter (defined in Section 2). This was further improved by Zhang & Sugiyama (2023), who proposed a computationally efficient algorithm with a regret bound of  $\tilde{O}(\sqrt{T})$ , thus achieving a  $\kappa$ -free bound (the leading term is free of  $\kappa$ ). However, both of these algorithms face challenges in real-world deployment due to infrastructural and practical constraints associated with updating the policy at every round.

Thus, the limited adaptivity framework was introduced to combat this challenge, wherein the algorithm could only undergo a limited number of policy switches. This framework consists of two paradigms: the first being the *Batched* Setting, where the batch lengths are predetermined and was first studied by Gao et al. (2019), who showed that  $\Omega(\log \log T)$  batches are necessary to obtain optimal minimax regret. The second setting is the *Rarely Switching* Setting, first introduced by Abbasi-Yadkori et al. (2011), where batch lengths are determined adaptively, based on a switching criterion, such as the determinant doubling trick, wherein the policy is updated every time the determinant of the information matrix doubles.

In the contextual setting, Ruan et al. (2021) used optimal designs to study the case where the arm sets themselves were generated stochastically, providing a bound of  $\tilde{O}(\sqrt{d \log dT})$  for the batched setting. This idea was then extended to the generalized linear setting by Sawarni et al. (2024), who proposed algorithms that could achieve  $\kappa$ -free regret in both the batched and rarely-switching settings (independent of  $\kappa$  in the leading term). However, to the best of our knowledge, the limited adaptivity framework has not yet been explored in the multinomial setting. The primary focus of this work is to propose optimal limited-adaptivity algorithms for the multinomial setting. We achieve this by extending the results of Sawarni et al. (2024) and Ruan et al. (2021) to the multinomial setting in the batched setting while preserving the regret bound of Zhang & Sugiyama (2023) in the first-order term. In the rarely-switching setting, we further build upon the work of Abbasi-Yadkori et al. (2011) and Sawarni et al. (2024) to adapt it for the multinomial case. This maintains the regret bound of Zhang & Sugiyama (2023) while also reducing the number of switches as compared to Sawarni et al. (2024).

### 2 Preliminaries

Notations: We denote all vectors with bold lower case letters, matrices with bold upper case letters, and sets with upper case calligraphic symbols. We write  $M \succeq 0$ , if matrix M is positive semidefinite (p.s.d). For a p.s.d matrix M, we define the norm of a vector x with respect to M as  $||x||_M = \sqrt{x^\top M x}$  and the spectral norm of M as  $||M||_2 = \sqrt{\lambda_{max} (M^\top M)}$  where  $\lambda_{max} (M)$ denotes the maximum eigenvalue of M. We denote the set  $\{1, \ldots, N\}$  as [N]. The Kronecker product of matrices  $A \in \mathbb{R}^{p \times q}$  and  $B \in \mathbb{R}^{r \times s}$  is defined as  $(A \otimes B)_{pr+v, qs+w} = A_{rs} \cdot B_{vw}$ , resulting in a  $pr \times qs$  matrix, where  $M_{ij}$  denotes the element of the matrix M present at the  $i^{th}$  row and the  $j^{th}$  column. Finally, we use  $\Delta(\mathcal{X})$  to denote the set of all probability distributions over  $\mathcal{X}$ . We use  $I_n$  to denote an identity matrix of dimension n, and we simply use I when the dimensions are clear from context.

**Multinomial Logistic Bandits**: In the Multinomial Logistic Bandit Setting, at each round t, the learner is presented with a set of arms  $\mathcal{X}_t \subseteq \mathbb{R}^d$ , and is expected to choose an arm  $\mathbf{x}_t \in \mathcal{X}_t$ . Based on the learner's choice, the environment provides an outcome  $y_t \in [K] \cup \{0\}^1$ . While choosing the arm at round t, the learner can utilize all prior information, which can be encoded in the filtration  $\mathcal{F}_t = \sigma(\mathcal{F}_0, \mathbf{x}_1, y_1, \dots, \mathbf{x}_{t-1}, y_{t-1})$ , where  $\mathcal{F}_0$  represents any prior information the learner had before starting the algorithm. The probability distribution over these K + 1 outcomes is modeled using a multinomial logistic function<sup>2</sup> as follows:

$$\mathbb{P}\left\{y_t = i \mid \boldsymbol{x}_t, \mathcal{F}_t\right\} = \begin{cases} \frac{\exp\left(\boldsymbol{x}_t^{\top} \boldsymbol{\theta}_i^*\right)}{1 + \sum\limits_{j=1}^{K} \exp\left(\boldsymbol{x}_t^{\top} \boldsymbol{\theta}_j^*\right)}, & 1 \le i \le K\\ \frac{1}{1 + \sum\limits_{j=1}^{K} \exp\left(\boldsymbol{x}_t^{\top} \boldsymbol{\theta}_j^*\right)}, & i = 0, \end{cases}$$

where  $\boldsymbol{\theta}^{\star} = \left(\boldsymbol{\theta}_{1}^{\star \top}, \dots, \boldsymbol{\theta}_{K}^{\star \top}\right)^{\top} \in \mathbb{R}^{dK}$  comprises the hidden optimal parameter vectors associated with each of the *K* outcomes. Based on the outcome  $y_t$ , the learner receives a reward  $\rho_{y_t} \geq 0$ . It is standard to set  $\rho_0 = 0$ . We assume that the reward vector  $\boldsymbol{\rho} = (\rho_1, \dots, \rho_K)^{\top}$  is fixed and known. We assume that  $||\boldsymbol{\theta}^{\star}||_2 \leq S$ ,  $||\boldsymbol{\rho}||_2 \leq R$ , and  $||\boldsymbol{x}||_2 \leq 1$ , for all  $\boldsymbol{x} \in \mathcal{X}_t$ , where *R* and *S* are fixed and

<sup>&</sup>lt;sup>1</sup>The outcome 0 indicates *no outcome*.

<sup>&</sup>lt;sup>2</sup>The multinomial logistic function is also referred to as the link function and would be used interchangeably throughout.

known beforehand. Note that when K = 1, the problem reduces to the binary logistic setting. For simplicity, we denote the probability of the  $i^{\text{th}}$  outcome  $\mathbb{P} \{y_t = i \mid x_t, \mathcal{F}_t\}$  as  $z_i(x_t, \theta^*)$  and denote the probability vector over the K outcomes as  $\mathbf{z}(\mathbf{x}_t, \theta^*) = (z_1(\mathbf{x}_t, \theta^*), \dots, z_K(\mathbf{x}_t, \theta^*))^{\top}$ . Then, it is easy to see that the expected reward of the learner at round t is given by  $\boldsymbol{\rho}^\top \mathbf{z}(\mathbf{x}_t, \theta^*)$ . The goal of the learner is to choose an arm  $\mathbf{x}_t, t \in [T]$  so as to minimize her regret, which can have different formulations based on the problem setting:

1. Stochastic Contextual setting : In this setting, at each time step, the feasible action sets are sampled from the same (unknown) distribution  $\mathcal{D}$ . Thus, the learner wishes to minimize her expected cumulative regret which is given by

$$R(T) = \mathbb{E}\left[\sum_{t=1}^{T} \left[\max_{\boldsymbol{x} \in \mathcal{X}_t} \boldsymbol{\rho}^\top \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}^\star) - \boldsymbol{\rho}^\top \boldsymbol{z}(\boldsymbol{x}_t, \boldsymbol{\theta}^\star)\right]\right]$$

Here, the expectation is over the distribution of the arm set  $\mathcal{D}$  and the randomness inherently present in the algorithm. In this setting, we assume that only M (fixed beforehand) policy updates can be made and the rounds at which these updates can happen need to be decided prior to starting the algorithm.

2. Adversarial Contextual setting : In this setting, there are no assumptions made on how the feature vectors of the arms are generated. Thus, allowing M policy updates, the algorithm can choose the rounds at which it updates its policy during the course of the algorithm. These dynamic updates are based on a simple switching criterion similar to the one presented in Abbasi-Yadkori et al. (2011). In this setting, the learner wishes to minimize her cumulative regret given by

$$R(T) = \sum_{t=1}^{T} \left[ \max_{\boldsymbol{x} \in \mathcal{X}_t} \boldsymbol{\rho}^\top \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}^\star) - \boldsymbol{\rho}^\top \boldsymbol{z}(\boldsymbol{x}_t, \boldsymbol{\theta}^\star) \right].$$

**Discussion on the Instance-Dependent Non-Linearity Parameter**  $\kappa$ **:** Several works on the binary logistic model and generalized linear model (Filippi et al., 2010; Faury et al., 2020) as well as the multinomial logistic model (Amani & Thrampoulidis, 2021; Zhang & Sugiyama, 2023) have mentioned the importance of an instance dependent, non-linearity parameter  $\kappa$ , and have stressed on the need to obtain regret guarantees independent of  $\kappa$  (at least in the leading term).  $\kappa$  was first defined for the binary logistic reward model setting Filippi et al. (2010). A natural extension to the multinomial logit setting was recently proposed in Amani & Thrampoulidis (2021). We use the same definition as Amani & Thrampoulidis (2021), i.e.,

$$\kappa = \sup \left\{ rac{1}{\lambda_{min}(\boldsymbol{A}(\boldsymbol{x}, \boldsymbol{ heta}))} : \boldsymbol{x} \in \mathcal{X}_1 \cup \ldots \cup \mathcal{X}_T, \boldsymbol{ heta} \in \Theta 
ight\},$$

where  $A(x, \theta) = \nabla z(x, \theta) = diag(z(x, \theta)) - z(x, \theta)z(x, \theta)^{\top}$ , is the gradient of the link function z with respect to the vector  $x^{\top}\theta$ . In Section 2, Faury et al. (2020), it was highlighted that that  $\kappa$  can grow exponentially in the instance parameters such as S and therefore regret proportional to  $\kappa$  could be detrimental when these parameters are large. In Section 3 of Amani & Thrampoulidis (2021), the authors show that  $\kappa$  in the multinomial setting also scales exponentially with the diameter of the parameter and action sets. We direct the reader to Section 3 of Amani & Thrampoulidis (2021) for a more elaborate discussion on the importance of  $\kappa$  in the multinomial setting.

**Optimal Design policies:** Optimal Experimental Designs are concerned with efficiently selecting the best data points so as to minimize the variance (or equivalently, maximize the information) of estimated parameters. For a set of points  $\mathcal{X} \subseteq \mathbb{R}^d$  and some distribution  $\pi$  defined on  $\mathcal{X}$ , The information matrix is defined as  $(\mathbb{E}_{x \sim \pi} x x^T)^{-1}$ . Several criteria are used to maximize the information, some of which are A-Criterion (minimize trace of the information matrix), E-Criterion (maximize the information matrix), and D-Criterion (maximize the determinant of the information matrix). One of the popular criteria used in bandit literature is the G-Optimal Design which is defined as follows:

**Definition 2.1. G-Optimal Design:** For a set  $\mathcal{X} \subseteq \mathbb{R}^d$ , the G-Optimal design  $\pi_G(\mathcal{X})$  is the solution to the following optimization problem:

$$\min_{\pi \in \Delta(\mathcal{X})} \max_{\boldsymbol{x} \in \mathcal{X}} \|\boldsymbol{x}\|_{V(\pi)^{-1}}, \quad \text{where} \quad \boldsymbol{V}(\pi) = \mathbb{E}_{\boldsymbol{x} \sim \pi} [\boldsymbol{x} \boldsymbol{x}^{\top}].$$

The General Equivalence Theorem (Kiefer & Wolfowitz, 1960; Lattimore & Szepesvári, 2020) establishes an equivalence between the G-Optimal and D-Optimal criteria. Specifically, it shows that for any set  $\mathcal{X} \subseteq \mathbb{R}^d$ , there exists a G-Optimal design  $\pi_G(\mathcal{X}) \in \Delta(\mathcal{X})$  such that:

$$\|\boldsymbol{x}\|_{V(\pi)^{-1}} \leq d \quad \forall \boldsymbol{x} \in \mathcal{X}.$$

Furthermore, if  $\mathcal{X}$  is a discrete set with finite cardinality, one can find a G-Optimal design in polytime with respect to  $|\mathcal{X}|$  such that the right-hand side can be relaxed to 2d (Lemma 3, Ruan et al. (2021)).

**Distributional Optimal design**: The extension of the G-Optimal design to the stochastic contextual setting worsens the bound on  $\|\boldsymbol{x}\|_{V(\pi)^{-1}}$ , i.e, in the worst case, the expected value of  $\|\boldsymbol{x}\|_{V(\pi)^{-1}}$  is upper bounded by  $d^2$ , where the expectation is over the arm set  $\mathcal{X}$ . To address this, Ruan et al. (2021) introduces the Distributional Optimal Design, formalized in the following result:

**Lemma 2.1.** (Theorem 5, Ruan et al. (2021)) Let  $\pi$  be the DISTRIBUTIONAL OPTIMAL DESIGN policy that has been learned from s independent samples  $\mathcal{X}_1, \ldots, \mathcal{X}_s \sim \mathcal{D}$ . Let V denote the expected design matrix,

$$oldsymbol{V} = \mathop{\mathbb{E}}\limits_{\mathcal{X}\sim\mathcal{D}} \mathop{\mathbb{E}}\limits_{oldsymbol{x}\sim\pi(\mathcal{X})} \left[oldsymbol{x}oldsymbol{x}^{ op} \mid \mathcal{X}
ight].$$

Then,

$$\mathbb{P}\left(\mathbb{E}_{\mathcal{X}\sim\mathcal{D}}\left[\max_{\boldsymbol{x}\in\mathcal{X}}\|\boldsymbol{x}\|_{\boldsymbol{V}^{-1}}\right] \leq \mathcal{O}(\sqrt{d\log d})\right) \geq 1 - \exp\left(\mathcal{O}(d^4\log^2 d - sd^{-1.2}\cdot 2^{-16})\right).$$

We utilize the **CoreLearning for Distributional G-Optimal Design** algorithm (Algorithm 3, Ruan et al. (2021)) to learn the distributional optimal design over a given set of context vectors. In this paper, we extend both the G-Optimal and Distributional Optimal Design frameworks to the multinomial logistic (MNL) setting by introducing *directionally scaled sets*. These sets are then used to construct the design policies employed in our batched algorithm.

#### **3** Batched Multinomial Contextual Bandit Algorithm: B-MNL-CB

In this section, we present our first algorithm, B-MNL-CB. This section is structured in the following manner: we introduce the algorithm and explain each step in detail. This is followed by a few salient remarks and the regret guarantee for the algorithm. We provide a proof sketch for this guarantee and guide the reader to the full proof in the appendix.

B-MNL-CB operates in the stochastic contextual setting (described in Section 2), building upon BATCHLINUCB-DG (Algorithm 5, Ruan et al. (2021)) and B-GLinCB (Algorithm 1, Sawarni et al. (2024)), both of which are batched algorithms. In the batched paradigm, the rounds at which the policy updates occur are fixed beforehand. We will refer to all the rounds between two consecutive policy updates as a *batch*. The horizon is divided into  $M = O(\log \log T)$  disjoint batches denoted by  $\{\mathcal{T}_{\beta}\}_{\beta=1}^{M}$ , and the lengths of the batches are denoted by  $\tau_{\beta} = |\mathcal{T}_{\beta}|$ .

The input to B-MNL-CB includes the number of batches M, the fixed (known) reward vector  $\rho$ , the known upper bound on  $||\theta^*||_2$  denoted by S, and the total number of rounds T. We denote the policy learned in each batch  $\beta$  by  $\pi_\beta$ , initializing  $\pi_0$  with the G-Optimal design. We also initialize  $\lambda$  to  $\sqrt{Kd \log T}$  and define the batch lengths  $\{\tau_\beta\}_{\beta=1}^M$  as follows:

$$\tau_{\beta} = \lfloor T^{1-2^{-\beta}} \rfloor \,\forall \beta \in [1, M]. \tag{1}$$

#### Algorithm 1 Batched Multinomial Contextual Bandit Algorithm: B-MNL-CB

1: **Input:** M,  $\rho$ , S, T2: Initialize  $\{\mathcal{T}_m\}_{m=1}^M$  as per 1,  $\lambda = \sqrt{Kd \log T}$ , and policy  $\pi_0$  as G-OPTIMAL DESIGN 3: for batches  $\beta \in [M]$  do 4: for each round  $t \in \mathcal{T}_{\beta}$  do Observe arm set  $\mathcal{X}_t$ 5: for j = 1 to  $\beta - 1$  do 6: Update arm set  $\mathcal{X}_t \leftarrow \mathrm{UL}_i(\mathcal{X}_t)$  (defined in 5) 7: 8: end for 9: Sample  $x_t \sim \pi_{\beta-1}(\mathcal{X}_t)$  and obtain outcome  $y_t$  along with the corresponding reward  $\rho_{y_t}$ . 10: end for Divide  $\mathcal{T}_{\beta}$  into two sets C and D such that  $|C| = |D|, C \cup D = \mathcal{T}_{\beta}$ , and  $C \cap D = \emptyset$ . Compute  $\hat{\theta}_{\beta} \leftarrow \arg \min \sum_{s \in C} \ell(\theta, \boldsymbol{x}_{s}, y_{s}), \boldsymbol{H}_{\beta} = \lambda \boldsymbol{I} + \sum_{s \in C} \frac{\boldsymbol{A}(\boldsymbol{x}_{t}, \hat{\theta}_{\beta}) \otimes \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top}}{B_{\beta}(\boldsymbol{x}_{t})}$ , and  $\pi_{\beta}$  using Algorithm 2 with the inputs  $(\beta, \{\mathcal{X}_{t}\}_{t \in D})$ 11: 12: 13: end for

We now provide a detailed explanation of the steps involved in the algorithm. In *Steps 3-13*, we iterate over all batches  $\beta \in [M]$  and rounds  $t \in \mathcal{T}_{\beta}$ . During batch  $\beta$  and round  $t \in \mathcal{T}_{\beta}$ , first, in *Step 5*, we obtain the set of feasible arms  $\mathcal{X}_t$  at round t. Then in *Steps 6-8*, we iterate over all the previous batches  $j \in [\beta - 1]$  to prune  $\mathcal{X}_t$  and retain only a subset of it via a *Successive Elimination* procedure described next.

#### 3.1 Successive Eliminations

For each prior batch  $j \in [\beta - 1]$ , we compute an upper confidence bound UCB $(j, x, \lambda)$  and a lower confidence bound LCB $(j, x, \lambda)$  as follows,

$$UCB(j, \boldsymbol{x}, \lambda) = \boldsymbol{\rho}^T \hat{\boldsymbol{\theta}}_j + \epsilon_1(j, \boldsymbol{x}, \lambda) + \epsilon_2(j, \boldsymbol{x}, \lambda),$$
(2)

$$LCB(j, \boldsymbol{x}, \lambda) = \boldsymbol{\rho}^T \hat{\boldsymbol{\theta}}_j - \epsilon_1(j, \boldsymbol{x}, \lambda) - \epsilon_2(j, \boldsymbol{x}, \lambda),$$
(3)

where the bonus terms  $\epsilon_1(j, \boldsymbol{x}, \lambda)$  and  $\epsilon_2(j, \boldsymbol{x}, \lambda)$  are defined as,

$$\epsilon_1(j,\boldsymbol{x},\boldsymbol{\lambda}) = \gamma(\boldsymbol{\lambda}) \|\boldsymbol{H}_j^{-\frac{1}{2}}(\boldsymbol{I} \otimes \boldsymbol{x})\boldsymbol{A}(\boldsymbol{x},\hat{\boldsymbol{\theta}}_j)\boldsymbol{\rho}\|_2, \\ \epsilon_2(j,\boldsymbol{x},\boldsymbol{\lambda}) = 3\gamma(\boldsymbol{\lambda})^2 \|\boldsymbol{\rho}\|_2 \|(\boldsymbol{I} \otimes \boldsymbol{x}^\top)\boldsymbol{H}_j^{-\frac{1}{2}}\|_2^2.$$
(4)

Here,  $\theta_j$  and  $H_j$  are the estimators (computed during *Steps 11,12* at the end of batch j) of the true parameter vector  $\theta^*$  and an optimal batch-level Hessian matrix  $H_j^*$  and  $\gamma(\lambda)$  is defined to be  $O(\sqrt{Kd \log T})$ . We provide more details on these in Section 3.2. In *Step 7*, for batch j, we eliminate a subset of  $\mathcal{X}_t$  using the upper and lower confidence bounds just defined. In particular, we eliminate all  $\boldsymbol{x} \in \mathcal{X}_t$  for which  $UCB(j, \boldsymbol{x}, \lambda) \leq \max_{\boldsymbol{x}'} LCB(j, \boldsymbol{x}', \lambda)$ . Thus, in *Step 7*,  $\mathcal{X}_t$  is updated to  $UL_j(\mathcal{X})$ , defined as,

$$\mathrm{UL}_{j}(\mathcal{X}) = \mathcal{X} \setminus \left\{ \boldsymbol{x} \in \mathcal{X} : \mathrm{UCB}(j, \boldsymbol{x}, \lambda) \leq \max_{\boldsymbol{y} \in \mathcal{X}} \mathrm{LCB}(j, \boldsymbol{y}, \lambda) \right\},\tag{5}$$

Following the successive eliminations over all prior batches  $j \in [\beta - 1]$ , in *Step 9*, we select an arm  $x_t$  from the pruned arm set according to the policy computed at the end of batch  $\beta - 1$  using Algorithm 2. The environment then returns the outcome  $y_t$  and the corresponding reward  $\rho_{y_t}$ . Details of the policy computation (Algorithm 2) are provided in Section 3.3. After completing all rounds in batch  $\beta$  (i.e.,  $\mathcal{T}_{\beta}$ ), we proceed to *Step 11*, where we partition these rounds equally into two sets, C and D. The set C is used to define a batch-level Hessian matrix  $H_{\beta}^{\star}$ , compute an estimator  $\hat{\theta}_{\beta}$  of  $\theta^{\star}$ , and construct a matrix  $H_{\beta}$  that estimates  $H_{\beta}^{\star}$  as described in the next section.

#### 3.2 Batch Level Hessian and Parameter Estimation

In batch  $\beta$ , we define a batch level Hessian matrix  $H_{\beta}^{\star} = \lambda I + \sum_{t \in C} A(x_t, \theta^{\star}) \otimes x_t x_t^{\top}$ , which is constructed using the set C. Since  $\theta^{\star}$  is unknown, we maintain an online proxy to estimate  $H_{\beta}^{\star}$ by calculating a scaled Hessian matrix  $H_{\beta} = \lambda I + \sum_{t \in C} \frac{A(x_t, \hat{\theta}_{\beta})}{B_{\beta}(x)} \otimes x_t x_t^{\top}$ . Here,  $B_{\beta}(x)$  is a normalizing factor which is obtained using the self-concordance properties of the link function and is given by:

$$B_{\beta}(\boldsymbol{x}) = \exp\left(\sqrt{6}\min\left\{\gamma(\lambda)\sqrt{\kappa} ||\boldsymbol{x}||_{\boldsymbol{V}_{\beta}^{-1}}, 2S\right\}\right),\tag{6}$$

where  $\gamma(\lambda) = \mathcal{O}(\sqrt{Kd \log T})$  is the confidence radius for the permissible set of  $\boldsymbol{\theta}$  and  $V_{\beta}$  is the design matrix given by  $V_{\beta} = \lambda I + \sum_{t \in C} \boldsymbol{x}_t \boldsymbol{x}_t^{\top}$ . Using the self-concordance properties of the link function, we can show that  $H_{\beta} \preccurlyeq H_{\beta}^{*}$ . The set *C* is also used to update the estimator  $\hat{\boldsymbol{\theta}}_{\beta}$ , which is done by minimizing the negative log likelihood  $\sum_{t \in C} \ell(\boldsymbol{\theta}, \boldsymbol{x}_t, y_t)$ , where  $\ell(\boldsymbol{\theta}, \boldsymbol{x}, y)$  is defined as,

$$\ell(\boldsymbol{\theta}, \boldsymbol{x}, y) = -\sum_{i=1}^{K} \mathbb{1}\left\{y = i\right\} \log \frac{1}{z_i(\boldsymbol{x}, \boldsymbol{\theta})} + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2,$$
(7)

Next, we explain how the policy is updated to  $\pi_{\beta}$  at the end of batch  $\beta$  using the rounds in set D.

#### 3.3 Policy calculation

Algorithm 2 Distributional Optimal Design for MNL bandits

- 1: Input Batch  $\beta$  and collection of arm sets  $\{X_j\}_j$
- 2: Create the sets  $\{F_i(\{\mathcal{X}_j\}_j,\beta)\}_{i=1}^K$  as defined in Equation 8.
- 3: Compute the distributional optimal design policy  $\pi_i$  for each of the sets  $F_i(\{\mathcal{X}_j\}_j, \beta)$ .
- 4: Compute the distributional optimal design policy  $\pi_0$  for the set  $\{\mathcal{X}_j\}_j$ .

5: **Return** 
$$\pi = \frac{1}{K+1} \sum_{i=0}^{K} \pi_i$$

To compute our final policy at the end of each batch, we utilize the idea of distributional optimal design, first introduced in Ruan et al. (2021) (See Section 2). Recently, Sawarni et al. (2024) used distributional optimal designs to develop limited adaptivity algorithms for stochastic contextual bandits for generalized linear bandits. A key step in their algorithm (Step 13 and Equation 4, Algorithm 1 in Sawarni et al. (2024)) involves scaling the arm set (after pruning using *successive eliminations*) with the derivative of the link function and a suitable normalization factor. Generalizing this idea to the MNL setting results in a matrix  $\tilde{X} = \frac{A(x, \hat{\theta}_t)^{\frac{1}{2}}}{B_{\beta}(x)} \otimes x$ . Note that the notion of distributional optimal designs introduced in Ruan et al. (2021) and used by Sawarni et al. (2024), applies only to vectors. Hence, in Algorithm 2, we construct several sets of vectors from  $\tilde{X}$  and learn the optimal design for each of these sets.

In *Step 12* of Algorithm 1, we invoke this algorithm (Algorithm 2) with inputs as the batch number  $\beta$  and the collection of all the pruned arm sets  $\{\mathcal{X}_t\}_{t\in D}$  (*Step 7*, Algorithm 1). We then create K different sets  $F_i(\{\mathcal{X}_t\}_{t\in D}, \beta)$  ( $i \in [K]$ ), which comprises of the arms in each arm set scaled by the  $i^{\text{th}}$  column of the gradient matrix. In particular,

$$F_i(\{\mathcal{X}_t\}_{t\in D},\beta) = \left\{ \left\{ \frac{\boldsymbol{A}(\boldsymbol{x},\hat{\boldsymbol{\theta}}_{\beta})^{\frac{1}{2}}}{\sqrt{B_{\beta}(\boldsymbol{x})}} \boldsymbol{e}_i \otimes \boldsymbol{x} : \boldsymbol{x} \in \mathcal{X}_t \right\} : t \in D \right\},\tag{8}$$

where  $e_i \in \mathbb{R}^K$  is the *i*<sup>th</sup> standard basis vector. We calculate the distributional optimal design for each of the sets  $F_i(\{\mathcal{X}_t\}_{t\in D}, \beta)$  using Algorithm 2 in Ruan et al. (2021). In such a case, it is easy to see that calculating the distributional optimal design over  $\tilde{X}$  can be done by calculating the distributional optimal designs for each of the sets  $F_i({\mathcal{X}_t}_{t\in D}, \beta)$ . We provide the proof for the same in Section 7.3. We also calculate the distributional optimal design over  ${\mathcal{X}_t}_{t\in D}$ . Finally, the policy returned is a convex combination (in this case, a uniform combination) over all the K + 1 designs that were calculated.

This completes our explanation of Algorithm 1. We provide a regret guarantee in Theorem 3.1.

**Remark 3.1.** A direct application of the scaling techniques introduced in Sawarni et al. (2024) for learning distributional optimal designs in the multinomial setting results in the creation of a scaled matrix. Since the notion of distributional optimal design introduced in Ruan et al. (2021) applies only to vectors, Algorithm 2 scales the original context vectors into K different sets and then learns the optimal designs for each of them.

**Remark 3.2.** Sawarni et al. (2024) introduces a warm-up round with length  $O(\kappa^{1/3})$ . Since  $\kappa$  can scale exponentially with several instance-dependent parameters, the warm-up round can result in a long exploration phase. Using the regret decomposition in Zhang & Sugiyama (2023), we can eliminate the dependence on  $\kappa$ , resulting in  $\kappa$ -free batch lengths, including the length of the warm-up round.

**Remark 3.3.** While Zhang & Sugiyama (2023) introduced a novel method of regret decomposition into the error terms (refer 4), a straightforward application to the limited adaptivity setting is not easy. Hence, with some additional insights, we incorporate their method into the batched setting while being able to match the leading term of their regret bound.

**Theorem 3.1.** (Regret of B-MNL-CB) With high probability, at the end of T rounds, the regret incurred by Algorithm 1 is bounded as  $R_T \leq R_1 + R_2$  where

$$R_1 = \tilde{O}\left(RS^{5/4}K^{5/2}d\sqrt{T}\right) \text{ and } R_2 = \tilde{O}\left(RS^{5/2}K^2d^2\kappa^{1/2}T^{1/4}\max\{e^{3S}K^{3/2}S^{-1},\kappa^{1/2}d\}\right).$$

#### **Proof Sketch:**

We know that the expected regret during batch  $\beta + 1$  is given by:

$$R_{\beta+1} = \mathbb{E}\left[\sum_{t\in\beta} \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_t^{\star}, \boldsymbol{\theta}^{\star}) - \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_t, \boldsymbol{\theta}^{\star})\right],$$

where  $x_t^{\star} = \underset{x \in \mathcal{X}_t}{\operatorname{arg\,max}} \rho^{\top} z(x, \theta^{\star})$  is the best arm at round t and the expectation is taken over the distribution of the arm set  $\mathcal{D}$ . Using ideas similar to Zhang & Sugiyama (2023), we can decompose the regret into

$$R(T) \leq 4 \sum_{t \in \beta} \left\{ \mathbb{E} \left[ \max_{\boldsymbol{x} \in \mathcal{X}_t} \epsilon_1(\beta, \boldsymbol{x}, \lambda) \right] + \mathbb{E} \left[ \max_{\boldsymbol{x} \in \mathcal{X}_t} \epsilon_2(\beta, \boldsymbol{x}, \lambda) \right] \right\},\$$

where  $\epsilon_1(\beta, \boldsymbol{x}, \lambda)$  and  $\epsilon_2(\beta, \boldsymbol{x}, \lambda)$  are as defined in 4. We proceed to bound each of these terms using the extension of distributional optimal design we introduced in Algorithm 2.

Directly extending the ideas of Ruan et al. (2021) and Sawarni et al. (2024) to construct the distributional optimal designs results in an attempt to learn the design for matrices  $\tilde{X}_{\beta} = \frac{A(\boldsymbol{x}, \hat{\theta}_{\beta})^{\frac{1}{2}}}{B_{\beta}(\boldsymbol{x})} \otimes \boldsymbol{x}$ . Hence, we create K different sets  $F_i(\mathcal{X})$  for all  $i \in [K]$  (defined in 8), such that

$$ilde{oldsymbol{X}}_eta ilde{oldsymbol{X}}_eta^ op = \sum_{i=1}^K \left\{ rac{oldsymbol{A}(oldsymbol{x}, \hat{oldsymbol{ heta}}_eta)^rac{1}{2}}{\sqrt{B_eta(oldsymbol{x})}} oldsymbol{e}_i \otimes oldsymbol{x} 
ight\} \left\{ rac{oldsymbol{A}(oldsymbol{x}, \hat{oldsymbol{ heta}}_eta)^rac{1}{2}}{\sqrt{B_eta(oldsymbol{x})}} oldsymbol{e}_i \otimes oldsymbol{x} 
ight\}^T.$$

Thus, learning the optimal design over  $\tilde{X}$  is equivalent to creating a convex combination of the designs learned over  $F_i(\mathcal{X})$  for all  $i \in [K]$ . This gives us a way of bounding the scaled Hessian

matrix  $H_{\beta}$  by the scaled Hessian matrices  $H^i_{\beta}$  constructed over  $F_i(\mathcal{X})$  for all  $i \in [K]$ . We then use methods similar to Sawarni et al. (2024) and Ruan et al. (2021) to obtain the bound on the regret for the batch  $\beta + 1$  as:

$$\begin{aligned} R_{\beta+1} &\leq +32RK\kappa^{1/2}d\gamma^2(\lambda)\left\{e^{3S}K^{3/2}S^{-1}\sqrt{\log(Kd)\log d} + 12\kappa^{1/2}d\right\}\left(\frac{\tau_{\beta+1}}{\tau_{\beta}}\right) \\ &+ 16RK^2\gamma(\lambda)\sqrt{d\log(Kd)}\left(\frac{\tau_{\beta+1}}{\sqrt{\tau_{\beta}}}\right) \end{aligned}$$

Finally, using the batch lengths defined in 1 and summing over all the M batches completes the proof. For the sake of brevity, we provide the complete proof in Section 7.

#### 4 Rarely Switching Multinomial Contextual Bandit Algorithm: RS-MNL

Algorithm 3 RS-MNL

1: Inputs:  $\rho, S, T$ 2: Initialize:  $H_1 = \lambda I$ ,  $\tau = 1$ ,  $\lambda := KdS^{-1/2}\log(T/\delta)$ ,  $\gamma := CS^{5/4}\sqrt{Kd\log(T/\delta)}$ 3: for t = 1, ..., T do Observe arm set  $\mathcal{X}_t$ 4: 5: if  $det(\boldsymbol{H}_t) > 2 det(\boldsymbol{H}_{\tau})$  then Set  $\tau = t$ 6: Update  $\hat{\boldsymbol{\theta}}_{\tau} \leftarrow \arg\min_{\boldsymbol{\theta}} \sum_{s \in [t-1]} \ell(\boldsymbol{\theta}, \boldsymbol{x}_s, y_s)$  and  $\boldsymbol{H}_t = \sum_{s \in [t-1]} \frac{\boldsymbol{A}(\boldsymbol{x}_s, \hat{\theta}\tau)}{B_{\tau}(\boldsymbol{x}_s)} \otimes \boldsymbol{x}_s \boldsymbol{x}_s^{\top} + \lambda \boldsymbol{I}_{Kd}$ 7: 8: end if Select  $\boldsymbol{x}_t = \arg \max \operatorname{UCB}(t, \tau, \boldsymbol{x})$ , observe  $y_t$ , and update  $\boldsymbol{H}_{t+1} \leftarrow \boldsymbol{H}_t + \frac{\boldsymbol{A}(\boldsymbol{x}_t, \hat{\theta}_{\tau})}{B_{\tau}(\boldsymbol{x}_t)} \otimes \boldsymbol{x}_t \boldsymbol{x}_t^{\top}$ 9: 10: end for

In this section, we present our second algorithm RS-MNL. We introduce the algorithm and explain the workings in a step-by-step fashion. We then mention a few salient remarks about our algorithm. We conclude with the regret guarantee of our algorithm, a proof sketch for the same, and guide the reader to the complete proof in the Appendix.

Our second algorithm, RS-MNL (Algorithm 3) operates in the Adversarial Contextual setting. In this setting, there are no assumptions on the generation of the feature vectors. RS-MNL also limits the number of policy updates in a rarely-switching fashion, i.e, the rounds where these updates are made are decided dynamically, based on a simple switching criterion, similar to the one used in Abbasi-Yadkori et al. (2011). While the algorithm is based on RS-GLinCB in Sawarni et al. (2024), a unique regret decomposition method allows for the removal of the warmup criterion, in turn, helping in the reduction in the number of switches made by the algorithm from  $O(\log^2 T)$ to  $O(\log T)$ . Further, we successfully remove the idea of *successive eliminations* based on the previous confidence regions and replace the idea with the maximization of the Upper Confidence Bound (UCB) of each arm.

The inputs to the algorithm are  $\rho$ , the fixed and known reward vector, S, the fixed and known upper bound on  $\|\theta\|_2$ , and T, the number of rounds for which the algorithm is played. In *Step 2*, we initialize the scaled Hessian matrix  $H_1$  to I,  $\lambda$  to  $KdS^{-1/2}\log(T/\delta)$ , and  $\gamma$  to  $CS^{5/4}\sqrt{Kd\log(T/\delta)}$ . Next, at every round  $t \in [T]$ , we receive the arm set  $\mathcal{X}_t$  in *Step 4*. During *Steps 5-8*, we check if the switching condition is met and update the policy accordingly.

#### 4.1 Switching Criterion and Policy Update:

We use  $\tau \leq t$  to denote the last time round at which a switch occurred for some round t. In *Step* 5, we check for the switching condition: if the determinant of the scaled Hessian matrix  $H_t = \lambda I + \sum_{s \in [t-1]} \frac{A(\boldsymbol{x}_s, \hat{\boldsymbol{\theta}}_\tau)}{B(\boldsymbol{x}_s)} \otimes \boldsymbol{x}_s \boldsymbol{x}_s^\top$  has increased by a constant factor (in this case, 2) as compared to

 $H_{\tau}$ . In case the switching condition is triggered, we set  $\tau = t$  in *Step 6* (since a switch was made in round t). We then compute  $\hat{\theta}_{\tau}$  by minimizing the negative log likelihood  $\sum_{s \in [t-1]} \ell(\theta, x_s, y_s)$  (see 7 for definition of  $\ell(\theta, x_s, y_s)$ ) over all previous rounds  $s \in [t-1]$ , and recompute the matrix  $H_t$  with respect to the newly calculated  $\hat{\theta}_{\tau}$  (*Step 7*). The switching criterion is similar to the one used in Abbasi-Yadkori et al. (2011) and helps to reduce the number of policy updates to  $O(\log T)$ .

#### 4.2 Arm Selection:

Next, in *Step 9*, we determine the arm  $x_t$  to be played based on the Upper Confidence Bound (UCB). The upper confidence bound UCB $(t, \tau, x)$  for an arm  $x \in \mathcal{X}_t$  with respect to the previous switching round  $\tau(\leq t)$  is defined as:

$$UCB(t, \tau, \boldsymbol{x}) = \boldsymbol{\rho}^T \hat{\boldsymbol{\theta}}_{\tau} + \epsilon_1(t, \tau, \boldsymbol{x}) + \epsilon_2(t, \tau, \boldsymbol{x}), \tag{9}$$

where the error terms  $\epsilon_1(t, \tau, x)$  and  $\epsilon_2(t, \tau, x)$  are defined as:

$$\epsilon_1(t,\tau,\boldsymbol{x}) = \sqrt{2}\gamma(\delta) \|\boldsymbol{H}_t^{-\frac{1}{2}}(\boldsymbol{I}\otimes\boldsymbol{x})\boldsymbol{A}(\boldsymbol{x},\hat{\boldsymbol{\theta}}_{\tau})\boldsymbol{\rho}\|_2, \ \epsilon_2(t,\tau,\boldsymbol{x}) = 6R\gamma(\delta)^2 \|(\boldsymbol{I}\otimes\boldsymbol{x}^{\top})\boldsymbol{H}_t^{-\frac{1}{2}}\|_2^2.$$
(10)

We then obtain the outcome  $y_t$ , which is sampled from  $z(x_t, \theta^*)$ , and receives the corresponding reward  $\rho_{y_t}$ . The algorithm then updates the scaled Hessian matrix  $H_{t+1}$ . In Theorem 4.1, we provide the regret guarantee for RS-MNL.

**Remark 4.1.** The goal of a rarely-switching algorithm is to reduce the number of switches (policy updates) that are made. Our algorithm successfully reduces the number of switches from  $O(\log^2 T)$  to  $O(\log T)$  due to the removal of the warm-up switching criterion. Additionally, the number of switches is independent of  $\kappa$ .

**Remark 4.2.** Similar to the batched setting, using the regret decomposition method introduced in Zhang & Sugiyama (2023) in the rarely-switching paradigm is non-trivial. We manage to extend their results to match the leading term of their regret bound while performing a switch  $O(\log T)$  times.

**Theorem 4.1.** With high probability, after T rounds, Algorithm 3 achieves the following regret:

$$R_T \le \tilde{O}\left(RK^{3/2}S^{5/4}d\sqrt{T}\right).$$

#### **Proof Sketch:**

The expression for total regret is given by

$$R(T) = \sum_{t=1}^{T} \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}^{\star}, \boldsymbol{\theta}^{\star}) - \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}),$$

where  $x_t^{\star} = \underset{x \in \mathcal{X}_t}{\arg \max} \rho^{\top} z(x, \theta^{\star})$  is the best arm at any given round t. Using a method similar to the one used in Zhang & Sugiyama (2023), we can upper bound the regret as

$$R(T) \leq 2\sum_{t=1}^{T} \left\{ \epsilon_1(t, \tau, \boldsymbol{x}_t) + \epsilon_2(t, \tau, \boldsymbol{x}_t) \right\},\,$$

where  $\epsilon_1(t, \tau, \boldsymbol{x}_t)$  and  $\epsilon_2(t, \tau, \boldsymbol{x}_t)$  are as defined in 10. We now wish to upper bound both the terms separately.

Bounding  $\epsilon_1(t, \tau, x_t)$  using the switching criterion in Abbasi-Yadkori et al. (2011) along with the selection rule in our algorithm can result in an exponential dependency in S. Sawarni et al. (2024)

was able to circumvent this exponential dependency by using an additional switching criterion, referred to as a *warmup criterion*. However, this results in the increase in the number of switches from  $O(\log T)$  to  $O(\log^2 T)$ . It also slows down the algorithm due to the successive eliminations done at each round (similar to the ones in Algorithm 1). Our algorithm gets rid of the exponential dependency from the first order term and the warm-up criterion by decomposing  $\epsilon_1(t, \tau, x_t)$  in an alternate manner, resulting in an improved runtime as well as  $O(\log T)$  switches.

We bound both  $\epsilon_1(t, \tau, \boldsymbol{x}_t)$  and  $\epsilon_2(t, \tau, \boldsymbol{x}_t)$  using an analysis similar to the one used for Theorem 3.1, where we attempt to upper bound the scaled Hessian matrix  $\boldsymbol{H}_t$  using the scaled Hessian matrices calculated over the K different scaled sets introduced in 2. Note that nowhere do these K different sets appear in the algorithm. They only serve to ease the analysis. Combining the bounds on each of the error terms finishes the proof. For the sake of brevity, we provide the complete proof in Section 8 of the Appendix.

## **5** Experiments

In this section, we compare our algorithm RS-MNL to several contextual logistic and MNL bandit algorithms<sup>3</sup>. We describe the experiments in detail below:



Experiment 1 (R(T) vs. T for the Logistic (K = 1) Setting): In this experiment, we compare our algorithm RS-MNL to several state-of-the-art contextual logistic bandit algorithms such as ada-OFU-ECOLog (Algorithm 2, Faury et al. (2022)), RS-GLinCB (Algorithm 2, Sawarni et al. (2024), OFUL-MLogB (Algorithm 2, Zhang & Sugiyama (2023)), and OFULog+ (Algorithm 1, Lee et al. (2024)). The dimension of the arms d is set to 3 and the number of outcomes K is set to 1, which reduces the problem to the logistic setting. The arm set  $\mathcal{X}$  is constructed by sampling 10 different arms from  $[-1, 1]^3$  and normalizing them to unit vectors. The optimal parameter  $\theta^*$  is chosen randomly from  $[-1, 1]^3$  and normalized so that  $\|\theta^*\| = S = 2$ . We run all the algorithms for  $T \in \{1000, 2450, 4500\}$  rounds and average the results over 10 different seeds (for sampling rewards). The results are plotted in Figure 1a. We see that RS-MNL is incredibly competitive with ada-OFU-ECOLog and OFULog+, while incurring much lower regret than RS-GLinCB and OFUL-MLogB. We showcase the results with two standard deviations in Section 10.

Experiment 2 (R(T) vs. T for K = 3): In this experiment, we compare our algorithm RS-MNL to OFUL-MLogB, the only algorithm that achieves an optimal ( $\kappa$ -free) regret while being computationally efficient for MNL Bandits (to the best of our knowledge). We set the number of outcomes K as 3 and the dimension of the arms d is set to 3. The arm set  $\mathcal{X}$  is constructed by sampling 10 different arms from  $[-1, 1]^3$  and normalizing them to unit vectors. The optimal parameter  $\theta^*$ is sampled from  $[-1, 1]^9$  (since  $\theta^* \in \mathbb{R}^{Kd}$  and normalized so that  $\|\theta^*\| = S = 2$ . The reward vector  $\rho$  is sampled from  $[0, 1]^3$  and normalized so that  $\|\rho\| = R = 2$ . We run both the algorithms for  $T \in \{1000, 2450, 4500\}$  rounds and average these results over 10 different seeds (for sampling rewards and  $\rho$ ). The results are plotted in 1b. We see that RS-MNL incurs much lower regret than OFUL-MLogB. We also showcase the results with two standard deviations in Section 10.

<sup>&</sup>lt;sup>3</sup>The code for the experiments can be found here.

**Experiment 3 (Number of Switches vs. T):** In this experiment, we plot the number of switches RS-MNL makes as a function of the number of rounds T. We assume that the instance is simulated in the same manner as **Experiment 2**. We run the algorithm for T = 5000 rounds and average over 10 different seeds. The results are shown in Figure 1c. We see that the number of switches made by RS-MNL exhibits a strong logarithmic dependence with  $t \in [T]$ . This is in agreement with Lemma 8.14, where we show that RS-MNL switches  $O(\log t)$  times, as compared to other algorithms, which switch (update) O(t) times.

#### 6 Conclusions and Future Work

In this paper, we present two algorithms B-MNL-CB and RS-MNL, for the multinomial logistic setting in the batched and rarely-switching paradigms, respectively. The batched setting involves fixing the policy update rounds at the start of the algorithm, while the rarely switching setting chooses the policy update rounds adaptively. Our first algorithm, B-MNL-CB manages to extend the notion of distributional optimal designs to the multinomial logit setting while being able to achieve an optimal regret of  $\tilde{O}(\sqrt{T})$  in  $\Omega(\log \log T)$  batches. Our second algorithm, RS-MNL, builds upon the rarelyswitching algorithm presented in Sawarni et al. (2024) and obtains an optimal regret of  $\tilde{O}(\sqrt{T})$ while being able to reduce the number of switches to  $O(\log T)$  using alternate ways of regret decomposition. The regret of our algorithms scales with the number of outcomes K as  $K^{7/2}$  and  $K^{5/2}$ respectively, which can be detrimental for problems with a large number of outcomes. We believe that this dependence on K can be further improved, which is an interesting line for future work.

### References

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In Advances in Neural Information Processing Systems 24 (NeurIPS), pp. 2312–2320, 2011.
- Marc Abeille, Louis Faury, and Clement Calauzenes. Instance-wise minimax-optimal algorithms for logistic bandits. In Arindam Banerjee and Kenji Fukumizu (eds.), *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130 of *Proceedings of Machine Learning Research*, pp. 3691–3699. PMLR, 13–15 Apr 2021. URL https://proceedings.mlr.press/v130/abeille21a.html.
- Sanae Amani and Christos Thrampoulidis. Ucb-based algorithms for multinomial logistic regression bandits, 2021. URL https://arxiv.org/abs/2103.11489.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. J. Mach. Learn. Res., 3(null):397–422, March 2003. ISSN 1532-4435.
- Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In Geoffrey Gordon, David Dunson, and Miroslav Dudík (eds.), Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, volume 15 of Proceedings of Machine Learning Research, pp. 208–214, Fort Lauderdale, FL, USA, 11–13 Apr 2011. PMLR. URL https://proceedings.mlr.press/v15/chulla.html.
- Louis Faury, Marc Abeille, Clément Calauzènes, and Olivier Fercoq. Improved optimistic algorithms for logistic bandits, 2020. URL https://arxiv.org/abs/2002.07530.
- Louis Faury, Marc Abeille, Kwang-Sung Jun, and Clément Calauzènes. Jointly efficient and optimal algorithms for logistic bandits, 2022. URL https://arxiv.org/abs/2201.01985.
- Sarah Filippi, Olivier Cappe, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The generalized linear case. In J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta (eds.), Advances in Neural Information Processing Systems, volume 23. Curran Associates, Inc., 2010. URL https://proceedings.neurips.cc/paper\_files/ paper/2010/file/c2626d850c80ea07e7511bbae4c76f4b-Paper.pdf.

- Zijun Gao, Yanjun Han, Zhimei Ren, and Zhengqing Zhou. Batched multi-armed bandits problem. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32, pp. 503–513. Curran Associates, Inc., 2019.
- International Stroke Trial Collaborative Group et al. The international stroke trial (ist): a randomised trial of aspirin, subcutaneous heparin, both, or neither among 19 435 patients with acute ischaemic stroke. *The Lancet*, 349(9065):1569–1581, 1997.
- Osama A. Hanna, Lin F. Yang, and Christina Fragouli. Efficient batched algorithm for contextual linear bandits with large action space via soft elimination. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2023. Curran Associates Inc.
- J. Kiefer and J. Wolfowitz. The equivalence of two extremum problems. Canadian Journal of Mathematics, 12:363–366, 1960. DOI: 10.4153/CJM-1960-030-4.
- Tor Lattimore and Csaba Szepesvári. Bandit Algorithms. Cambridge University Press, 2020.
- Junghyun Lee, Se-Young Yun, and Kwang-Sung Jun. Improved regret bounds of (multinomial) logistic bandits via regret-to-confidence-set conversion, 2024. URL https://arxiv.org/abs/2310.18554.
- L. Li, Y. Lu, and D. Zhou. Provably optimal algorithms for generalized linear contextual bandits. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 2071–2080. JMLR.org, 2017.
- Yufei Ruan, Jiaqi Yang, and Yuan Zhou. Linear bandits with limited adaptivity and learning distributional optimal design, 2021. URL https://arxiv.org/abs/2007.01980.
- Ayush Sawarni, Nirjhar Das, Siddharth Barman, and Gaurav Sinha. Generalized linear bandits with limited adaptivity, 2024. URL https://arxiv.org/abs/2404.06831.
- Yu-Jie Zhang and Masashi Sugiyama. Online (multinomial) logistic bandit: Improved regret and constant computation cost. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 29741–29782. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper\_files/paper/2023/file/ 5ef04392708bb2340cb9b7da41225660-Paper-Conference.pdf.

## **Supplementary Materials**

The following content was not necessarily subject to peer review.

Throughout the appendix, for a matrix A, we shall define  $\lambda_{max}(A)$  and  $\lambda_{min}(A)$  as the maximum and minimum eigenvalue of A respectively. Further, the norm of a matrix A is defined as  $\|A\|_2^2 =$  $\lambda_{max} (\mathbf{A}^{\top} \mathbf{A}).$ 

Without loss of generality, we also assume that  $\kappa$ , K, d, R, S, and T are greater than 1 throughout the appendix.

#### 7 Batched Multinomial Contextual Bandit Algorithm: B-MNL-CB

#### 7.1 Notations

× т

We first list a few matrices, vectors, and scalars that are commonly used throughout this section:

1. 
$$V_{\beta} = \lambda I_{d \times d} + \sum_{t \in \mathcal{T}_{\beta}} x_t x_t^{\top}$$
  
2.  $\tilde{V}_{\beta} = I_{K \times K} \otimes V_{\beta}$   
3.  $A(x, \theta) = \operatorname{diag}(z(x, \theta)) - z(x, \theta)z(x, \theta)^{\top}$   
4.  $M(x, \theta_1, \theta_2) = \int_0^1 A(x, v\theta_1 + (1 - v)\theta_2) \, dv$   
5.  $H_{\beta}^* := \lambda I_{Kd \times Kd} + \sum_{t \in \mathcal{T}_{\beta}} A(x_t, \theta^*) \otimes x_t x_t^{\top}$   
6.  $\gamma(\lambda) = 12S\sqrt{\log T + Kd} + 8S\lambda^{-1/2}(\log T + Kd) + 2S^{3/2}\lambda^{1/2}$   
7.  $B_{\beta}(x) = \exp\left(\sqrt{6}\min\left\{\kappa^{1/2}\gamma(\delta) ||x||_{V_{\beta}^{-1}}, 2S\right\}\right)$   
8.  $H_{\beta} = \lambda I_{Kd \times Kd} + \sum_{t \in \mathcal{T}_{\beta}} \frac{A(x_t, \hat{\theta}_{\beta})}{B_{\beta}(x_t)} \otimes x_t x_t^{\top}$   
9.  $\tilde{X}_{\beta} = \frac{A(x, \hat{\theta}_{\beta})^{1/2}}{\sqrt{B_{\beta}(x)}} \otimes x$   
10.  $\tilde{x}_{\beta}^{(i)} = \frac{A(x, \hat{\theta}_{\beta})^{1/2}}{\sqrt{B_{\beta}(x)}} e_i \otimes x$   
11.  $m_s = (\mathbbm{1}\{y_s = 1\}, \dots, \mathbbm{1}\{y_s = K\})^{\top}$ 

We now present the regret upper bound for B-MNL-CB by restating Theorem 3.1:

**Theorem 7.1.** (Regret of B-MNL-CB) With high probability, at the end of T rounds, the regret incurred by Algorithm 1 is bounded above by  $R_T$  where

$$R_T \le \tilde{\mathcal{O}}\left(RS^{5/4}K^{5/2}d\sqrt{T} + RS^{5/2}K^2d^2\kappa^{1/2}T^{1/4}\max\{e^{3S}K^{3/2}S^{-1}, \kappa^{1/2}d\}\right)$$

*Proof.* From Lemma 7.10, we have an upper bound for the regret incurred for any round  $t \in \mathcal{T}_{\beta+1}$ . Thus, the regret incurred in batch  $\beta + 1$  is given by:

$$R_{\beta+1} \leq 16RK^2\gamma(\lambda)\sqrt{d\log(Kd)} \left(\frac{\tau_{\beta+1}}{\sqrt{\tau_{\beta}}}\right) + 32RK\kappa^{1/2}d\gamma^2(\lambda) \left\{e^{3S}K^{3/2}S^{-1}\sqrt{\log(Kd)\log d} + 12\kappa^{1/2}d\right\} \left(\frac{\tau_{\beta+1}}{\tau_{\beta}}\right)$$

Choosing the batch lengths as  $\tau_{\beta} = T^{1-2^{-\beta}}$  results in the following observation (Hanna et al., 2023; Gao et al., 2019):

$$\frac{\tau_{\beta+1}}{\sqrt{\tau_{\beta}}} \le 2\sqrt{T} \qquad \qquad \frac{\tau_{\beta+1}}{\tau_{\beta}} \le T^{\frac{1}{4}}$$

Thus, the regret incurred in batch  $\beta + 1$  is bounded by:

$$R_{\beta+1} \le 32RK^2\gamma(\lambda)\sqrt{d\log(Kd)}\sqrt{T} + 32RK\kappa^{1/2}d\gamma^2(\lambda)\left\{e^{3S}K^{3/2}S^{-1}\sqrt{\log(Kd)\log d} + 12\kappa^{1/2}d\right\}T^{1/4}d\gamma^2(\lambda)\left\{e^{3S}K^{3/2}S^{-1}\sqrt{\log(Kd)\log d} + 12\kappa^{1/2}d\right\}T^{1/4}d\gamma^2(\lambda)d\gamma^2(\lambda)\left\{e^{3S}K^{3/2}S^{-1}\sqrt{\log(Kd)\log d} + 12\kappa^{1/2}d\right\}T^{1/4}d\gamma^2(\lambda)$$

We now trivially upper bound the regret for  $\mathcal{T}_1$  as  $R\tau_1 = R\sqrt{T}$ . Thus, adding the regret incurred in each batch over all batches  $\beta \in [1, \log \log T + 1]$  results in:

$$\begin{aligned} R_T &\leq \left( 32RK^2 \gamma(\lambda) \sqrt{d \log(Kd)} + R \right) \sqrt{T} \log \log T \\ &+ 32RK\kappa^{1/2} d\gamma^2(\lambda) \left\{ e^{3S} K^{3/2} S^{-1} \sqrt{\log(Kd) \log d} + 12\kappa^{1/2} d \right\} T^{1/4} \log \log T \end{aligned}$$

From Lemma 7.1, setting  $\lambda = S^{-1/2} K d \log T$  along with the fact that  $K d + \log T \leq K d \log T$  results in  $\gamma(\lambda) \leq 22S^{5/4} \sqrt{K d \log T}$ . Substituting the value of  $\gamma(\lambda)$  gives us:

$$R_T \le \left(704S^{5/4}RK^{5/2}d\sqrt{\log T\log(Kd)} + R\right)\sqrt{T}\log\log T + 14784RS^{5/2}K^2d^2\kappa^{1/2}\left\{e^{3S}K^{3/2}S^{-1}\sqrt{\log(Kd)\log d} + 12\kappa^{1/2}d\right\}T^{1/4}\log^2 T\log\log T$$

This concludes the proof.

#### 7.2 Supporting Lemmas for Theorem 7.1

**Lemma 7.1.** For batch  $\beta$ , denoted by  $\mathcal{T}_{\beta}$ , let  $\{x_1, \ldots, x_{\tau_{\beta}}\}$  be a set of i.i.d arms and  $\{r_1, \ldots, r_{\tau_{\beta}}\}$  be the corresponding rewards associated with these arms, where  $\tau_{\beta} = |\mathcal{T}_{\beta}|$ . Define  $\hat{\theta}_{\beta}$  to be the *MLE* estimate for this batch, i.e

$$\hat{\boldsymbol{\theta}}_{\beta} = \arg\min_{\boldsymbol{\theta}} \sum_{s \in \mathcal{T}_{\beta}} \sum_{i=1}^{K} \mathbb{1}\{y_s = i\} \log z_i(\boldsymbol{x}_s, \boldsymbol{\theta}) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

Let the optimal Hessian matrix for batch  $\beta$ ,  $H_{\beta}^{\star}$ , be defined as in Section 7.1. Then, with probability greater than  $1 - \frac{1}{T^2}$ , we have:

$$\left\| \boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta} \right\|_{\boldsymbol{H}_{\beta}^{\star}} \leq 12S\sqrt{\log T + Kd} + 8S\lambda^{-1/2}(\log T + Kd) + 2S^{3/2}\lambda^{1/2}$$

*Proof.* For a batch  $\beta$ , we define the following quantity:

$$oldsymbol{G}_eta(oldsymbol{ heta}_1,oldsymbol{ heta}_2) = \sum_{t\in\mathcal{T}_eta}oldsymbol{M}(oldsymbol{x},oldsymbol{ heta}_1,oldsymbol{ heta}_2)\otimesoldsymbol{x}_toldsymbol{x}_t^ op+\lambdaoldsymbol{I}_{Kd imes Kd}$$

Then,

$$\begin{split} \left\| \boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta} \right\|_{\boldsymbol{H}^{\star}_{\beta}} &\leq \sqrt{1 + 2S} \left\| \boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta} \right\|_{\boldsymbol{G}_{\beta}(\boldsymbol{\theta}^{\star}, \hat{\boldsymbol{\theta}}_{\beta})} \\ &\leq \sqrt{1 + 2S} \left\| S_{\beta}(\boldsymbol{\theta}^{\star}, \hat{\boldsymbol{\theta}}_{\beta}) \left( \boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta} \right) \right\|_{\boldsymbol{G}^{-1}_{\beta}(\boldsymbol{\theta}^{\star}, \hat{\boldsymbol{\theta}}_{\beta})} \\ &\leq \sqrt{1 + 2S} \left\| \sum_{t \in \mathcal{T}_{\beta}} \left[ \boldsymbol{M}(\boldsymbol{x}, \boldsymbol{\theta}^{\star}, \hat{\boldsymbol{\theta}}_{\beta}) \otimes \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top} + \boldsymbol{I}_{Kd \times Kd} \right] \left( \boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta} \right) \right\|_{\boldsymbol{G}^{-1}_{\beta}(\boldsymbol{\theta}^{\star}, \hat{\boldsymbol{\theta}}_{\beta})} \\ &\stackrel{(ii)}{\leq} \sqrt{1 + 2S} \left\| \sum_{t \in \mathcal{T}_{\beta}} \left[ \boldsymbol{M}(\boldsymbol{x}, \boldsymbol{\theta}^{\star}, \hat{\boldsymbol{\theta}}_{\beta}) \otimes \boldsymbol{x}_{t}^{\top} \right] \left( \boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta} \right) \otimes \boldsymbol{x}_{t} + \lambda \left( \boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta} \right) \right\|_{\boldsymbol{G}^{-1}_{\beta}(\boldsymbol{\theta}^{\star}, \hat{\boldsymbol{\theta}}_{\beta})} \\ &\stackrel{(iii)}{\leq} \sqrt{1 + 2S} \left\| \sum_{t \in \mathcal{T}_{\beta}} \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{z}(\boldsymbol{x}_{t}, \hat{\boldsymbol{\theta}}_{\beta}) \right] \otimes \boldsymbol{x}_{t} - \lambda \hat{\boldsymbol{\theta}}_{\beta} \right\|_{\boldsymbol{G}^{-1}_{\beta}(\boldsymbol{\theta}^{\star}, \hat{\boldsymbol{\theta}}_{\beta})} \\ &\stackrel{(iv)}{\leq} (1 + 2S) \left\| \sum_{t \in \mathcal{T}_{\beta}} \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{z}(\boldsymbol{x}_{t}, \hat{\boldsymbol{\theta}}_{\beta}) \right] \otimes \boldsymbol{x}_{t} - \lambda \hat{\boldsymbol{\theta}}_{\beta} \right\|_{\boldsymbol{H}^{\star}_{\beta}^{-1}} \\ &+ \sqrt{\lambda(1 + 2S)} \left\| \boldsymbol{\theta}^{\star} \right\|_{2} \\ &\stackrel{(v)}{\leq} 3S \left\| \sum_{t \in \mathcal{T}_{\beta}} \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s} \right] \otimes \boldsymbol{x}_{t} \right\|_{\boldsymbol{H}^{\star}_{\beta}^{-1}} \end{aligned}$$

where (i) follows from Lemma 9.2, (ii) follows from Mixed Product Property, (iii) follows from the Mean value Theorem and the triangle inequality, (iv) follows from the fact that  $G_{\beta} \succeq \lambda I$  and Lemma 9.2, and (v) follows from Lemma 9.3 and the fact that  $||\theta^*||_2 \leq S$ .

Now, consider the following term:

$$\left\| \sum_{t \in \mathcal{T}_{\beta}} \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s} \right] \otimes \boldsymbol{x}_{t} \right\|_{\boldsymbol{H}_{\beta}^{\star}^{-1}} = \left\| \left\| \sum_{t \in \mathcal{T}_{\beta}} \boldsymbol{H}_{\beta}^{\star}^{-1/2} \left( \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s} \right] \otimes \boldsymbol{x}_{t} \right) \right\|_{2} \right\|_{2} \\ = \max_{\boldsymbol{y} \in \mathcal{B}_{2}(Kd)} \left\langle \boldsymbol{y}, \sum_{t \in \mathcal{T}_{\beta}} \boldsymbol{H}_{\beta}^{\star}^{-1/2} \left( \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s} \right] \otimes \boldsymbol{x}_{t} \right) \right\rangle$$

where  $\mathcal{B}_2(Kd)$  represents the Kd-dimensional unit ball with respect to the  $\ell_2$  norm. We construct an  $\epsilon$ -net for this unit ball, denoted as  $C_{\epsilon}$ . For any  $\boldsymbol{y} \in \mathcal{B}_2(Kd)$ , we define  $\boldsymbol{y}_{\epsilon} = \underset{\boldsymbol{x} \in C_{\epsilon}}{\arg \min} ||\boldsymbol{y} - \boldsymbol{x}||_2$ , then,

$$\begin{split} \left\| \sum_{t \in \mathcal{T}_{\beta}} \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s} \right] \otimes \boldsymbol{x}_{t} \right\|_{\boldsymbol{H}_{\beta}^{\star}^{-1}} &= \max_{\boldsymbol{y} \in \mathcal{B}_{2}(Kd)} \left\langle \boldsymbol{y}, \sum_{t \in \mathcal{T}_{\beta}} \boldsymbol{H}_{\beta}^{\star - 1/2} \left( \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s} \right] \otimes \boldsymbol{x}_{t} \right) \right\rangle \\ &= \max_{\boldsymbol{y} \in \mathcal{B}_{2}(Kd)} \left\langle (\boldsymbol{y} - \boldsymbol{y}_{\epsilon}) + \boldsymbol{y}_{\epsilon}, \sum_{t \in \mathcal{T}_{\beta}} \boldsymbol{H}_{\beta}^{\star - 1/2} \left( \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s} \right] \otimes \boldsymbol{x}_{t} \right) \right\rangle \end{split}$$

Thus, an application of the Cauchy-Schwarz inequality along with the fact that  $||y - y_{\epsilon}||_2 \le \epsilon$  gives us

$$\left|\left|\sum_{t\in\mathcal{T}_{\beta}}\left[\boldsymbol{z}(\boldsymbol{x_{t}},\boldsymbol{\theta}^{\star})-\boldsymbol{m}_{s}\right]\otimes\boldsymbol{x}_{t}\right|\right|_{\boldsymbol{H}_{\beta}^{\star}^{-1}}\leq\frac{1}{1-\epsilon}\left\langle\boldsymbol{y}_{\epsilon},\sum_{t\in\mathcal{T}_{\beta}}\boldsymbol{H}_{\beta}^{\star}^{-\frac{1}{2}}\left(\left[\boldsymbol{z}(\boldsymbol{x_{t}},\boldsymbol{\theta}^{\star})-\boldsymbol{m}_{s}\right]\otimes\boldsymbol{x}_{t}\right)\right\rangle$$

The above term can be bounded using the Bernstein Inequality (Lemma 9.4), which has been done in Lemma 7.2. We note that  $|C_{\epsilon}| \leq \left(\frac{2}{\epsilon}\right)^{Kd}$ . We now set  $\epsilon = 0.5$  and  $\delta = (T^2|C_{\epsilon}|)^{-1}$  and then perform a union bound over  $C_{\epsilon}$ . We get that with probability greater than  $1 - \frac{1}{T^2}$ , we have:

$$\left\| \left\| \sum_{t \in \mathcal{T}_{\beta}} \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s} \right] \otimes \boldsymbol{x}_{t} \right\|_{\boldsymbol{H}_{\beta}^{\star}^{-1}} \leq 2 \left( \sqrt{2 \log \left( T^{2} 4^{Kd} \right)} + \frac{4}{3} \lambda^{-1/2} \log \left( T^{2} 4^{Kd} \right) \right) \\ \leq 4 \sqrt{\log T + Kd} + \frac{8}{3} \lambda^{-1/2} (\log T + Kd)$$

Substituting this into the original bound finishes the proof.

**Lemma 7.2.** Let y be a fixed vector with  $||y||_2 \le 1$ , then, with probability at least  $1 - \delta$ 

$$\sum_{t \in \mathcal{T}_{\beta}} \left[ \boldsymbol{y}^{\top} \boldsymbol{H}_{\beta}^{\star - \frac{1}{2}} \left[ \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s} \right] \otimes \boldsymbol{x}_{t} \right] \leq \sqrt{2 \log \frac{1}{\delta}} + \frac{4}{3\sqrt{\lambda}} \log \frac{1}{\delta}$$

*Proof.* Denote  $\varphi_t = \boldsymbol{y}^\top \boldsymbol{H}_{\beta}^{\star^{-\frac{1}{2}}} ([\boldsymbol{z}(\boldsymbol{x}_t, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_s] \otimes \boldsymbol{x}_t)$ . From Lemma 9.5, we have that  $\mathbb{E}[\varphi_t] = 0$ .

Also,

$$\begin{split} \mathbb{V}\left[\varphi_{t}\right] &= \mathbb{E}\left[\varphi_{t}^{2}\right] - \mathbb{E}\left[\varphi_{t}\right]^{2} \stackrel{(i)}{=} \mathbb{E}\left[\varphi_{t}\varphi_{t}^{\top}\right] \\ &= \mathbb{E}\left[\boldsymbol{y}^{\top}\boldsymbol{H}_{\beta}^{\star} \stackrel{-\frac{1}{2}}{\left(\left[\boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s}\right] \otimes \boldsymbol{x}_{t}\right)\left(\left[\boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s}\right] \otimes \boldsymbol{x}_{t}\right)^{\top}\boldsymbol{H}_{\beta}^{\star} \stackrel{-\frac{1}{2}}{\boldsymbol{y}}\right] \\ \stackrel{(ii)}{=} \boldsymbol{y}^{\top}\boldsymbol{H}_{\beta}^{\star} \stackrel{-\frac{1}{2}}{\mathbb{E}}\left[\left[\boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s}\right]\left[\boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s}\right]^{\top} \otimes \boldsymbol{x}_{t}\boldsymbol{x}_{t}^{\top}\right]\boldsymbol{H}_{\beta}^{\star} \stackrel{-\frac{1}{2}}{\boldsymbol{y}} \\ &= \boldsymbol{y}^{\top}\boldsymbol{H}_{\beta}^{\star} \stackrel{-\frac{1}{2}}{\left(\mathbb{E}\left[\left[\boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s}\right]\left[\boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{\star}) - \boldsymbol{m}_{s}\right]^{\top}\right] \otimes \boldsymbol{x}_{t}\boldsymbol{x}_{t}^{\top}\right)\boldsymbol{H}_{\beta}^{\star} \stackrel{-\frac{1}{2}}{\boldsymbol{y}} \\ \stackrel{(iii)}{=} \boldsymbol{y}^{\top}\boldsymbol{H}_{\beta}^{\star} \stackrel{-\frac{1}{2}}{\left(\boldsymbol{A}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{\star}) \otimes \boldsymbol{x}_{t}\boldsymbol{x}_{t}^{\top}\right)\boldsymbol{H}_{\beta}^{\star} \stackrel{-\frac{1}{2}}{\boldsymbol{y}} \stackrel{(iv)}{=} \boldsymbol{y}^{\top}\boldsymbol{H}_{\beta}^{\star} \stackrel{-\frac{1}{2}}{\left(\boldsymbol{H}_{\beta}^{\star} - \lambda\boldsymbol{I}\right)\boldsymbol{H}_{\beta}^{\star} \stackrel{-\frac{1}{2}}{\boldsymbol{y}} \\ &\leq \boldsymbol{y}^{\top}\boldsymbol{y} \leq 1 \end{split}$$

where (i) follows from the fact that  $\varphi_t$  is a scalar and  $\mathbb{E}[\varphi_t] = 0$ , (ii) follows from the fact that  $(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top$  and the mixed-product property of the Kronecker Product, (iii) follows from Lemma 9.5, and (iv) follows from the definition of  $\mathbf{H}_{\beta}^{\star}$ .

Finally, we note that

$$\begin{split} |\varphi_t - \mathbb{E}\left[\varphi_t\right]| &= |\varphi_t| = \left| \boldsymbol{y}^\top \boldsymbol{H}_{\beta}^{\star - \frac{1}{2}} \left( \left[ \boldsymbol{z}(\boldsymbol{x}_t, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_s \right] \otimes \boldsymbol{x}_t \right) \right| \stackrel{(i)}{\leq} ||\boldsymbol{y}||_2 \left| \left| \boldsymbol{H}_{\beta}^{\star - \frac{1}{2}} \left( \left[ \boldsymbol{z}(\boldsymbol{x}_t, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_s \right] \otimes \boldsymbol{x}_t \right) \right| \right|_2 \\ &\stackrel{(ii)}{\leq} \left| \left| \boldsymbol{H}_{\beta}^{\star - \frac{1}{2}} \right| \left| \left| \left| \left( \boldsymbol{z}(\boldsymbol{x}_t, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_s \right) \otimes \boldsymbol{x}_t \right| \right|_2 \stackrel{(iii)}{\leq} \frac{1}{\sqrt{\lambda}} \left| \left| \boldsymbol{z}(\boldsymbol{x}_t, \boldsymbol{\theta}^{\star}) - \boldsymbol{m}_s \right| \right|_2 ||\boldsymbol{x}_t||_2 \\ &\stackrel{(iv)}{\leq} \frac{1}{\sqrt{\lambda}} \left( ||\boldsymbol{z}(\boldsymbol{x}_t, \boldsymbol{\theta}^{\star})||_2 + ||\boldsymbol{m}_s||_2 \right) \stackrel{(v)}{\leq} \frac{2}{\sqrt{\lambda}} \end{split}$$

where (i) follows from Cauchy-Schwarz, (ii) follows from the fact that  $||\boldsymbol{y}||_2 \leq 1$  and  $||\boldsymbol{Ax}||_2 \leq ||\boldsymbol{A}|| ||\boldsymbol{x}||_2$ , (iii) follows from  $\boldsymbol{H}^{\star}_{\beta} \succeq \lambda \boldsymbol{I}$  and the fact that  $||\boldsymbol{a} \otimes \boldsymbol{b}||_2 = ||\boldsymbol{a}||_2 ||\boldsymbol{b}||_2$ , (iv) follows from  $||\boldsymbol{x}||_2 \leq 1$  and uses the triangle inequality, and (v) follows from the fact  $||\boldsymbol{z}(\boldsymbol{x},\boldsymbol{\theta})||_2 \leq ||\boldsymbol{z}(\boldsymbol{x},\boldsymbol{\theta})||_2 \leq ||\boldsymbol{z}(\boldsymbol{x},\boldsymbol{\theta})||_1 \leq 1$ .

Substituting v = 1 and  $b = \frac{2}{\sqrt{\lambda}}$  in Lemma 9.4 finishes the proof.

**Lemma 7.3.** Let  $V_{\beta}$  and  $H_{\beta}^{\star}$  be the design and optimal Hessian matrices defined as in Section 7.1. Then, we have that

$$V_{\beta} \preccurlyeq \kappa H_{\beta}^{\star}$$

*Proof.* From the definition of  $\kappa$ , we know that  $A(x, \theta) \geq \frac{1}{\kappa}I$ .

Hence, using the fact that  $\kappa > 1$ , we can say that

$$\begin{split} \tilde{\boldsymbol{V}}_{\beta} &= \boldsymbol{I}_{K \times K} \otimes \boldsymbol{V}_{\beta} = \boldsymbol{I}_{K \times K} \otimes \left( \lambda \boldsymbol{I}_{d \times d} + \sum_{t \in \mathcal{T}_{\beta}} \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top} \right) = \lambda \boldsymbol{I}_{K d \times K d} + \boldsymbol{I}_{K \times K} \otimes \sum_{t \in \mathcal{T}_{\beta}} \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top} \\ &\preccurlyeq \lambda \boldsymbol{I}_{K d \times K d} + \kappa \sum_{t \in \mathcal{T}_{\beta}} \boldsymbol{A}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) \otimes \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top} \preccurlyeq \kappa \boldsymbol{H}_{\beta}^{\star} \end{split}$$

**Lemma 7.4.** Let  $H_{\beta}^{*}$  and  $H_{\beta}$  be the optimal and proxy Hessian matrices in batch  $\beta$  as defined in Section 7.1. Then, we have that

$$H_{\beta} \preccurlyeq H_{\beta}^{\star}$$

*Proof.* From Lemma 9.1, we have that

$$oldsymbol{A}(oldsymbol{x}, \hat{oldsymbol{ heta}}_eta) \preccurlyeq oldsymbol{A}(oldsymbol{x}, oldsymbol{ heta}^{\star}) \exp\left(\sqrt{6} \left| \left| (oldsymbol{I} \otimes oldsymbol{x}^{ op}) (oldsymbol{ heta}^{\star} - \hat{oldsymbol{ heta}}_eta) 
ight| 
ight|_2 
ight)$$

We can bound  $\left| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta}) \right| \right|_2$  as follows:

$$\begin{split} \left| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta}) \right| \right|_{2} \stackrel{(i)}{\leq} 2S \left| \left| \boldsymbol{I} \otimes \boldsymbol{x}^{\top} \right| \right|_{2} \stackrel{(ii)}{=} 2S \sqrt{\lambda_{max} \left( (\boldsymbol{I} \otimes \boldsymbol{x}) (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) \right)} \\ \stackrel{(iii)}{=} 2S \sqrt{\lambda_{max} \left( \boldsymbol{I} \otimes \boldsymbol{x} \boldsymbol{x}^{\top} \right)} \stackrel{(iv)}{\leq} 2S \end{split}$$

where (i) uses the sub-multiplicativity of the norm, a triangle inequality, and the fact that  $||\boldsymbol{\theta}^{\star}||_{2} \leq S$ , (ii) uses the definition of the norm, i.e.,  $||\boldsymbol{A}||_{2} = \sqrt{\lambda_{max} (\boldsymbol{A}^{\top} \boldsymbol{A})}$ , (iii) follows from the Mixed-Product property of Kronecker Products, and (iv) follows from the fact that  $\lambda (\boldsymbol{A} \otimes \boldsymbol{B}) = \lambda(\boldsymbol{A})\lambda(\boldsymbol{B})$  and since  $\boldsymbol{x}\boldsymbol{x}^{\top}$  is a rank-one matrix, the only eigenvalues are  $||\boldsymbol{x}||_{2}^{2}$  and 0, and  $0 \leq ||\boldsymbol{x}||_{2} \leq 1$ .

We can also bound  $\left| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}_{\beta}}) \right| \right|_2$  as follows:

$$\begin{split} \left| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta}) \right| \right|_{2} &= \left| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) \boldsymbol{H}_{\beta}^{\star - 1/2} \boldsymbol{H}_{\beta}^{\star 1/2} (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta}) \right| \right|_{2} \overset{(ii)}{\leq} \left| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) \boldsymbol{H}_{\beta}^{\star - 1/2} \right| \right|_{2} \left| \left| \boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta} \right| \right|_{\boldsymbol{H}_{\beta}^{\star}} \\ & \stackrel{(ii)}{\leq} \kappa^{1/2} \gamma(\lambda) \left| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) \kappa^{1/2} \tilde{\boldsymbol{V}}_{\beta}^{-1/2} \right| \right|_{2} \overset{(iii)}{=} \kappa^{1/2} \gamma(\lambda) \sqrt{\lambda_{max} \left( \tilde{\boldsymbol{V}}_{\beta}^{-1/2} (\boldsymbol{I} \otimes \boldsymbol{x}) (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) \tilde{\boldsymbol{V}}_{\beta}^{-1/2} \right)} \\ & \stackrel{(iv)}{=} \kappa^{1/2} \gamma(\lambda) \sqrt{\lambda_{max} \left( (\boldsymbol{I} \otimes \boldsymbol{V}_{\beta}^{-1/2}) (\boldsymbol{I} \otimes \boldsymbol{x}) (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{I} \otimes \boldsymbol{V}_{\beta}^{-1/2}) \right)} \\ & \stackrel{(v)}{=} \kappa^{1/2} \gamma(\lambda) \sqrt{\lambda_{max} \left( \boldsymbol{I} \otimes \boldsymbol{V}_{\beta}^{-1/2} \boldsymbol{x} \boldsymbol{x}^{\top} \boldsymbol{V}_{\beta}^{-1/2} \right)} \overset{(vi)}{=} \kappa^{1/2} \gamma(\lambda) \left| |\boldsymbol{x}| |_{\boldsymbol{V}_{\beta}^{-1}} \right| \end{split}$$

where (i) follows from the sub-multiplicativity of the norm, (ii) follows from Lemma 7.1 and Lemma 7.3, (iii) follows from the definition of the norm, (iv) follows from the definition of  $\tilde{V}_{\beta}$  and the fact that  $(\mathbf{A} \otimes \mathbf{B})^n = \mathbf{A}^n \otimes \mathbf{B}^n$ , (v) follows from the Mixed-Product property, and (vi) follows from  $\lambda(\mathbf{A} \otimes \mathbf{B}) = \lambda(\mathbf{A})\lambda(\mathbf{B})$ .

Thus, we can say that  $\left| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\beta}) \right| \right|_{2} \leq \min \left\{ \gamma(\lambda) \kappa^{1/2} ||\boldsymbol{x}||_{\boldsymbol{V}_{\beta}^{-1}}, 2S \right\}.$ Define  $B_{\beta}(\boldsymbol{x}) = \exp \left( \sqrt{6} \min \left\{ \gamma(\lambda) \kappa^{1/2} ||\boldsymbol{x}||_{\boldsymbol{V}_{\beta}^{-1}}, 2S \right\} \right).$  Then,  $\boldsymbol{A}(\boldsymbol{x}, \hat{\boldsymbol{\theta}}_{\beta}) \preccurlyeq \boldsymbol{A}(\boldsymbol{x}, \boldsymbol{\theta}^{\star}) B_{\beta}(\boldsymbol{x}).$ Hence, we can say,

$$\boldsymbol{H}_{\beta} = \lambda \boldsymbol{I}_{Kd \times Kd} + \sum_{t \in \beta} \frac{\boldsymbol{A}(\boldsymbol{x}_{t}, \hat{\boldsymbol{\theta}}_{\beta})}{\boldsymbol{B}_{\beta}(\boldsymbol{x}_{t})} \otimes \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top} \preccurlyeq \lambda \boldsymbol{I}_{Kd \times Kd} + \sum_{t \in \beta} \boldsymbol{A}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) \otimes \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top} = \boldsymbol{H}_{\beta}^{\star}$$

Lemma 7.5. (Proposition 1, Zhang & Sugiyama (2023)) For any arm x, we have that,

$$\left| \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}^{\star}) - \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}_{j}) \right| \leq \epsilon_{1}(j, \boldsymbol{x}, \lambda) + \epsilon_{2}(j, \boldsymbol{x}, \lambda)$$

where

$$\epsilon_1(j, \boldsymbol{x}, \lambda) = \gamma(\lambda) \left\| \left| \boldsymbol{H}_j^{-1/2}(\boldsymbol{I} \otimes \boldsymbol{x}) \boldsymbol{A}(\boldsymbol{x}, \boldsymbol{\theta}_j) \boldsymbol{\rho} \right| \right\|_2 \text{ and } \epsilon_2(j, \boldsymbol{x}, \lambda) = 3R\gamma(\lambda)^2 \left\| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^\top) \boldsymbol{H}_j^{-1/2} \right| \right\|_2^2$$

*Proof.* We provide the proof for the sake of completeness:

$$\begin{aligned} \left| \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}^{\star}) - \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}_{j}) \right| &= \left| \sum_{i=1}^{K} \rho_{i} \left[ z_{i}(\boldsymbol{x}, \boldsymbol{\theta}^{\star}) - z_{i}(\boldsymbol{x}, \boldsymbol{\theta}_{j}) \right] \right| \\ &= \left| \sum_{i=1}^{K} \rho_{i} \nabla z_{i}(\boldsymbol{x}, \boldsymbol{\theta}_{j})^{\top} \left[ (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j}) \right] + \sum_{i=1}^{K} \rho_{i} \left| \left| (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j}) \right| \right|_{\boldsymbol{Z}_{i}} \right| \\ &\leq \left| \boldsymbol{\rho}^{\top} \boldsymbol{A}(\boldsymbol{x}, \boldsymbol{\theta}_{j}) (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j}) \right| + \left| \sum_{i=1}^{K} \rho_{i} \left| \left| (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j}) \right| \right|_{\boldsymbol{Z}_{i}} \right| \end{aligned}$$

where

$$\boldsymbol{Z}_{i} = \int_{0}^{1} (1-v) \nabla^{2} z_{i}(\boldsymbol{x}, v\boldsymbol{\theta}^{\star} + (1-v)\boldsymbol{\theta}_{j}) \, \mathrm{d}v$$

Beginning with the first term :

$$\begin{aligned} \left| \boldsymbol{\rho}^{\top} \boldsymbol{A}(\boldsymbol{x}, \boldsymbol{\theta}_{j}) (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j}) \right| &= \left| \boldsymbol{\rho}^{\top} \boldsymbol{A}(\boldsymbol{x}, \boldsymbol{\theta}_{j}) (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) \boldsymbol{H}_{j}^{\star^{-1/2}} \boldsymbol{H}_{j}^{\star^{-1/2}} \boldsymbol{H}_{j}^{\star^{-1/2}} (\boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j}) \right| \\ & \stackrel{(i)}{\leq} \left| \left| \boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j} \right| \right|_{\boldsymbol{H}_{j}^{\star}} \left\| \left| \boldsymbol{\rho}^{\top} \boldsymbol{A}(\boldsymbol{x}, \boldsymbol{\theta}_{j}) (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) \boldsymbol{H}_{j}^{\star^{-1/2}} \right| \right|_{2} \\ & \leq \gamma(\lambda) \left\| \left| \boldsymbol{H}_{j}^{\star^{-1/2}} (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}) \boldsymbol{A}(\boldsymbol{x}, \boldsymbol{\theta}_{j}) \boldsymbol{\rho} \right| \right\|_{2} \\ & \stackrel{(ii)}{\leq} \gamma(\lambda) \left\| \left| \boldsymbol{H}_{j}^{-1/2} (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}) \boldsymbol{A}(\boldsymbol{x}, \boldsymbol{\theta}_{j}) \boldsymbol{\rho} \right| \right\|_{2} \end{aligned}$$

where (i) follows from the sub-multiplicativity of the norm and (ii) is due to Lemma 7.4.

For the second term, for some  $k \in [1, K]$ , we make the following observation:

$$\boldsymbol{Z}_{k} = \int_{0}^{1} (1-v) \nabla^{2} z_{k}(\boldsymbol{x}, v\boldsymbol{\theta}^{\star} + (1-v)\boldsymbol{\theta}_{j}) \, \mathrm{d}v \preccurlyeq 3\boldsymbol{I} \int_{0}^{1} (1-v) \, \mathrm{d}v \preccurlyeq 3\boldsymbol{I}$$

Thus, we have:

$$\begin{aligned} \left| \sum_{i=1}^{K} \rho_{i} \left| \left| (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j}) \right| \right|_{\boldsymbol{Z}_{i}}^{2} \right| &\leq \left| \sum_{i=1}^{K} 3\rho_{i} \left| \left| (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j}) \right| \right|_{2}^{2} \right| \\ &\leq 3R \left| \left| (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) \boldsymbol{H}_{j}^{\star - 1/2} \boldsymbol{H}_{j}^{\star - 1/2} \boldsymbol{H}_{j}^{\star - 1/2} (\boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j}) \right| \right|_{2}^{2} \\ &\leq 3R \left| \left| \boldsymbol{\theta}^{\star} - \boldsymbol{\theta}_{j} \right| \right|_{\boldsymbol{H}_{j}^{\star}}^{2} \left| \left| (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) \boldsymbol{H}_{j}^{\star - 1/2} \right| \right|_{2}^{2} \\ &\leq 3R \gamma(\lambda)^{2} \left| \left| (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) \boldsymbol{H}_{j}^{\star - 1/2} \right| \right|_{2}^{2} \\ &\leq 3R \gamma(\lambda)^{2} \left| \left| (\boldsymbol{I}_{K \times K} \otimes \boldsymbol{x}^{\top}) \boldsymbol{H}_{j}^{- 1/2} \right| \right|_{2}^{2} \end{aligned} \end{aligned}$$

**Lemma 7.6.** Let  $x_t^*$  be the optimal arm at round t, i.e  $x_t^* = \arg \max_{x \in \mathcal{X}_t} \rho^\top z(x, \theta^*)$ . Then, the optimal arm never gets eliminated in any round.

Proof. From Lemma 7.5, we know that

$$\left| \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}^{\star}) - \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}_{j}) \right| \leq \epsilon_{1}(j, \boldsymbol{x}, \lambda) + \epsilon_{2}(j, \boldsymbol{x}, \lambda)$$

Also, from Algorithm 1, we have the definitions of  $UCB(j, x, \lambda)$  and  $LCB(j, x, \lambda)$  as:

$$UCB(j, \boldsymbol{x}, \lambda) = \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}_j) + \epsilon_1(j, \boldsymbol{x}, \lambda) + \epsilon_2(j, \boldsymbol{x}, \lambda)$$
$$LCB(j, \boldsymbol{x}, \lambda) = \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}_j) - \epsilon_1(j, \boldsymbol{x}, \lambda) - \epsilon_2(j, \boldsymbol{x}, \lambda)$$

From Algorithm 1, we know that an arm  $\boldsymbol{x} \in \mathcal{X}_t$  gets eliminated if  $UCB(j, \boldsymbol{x}, \lambda) \leq \max_{\boldsymbol{y} \in \mathcal{X}_t} LCB(j, \boldsymbol{y}, \lambda)$ . Thus, showing that  $UCB(j, \boldsymbol{x}_t^{\star}, \lambda) \geq \max_{\boldsymbol{y} \in \mathcal{X}_t} LCB(j, \boldsymbol{y}, \lambda)$  accounts to showing that  $\boldsymbol{x}_t^{\star}$  never gets eliminated.

We assume that  $\arg \max_{y \in \mathcal{X}_t} \text{LCB}(j, y, \lambda) = y$ . Then, for any arm  $x \in \mathcal{X}_t$ , we have that

$$\begin{aligned} \mathsf{LCB}(j, \boldsymbol{x}, \lambda) &\leq \max_{\boldsymbol{y} \in \mathcal{X}_{t}} \mathsf{LCB}(j, \boldsymbol{y}, \lambda) \\ &= \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{y}, \boldsymbol{\theta}_{j}) - \epsilon_{1}(j, \boldsymbol{y}, \lambda) - \epsilon_{2}(j, \boldsymbol{y}, \lambda) \\ &\leq \left[ \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{y}, \boldsymbol{\theta}^{\star}) + \epsilon_{1}(j, \boldsymbol{y}, \lambda) + \epsilon_{2}(j, \boldsymbol{y}, \lambda) \right] - \epsilon_{1}(j, \boldsymbol{y}, \lambda) - \epsilon_{2}(j, \boldsymbol{y}, \lambda) \\ &= \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{y}, \boldsymbol{\theta}^{\star}) \\ &\leq \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}^{\star}, \boldsymbol{\theta}^{\star}) \\ &\leq \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}^{\star}, \boldsymbol{\theta}_{j}) + \epsilon_{1}(j, \boldsymbol{x}_{t}^{\star}, \lambda) + \epsilon_{2}(j, \boldsymbol{x}_{t}^{\star}, \lambda) \\ &= \mathsf{UCB}(j, \boldsymbol{x}_{t}^{\star}, \lambda) \end{aligned}$$

where (*i*) follows from Lemma 7.5, (*ii*) follows from the fact that  $x_t^* = \underset{y \in \mathcal{X}_t}{\arg \max} \rho^\top z(y, \theta^*)$ , and (*iii*) again follows from Lemma 7.5.

**Lemma 7.7.** Let  $B_{\beta}(\mathbf{x})$  be as defined in Section 7.1. Then, we have that

$$\sqrt{B_{\beta}(\boldsymbol{x})} \leq \frac{1}{2} e^{3S} \gamma(\lambda) \kappa^{1/2} S^{-1} \|\boldsymbol{x}\|_{\boldsymbol{V}_{\beta}^{-1}} + 1$$

Proof.

$$\sqrt{B_{\beta}(\boldsymbol{x})} = \exp\left(\sqrt{6}\min\left\{S, \frac{1}{2}\gamma(\lambda)\kappa^{1/2}\|\boldsymbol{x}\|_{\boldsymbol{V}_{\beta}^{-1}}\right\}\right) \leq \frac{1}{2}e^{3S}\gamma(\lambda)\kappa^{1/2}S^{-1}\|\boldsymbol{x}\|_{\boldsymbol{V}_{\beta}^{-1}} + 1$$

where the inequality follows from Lemma 9.6 by choosing  $\min\{2S, \gamma(\lambda)\kappa^{1/2} \|\boldsymbol{x}\|_{\boldsymbol{V}_{\beta}^{-1}}\} = \gamma(\lambda)\kappa^{1/2} \|\boldsymbol{x}\|_{\boldsymbol{V}_{\beta}^{-1}}$  and  $M = \sqrt{6}S$ .

**Lemma 7.8.** Let  $\epsilon_1(\beta, \boldsymbol{x}, \lambda)$  be as defined in Lemma 7.5. Then, we have

$$\mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}}\left[\max_{\boldsymbol{x}\in\mathcal{X}}\epsilon_{1}(\beta,\boldsymbol{x},\lambda)\right] \leq \frac{8R\kappa^{1/2}K^{5/2}de^{3S}\gamma(\lambda)^{2}S^{-1}\sqrt{\log Kd\log d}}{\tau_{\beta}} + \frac{4RK^{2}d^{1/2}\gamma(\lambda)\sqrt{\log(Kd)}}{\sqrt{\tau_{\beta}}}$$

Proof.

$$\begin{split} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \epsilon_{1}(\beta,\boldsymbol{x},\lambda) \right] &= \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \gamma(\lambda) \left\| \left| \boldsymbol{H}_{\beta}^{-1/2}(\boldsymbol{I}\otimes\boldsymbol{x})\boldsymbol{A}(\boldsymbol{x},\hat{\boldsymbol{\theta}}_{\beta})\boldsymbol{\rho} \right| \right|_{2} \right] \\ & \stackrel{(i)}{\leq} \gamma(\lambda) \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left\| \left| \boldsymbol{H}_{\beta}^{-1/2}(\boldsymbol{I}\otimes\boldsymbol{x})\boldsymbol{A}(\boldsymbol{x},\hat{\boldsymbol{\theta}}_{\beta})^{1/2} \right| \right|_{2} \left\| \boldsymbol{\rho} \right\|_{\boldsymbol{A}(\boldsymbol{x},\hat{\boldsymbol{\theta}}_{\beta})} \right] \\ & \stackrel{(ii)}{\leq} R\gamma(\lambda) \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left\| \left| \boldsymbol{A}(\boldsymbol{x},\hat{\boldsymbol{\theta}}_{\beta})^{1/2}(\boldsymbol{I}\otimes\boldsymbol{x}^{\top})\boldsymbol{H}_{\beta}^{-1/2} \right| \right|_{2} \right] \\ & \stackrel{(iii)}{\leq} R\gamma(\lambda) \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left\| \left| \sqrt{B_{\beta}(\boldsymbol{x})} \tilde{\boldsymbol{X}}_{\beta}^{\top} \boldsymbol{H}_{\beta}^{-1/2} \right| \right|_{2} \right] \\ & \stackrel{(iv)}{\leq} 4R\gamma(\lambda) K^{2} \sqrt{\frac{d\log Kd}{\tau_{\beta}}} \left\{ \frac{1}{2} e^{3S} \gamma(\lambda) \kappa^{1/2} S^{-1} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left\| \boldsymbol{x} \right\|_{V_{\beta}^{-1}} \right] + 1 \right\} \\ & \stackrel{(v)}{\leq} 4R\gamma(\lambda) K^{2} \sqrt{\frac{d\log Kd}{\tau_{\beta}}} \left\{ 2e^{3S} \gamma(\lambda) \kappa^{1/2} S^{-1} \sqrt{\frac{Kd\log d}{\tau_{\beta}}} + 1 \right\} \\ & \leq \frac{8R\kappa^{1/2} K^{5/2} de^{3S} \gamma(\lambda)^{2} S^{-1} \sqrt{\log Kd\log d}}{\tau_{\beta}}} + \frac{4RK^{2} d^{1/2} \gamma(\lambda) \sqrt{\log(Kd)}}{\sqrt{\tau_{\beta}}} \end{split}$$

where (i) follows from  $||Ax||_2 \leq ||A||_2 ||x||_2$ , (ii) follows from the fact that  $A(x, \theta) \preccurlyeq I$ , (iii) follows from the definition of  $\tilde{X}$ , (iv) follows from Lemma 7.7, the fact that  $\max \{ab\} \leq \max \{a\} \max \{b\}$ , and Lemma 7.18, and (v) follows from Lemma 7.17.

**Lemma 7.9.** Let  $\epsilon_2(\beta, \boldsymbol{x}, \lambda)$  be as defined in Lemma 7.5. Then, we have

$$\mathbb{E}_{\mathcal{X} \sim \mathcal{D}_{\beta+1}} \left[ \epsilon_2(\beta, \boldsymbol{x}, \lambda) \right] \le \frac{96R\kappa\gamma(\lambda)^2}{\tau_\beta} K d^2$$

Proof. Recall from Lemma 7.5, in one of the intermediate steps, we have that

$$\epsilon_2(\beta, \boldsymbol{x}, \lambda) = 3R\gamma(\lambda)^2 \left| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^\top) \boldsymbol{H}_{\beta}^{\star - 1/2} \right| \right|_2^2$$

Thus, we have

$$\begin{split} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \epsilon_{2}(\beta, \boldsymbol{x}, \lambda) \right] &= \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} 3R\gamma(\lambda)^{2} \left| \left| (\boldsymbol{I}\otimes\boldsymbol{x}^{\top})\boldsymbol{H}_{\beta}^{\star - 1/2} \right| \right|_{2}^{2} \right] \\ &= 3R\gamma(\lambda)^{2} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left| \left| (\boldsymbol{I}\otimes\boldsymbol{x}^{\top})\boldsymbol{H}_{\beta}^{\star - 1/2} \right| \right|_{2}^{2} \right] \\ &\stackrel{(i)}{\leq} 3R\kappa\gamma(\lambda)^{2} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left| \left| (\boldsymbol{I}\otimes\boldsymbol{x}^{\top})\tilde{\boldsymbol{V}}_{\beta}^{-1/2} \right| \right|_{2}^{2} \right] \\ &\stackrel{(ii)}{\leq} 3R\kappa\gamma(\lambda)^{2} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left| \left| (\boldsymbol{I}\otimes\boldsymbol{x}^{\top})(\boldsymbol{I}\otimes\boldsymbol{V}_{\beta}^{-1/2}) \right| \right|_{2}^{2} \right] \\ &\stackrel{(iii)}{\leq} 3R\kappa\gamma(\lambda)^{2} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left| \left| (\boldsymbol{I}\otimes\boldsymbol{x}^{\top})(\boldsymbol{I}\otimes\boldsymbol{V}_{\beta}^{-1/2}) \right| \right|_{2}^{2} \right] \\ &\stackrel{(iii)}{\leq} 3R\kappa\gamma(\lambda)^{2} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left| \left| \boldsymbol{x} \right| \right|_{\boldsymbol{V}_{\beta}^{-1}}^{2} \right] \stackrel{(iv)}{\leq} \frac{48R\kappa\gamma(\lambda)^{2}}{\tau_{\beta}} (K+1)d^{2} \leq \frac{96R\kappa\gamma(\lambda)^{2}}{\tau_{\beta}} Kd^{2} \end{split}$$

where (*i*) follows from Lemma 7.3, (*ii*) follows from the definition of  $\tilde{V}_{\beta}$ , (*iii*) follows from the Mixed-Product Property and the fact that  $\lambda(\boldsymbol{A} \otimes \boldsymbol{B}) = \lambda(\boldsymbol{A})\lambda(\boldsymbol{B})$ , and (*iv*) follows from Lemma 7.16.

**Lemma 7.10.** Let t be a time round in batch  $\beta + 1$ , i.e  $t \in T_{\beta}$ . Then, the expected regret incurred at round t, denoted as  $R_t$  can be bounded as:

$$R_t \le \frac{32RK\kappa^{1/2}d\gamma(\lambda)^2}{\tau_{\beta}} \left\{ e^{3S}K^{3/2}S^{-1}\sqrt{\log(Kd)\log d} + 12\kappa^{1/2}d \right\} + \frac{16RK^2d^{1/2}\gamma(\lambda)\sqrt{\log(Kd)}}{\sqrt{\tau_{\beta}}}$$

Proof. Using Lemma 7.5,

 $\rho^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}^{\star}, \boldsymbol{\theta}^{\star}) - \rho^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) \leq \rho^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}^{\star}, \boldsymbol{\theta}_{\beta}) - \rho^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}_{\beta}) + \epsilon_{1}(\beta, \boldsymbol{x}_{t}^{\star}, \lambda) + \epsilon_{2}(\beta, \boldsymbol{x}_{t}^{\star}, \lambda) + \epsilon_{1}(\beta, \boldsymbol{x}_{t}, \lambda) + \epsilon_{2}(\beta, \boldsymbol{x}_{t}, \lambda) +$ 

$$\boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_t, \boldsymbol{\theta}_{\beta}) + \epsilon_1(\beta, \boldsymbol{x}_t, \lambda) + \epsilon_2(\beta, \boldsymbol{x}_t, \lambda) \geq \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_t^{\star}, \boldsymbol{\theta}_{\beta}) - \epsilon_1(\beta, \boldsymbol{x}_t^{\star}, \lambda) - \epsilon_2(\beta, \boldsymbol{x}_t^{\star}, \lambda)$$

Thus, we get

$$\rho^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}^{\star}, \boldsymbol{\theta}^{\star}) - \rho^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) \leq 2\epsilon_{1}(\beta, \boldsymbol{x}_{t}, \lambda) + 2\epsilon_{2}(\beta, \boldsymbol{x}_{t}, \lambda) + 2\epsilon_{1}(\beta, \boldsymbol{x}_{t}^{\star}, \lambda) + 2\epsilon_{2}(\beta, \boldsymbol{x}_{t}^{\star}, \lambda) \\ \leq 4 \max_{\boldsymbol{x} \in \mathcal{X}} \epsilon_{1}(\beta, \boldsymbol{x}, \lambda) + 4 \max_{\boldsymbol{x} \in \mathcal{X}} \epsilon_{2}(\beta, \boldsymbol{x}, \lambda)$$

Taking an expectation on both sides, we get

$$\mathbb{E}_{\mathcal{X} \sim \mathcal{D}_{\beta+1}} \left[ \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}^{\star}, \boldsymbol{\theta}^{\star}) - \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) \right] \leq 4 \left( \mathbb{E}_{\mathcal{X} \sim \mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x} \in \mathcal{X}} \epsilon_{1}(\beta, \boldsymbol{x}, \lambda) + \max_{\boldsymbol{x} \in \mathcal{X}} \epsilon_{2}(\beta, \boldsymbol{x}, \lambda) \right] \right)$$

$$\leq \frac{32RK\kappa^{1/2}d\gamma(\lambda)^{2}}{\tau_{\beta}} \left\{ e^{3S}K^{3/2}S^{-1}\sqrt{\log(Kd)\log d} + 12\kappa^{1/2}d \right\} + \frac{16RK^{2}d^{1/2}\gamma(\lambda)\sqrt{\log(Kd)}}{\sqrt{\tau_{\beta}}}$$

which follows from Lemma 7.8 and Lemma 7.9.

#### 7.3 Supporting Results on Optimal Designs for 7

Recall from Section 7.1,

$$ilde{oldsymbol{X}}_eta = rac{oldsymbol{A}(oldsymbol{x},\hat{oldsymbol{ heta}}_eta)^rac{1}{2}}{\sqrt{B_eta(oldsymbol{x})}}\otimesoldsymbol{x}$$

Also, recall that at each round  $t \in [T]$ , the feasible set of context vectors  $\mathcal{X}_t$  is being sampled from some distribution  $\mathcal{D}$ . For a given batch  $\beta$ , we denote  $\mathcal{D}_{\beta}$  to be the distribution of the pruned arm-sets post the successive elimination procedure (Section 3.1). Thus, we have that  $\mathcal{D}_{\beta+1} \subset \mathcal{D}_{\beta}$ .

We now define K different partitions of  $\tilde{X}_{\beta}$  as follows:

$$ilde{oldsymbol{x}}_eta^{(i)} = rac{oldsymbol{A}(oldsymbol{x},oldsymbol{ heta}_eta)^rac{1}{2}}{\sqrt{B_eta(oldsymbol{x})}}oldsymbol{e}_i\otimesoldsymbol{x}$$

where  $i \in [K]$  and  $e_i$  is the K-dimensional standard basis vector. We first show a few relations between  $\tilde{X}_{\beta}$  and  $\tilde{x}_{\beta}^{(i)}$ :

**Lemma 7.11.** Let  $ilde{X}_{eta}$  and  $ilde{x}_{eta}^{(i)}$  be defined as above. Then, we have

$$ilde{oldsymbol{X}}_{eta} ilde{oldsymbol{X}}_{eta}^{ op} = \sum_{i=1}^{K} ilde{oldsymbol{x}}_{eta}^{(i)} ilde{oldsymbol{x}}_{eta}^{(i)}^{ op}$$

Proof.

$$egin{aligned} &\sum_{i=1}^{K} ilde{m{x}}_{eta}^{(i)} ilde{m{x}}_{eta}^{(i)} &^{ op} &= \sum_{i=1}^{K} \left( rac{m{A}(m{x}, \hat{m{ heta}}_{eta})^{rac{1}{2}}}{\sqrt{B_{eta}(m{x})}} m{e}_i \otimes m{x} 
ight) \left( m{e}_i^{ op} rac{m{A}(m{x}, \hat{m{ heta}}_{eta})^{rac{1}{2}}}{\sqrt{B_{eta}(m{x})}} \otimes m{x}^{ op} 
ight) \ &= rac{1}{B_{eta}(m{x})} \sum_{i=1}^{K} m{A}(m{x}, \hat{m{ heta}}_{eta})^{rac{1}{2}} m{e}_i m{e}_i^{ op} m{A}(m{x}, \hat{m{ heta}}_{eta})^{rac{1}{2}} \otimes m{x}^{ op} 
ight) \ &= rac{1}{B_{eta}(m{x})} \sum_{i=1}^{K} m{A}(m{x}, \hat{m{ heta}}_{eta})^{rac{1}{2}} m{e}_i m{e}_i^{ op} m{A}(m{x}, \hat{m{ heta}}_{eta})^{rac{1}{2}} \otimes m{x} m{x}^{ op} \ &= rac{1}{B_{eta}(m{x})} m{A}(m{x}, \hat{m{ heta}}_{m{ heta}})^{rac{1}{2}} \left(\sum_{i=1}^{K} m{e}_i m{e}_i^{ op} 
ight) m{A}(m{x}, \hat{m{ heta}}_{m{ heta}})^{rac{1}{2}} \otimes m{x} m{x}^{ op} \ &= rac{m{A}(m{x}, \hat{m{ heta}}_{m{ heta}})^{rac{1}{2}} \left(\sum_{i=1}^{K} m{e}_i m{e}_i^{ op} 
ight) m{A}(m{x}, \hat{m{ heta}}_{m{ heta}})^{rac{1}{2}} \otimes m{x} m{x}^{ op} \ &= rac{m{A}(m{x}, \hat{m{ heta}}_{m{ heta}}) \otimes m{x} m{x}^{ op} = m{X}_{m{ heta}} m{X}_{m{ heta}}^{ op} \end{split}$$

where we use the fact that  $\sum_{i=1}^{K} e_i e_i^{\top} = I_{K \times K}$ .

**Lemma 7.12.** Let  $M \in \mathbb{R}^{Kd}$  be any positive-semidefinite matrix. Then,

$$\lambda_{max}\left(\tilde{\boldsymbol{X}}_{\beta}^{\top}\boldsymbol{M}\tilde{\boldsymbol{X}}_{\beta}\right) \leq \sum_{i=1}^{K}\left|\left|\tilde{\boldsymbol{x}}_{\beta}^{(i)}\right|\right|_{\boldsymbol{M}}^{2}$$

Proof.

$$\begin{split} \lambda_{max} \left( \tilde{\boldsymbol{X}}_{\beta}^{\top} \boldsymbol{M} \tilde{\boldsymbol{X}}_{\beta} \right) &\stackrel{(i)}{=} \lambda_{max} \left( \tilde{\boldsymbol{X}}_{\beta} \tilde{\boldsymbol{X}}_{\beta}^{\top} \boldsymbol{M} \right) \stackrel{(ii)}{=} \lambda_{max} \left( \sum_{i=1}^{K} \tilde{\boldsymbol{x}}_{\beta}^{(i)} \tilde{\boldsymbol{x}}_{\beta}^{(i) \top} \boldsymbol{M} \right) \\ &\stackrel{(iii)}{\leq} \sum_{i=1}^{K} \lambda_{max} \left( \tilde{\boldsymbol{x}}_{\beta}^{(i)} \tilde{\boldsymbol{x}}_{\beta}^{(i) \top} \boldsymbol{M} \right) \stackrel{(iv)}{=} \sum_{i=1}^{K} \lambda_{max} \left( \tilde{\boldsymbol{x}}_{\beta}^{(i) \top} \boldsymbol{M} \tilde{\boldsymbol{x}}_{\beta}^{(i)} \right) = \sum_{i=1}^{K} \left\| \tilde{\boldsymbol{x}}_{\beta}^{(i)} \right\|_{\boldsymbol{M}}^{2} \end{split}$$

where (i) follows from the cyclic property of eigenvalues, (ii) follows from Lemma 7.11, (iii) follows from the fact that  $\lambda_{max} (\mathbf{A} + \mathbf{B}) \leq \lambda_{max} (\mathbf{A}) + \lambda_{max} (\mathbf{B})$ , and (iv) again follows from the cyclic property of eigenvalues.

We first redefine the Distributional Optimal Design (Definition ??) for a set  $\mathcal{X}$ .

٦

$$\pi(\mathcal{X}) = \begin{cases} \pi_G(\mathcal{X}) & \text{w.p.}\frac{1}{2} \\ \pi_{M_i}^S(\mathcal{X}) & \text{w.p.}\frac{p_i}{2} \end{cases}$$

where  $\pi_G$  is the G-optimal design and  $\pi_{M_i}^S$  represents the Softmax Policy with respect to  $M_i$ . We refer the reader to Definition ?? for more details.

We now define a few notations regarding some of the information and design matrices used throughout this section.

1. 
$$\mathbb{I}_{\mathcal{D}}^{\lambda}(\pi) = \underset{\mathcal{X}\sim\mathcal{D}}{\mathbb{E}} \left[ \underset{\boldsymbol{x}\sim\pi(\mathcal{X})}{\mathbb{E}} \tilde{\boldsymbol{X}}_{\beta} \tilde{\boldsymbol{X}}_{\beta}^{\top} \right]$$
  
2. 
$$\mathbb{W}_{\mathcal{D}}^{(i)}(\pi) = \underset{\mathcal{X}\sim\mathcal{D}}{\mathbb{E}} \left[ \underset{\boldsymbol{x}\sim\pi(\mathcal{X})}{\mathbb{E}} \tilde{\boldsymbol{x}}_{\beta}^{(i)} \tilde{\boldsymbol{x}}_{\beta}^{(i)} \right]$$
  
3. 
$$\mathbb{W}_{\mathcal{D}}^{(0)}(\pi) = \underset{\mathcal{X}\sim\mathcal{D}}{\mathbb{E}} \left[ \underset{\boldsymbol{x}\sim\pi(\mathcal{X})}{\mathbb{E}} \boldsymbol{x} \boldsymbol{x}^{\top} \right]$$

г

Suppose Algorithm 2 is called with the inputs  $\beta$  and  $S_{\beta}$ , where  $\beta$  is the current batch index. Then, the policy returned by the algorithm is denoted by  $\pi_{\beta}$ , where

$$\pi_{\beta} = \frac{1}{K+1} \left( \sum_{i=0}^{K} \pi_{\beta,i} \right)$$

where  $\pi_{\beta,0}$  and  $\pi_{\beta,i}$   $i \in [K]$  represents the Distributional Optimal Design learned over  $S_{\beta}$  and  $F_i(S_{\beta}, \beta), i \in [K]$  respectively. Here  $F_i$  is as defined in Equation 8.

We now state a few results that relate the design matrices H and V as well as the matrices  $|\mathbb{I}|$  and  $\mathbb{W}$ .

**Lemma 7.13.** (Lemma A.16, Sawarni et al. (2024)) Let  $V_{\beta}$  and  $H_{\beta}$  as defined in Section 7.1 and  $W_{\mathcal{D}}^{(0)}(\pi_{\beta})$  and  $I_{\mathcal{D}}^{\lambda}(\pi_{\beta})$  be as defined in Section 7.3.

Then, with probability at least  $1 - \frac{1}{T^2}$ , we have that

$$\begin{aligned} \boldsymbol{V}_{\beta} &\succeq \frac{\tau_{\beta}}{8} \mathbb{W}_{\mathcal{D}}^{(0)}(\pi_{\beta}) \\ \boldsymbol{H}_{\beta} &\succeq \frac{\tau_{\beta}}{8} \mathbb{I}_{\mathcal{D}}^{\lambda}(\pi_{\beta}) \end{aligned}$$

**Lemma 7.14.** For all  $i \in [0, K]$ , we have that

$$(K+1)\mathbb{I}_{\mathcal{D}}^{\lambda}(\pi_{\beta}) \succeq \mathbb{I}_{\mathcal{D}}^{\lambda}(\pi_{\beta,i})$$

Proof.

$$\mathbb{I}_{\mathcal{D}}^{\lambda}(\pi_{\beta}) = \mathbb{E}_{\mathcal{X}\sim\mathcal{D}}\left[\mathbb{E}_{\boldsymbol{x}\sim\pi_{\beta}(\mathcal{X})}\tilde{\boldsymbol{X}}_{\beta}\tilde{\boldsymbol{X}}_{\beta}^{\top}\right] \stackrel{(i)}{\succcurlyeq} (K+1)^{-1} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}}\left[\sum_{i=0}^{K} \mathbb{E}_{\boldsymbol{x}\sim\pi_{\beta,i}(\mathcal{X})}\tilde{\boldsymbol{X}}_{\beta}\tilde{\boldsymbol{X}}_{\beta}^{\top}\right]$$
$$\approx (K+1)^{-1} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}}\left[\mathbb{E}_{\boldsymbol{x}\sim\pi_{\beta,i}(\mathcal{X})}\tilde{\boldsymbol{X}}_{\beta}\tilde{\boldsymbol{X}}_{\beta}^{\top}\right] = (K+1)^{-1}\mathbb{I}_{\mathcal{D}}^{\lambda}(\pi_{\beta,i})$$

where (i) follows from the definition of  $\pi_{\beta}$ .

**Lemma 7.15.** For all  $i \in [K]$ , we have that

$$\mathbb{I}^{\lambda}_{\mathcal{D}}(\pi) \succcurlyeq \mathbb{W}^{(i)}_{\mathcal{D}}(\pi)$$

Proof.

$$\mathbb{I}_{\mathcal{D}}^{\lambda}(\pi) = \mathbb{E}_{\mathcal{X}\sim\mathcal{D}}\left[\mathbb{E}_{\boldsymbol{x}\sim\pi(\mathcal{X})}\tilde{\boldsymbol{X}}_{\beta}\tilde{\boldsymbol{X}}_{\beta}^{\top}\right] \stackrel{(i)}{=} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}}\left[\mathbb{E}_{\boldsymbol{x}\sim\pi(\mathcal{X})}\sum_{i=1}^{K}\tilde{\boldsymbol{x}}_{\beta}^{(i)}\tilde{\boldsymbol{x}}_{\beta}^{(i)}^{\top}\right] \succcurlyeq \mathbb{E}_{\mathcal{X}\sim\mathcal{D}}\left[\mathbb{E}_{\boldsymbol{x}\sim\pi(\mathcal{X})}\tilde{\boldsymbol{x}}_{\beta}^{(i)}\tilde{\boldsymbol{x}}_{\beta}^{(i)}^{\top}\right] = \mathbb{W}_{\mathcal{D}}^{(i)}(\pi)$$

Using the lemmas stated above, we now derive a few results.

**Lemma 7.16.** Let  $V_{\beta}$  be as defined in Section 7.1 and  $\tau_{\beta}$  be the length of the  $\beta$  batch, i.e  $|\mathcal{T}_{\beta} = \tau_{\beta}$ . Then, we have

$$\mathbb{E}_{\mathcal{X} \sim \mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x} \in \mathcal{X}} ||\boldsymbol{x}||_{\boldsymbol{V}_{\beta}^{-1}}^{2} \right] \leq \frac{16}{\tau_{\beta}} (K+1) d^{2}$$

Proof.

$$\mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}}\left[\max_{\boldsymbol{x}\in\mathcal{X}}||\boldsymbol{x}||_{V_{\beta}^{-1}}^{2}\right] \stackrel{(i)}{\leq} \frac{8}{\tau_{\beta}} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}}\left[\max_{\boldsymbol{x}\in\mathcal{X}}||\boldsymbol{x}||_{\mathbb{W}_{\mathcal{D}_{\beta}}^{(0)-1}(\pi_{\beta})}^{2}\right] \\
\stackrel{(ii)}{\leq} \frac{8}{\tau_{\beta}}(K+1) \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}}\left[\max_{\boldsymbol{x}\in\mathcal{X}}||\boldsymbol{x}||_{\mathbb{W}_{\mathcal{D}_{\beta}}^{(0)-1}(\pi_{\beta},0)}^{2}\right] \\
\stackrel{(iii)}{\leq} \frac{16}{\tau_{\beta}}(K+1) \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}}\left[\max_{\boldsymbol{x}\in\mathcal{X}}||\boldsymbol{x}||_{\mathbb{W}_{\mathcal{D}_{\beta}}^{(0)-1}(\pi_{G})}^{2}\right] \\
\stackrel{(iv)}{\leq} \frac{16}{\tau_{\beta}}(K+1)d^{2}$$

where (*i*) follows from Lemma 7.13, (*ii*) follows from Lemma 7.14 and the fact that  $\mathcal{D}_{\beta+1} \subset \mathcal{D}_{\beta}$ and hence,  $\mathbb{E}_{\mathcal{D}_{\beta+1}} \leq \mathbb{E}_{\mathcal{D}_{\beta}}$ , (*iii*) follows from the definition of  $\pi_{\beta,0}$  and uses the fact that  $\pi_{\beta,0} \geq \frac{\pi_G}{2}$ , and (*iv*) follows from Lemma 9.8.

**Lemma 7.17.** Let  $V_{\beta}$  be as defined in Section 7.1 and  $\tau_{\beta}$  be the length of the  $\beta$  batch, i.e  $|\mathcal{T}_{\beta}| = \tau_{\beta}$ . Then, we have

$$\mathbb{E}_{\mathcal{X} \sim \mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x} \in \mathcal{X}} ||\boldsymbol{x}||_{\boldsymbol{V}_{\beta}^{-1}} \right] \le 4\sqrt{\frac{Kd\log d}{\tau_{\beta}}}$$

Proof.

$$\mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} ||\boldsymbol{x}||_{\boldsymbol{V}_{\beta}^{-1}} \right] \stackrel{(i)}{\leq} \sqrt{\frac{8}{\tau_{\beta}}} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} ||\boldsymbol{x}||_{\mathbb{W}_{\mathcal{D}_{\beta}}^{(0)-1}(\pi_{\beta})} \right] \\
\stackrel{(ii)}{\leq} \sqrt{\frac{8}{\tau_{\beta}}(K+1)} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} ||\boldsymbol{x}||_{\mathbb{W}_{\mathcal{D}_{\beta}}^{(0)-1}(\pi_{\beta,0})} \right] \\
\stackrel{(iii)}{\leq} \sqrt{\frac{8}{\tau_{\beta}}(K+1)d\log d} \\
\leq 4\sqrt{\frac{Kd\log d}{\tau_{\beta}}}$$

where (i) follows from Lemma 7.13 and the fact that  $\mathcal{D}_{\beta+1} \subset \mathcal{D}_{\beta}$ , (ii) follows in a similar manner as Lemma 7.14, and (iii) follows from Lemma 9.7.

**Lemma 7.18.** Let  $\tilde{X}_{\beta}$  and  $H_{\beta}$  be as defined in Section 7.1. Denote  $\tau_{\beta} = |\mathcal{T}_{\beta}|$ . Then, we have that

$$\mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}}\left[\max_{\boldsymbol{x}\in\mathcal{X}}\left\|\left|\tilde{\boldsymbol{X}}_{\beta}^{\top}\boldsymbol{H}_{\beta}^{-1/2}\right|\right|_{2}\right] \leq 4K^{2}\sqrt{\frac{d\log(Kd)}{\tau_{\beta}}}$$

Proof.

$$\begin{split} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta+1}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left\| \left| \tilde{\boldsymbol{X}}_{\beta}^{\top} \boldsymbol{H}_{\beta}^{-1/2} \right\|_{2} \right] \stackrel{(i)}{\leq} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \sqrt{\lambda_{max} \left( \boldsymbol{H}_{\beta}^{-1/2} \tilde{\boldsymbol{X}}_{\beta} \tilde{\boldsymbol{X}}_{\beta}^{\top} \boldsymbol{H}_{\beta}^{-1/2} \right)} \right] \\ \stackrel{(ii)}{=} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \sqrt{\lambda_{max} \left( \tilde{\boldsymbol{X}}_{\beta}^{\top} \boldsymbol{H}_{\beta}^{-1} \tilde{\boldsymbol{X}}_{\beta} \right)} \right] \\ \stackrel{(iii)}{\leq} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \sqrt{\sum_{i=1}^{K} \left\| \left| \tilde{\boldsymbol{x}}_{\beta}^{(i)} \right\|_{\mathbf{H}_{\beta}^{-1}}^{2}} \right] \\ \stackrel{(iv)}{\leq} \sqrt{\frac{8}{\tau_{\beta}}} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \sqrt{\sum_{i=1}^{K} \left\| \left| \tilde{\boldsymbol{x}}_{\beta}^{(i)} \right\|_{\mathbb{T}_{D_{\beta}}^{2}(\pi_{\beta})}^{2}} \right] \\ \stackrel{(v)}{\leq} \sqrt{\frac{8}{\tau_{\beta}} (K+1)} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \sqrt{\sum_{i=1}^{K} \left\| \left| \tilde{\boldsymbol{x}}_{\beta}^{(i)} \right\|_{\mathbb{T}_{D_{\beta}}^{2}(\pi_{\beta},i)}^{2}} \right] \\ \stackrel{(vi)}{\leq} \sqrt{\frac{8}{\tau_{\beta}} (K+1)} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \sqrt{\sum_{i=1}^{K} \left\| \left| \tilde{\boldsymbol{x}}_{\beta}^{(i)} \right\|_{\mathbb{W}_{D_{\beta}}^{(i)}(\pi_{\beta},i)}^{2}} \right] \\ \stackrel{(vii)}{\leq} \sqrt{\frac{8}{\tau_{\beta}} (K+1)} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \sum_{i=1}^{K} \left\| \left| \tilde{\boldsymbol{x}}_{\beta}^{(i)} \right\|_{\mathbb{W}_{D_{\beta}}^{(i)}(\pi_{\beta},i)}^{2}} \right] \\ \stackrel{(vii)}{\leq} \sqrt{\frac{8}{\tau_{\beta}} (K+1)} \sum_{i=1}^{K} \mathbb{E}_{\mathcal{X}\sim\mathcal{D}_{\beta}} \left[ \max_{\boldsymbol{x}\in\mathcal{X}} \left\| \left| \tilde{\boldsymbol{x}}_{\beta}^{(i)} \right\|_{\mathbb{W}_{D_{\beta}}^{(i)}(\pi_{\beta},i)} \right] \\ \stackrel{(viii)}{\leq} K \sqrt{\frac{8}{\tau_{\beta}} (K+1)} K d\log(Kd) \leq 4K^{2} \sqrt{\frac{d\log(Kd)}{\tau_{\beta}}} \end{split} \right]$$

where (i) follows from the definition of the norm  $||\mathbf{A}||_2 = \sqrt{\lambda_{max} (\mathbf{A}^\top \mathbf{A})}$  and the fact that  $\mathcal{D}_{\beta+1} \subset \mathcal{D}_{\beta}$ , (ii) follows from the cyclic property of eigenvalues, (iii) follows from Lemma 7.12, (iv) follows from Lemma 7.13, (v) follows from Lemma 7.14, (vi) follows from Lemma 7.15, (vii) uses the fact that for  $\{a_i\}_{i=1}^N$ ,  $\sqrt{\sum_{i=1}^N a_i^2} \leq \sqrt{\left(\sum_{i=1}^N a_i\right)^2} = \sum_{i=1}^N a_i$ , (viii) uses the linearity of expectations and the fact that  $\max_x [f(x) + g(x)] \leq \max_x f(x) + \max_x g(x)$ , and (ix) follows from Lemma 9.7.

## 8 Rarely Switching Multinomial Contextual Bandit Algorithm: RS-MNL

#### 8.1 Notations

We first define a few matrices, vectors, and scalars that are used throughout this section (here,  $e_i$  denotes the  $i^{th}$ -standard basis vector):

1. 
$$V_t = \lambda I_{d \times d} + \sum_{s \in [t]} x_s x_s^{\top}$$
  
2.  $\tilde{V}_t = I_{K \times K} \otimes V_t$   
3.  $A(x, \theta) = \operatorname{diag}(z(x, \theta)) - z(x, \theta)z(x, \theta)^{\top}$   
4.  $M(x, \theta_1, \theta_2) = \int_0^1 A(x, v\theta_1 + (1 - v)\theta_2) dv$   
5.  $H_t^{\star} = \lambda I_{Kd \times Kd} + \sum_{s \in [t]} A(x_s, \theta^{\star}) \otimes x_s x_t^{\top}$   
6.  $\gamma(\delta) = CS^{5/4} \sqrt{Kd \log(T/\delta)}$   
7.  $B_t(x) = \exp\left(\sqrt{6} \min\left\{2\kappa^{1/2}\gamma(\delta) \|x\|_{V_t^{-1}}, 2S\right\}\right)$   
8.  $H_t(\theta) = \lambda I_{Kd \times Kd} + \sum_{s \in [t]} \frac{A(x_s, \theta)}{B_s(x_s)} \otimes x_s x_s^{\top}$   
9.  $H_t(\theta) = \lambda I_{Kd \times Kd} + \sum_{s \in [t]} \frac{A(x_s, \theta)}{B_s(x_s)} \otimes x_s x_s^{\top}$   
10.  $\tilde{X}_t(\theta) = \frac{A(x_t, \theta)^{\frac{1}{2}}}{\sqrt{B_t(x_t)}} e_i \otimes x_t$ 

12. 
$$\boldsymbol{H}_{t}^{i}(\boldsymbol{\theta}) = \sum_{s \in [t]} \tilde{\boldsymbol{x}}_{s}^{(i)}(\boldsymbol{\theta}) \tilde{\boldsymbol{x}}_{s}^{(i)}(\boldsymbol{\theta})^{\top} + \lambda \boldsymbol{I}$$

We now present the regret upper bound for RS-MNL by restating Theorem 4.1.

**Theorem 8.1.** With high probability, the regret incurred by Algorithm 3 is bounded above by  $R_T$  where:

$$R_T \le CRK^{3/2}S^{5/4}(\log T\log(T/\delta))^{1/2}d\sqrt{T} + CRK^2d^2S^{5/2}\log T\log(T/\delta)\kappa^{1/2}e^{2S}(e^S + K\kappa^{1/2})$$

*Proof.* For any round  $t \in [T]$ , let  $\tau_t \leq t$  denote the last round at which a switch was made. Then, using the value of  $\gamma(\delta)$  alongside Lemma 8.8, Lemma 8.12, and Lemma 8.13, we get:

$$\begin{split} R(T) &\leq \sum_{t \in [T]} |\boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}^{\star}, \boldsymbol{\theta}^{\star}) - \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star})| \leq \sum_{t \in [T]} 2\epsilon_{1}(t, \tau_{t}, \boldsymbol{x}_{t}) + 2\epsilon_{2}(t, \tau_{t}, \boldsymbol{x}_{t}) \\ &\leq 4RKd^{1/2}(\log T)^{1/2}\gamma(\delta)\sqrt{T} + 8RKd\log T\kappa^{1/2}e^{3S}\gamma(\delta)^{2} + 24dRK^{2}e^{2S}\kappa\gamma(\delta)^{2}\log T \\ &\leq CRK^{3/2}S^{5/4}(\log T\log(T/\delta))^{1/2}d\sqrt{T} + CRK^{2}d^{2}S^{5/2}\log T\log(T/\delta)\kappa^{1/2}e^{2S}(e^{S} + K\kappa^{1/2}) \\ & \Box \end{split}$$

#### 8.2 Supporting Lemmas for 8

**Lemma 8.1.** Let  $\{x_1, \ldots, x_{\tau}\}$  be a set of arms and  $\{r_1, \ldots, r_{\tau}\}$  be the set of corresponding rewards associated with the arms. Define  $\hat{\theta}_{\tau}$  be the MLE estimate calculated using this set of arms and rewards, *i.e* 

$$\hat{\boldsymbol{\theta}}_{\tau} = \arg\min_{\boldsymbol{\theta}} \sum_{s \in [\tau]} \sum_{i=1}^{K} \mathbb{1}\{y_s = i\} \log z_i(\boldsymbol{x}_s, \boldsymbol{\theta}) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

Let  $H_{\tau}^{\star}$  be as defined in Section 8.1. Then, with high probability, and the choice of  $\lambda = KdS^{-1/2}\log(T/\delta)$ , we have that

$$\|\hat{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}^{\star}\|_{\boldsymbol{H}_{\tau}^{\star}} \leq CS^{5/4}\sqrt{Kd\log(T/\delta)}$$

*Proof.* We define  $G_{\tau}(\theta_1, \theta_2)$  as:

$$G_{ au}(oldsymbol{ heta}_1,oldsymbol{ heta}_2) = \sum_{t\in[ au]} M(oldsymbol{x}_t,oldsymbol{ heta}_1,oldsymbol{ heta}_2) \otimes oldsymbol{x}_toldsymbol{x}_t^ op + \lambda oldsymbol{I}$$

where  $M(\boldsymbol{x}, \boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$  is as defined in Section 8.1. Thus, from Lemma 9.2, we have that

$$(1+2S)^{-1}\boldsymbol{G}_{\tau} \succcurlyeq \boldsymbol{H}_{\tau}^{\star}$$

Thus, we have

$$\begin{split} \left\| \hat{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}^{*} \right\|_{\boldsymbol{H}_{\tau}^{*}} &\leq \sqrt{1+2S} \left\| \hat{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}^{*} \right\|_{\boldsymbol{G}_{\tau}(\hat{\boldsymbol{\theta}}_{\tau},\boldsymbol{\theta}^{*})} \\ &\leq \sqrt{1+2S} \left\| \left[ \sum_{t \in [\tau]} \boldsymbol{M}(\hat{\boldsymbol{\theta}}_{\tau},\boldsymbol{\theta}^{*}) \otimes \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top} + \lambda \boldsymbol{I}_{Kd \times Kd} \right] \left( \hat{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}^{*} \right) \right\|_{\boldsymbol{G}_{\tau}^{-1}(\hat{\boldsymbol{\theta}}_{\tau},\boldsymbol{\theta}^{*})} \\ & \stackrel{(i)}{\leq} \sqrt{1+2S} \left\| \sum_{t \in [\tau]} \left[ \boldsymbol{M}(\boldsymbol{x},\boldsymbol{\theta}^{*},\hat{\boldsymbol{\theta}}_{\tau}) \otimes \boldsymbol{x}_{t}^{\top} \right] \left( \boldsymbol{\theta}^{*} - \hat{\boldsymbol{\theta}}_{\tau} \right) \otimes \boldsymbol{x}_{t} + \lambda \left( \hat{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}^{*} \right) \right\|_{\boldsymbol{G}_{\tau}^{-1}(\hat{\boldsymbol{\theta}}_{\tau},\boldsymbol{\theta}^{*})} \\ & \stackrel{(iii)}{\leq} \sqrt{1+2S} \left\| \sum_{t \in [\tau]} \left[ \boldsymbol{X}(\boldsymbol{x},\boldsymbol{\theta}^{*},\hat{\boldsymbol{\theta}}_{\tau}) \otimes \boldsymbol{x}_{t}^{\top} \right] \left( \boldsymbol{\theta}^{*} - \hat{\boldsymbol{\theta}}_{\tau} \right) \otimes \boldsymbol{x}_{t} + \lambda \left( \hat{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}^{*} \right) \right\|_{\boldsymbol{G}_{\tau}^{-1}(\boldsymbol{\theta}_{1},\boldsymbol{\theta}_{2})} \\ & \stackrel{(iii)}{\leq} \sqrt{1+2S} \left\| \sum_{t \in [\tau]} \left[ \boldsymbol{z}(\boldsymbol{x}_{t},\hat{\boldsymbol{\theta}}_{\tau}) - \boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{*}) \right] \otimes \boldsymbol{x}_{t} + \lambda \left( \hat{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}^{*} \right) \right\|_{\boldsymbol{G}_{\tau}^{-1}(\boldsymbol{\theta}_{1},\boldsymbol{\theta}_{2})} \\ & \stackrel{(iv)}{\leq} \sqrt{1+2S} \left\| \sum_{t \in [\tau]} \left[ \boldsymbol{m}_{t} - \boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{*}) \right] \otimes \boldsymbol{x}_{t} - \lambda \boldsymbol{\theta}^{*} \right\|_{\boldsymbol{G}_{\tau}^{-1}(\boldsymbol{\theta}_{1},\boldsymbol{\theta}_{2})} \\ & \stackrel{(v)}{\leq} \sqrt{1+2S} \left\| \sum_{t \in [\tau]} \left[ \boldsymbol{m}_{t} - \boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{*}) \right] \otimes \boldsymbol{x}_{t} \right\|_{\boldsymbol{H}_{\tau}^{*}^{-1}} \\ & + \lambda \sqrt{1+2S} \left\| \boldsymbol{\theta}^{*} \right\|_{\boldsymbol{G}_{\tau}^{-1}(\boldsymbol{\theta}_{1},\boldsymbol{\theta}_{2})} \\ & \stackrel{(vi)}{\leq} (1+2S) \left\| \sum_{t \in [\tau]} \left[ \boldsymbol{m}_{t} - \boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{*}) \right] \otimes \boldsymbol{x}_{t} \right\|_{\boldsymbol{H}_{\tau}^{*}^{-1}} \\ & + \lambda \sqrt{1+2S} \left\| \boldsymbol{\theta}^{*} \right\|_{\boldsymbol{G}_{\tau}^{-1}(\boldsymbol{\theta}_{1},\boldsymbol{\theta}_{2})} \\ & \stackrel{(vi)}{\leq} 3S \left\| \sum_{t \in [\tau]} \boldsymbol{\epsilon}_{t} \otimes \boldsymbol{x}_{t} \right\|_{\boldsymbol{H}_{\tau}^{*}^{-1}} \\ & + \sqrt{3}\lambda^{1/2}S^{3/2} \end{aligned} \right\}$$

where (i) follows from Lemma 9.2, (ii) follows from Mixed Product Property, (iii) follows from the Mean value Theorem, (iv) from Lemma 9.3, (v) follows from Cauchy-Schwarz, and (vi) follows from the fact that  $G_{\tau} \geq \lambda I$  and  $||\boldsymbol{\theta}||_2 \leq S$ .

Note that  $\epsilon_t = m_t - z(x_t, \theta^*)$  and since  $\mathbb{E}[m_t] = z(x_t, \theta^*)$ , we get  $\mathbb{E}[\epsilon_t \epsilon_t^\top] = A(x_t, \theta^*)$ . Also, note that  $\|\epsilon_t\|_1 \le \|m_t\|_1 + \|z(x_t, \theta^*)\|_1 \le 2$ . Thus, using Lemma 9.10, we get

$$\left\| \hat{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}^{\star} \right\|_{\boldsymbol{H}_{\tau}^{\star}} \leq 3S \left( \frac{\sqrt{\lambda}}{4} + \frac{4}{\sqrt{\lambda}} \log \left( \frac{\det \boldsymbol{H}_{\tau}^{1/2}}{\delta \lambda^{\frac{dK}{2}}} \right) + \frac{4}{\sqrt{\lambda}} K d \log 2 \right) + 2S^{3/2} \lambda^{1/2}$$

where  $\boldsymbol{H}_{\tau} = \lambda \boldsymbol{I} + \sum_{t \in \tau} A(\boldsymbol{x}_t, \boldsymbol{\theta}^{\star}) \otimes \boldsymbol{x}_t \boldsymbol{x}_t^{\top}$ .

We can calculate det  $H_{\tau}$  as follows:

$$\det \boldsymbol{H}_{\tau} \stackrel{(i)}{\leq} \left(\frac{\operatorname{trace} \boldsymbol{H}_{\tau}}{Kd}\right)^{Kd} \\ \leq \left(\frac{\operatorname{trace} \lambda \boldsymbol{I} + \operatorname{trace} \sum_{t \in \tau} A(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) \otimes \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top}}{Kd}\right)^{Kd} \\ \stackrel{(ii)}{\leq} \left(\frac{\lambda Kd + \tau \|\boldsymbol{x}_{t}\|_{2}^{2}}{Kd}\right)^{Kd} \\ \stackrel{(iii)}{\leq} \lambda^{Kd} \left(1 + \frac{\tau}{\lambda Kd}\right)^{Kd}$$

where (i) follows from Lemma 9.11, (ii) follows from the fact that tr  $(\mathbf{A} \otimes \mathbf{B}) = \sum \lambda(\mathbf{A})\lambda(\mathbf{B})$  and the fact that  $\mathbf{A}(\mathbf{x}, \boldsymbol{\theta}^{\star}) \preccurlyeq \mathbf{I}$  and the only non-zero eigenvalue of  $\mathbf{x}_t \mathbf{x}_t^{\top}$  is  $\|\mathbf{x}_t\|_2^2$ , and (iii) follows since  $\|\mathbf{x}\| \le 1$ .

Thus, we have

$$\begin{split} \left| \left| \hat{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}^{\star} \right| \right|_{\boldsymbol{H}_{\tau}^{\star}} &\leq 3S \left( \frac{\sqrt{\lambda}}{4} + \frac{4}{\sqrt{\lambda}} \log \left( \frac{\left(1 + \frac{\tau}{\lambda K d}\right)^{\frac{Kd}{2}}}{\delta} \right) + \frac{4}{\sqrt{\lambda}} K d \log 2 \right) + 2S^{3/2} \lambda^{1/2} \\ &= 3S \left( \frac{\sqrt{\lambda}}{4} + \frac{2Kd}{\sqrt{\lambda}} \log \left( 1 + \frac{\tau}{\lambda d} \right) + \frac{4}{\sqrt{\lambda}} \log \frac{1}{\delta} + \frac{4}{\sqrt{\lambda}} K d \log 2 \right) + 2S^{3/2} \lambda^{1/2} \end{split}$$

Finally, by setting  $\lambda = K dS^{-1/2} \log(T/\delta)$  and simplifying the constants, we get that for some appropriately tuned constant C

$$\left\| \hat{\boldsymbol{\theta}}_{\tau} - \boldsymbol{\theta}^{\star} \right\|_{\boldsymbol{H}_{\tau}^{\star}} \leq C S^{5/4} \sqrt{K d \log(T/\delta)}$$

From here on, we shall use the notation  $\gamma(\delta) = CS^{5/4}\sqrt{Kd\log(T/\delta)}$ . Lemma 8.2. Let  $\tilde{V}_t$  and  $H_t^*$  be defined as in Section 8.1. Then, for any round  $t \in [T]$ , we have that

$$V_t \preccurlyeq \kappa H_t^{\star}$$

*Proof.* From the definition of  $\kappa$ , we have  $A(x, \theta) \geq \frac{1}{\kappa}I$ . Hence, using the fact that  $\kappa \geq 1$ , we have

$$\begin{split} ilde{oldsymbol{V}_t} &= oldsymbol{I}_{K imes K} \otimes oldsymbol{V}_t = oldsymbol{I}_{K imes K} \otimes o$$

**Lemma 8.3.** Let  $\tilde{V}_t$  and  $H_t(\theta)$  be defined as in Section 8.1. Then, for any round  $t \in [T]$ , we have that

$$V_t \preccurlyeq \kappa H_t(\theta)$$

*Proof.* From the definition of  $\kappa$ , we have  $A(x, \theta) \geq \frac{1}{\kappa}I$ . Hence, using the fact that  $\kappa \geq 1$ , we have

$$\begin{split} \tilde{\boldsymbol{V}}_{t} &= \boldsymbol{I}_{K \times K} \otimes \boldsymbol{V}_{t} = \boldsymbol{I}_{K \times K} \otimes \left( \lambda \boldsymbol{I}_{d \times d} + \sum_{s \in [t]} \boldsymbol{x}_{s} \boldsymbol{x}_{s}^{\top} \right) \\ &= \lambda \boldsymbol{I}_{K d \times K d} + \boldsymbol{I}_{K \times K} \otimes \sum_{s \in [t]} \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top} \\ &\preccurlyeq \kappa \lambda \boldsymbol{I}_{K d \times K d} + \kappa \sum_{s \in [t]} \boldsymbol{A}(\boldsymbol{x}_{t}, \boldsymbol{\theta}) \otimes \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top} \\ &\preccurlyeq \kappa \lambda \boldsymbol{I}_{K d \times K d} + \kappa \sum_{s \in [t]} \frac{\boldsymbol{A}(\boldsymbol{x}_{t}, \boldsymbol{\theta})}{B_{s}(\boldsymbol{x}_{s})} \otimes \boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\top} \\ &\preccurlyeq \kappa \boldsymbol{H}_{t}(\boldsymbol{\theta}) \end{split}$$

where the second to last inequality follows since  $B_t(x) \ge 1$ .

**Lemma 8.4.** Let  $1, \tau_1, \ldots, \tau_m$  be the rounds at which a switch occurs, i.e det  $H_{\tau_{i+1}}(\hat{\theta}_{\tau_i}) \geq 2 \det H_{\tau_i}(\hat{\theta}_{\tau_i}) \forall i \in [m]$ . Let  $H_t(\theta)$  and  $H_t^*$  be defined as in Section 8.1. Then, for all  $i \in [m]$ , we have that

$$oldsymbol{H}_{ au_i}(\hat{oldsymbol{ heta}}_{ au_i}) \preccurlyeq oldsymbol{H}_{ au_i}^\star$$

*Proof.* From Lemma 9.1, for some x such that  $||x|| \le 1$  and some  $\tau \in \{\tau_1, \ldots, \tau_m\}$ , we have that

$$A(\boldsymbol{x}, \hat{\boldsymbol{\theta}}_{\tau}) \preccurlyeq A(\boldsymbol{x}, \boldsymbol{\theta}^{\star}) \exp\left(\sqrt{6} \left\| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\tau}) \right\|_{2} \right)$$

Now, we can bound  $\left\| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\tau}) \right\|_2$  as follows:

$$\begin{split} \left\| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\tau_t}) \right\|_2 \stackrel{(i)}{\leq} 2S \| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) \|_2 \stackrel{(ii)}{=} 2S \sqrt{\lambda_{max} \left( (\boldsymbol{I} \otimes \boldsymbol{x}) (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) \right)} \\ \stackrel{(iii)}{=} 2S \sqrt{\lambda_{max} \left( \boldsymbol{I} \otimes \boldsymbol{x} \boldsymbol{x}^{\top} \right)} \stackrel{(iv)}{\leq} 2S \end{split}$$

where (i) uses Cauchy-Schwarz inequality and the fact that  $\|\boldsymbol{\theta}\|_2 \leq S$ , (ii) uses the definition of the norm as  $\|\boldsymbol{A}\|_2 = \sqrt{\lambda_{max} (\boldsymbol{A}^\top \boldsymbol{A})}$ , (iii) follows from the mixed product property of tensor products, and (iv) follows from the fact that  $\lambda_{max} (\boldsymbol{A} \otimes \boldsymbol{B}) = \lambda_{max} (\boldsymbol{A}) \lambda_{max} (\boldsymbol{B})$  and  $\lambda_{max} (\boldsymbol{x}\boldsymbol{x}^\top) = \|\boldsymbol{x}\|_2^2 \leq 1$ .

We can also bound  $\left\| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top})(\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\tau}) \right\|_{2}$  in the following way (note that the *d*-dimensional unit ball is represented as  $\mathcal{B}_{2}(d)$ ):

$$\begin{split} \|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})(\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\tau})\|_{2} &= \|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})\boldsymbol{H}_{\tau}^{\star^{-1/2}}\boldsymbol{H}_{\tau}^{\star^{1/2}}(\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\tau})\|_{2} \\ &\stackrel{(i)}{\leq} \|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})\boldsymbol{H}_{\tau}^{\star^{-1/2}}\|_{2}\gamma(\delta) \\ &\stackrel{(ii)}{\leq} \kappa^{1/2}\|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})\tilde{\boldsymbol{V}}_{\tau}^{-1/2}\|_{2}\gamma(\delta) \\ &\stackrel{(iii)}{\leq} \kappa^{1/2}\|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})(\boldsymbol{I} \otimes \boldsymbol{V}_{\tau}^{-1/2})\|_{2}\gamma(\delta) \\ &\stackrel{(iv)}{=} \kappa^{1/2}\sqrt{\lambda_{max}\left((\boldsymbol{I} \otimes \boldsymbol{V}_{\tau}^{-1/2})(\boldsymbol{I} \otimes \boldsymbol{x})(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})(\boldsymbol{I} \otimes \boldsymbol{V}_{\tau}^{-1/2})\right)}\gamma(\delta) \\ &\stackrel{(v)}{=} \kappa^{1/2}\gamma(\delta)|\boldsymbol{x}\|_{\boldsymbol{V}_{\tau}^{-1}} \\ &\leq 2\kappa^{1/2}\gamma(\delta)|\boldsymbol{x}\|_{\boldsymbol{V}_{\tau}^{-1}} \end{split}$$

where (i) is obtained from the fact that  $\|Ax\|_2 \leq \|A\|_2 \|x\|_2$  and from Lemma 8.1, (ii) follows from Lemma 8.2, (iii) is obtained from the definition of  $\tilde{V}$  and the fact that  $(A \otimes B)^n = A^n \otimes B^n$ , (iv) follows from the definition of the norm, i.e,  $\|A\|_2 = \sqrt{\lambda_{max} (A^{\top}A)}$ , and (v) follows from the cyclic property of eigenvalues and the fact that  $\lambda_{max} (A \otimes B) = \lambda_{max} (A) \lambda_{max} (B)$ .

Thus, by combining both bounds, we obtain

$$A(\boldsymbol{x}, \hat{\boldsymbol{\theta}}_{\tau}) \preccurlyeq A(\boldsymbol{x}, \boldsymbol{\theta}^{\star}) \exp\left(\sqrt{6} \min\left\{\sqrt{2}\kappa^{1/2}\gamma(\delta) |\boldsymbol{x}\|_{\boldsymbol{V}_{\tau}^{-1}}, 2S\right\}\right)$$

Let  $B_{\tau}(\boldsymbol{x})$  denote the value  $\exp\left(\sqrt{6}\min\left\{\sqrt{2\kappa^{1/2}\gamma(\delta)}|\boldsymbol{x}\|_{\boldsymbol{V}_{\tau}^{-1}},2S\right\}\right)$ . Then, we have that

$$\boldsymbol{H}_{\tau}^{\star} = \lambda \boldsymbol{I} + \sum_{s \in [\tau]} \boldsymbol{A}(\boldsymbol{x}_{s}, \boldsymbol{\theta}^{\star}) \otimes \boldsymbol{x}_{s} \boldsymbol{x}_{s} \succcurlyeq \lambda \boldsymbol{I} + \sum_{s \in [\tau]} \frac{\boldsymbol{A}(\boldsymbol{x}_{s}, \hat{\boldsymbol{\theta}}_{\tau})}{B_{\tau}(\boldsymbol{x}_{s})} \otimes \boldsymbol{x}_{s} \boldsymbol{x}_{s} = \boldsymbol{H}_{\tau}(\hat{\boldsymbol{\theta}}_{\tau})$$

г		

**Lemma 8.5.** For time round t, let  $\tau_t \leq t$  be the last time round at which a switch occurred, i.e det  $H_t(\hat{\theta}_{\tau_t}) \leq 2$  det  $H_{\tau_t}(\hat{\theta}_{\tau_t})$ . Let  $H_t(\theta)$  and  $H_t^*$  be defined as in Section 8.1.

$$\boldsymbol{H}_t(\hat{\boldsymbol{ heta}}_{ au_t}) \preccurlyeq \boldsymbol{H}_t^\star$$

*Proof.* Similar to Lemma 8.4 for some x such that  $||x|| \le 1$ , we have that

$$A(\boldsymbol{x}, \hat{\boldsymbol{\theta}}_{\tau_t}) \preccurlyeq A(\boldsymbol{x}, \boldsymbol{\theta}^{\star}) \exp\left(\sqrt{6} \left\| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\tau_t}) \right\|_2 \right)$$

Now, we can bound  $\left\| (\mathbf{I} \otimes \mathbf{x}^{\top}) (\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\tau_t}) \right\|_2$  in two different ways: the first way results in 2*S*, following the same method as Lemma 8.4. We can also bound it in the following way:

$$\begin{split} \|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})(\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\tau_{t}})\|_{2} &= \|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})\boldsymbol{H}_{\tau_{t}}^{\star}{}^{-1/2}\boldsymbol{H}_{\tau_{t}}^{\star}{}^{1/2}(\boldsymbol{\theta}^{\star} - \hat{\boldsymbol{\theta}}_{\tau_{t}})\|_{2} \\ &\stackrel{(i)}{\leq} \|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})\boldsymbol{H}_{\tau_{t}}^{\star}{}^{-1/2}\|_{2}\gamma(\delta) \\ &\stackrel{(iii)}{\leq} \|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})\boldsymbol{H}_{\tau_{t}}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1/2}\|_{2}\gamma(\delta) \\ &\stackrel{(iii)}{\leq} \sqrt{2}\|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})\boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1/2}\|_{2}\gamma(\delta) \\ &\stackrel{(iv)}{\leq} \sqrt{2}\kappa^{1/2}\|(\boldsymbol{I} \otimes \boldsymbol{x}^{\top})\tilde{\boldsymbol{V}}^{-1/2}\|_{2}\gamma(\delta) \\ &\leq 2\kappa^{1/2}\gamma(\delta)|\boldsymbol{x}\|_{\boldsymbol{V}_{t}^{-1}} \end{split}$$

where (i) is obtained from the fact that  $\|Ax\|_2 \leq \|A\|_2 \|x\|_2$  and from Lemma 8.1, (ii) follows from Lemma 8.4, (iii) follows from the combination of Lemma 9.13 and the fact that det  $H_t(\hat{\theta}_{\tau_t}) \leq$ 2 det  $H_{\tau_t}(\hat{\theta}_{\tau_t})$ , (iv) follows from Lemma 8.3, and (v) follows from the same steps used in Lemma 8.6.

Combining the bounds in the same manner as Lemma 8.4 finishes the proof.

**Lemma 8.6.** For time round t, let  $\tau_t \leq t$  be the last time round at which a switch occurred. Let  $H_t^{(i)}(\hat{\theta}_{\tau_t})$  and  $H_t(\hat{\theta}_{\tau_t})$  be defined as in Section 8.1. Then, we have

$$\boldsymbol{H}_{t}^{(i)}(\hat{\boldsymbol{ heta}}_{ au_{t}}) \preccurlyeq \boldsymbol{H}_{t}(\hat{\boldsymbol{ heta}}_{ au_{t}})$$

Proof. We have:

$$\begin{split} \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}}) &= \lambda \boldsymbol{I} + \sum_{s \in [t]} \tilde{\boldsymbol{X}}_{s}(\hat{\boldsymbol{\theta}}_{\tau_{t}}) \tilde{\boldsymbol{X}}_{s}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{\top} \\ &\stackrel{(i)}{=} \lambda \boldsymbol{I} + \sum_{s \in [t]} \sum_{i=1}^{K} \tilde{\boldsymbol{x}}_{s}^{(i)}(\hat{\boldsymbol{\theta}}_{\tau_{t}}) \tilde{\boldsymbol{x}}_{s}^{(i)}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{\top} \\ & \succcurlyeq \lambda \boldsymbol{I} + \sum_{s \in [t]} \tilde{\boldsymbol{x}}_{s}^{(i)}(\hat{\boldsymbol{\theta}}_{\tau_{t}}) \tilde{\boldsymbol{x}}_{s}^{(i)}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{\top} \\ &= \boldsymbol{H}_{\tau}^{i}(\hat{\boldsymbol{\theta}}_{\tau_{t}}) \end{split}$$

where (i) follows from Lemma 7.11.

**Lemma 8.7.** Let  $\tau_t \leq t$  be the last time round at which a switch was made. In other words, det  $H_t(\hat{\theta}_{\tau_t}) \leq 2$  det  $H_{\tau_t}(\hat{\theta}_{\tau_t})$ . Then, for any arm  $\boldsymbol{x}$ , we have that,

$$\left| \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}, \boldsymbol{\theta}^{\star}) - \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}, \hat{\boldsymbol{\theta}}_{\tau}) \right| \leq \epsilon_1(t, \tau_t, \boldsymbol{x}) + \epsilon_2(t, \tau_t, \boldsymbol{x})$$

where

$$\epsilon_{1}(t,\tau_{t},\boldsymbol{x}) = \sqrt{2}\gamma(\delta) \left\| \left| \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1/2} (\boldsymbol{I} \otimes \boldsymbol{x}) \boldsymbol{A}(\boldsymbol{x},\hat{\boldsymbol{\theta}}_{\tau_{t}}) \boldsymbol{\rho} \right\|_{2} \\ \epsilon_{2}(t,\tau_{t},\boldsymbol{x}) = 6R\gamma(\delta)^{2} \left\| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^{\top}) \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1/2} \right\|_{2}^{2} \right\|_{2}$$

*Proof.* The proof follows on the same lines as that of Lemma 7.5 and uses the fact that  $\frac{\det \boldsymbol{H}_{\tau_t}(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1}}{\det \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1}} \leq 2$  combined with Lemma 9.13 to convert  $\boldsymbol{H}_{\tau_t}(\hat{\boldsymbol{\theta}}_{\tau_t})$  to  $\boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})$ .

**Lemma 8.8.** Let  $\tau_t \leq t$  be the last time step at which a switch was made. Let  $\epsilon_1(t, \tau_t, \mathbf{x})$  and  $\epsilon_2(t, \tau_t, \mathbf{x})$  be as defined in Lemma 8.7. Then, the regret at time step t can be bounded as

$$\left|\boldsymbol{\rho}^{\top}\boldsymbol{z}(\boldsymbol{x}^{*},\boldsymbol{\theta}^{\star})-\boldsymbol{\rho}^{\top}\boldsymbol{z}(\boldsymbol{x}_{t},\boldsymbol{\theta}^{\star})\right|\leq 2\epsilon_{1}(t,\tau_{t},\boldsymbol{x}_{t})+2\epsilon_{2}(t,\tau_{t},\boldsymbol{x}_{t})$$

Proof.

$$\begin{aligned} \left| \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}^{\star}, \boldsymbol{\theta}^{\star}) - \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}, \boldsymbol{\theta}^{\star}) \right| & \stackrel{(i)}{\leq} \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}^{\star}, \hat{\boldsymbol{\theta}}_{\tau_{t}}) + \epsilon_{1}(t, \tau_{t}, \boldsymbol{x}^{\star}) + \epsilon_{2}(t, \tau_{t}, \boldsymbol{x}^{\star}) - \boldsymbol{\rho}^{\top} \boldsymbol{z}(\boldsymbol{x}_{t}, \hat{\boldsymbol{\theta}}_{\tau_{t}}) + \epsilon_{1}(t, \tau_{t}, \boldsymbol{x}_{t}) + \epsilon_{2}(t, \tau_{t}, \boldsymbol{x}_{t}) \\ & \stackrel{(ii)}{\leq} 2\epsilon_{1}(t, \tau_{t}, \boldsymbol{x}_{t}) + 2\epsilon_{2}(t, \tau_{t}, \boldsymbol{x}_{t}) \end{aligned}$$

where (i) follows from Lemma 8.7 and (ii) follows from the fact that  $\boldsymbol{x}_t = \arg \max_{\boldsymbol{x} \in \mathcal{X}} \text{UCB}(t, \tau_t, \boldsymbol{x}) = \arg \max_{\boldsymbol{x} \in \mathcal{X}} \left[ \boldsymbol{\rho}^\top \boldsymbol{z}(\boldsymbol{x}, \hat{\boldsymbol{\theta}}_{\tau_t}) + \epsilon_1(t, \tau_t, \boldsymbol{x}) + \epsilon_2(t, \tau_t, \boldsymbol{x}) \right]$  and hence, gets selected at round t.

**Lemma 8.9.** Let  $B_t(x)$  be as defined in Section 8.1. Then, we have that

$$\sqrt{B_t(\boldsymbol{x})} \le e^{3S} \kappa^{1/2} \gamma(\delta) \|\boldsymbol{x}\|_{\boldsymbol{V}_t^{-1}} + 1$$

Proof.

$$\begin{split} \sqrt{B_t(\boldsymbol{x})} &= \exp\left(\sqrt{6}\min\left\{\kappa^{1/2}\gamma(\delta)\|\boldsymbol{x}\|_{\boldsymbol{V}_t^{-1}}, S\right\}\right)\\ &\stackrel{(i)}{\leq} e^{3S}\kappa^{1/2}\gamma(\delta)\|\boldsymbol{x}\|_{\boldsymbol{V}_t^{-1}} + 1 \end{split}$$

where (i) follows from Lemma 9.6 by choosing  $\min \left\{ \kappa^{1/2} \gamma(\delta) \| \boldsymbol{x} \|_{\boldsymbol{V}_t^{-1}}, S \right\} = \kappa^{1/2} \gamma(\delta) \| \boldsymbol{x} \|_{\boldsymbol{V}_t^{-1}}$ and  $M = \sqrt{6}S$ .

**Lemma 8.10.** Let  $\tilde{X}_{\tau}(\theta)$  and  $\tilde{x}_{\tau}^{(i)}(\theta)$  be defined as in Section 8.1. Then, we have

$$\tilde{\boldsymbol{X}}_{\tau}(\boldsymbol{\theta})\tilde{\boldsymbol{X}}_{\tau}(\boldsymbol{\theta})^{\top} = \sum_{i=1}^{K} \tilde{\boldsymbol{x}}_{\tau}^{(i)}(\boldsymbol{\theta})\tilde{\boldsymbol{x}}_{\tau}^{(i)}(\boldsymbol{\theta})^{\top}$$

*Proof.* The proof follows on the same lines as Lemma 7.11.

**Lemma 8.11.** Let  $M \in \mathbb{R}^{Kd}$  be any positive-semidefinite matrix. Then,

$$\lambda_{max}\left( ilde{oldsymbol{X}}_{ au}(oldsymbol{ heta})^{ op}oldsymbol{M} ilde{oldsymbol{X}}_{ au}(oldsymbol{ heta})
ight) \leq \sum_{i=1}^{K}\left|\left| ilde{oldsymbol{x}}_{ au}^{(i)}(oldsymbol{ heta})
ight|
ight|_{oldsymbol{M}}^{2}$$

*Proof.* The proof follows on the same lines as Lemma 7.12.

**Lemma 8.12.** Let  $\epsilon_1(t, \tau, x)$  be as defined in Lemma 8.7, and  $\tau_t$  be the last switching round before round t. Then, we have that

$$\sum_{t \in [T]} \epsilon_1(t, \tau_t, \boldsymbol{x}_t) \le 8RKd \log T \kappa^{1/2} e^{3S} \gamma(\delta)^2 + 4RKd^{1/2} (\log T)^{1/2} \gamma(\delta) \sqrt{T}$$

Proof.

$$\begin{split} \sum_{t\in[T]} \epsilon_1(t,\tau_t,\boldsymbol{x}_t) &= \sqrt{2}\gamma(\delta) \sum_{t\in[T]} \left| \left| \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} (\boldsymbol{I}\otimes\boldsymbol{x}_t) \boldsymbol{A}(\boldsymbol{x}_t,\hat{\boldsymbol{\theta}}_{\tau_t}) \boldsymbol{\rho} \right| \right|_2 \\ & \stackrel{(i)}{\leq} \sqrt{2}\gamma(\delta) \sum_{t\in[T]} \left| \left| \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} (\boldsymbol{I}\otimes\boldsymbol{x}_t) \boldsymbol{A}(\boldsymbol{x}_t,\hat{\boldsymbol{\theta}}_{\tau_t})^{1/2} \right| \right|_2 ||\boldsymbol{\rho}||_{\boldsymbol{A}(\boldsymbol{x}_t,\hat{\boldsymbol{\theta}}_{\tau_t})} \\ & \leq \sqrt{2}R\gamma(\delta) \sum_{t\in[T]} \left| \left| \boldsymbol{A}(\boldsymbol{x}_t,\hat{\boldsymbol{\theta}}_{\tau_t})^{1/2} (\boldsymbol{I}\otimes\boldsymbol{x}_t^{\top}) \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} \right| \right|_2 \\ & \stackrel{(ii)}{\leq} \sqrt{2}R\gamma(\delta) \sum_{t\in[T]} \left| \left| \sqrt{B_t(\boldsymbol{x}_t)} \tilde{\boldsymbol{X}}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{\top} \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} \right| \right|_2 \\ & \stackrel{(iii)}{\leq} \sqrt{2}R\gamma(\delta) \sum_{t\in[T]} \left| \left| \tilde{\boldsymbol{X}}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{\top} \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} \right| \right|_2 \left\{ e^{3S} \kappa^{1/2} \gamma(\delta) \|\boldsymbol{x}_t\|_{\boldsymbol{V}_t^{-1}} + 1 \right\} \end{split}$$

where (i) follows from  $||\mathbf{A}\mathbf{x}||_2 \leq ||\mathbf{A}||_2 ||\mathbf{x}||_2$ , (ii) follows from the definition of  $\tilde{\mathbf{X}}(\boldsymbol{\theta})$ , and (iii) follows from Lemma 8.9.

We now bound the term  $\sum_{t \in [T]} \left\| \tilde{X}_t(\hat{\theta}_{\tau_t})^\top H_t(\hat{\theta}_{\tau_t})^{-1/2} \right\|_2$ :

$$\begin{split} \sum_{t\in[T]} \left| \left| \tilde{\boldsymbol{X}}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}}) \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1/2} \right| \right|_{2} &= \sum_{t\in[T]} \sqrt{\lambda_{max} \left( \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1/2} \tilde{\boldsymbol{X}}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{\top} \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1/2} \right)} \\ &= \sum_{t\in[T]} \sqrt{\lambda_{max} \left( \tilde{\boldsymbol{X}}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{\top} \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1} \tilde{\boldsymbol{X}}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}}) \right)} \\ & \stackrel{(i)}{=} \sum_{t\in[T]} \sqrt{\sum_{i=1}^{K} \left| \left| \tilde{\boldsymbol{x}}_{t}^{(i)}(\hat{\boldsymbol{\theta}}_{\tau_{t}}) \right| \right|_{H_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1}}^{2}} \\ & \stackrel{(ii)}{\leq} \sum_{t\in[T]} \sqrt{\sum_{i=1}^{K} \left| \left| \tilde{\boldsymbol{x}}_{t}^{(i)}(\hat{\boldsymbol{\theta}}_{\tau_{t}}) \right| \right|_{H_{t}^{i}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1}}^{2}} \\ & \stackrel{(iii)}{\leq} \sqrt{T} \sqrt{\sum_{t\in[T]} \sum_{i=1}^{K} \left| \left| \tilde{\boldsymbol{x}}_{t}^{(i)}(\hat{\boldsymbol{\theta}}_{\tau_{t}}) \right| \right|_{H_{t}^{i}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1}}^{2}} \\ & \stackrel{(iv)}{\leq} 2K \sqrt{dT \log T} \end{split}$$

where (i) follows from Lemma 8.11, (ii) follows from Lemma 8.6, (iii) follows from Cauchy-Schwarz, and (iv) follows from Lemma 9.12 and the fact that  $||\tilde{\boldsymbol{x}}^{(i)}(\boldsymbol{\theta})||_2 \leq ||\boldsymbol{A}(\boldsymbol{x},\boldsymbol{\theta})||_2 ||\boldsymbol{x}||_2 \leq 1$ .

We also bound the term 
$$\sum_{t \in [T]} \left| \left| \tilde{X}_t(\hat{\theta}_{\tau_t})^\top H_t(\hat{\theta}_{\tau_t})^{-1/2} \right| \right|_2 \|x_t\|_{V_t^{-1}}$$
 as follows:

$$\sum_{t \in [T]} \left\| \tilde{\boldsymbol{X}}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{\top} \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1/2} \right\|_{2} \|\boldsymbol{x}_{t}\|_{\boldsymbol{V}_{t}^{-1}} \stackrel{(i)}{\leq} \sqrt{\sum_{t \in [T]} \left\| \tilde{\boldsymbol{X}}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{\top} \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{t}})^{-1/2} \right\|_{2}^{2}} \sqrt{\sum_{t \in [T]} \|\boldsymbol{x}_{t}\|_{\boldsymbol{V}_{t}^{-1}}^{2}} \stackrel{(ii)}{\leq} 2K\sqrt{d\log T} \sqrt{\sum_{t \in [T]} \|\boldsymbol{x}_{t}\|_{\boldsymbol{V}_{t}^{-1}}^{2}} \stackrel{(ii)}{\leq} 4Kd\log T$$

where (i) follows from Cauchy-Schwarz, (ii) follows from the same steps used to bound  $\sum_{t \in [T]} \left\| \tilde{\boldsymbol{X}}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^\top \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} \right\|_2 \text{ above, and } (iii) \text{ follows from Lemma 9.12.}$ 

Substituting back, we get

$$\sum_{t \in [T]} \epsilon_1(t, \tau_t, \boldsymbol{x}_t) \le 4\sqrt{2}RKd\log T\kappa^{1/2}e^{3S}\gamma(\delta)^2 + 2\sqrt{2}RKd^{1/2}(\log T)^{1/2}\gamma(\delta)\sqrt{T}$$
$$\le 8RKd\log T\kappa^{1/2}e^{3S}\gamma(\delta)^2 + 4RKd^{1/2}(\log T)^{1/2}\gamma(\delta)\sqrt{T}$$

**Lemma 8.13.** Let  $\epsilon_2(t, \tau, x)$  be as defined in Lemma 8.7, and  $\tau_t$  be the last switching round before round t. Then, we have that

$$\sum_{t \in [T]} \epsilon_2(t, \tau_t, \boldsymbol{x}_t) \le 24 dR K^2 e^{2S} \kappa \gamma(\delta)^2 \log T$$

Proof.

$$\begin{split} \sum_{t\in[T]} \epsilon_2(t,\tau,\boldsymbol{x}_t) &= 6R\gamma(\delta)^2 \sum_{t\in[T]} \left| \left| (\boldsymbol{I} \otimes \boldsymbol{x}^\top) \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} \right| \right|_2^2 \\ &\stackrel{(i)}{=} 6R\gamma(\delta)^2 \sum_{t\in[T]} \left| \left| \boldsymbol{A}(\boldsymbol{x}_t, \hat{\boldsymbol{\theta}}_{\tau})^{-1/2} \right| \right|_2 \left| \left| \tilde{\boldsymbol{X}}_t(\hat{\boldsymbol{\theta}}_{\tau_t}) \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} \right| \right|_2^2 B_t(\boldsymbol{x}_t) \\ &\stackrel{(ii)}{\leq} 6R\gamma(\delta)^2 e^{2S} \sum_{t\in[T]} \left| \left| \boldsymbol{A}(\boldsymbol{x}_t, \hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} \right| \right|_2^2 \left| \left| \tilde{\boldsymbol{X}}_t(\hat{\boldsymbol{\theta}}_{\tau_t}) \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} \right| \right|_2^2 \\ &\stackrel{(iii)}{\leq} 6R\gamma(\delta)^2 e^{2S} \kappa \sum_{t\in[T]} \left| \left| \tilde{\boldsymbol{X}}_t(\hat{\boldsymbol{\theta}}_{\tau_t}) \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_t})^{-1/2} \right| \right|_2^2 \\ &\stackrel{(iv)}{\leq} 24dRK^2 e^{2S} \kappa \gamma(\delta)^2 \log T \end{split}$$

where (i) follows from the definition of  $\tilde{X}$  and Lemma 8.6, (ii) follows from the definition of  $B_t(x)$ , (iii) follows from the fact that  $A(x, \theta) \leq \frac{1}{\kappa}I$ , and (iv) follows from the methods used in Lemma 8.12.

**Lemma 8.14.** Let Algorithm 3 be run for t rounds. Then, the switching criterion is triggered a maximum of  $dK \log(1 + \frac{t}{d\lambda})$  times.

*Proof.* Let  $\tau_0, \tau_1, \ldots, \tau_m \in [1, t]$  be the time steps at which the switching criterion in Algorithm 3 is triggered, i.e., 2 det  $H_{\tau_i}(\hat{\theta}_{\tau_i}) \leq \det H_{\tau_{i+1}}(\hat{\theta}_{\tau_i})$  for  $i \in [m-1]$ , and  $\tau_m = t$ . Note that  $H_{\tau_0}(\theta) = \lambda I_{Kd \times Kd}$ .

$$\frac{\det \boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_{m-1}})}{\det \boldsymbol{H}_{\tau_0}(\boldsymbol{\theta})} = \frac{\det \boldsymbol{H}_{\tau_m}(\hat{\boldsymbol{\theta}}_{\tau_{m-1}})}{\det \boldsymbol{H}_{\tau_{m-1}}(\hat{\boldsymbol{\theta}}_{\tau_{m-1}})} \times \frac{\det \boldsymbol{H}_{\tau_{m-1}}(\hat{\boldsymbol{\theta}}_{\tau_{m-2}})}{\det \boldsymbol{H}_{\tau_{m-2}}(\hat{\boldsymbol{\theta}}_{\tau_{m-2}})} \times \ldots \times \frac{\det \boldsymbol{H}_{\tau_1}(\hat{\boldsymbol{\theta}}_{\tau_0})}{\det \boldsymbol{H}_{\tau_0}(\boldsymbol{\theta})} \\ \ge 2^m$$

and hence, det  $H_t(\hat{\theta}_{\tau_{m-1}}) \ge 2^m \lambda^{Kd}$  since det  $H_1 = \lambda^{Kd}$ . Also, we can say that:

$$\det \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{m-1}}) \stackrel{(i)}{\leq} \left(\frac{\operatorname{trace} \boldsymbol{H}_{t}(\hat{\boldsymbol{\theta}}_{\tau_{m-1}})}{Kd}\right)^{Kd}$$

$$\stackrel{(ii)}{\leq} \left(\frac{\sum_{i \in [K]} \operatorname{trace} \boldsymbol{H}_{t}^{i}(\hat{\boldsymbol{\theta}}_{\tau_{m-1}})}{Kd}\right)^{Kd}$$

$$\stackrel{(iii)}{\leq} \left(\frac{\lambda Kd + \sum_{i \in [K]} \sum_{s \in [t]} \|\tilde{\boldsymbol{x}}_{s}^{(i)}(\hat{\boldsymbol{\theta}}_{\tau_{m-1}})\|_{2}^{2}}{Kd}\right)^{Kd}$$

$$\stackrel{(iv)}{\leq} \left(\lambda + \frac{t}{d}\right)^{Kd}$$

Here (i) follows from Lemma 9.11, (ii) follows from Lemma 8.6 alongside the linearity of the trace operator, (iii) follows from the definition of  $H_t^i(\theta)$  and the fact that the only non-zero eigenvalue of  $xx^{\top}$  is  $||x||_2^2$ , and (iv) is due to the fact that  $||\tilde{x}_t^{(i)}(\theta)||_2 \leq ||A(x_t, \theta)|| \leq 1$ . Thus, we have

$$2^m \lambda^{Kd} \leq \det(\boldsymbol{H}_t(\hat{\boldsymbol{\theta}}_{\tau_{m-1}}) \leq \left(\lambda + \frac{t}{d}\right)^{Kd}$$

and hence,  $2^m \leq \left(1 + \frac{t}{\lambda d}\right)^{Kd}$ . This finishes the proof.

### 9 General Lemmas and Results

**Lemma 9.1.** (Self-Concordance) Let  $A(x, \theta) = \nabla z(x, \theta)$ . Then,  $A(x, \theta)$  is (M, v)-generalized self-concordant with v = 1 and  $M = \sqrt{6}$ . In other words, for any given  $x_1, x_2, \theta_1, \theta_2$ , denote  $A_1 = A(x_1, \theta_1)$  and  $A_2 = A(x_2, \theta_2)$ . Then, we have

$$\boldsymbol{A}_{2}\exp\left(-\sqrt{6}\left|\left|\left(\boldsymbol{I}\otimes\boldsymbol{x}_{1}^{\top}\right)\boldsymbol{\theta}_{1}-\left(\boldsymbol{I}\otimes\boldsymbol{x}_{2}^{\top}\right)\boldsymbol{\theta}_{2}\right|\right|_{2}\right)\preccurlyeq\boldsymbol{A}_{1}\preccurlyeq\boldsymbol{A}_{2}\exp\left(\sqrt{6}\left|\left|\left(\boldsymbol{I}\otimes\boldsymbol{x}_{1}^{\top}\right)\boldsymbol{\theta}_{1}-\left(\boldsymbol{I}\otimes\boldsymbol{x}_{2}^{\top}\right)\boldsymbol{\theta}_{2}\right|\right|_{2}\right)$$

**Lemma 9.2.** (Lemma 13, Amani & Thrampoulidis (2021)) Let  $\beta = \{t_1, \ldots, t_N\}$  be a set of time indices and define

$$oldsymbol{G}_eta(oldsymbol{ heta}_1,oldsymbol{ heta}_2) = \sum_{t\ineta}oldsymbol{M}(oldsymbol{x},oldsymbol{ heta}_1,oldsymbol{ heta}_2)\otimesoldsymbol{x}_toldsymbol{x}_t^ op+\lambdaoldsymbol{I}_{Kd imes Kd}$$

and

$$oldsymbol{H}^{\star}_{eta} = \sum_{t\ineta}oldsymbol{A}(oldsymbol{x}_t,oldsymbol{ heta}^{\star})\otimesoldsymbol{x}_toldsymbol{x}_t^{ op} + \lambdaoldsymbol{I}_{Kd imes Kd}$$

where

$$\boldsymbol{M}(\boldsymbol{x},\boldsymbol{\theta}_1,\boldsymbol{\theta}_2) = \int_0^1 \boldsymbol{A}(\boldsymbol{x},v\boldsymbol{\theta}_1 + (1-v)\boldsymbol{\theta}_2) \, \mathrm{d}v$$

Then,

$$\boldsymbol{G}_{\beta}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\star}) \succcurlyeq \frac{1}{1+2S} \boldsymbol{H}_{\beta}^{\star}$$

Lemma 9.3. Define the log-likelihood function as follows:

$$\mathcal{L}_t(\boldsymbol{\theta}) = \sum_{s=1}^{t-1} \sum_{i=1}^K \mathbb{1}\left\{y_s = i\right\} \log \frac{1}{\boldsymbol{z}_i(\boldsymbol{x}_s, \boldsymbol{\theta})} + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$$

Let  $\hat{\theta}$  be the MLE of  $\theta^*$ , i.e.,  $\hat{\theta} = \arg\min_{\theta} \mathcal{L}_t(\theta)$ , then

$$\sum_{s=1}^{t-1} oldsymbol{z}(oldsymbol{x}_s, \hat{oldsymbol{ heta}}) \otimes oldsymbol{x}_s + \lambda \hat{oldsymbol{ heta}} = \sum_{s=1}^{t-1} oldsymbol{m}_s \otimes oldsymbol{x}_s$$

where  $\boldsymbol{m}_s = (\mathbb{1}\{y_s = 1\}, \dots, \mathbb{1}\{y_s = K\})^\top$  is the user-response vector.

*Proof.* For the sake of convenience, define the loss incurred at round t (without the regularization term) as

$$l_t(\boldsymbol{\theta}) = \sum_{i=1}^{K} \mathbb{1}\left\{y_s = i\right\} \log \frac{1}{\boldsymbol{z}_i(\boldsymbol{x}_s, \boldsymbol{\theta})}$$

Then, it is easy to see that

$$\begin{aligned} \frac{\partial l_t(\boldsymbol{\theta})}{\partial \theta_m} &= -\sum_{i=1}^K \mathbb{1}\left\{y_s = i\right\} \frac{1}{\boldsymbol{z}_i(\boldsymbol{x}_s, \boldsymbol{\theta})} \frac{\partial \boldsymbol{z}_i(\boldsymbol{x}_s, \boldsymbol{\theta})}{\partial \theta_m} \\ &= -\sum_{i=1}^K \mathbb{1}\left\{y_s = i\right\} \frac{1}{\boldsymbol{z}_i(\boldsymbol{x}_s, \boldsymbol{\theta})} \left[\mathbb{1}\left\{i = m\right\} \boldsymbol{z}_i(\boldsymbol{x}_s, \boldsymbol{\theta}) - \boldsymbol{z}_i(\boldsymbol{x}_s, \boldsymbol{\theta}) \boldsymbol{z}_m(\boldsymbol{x}_s, \boldsymbol{\theta})\right] \otimes \boldsymbol{x}_s \\ &= \left[\mathbb{1}\left\{y_s = m\right\} - \boldsymbol{z}_m(\boldsymbol{x}_s, \boldsymbol{\theta})\right] \otimes \boldsymbol{x}_s \end{aligned}$$

and hence,

$$abla l_t(oldsymbol{ heta}) = [oldsymbol{m}_s - oldsymbol{z}(oldsymbol{x}_s,oldsymbol{ heta})] \otimes oldsymbol{x}_s$$

Since  $\hat{\theta} = \underset{\theta}{\operatorname{arg\,min}} \mathcal{L}_t(\theta)$ , we have that  $\nabla \mathcal{L}_t(\hat{\theta}) = \underset{\theta}{\operatorname{arg\,min}} \sum_{s=1}^{t-1} l_s(\hat{\theta}) + \lambda \hat{\theta} = 0$ , which results in the claim.

**Lemma 9.4.** (Bernstein's Inequality) Let  $X_1, \ldots, X_n$  be a sequence of independent random variables with  $|X_t - \mathbb{E}[X_t]| \leq b$ . Let  $S = \sum_{t=1}^n (X_t - \mathbb{E}[X_t])$  and  $v = \sum_{t=1}^n \mathbb{V}[X_t]$ . Then, for any  $\delta \in [0, 1]$ , we have  $\mathbb{P}\left\{S \geq \sqrt{2v \log \frac{1}{\delta}} + \frac{2b}{3} \log \frac{1}{\delta}\right\} \leq \delta$ 

**Lemma 9.5.** Let  $m_s = (\mathbb{1} \{y_s = 1\}, \dots, \mathbb{1} \{y_s = K\})$  be the user-response vector as defined in Section 7.1. Then,

$$\mathbb{E}[\boldsymbol{m}_s] = \boldsymbol{z}(\boldsymbol{x}_s, \boldsymbol{\theta}^{\star}) \text{ and } \mathbb{E}\left[\boldsymbol{m}_s \boldsymbol{m}_s^{\top}\right] = diag(\boldsymbol{z}(\boldsymbol{x}_s, \boldsymbol{\theta}^{\star}))$$

*Proof.* Since  $m_s = (1 \{ y_s = 1 \}, \dots, 1 \{ y_s = K \})$ , we have

$$\mathbb{E}[\boldsymbol{m}_s] = (\mathbb{E}[\mathbb{1}\{y_s = 1\}], \dots, \mathbb{E}[\mathbb{1}\{y_s = K\}]) = (z_1(\boldsymbol{x}_s, \boldsymbol{\theta}^{\star}), \dots, z_K(\boldsymbol{x}_s, \boldsymbol{\theta}^{\star})) = \boldsymbol{z}(\boldsymbol{x}_s, \boldsymbol{\theta}^{\star})$$

For the second part, note that

$$\begin{bmatrix} \boldsymbol{m}_s \boldsymbol{m}_s^\top \end{bmatrix}_{i,j} = \mathbb{1} \{ y_s = i \} \mathbb{1} \{ y_s = j \} = \begin{cases} \mathbb{1} \{ y_s = i \} & i = j \\ 0 & i \neq j \end{cases}$$

Thus, we have

$$\mathbb{E}\left[\boldsymbol{m}_{s}\boldsymbol{m}_{s}^{\top}\right] = \mathbb{E}\left[\operatorname{diag}\left(\mathbbm{1}\left\{y_{s}=1\right\},\ldots,\mathbbm{1}\left\{y_{s}=K\right\}\right)\right] = \operatorname{diag}\left(\mathbb{E}\left[\mathbbm{1}\left\{y_{s}=1\right\}\right],\ldots,\mathbb{E}\left[\mathbbm{1}\left\{y_{s}=K\right\}\right]\right) \\ = \operatorname{diag}\left(z_{1}(\boldsymbol{x}_{s},\boldsymbol{\theta}^{\star}),\ldots,z_{K}(\boldsymbol{x}_{s},\boldsymbol{\theta}^{\star})\right) = \operatorname{diag}(\boldsymbol{z}(\boldsymbol{x}_{s},\boldsymbol{\theta}^{\star}))$$

**Lemma 9.6.** (*Claim A.8, Sawarni et al.* (2024)) For any  $x \in [0, M]$ ,

$$e^x \le e^M\left(\frac{x}{M}\right) + 1$$

Lemma 9.7. (Theorem 5, Ruan et al. (2021)) Let  $\pi$  represent the G-Optimal Distributional Design learnt from  $\mathcal{X}_1 \dots \mathcal{X}_s \stackrel{i.i.d}{\sim} \mathcal{D}$  and let  $\mathbf{W}$  be the expected data matrix, i.e.  $\mathbf{W} = \lambda \mathbf{I} + \mathbb{E}_{\mathbf{X} \sim \mathcal{D}} \left[ \mathbb{E}_{\mathbf{x} \sim \pi(\mathcal{X})} \mathbf{x} \mathbf{x}^\top | \mathcal{X} \right]$ , then, we have  $P \left\{ \mathbb{E}_{\mathcal{X} \sim \mathcal{D}} \left[ \max_{\mathbf{x} \in \mathcal{X}} ||\mathbf{x}||_{\mathbf{W}^{-1}} \right] \leq O(\sqrt{d \log d}) \right\} \geq 1 - \exp\left(O(d^4 \log^2 d) - sd^{-12}2^{-16}\right)$ 

**Lemma 9.8.** (Lemma 4, Ruan et al. (2021)) Let  $\pi_G$  represent the G-Optimal design and define the design matrix  $\mathbf{W}_G = \lambda \mathbf{I} + \mathop{\mathbb{E}}_{\mathcal{X} \sim \mathcal{D}} \left[ \mathop{\mathbb{E}}_{\boldsymbol{x} \in \pi_G(\mathcal{X})} \boldsymbol{x} \boldsymbol{x}^\top \mid \mathcal{X} \right]$ , then we have $\mathop{\mathbb{E}}_{\mathcal{X} \sim \mathcal{D}} \left[ \max_{\boldsymbol{x} \in \mathcal{X}} ||\boldsymbol{x}||_{\boldsymbol{W}_G^{-1}}^2 \right] \leq O(d^2)$ 

**Lemma 9.9.** (Lemma A.15, Sawarni et al. (2024), Ruan et al. (2021)) Let  $x_1 \dots x_n \sim \mathcal{D}$  be vectors with  $||x||_2 \leq 1$ , then

$$\mathbb{P}\left\{3\epsilon N\boldsymbol{I} + \sum_{i=1}^{n} \boldsymbol{x}_{i}\boldsymbol{x}_{i}^{\top} \succeq \frac{n}{8} \mathop{\mathbb{E}}_{\boldsymbol{x} \sim \mathcal{D}} \left[\boldsymbol{x}\boldsymbol{x}^{\top}\right]\right\} \geq 1 - 2d\exp\left(-\frac{\epsilon n}{8}\right)$$

**Lemma 9.10.** (Lemma 6, Zhang & Sugiyama (2023)) Let  $\{\mathcal{F}_t\}_{t=1}^{\infty}$  be a filteration and  $\{x_t\}_{t=1}^{\infty}$  be a stochastic process in  $\mathcal{B}_2(d) = \{x \in \mathbb{R}^d \mid ||x_t||_2 \leq 1\}$  such that  $x_t$  is  $\mathcal{F}_t$ -measurable. Let  $\{\epsilon_t\}_{t=1}^{\infty}$  be a martingale difference sequence such that  $\epsilon_t$  is  $\mathcal{F}_{t+1}$ -measurable. Assume that conditioned on  $\mathcal{F}_t$ , we have  $||\epsilon_t||_1 \leq 2$  almost surely, and is denoted by  $\eta_t = \mathbb{E}\left[\epsilon_t \epsilon_t^\top \mid \mathcal{F}_t\right]$ . Let  $\lambda > 0$  and for any  $t \geq 1$ , define

$$oldsymbol{S}_t = \sum_{s=1}^{t-1} oldsymbol{\epsilon}_s \otimes oldsymbol{x}_s$$
 and  $oldsymbol{H}_t = \lambda oldsymbol{I}_{dK imes dK} + \sum_{s=1}^{t-1} oldsymbol{\eta}_s \otimes oldsymbol{x}_s oldsymbol{x}_s^ op$ 

Then, for any  $\delta \in (0, 1)$ , we have

$$\mathbb{P}\left\{\exists t > 1, ||\boldsymbol{S}_t||_{\boldsymbol{H}_t^{-1}} \geq \frac{\sqrt{\lambda}}{4} + \frac{4}{\sqrt{\lambda}}\log\left(\frac{\det \boldsymbol{H}_t^{1/2}}{\delta\lambda^{\frac{dK}{2}}}\right) + \frac{4}{\sqrt{\lambda}}Kd\log 2\right\} \leq \delta$$

**Lemma 9.11.** (Determinant-Trace Inequality) Let the determinant and trace of a p.s.d matrix  $A \in \mathbb{R}^{d \times d}$  be denoted by det A and trace A. Then, we have

$$det \mathbf{A} \le \left(\frac{trace \mathbf{A}}{d}\right)^d$$

*Proof.* Let the eigenvalues of A be denoted by  $\lambda(A) \ge 0$  since  $A \ge 0$ . Then, we know, det  $A = \prod \lambda(A)$  and trace  $A = \sum \lambda(A)$ . Thus, applying the inequality for arithmetic means and geometric means, we get that

$$\left(\prod \lambda(\boldsymbol{A})\right)^{1/d} \leq \frac{\sum \lambda(\boldsymbol{A})}{d} \implies \det \boldsymbol{A} \leq \left(\frac{\operatorname{trace} \boldsymbol{A}}{d}\right)^{d}$$

**Lemma 9.12.** (Elliptical Potential Lemma, Lemma 11, Abbasi-Yadkori et al. (2011)) Let  $\{\boldsymbol{x}_s\}_{s=1}^t$ represent a set of vectors in  $\mathbb{R}^d$  and let  $||\boldsymbol{x}_s||_2 \leq L$ . Let  $\boldsymbol{V}_s = \lambda \boldsymbol{I}_{d \times d} + \sum_{m=1}^{s-1} x_m x_m^{\top}$ . Then, for  $\lambda \geq 1$ 

$$\sum_{s=1}^{t} ||\boldsymbol{x}_{s}||_{\boldsymbol{V}_{s}^{-1}}^{2} \leq 2d \log \left(1 + \frac{tL^{2}}{\lambda d}\right) \leq 4d \log(tL^{2})$$

**Lemma 9.13.** (*Lemma 12, Abbasi-Yadkori et al. (2011)*) If  $A \succcurlyeq B \succcurlyeq 0$ , then

$$\sup_{\boldsymbol{x}\neq\boldsymbol{0}}\frac{\boldsymbol{x}^{\top}\boldsymbol{A}\boldsymbol{x}}{\boldsymbol{x}^{\top}\boldsymbol{B}\boldsymbol{x}}\leq\frac{\det\left(\boldsymbol{A}\right)}{\det\left(\boldsymbol{B}\right)}$$

#### **10** Additional Experiments

In this section, we supplement the experiments from Section 5 (in particular, **Experiment 1** and **Experiment 2**).



**Experiment 1** (R(T) vs. T for the Logistic (K = 1) Setting): In this experiment, we use the same instance as in Experiment 1 (Section 5) and average the regret over 10 different seeds for sampling rewards. The averaged results with two standard deviations can be found in Figure 2a.

**Experiment 2** (R(T) vs. T for K = 3): In this experiment, we use the same instance as in **Experiment 2** (Section 5) and average the regret over 10 different seeds for sampling rewards. The averaged results with two standard deviations are reported in Figure 2b.