

---

# Preference-based Reinforcement Learning beyond Pairwise Comparisons: Benefits of Multiple Options

---

**Joongkyu Lee**  
Seoul National University  
Seoul, South Korea  
jklee0717@snu.ac.kr

**Seouh-won Yi**  
Seoul National University  
Seoul, South Korea  
uniqueseouh@snu.ac.kr

**Min-hwan Oh**  
Seoul National University  
Seoul, South Korea  
minoh@snu.ac.kr

## Abstract

We study online preference-based reinforcement learning (PbRL) with the goal of improving sample efficiency. While a growing body of theoretical work has emerged—motivated by PbRL’s recent empirical success, particularly in aligning large language models (LLMs)—most existing studies focus only on pairwise comparisons. A few recent works [94, 49, 77] have explored using multiple comparisons and ranking feedback, but their performance guarantees fail to improve—and can even deteriorate—as the feedback length increases, despite the richer information available. To address this gap, we adopt the Plackett–Luce (PL) model for ranking feedback over action subsets and propose M-AUP0, an algorithm that selects multiple actions by maximizing the average uncertainty within the offered subset.

We prove that M-AUP0 achieves a suboptimality gap of  $\tilde{O}\left(\frac{d}{T}\sqrt{\sum_{t=1}^T \frac{1}{|S_t|}}\right)$ , where  $T$  is the total number of rounds,  $d$  is the feature dimension, and  $|S_t|$  is the size of the subset at round  $t$ . This result shows that larger subsets directly lead to improved performance and, notably, the bound avoids the exponential dependence on the unknown parameter’s norm, which was a fundamental limitation in most previous works. Moreover, we establish a near-matching lower bound of  $\Omega\left(\frac{d}{K\sqrt{T}}\right)$ , where  $K$  is the maximum subset size. To the best of our knowledge, this is the first theoretical result in PbRL with ranking feedback that explicitly shows improved sample efficiency as a function of the subset size.

## 1 Introduction

The framework of *Preference-based Reinforcement Learning* (PbRL) [12, 81, 82, 70] was introduced to address the difficulty of designing effective reward functions, which often demands substantial and complex engineering effort [80, 82]. PbRL has been successfully applied in diverse domains, including robot training, stock prediction, recommender systems, and clinical trials [28, 64, 17, 36, 52]. Notably, PbRL also serves as a foundational framework for Reinforcement Learning from Human Feedback (RLHF) when feedback is provided in the form of preferences rather than explicit scalar rewards. This preference-based approach has proven highly effective in aligning Large Language Models (LLMs) with human values and preferences [17, 57, 61].

Given its practical success, the field has also seen significant theoretical advances [16, 47, 70, 94, 87, 92, 91, 84, 72, 51, 13, 63, 21, 18, 49, 75, 71, 77, 86, 14, 37]. However, despite this progress, most existing models remain limited to handling only *pairwise* comparison feedback. A few works [94, 49, 77] explore the more general setting of *multiple* comparisons, offering a strict extension beyond the pairwise case. Zhu et al. [94] study the offline setting, where a dataset of questions (or contexts) along with corresponding ranking feedback over  $K$  answers (or actions), labeled by human annotators,

Table 1:  $T$  denotes the number of rounds (or the number of data points in the offline setting),  $K$  is the (maximum) size of the offered action set (i.e., *assortment*), and  $d$  is the feature dimension, and  $1/\kappa = \mathcal{O}(e^B)$ .  $\rho$  represents the unknown context distribution. Here,  $\tilde{\mathcal{O}}$  hides logarithmic factors and polynomial dependencies on  $B$ . “Sq. Pred. Error” refers to the squared prediction error.

	Setting	Context	Assortment	Measure	Result
Zhu et al. [94]	Offline	Accessible $\mathcal{X}$	Given	Suboptimality	$\tilde{\mathcal{O}}\left(\frac{K^2}{\kappa}\sqrt{\frac{d}{T}}\right)$
Mukherjee et al. [49]	Online	Accessible $\mathcal{X}$	Given	Pred. Error	$\tilde{\mathcal{O}}\left(\frac{K^3 d}{\kappa\sqrt{T}}\right)$
Thekumparampil et al. [77]	Online	No context	Select $K$	Pred. Error	$\tilde{\mathcal{O}}\left(\frac{K^3 d}{\kappa\sqrt{T}}\right)$
<b>This work</b> (Theorem 1, E.1)	Online	Sampled $x \sim \rho$	Select $\leq K$	Suboptimality	$\tilde{\mathcal{O}}\left(\frac{d}{T}\sqrt{\sum_{t=1}^T \frac{1}{ S_t }}\right)$
<b>This work</b> (Theorem 2)	Lower Bound	Sampled $x \sim \rho$	Select $\leq K$	Suboptimality	$\Omega\left(\frac{d}{K\sqrt{T}}\right)$

is available. Mukherjee et al. [49] investigate the online learning-to-rank problem [60], where a dataset of questions with  $K$  candidate answers is provided, but no feedback is initially available. Thekumparampil et al. [77] consider a context-free setting (i.e., a singleton context), and the goal is to learn the ranking of  $N \geq K$  answers based on ranking feedback obtained from subsets of size  $K$ . However, all of their theoretical performance guarantees fail to show that using multiple comparisons provides any advantage over the pairwise setting (see Table 1). This is counterintuitive, as ranking feedback is inherently more informative than pairwise feedback. Specifically, since a ranking over  $K$  actions provides  $\binom{K}{2}$  pairwise comparisons, it should, in principle, enable faster learning and lead to stronger performance guarantees. Thus, the following fundamental question remains open:

*Can we design an algorithm that achieves a strictly better theoretical guarantee under multiple-option feedback compared to the pairwise comparisons in the online PbRL setting?*

In this paper, we assume that the ranking feedback follows the Plackett-Luce (PL) model [59, 45], where, in each round, the learner receives ranking feedback over a subset of up to  $K$  actions.

## 1.1 Main Contributions.

**Improved sample efficiency via larger subsets:** We propose M-AUPD, a novel algorithm for online PbRL with PL ranking feedback, which selects action subsets by maximizing *average uncertainty*, thereby explicitly exploiting the richer information available from ranking feedback. We establish a suboptimality gap of  $\tilde{\mathcal{O}}\left(\frac{d}{T}\sqrt{\sum_{t=1}^T \frac{1}{|S_t|}}\right)$ , where  $|S_t|$  is the size of the action subset offered at round  $t$ . This result provides the first rigorous theoretical guarantee that larger subsets directly improve sample efficiency. To the best of our knowledge, this is the first theoretical work in PbRL that explicitly demonstrates performance improvements as a function of the subset size  $|S_t|$ .

**Free of  $\mathcal{O}(e^B)$  dependency:** Our result eliminates the exponential dependence on the parameter norm bound,  $\mathcal{O}(e^B)$ , by employing novel matrix concentration inequalities for the Hessian matrix  $H_t$  (Lemma D.2 and E.1). This represents a significant improvement over most prior works, where performance guarantees typically depend on  $\mathcal{O}(e^B)$  [65, 70, 94, 87, 92, 18, 86, 77, 37]. Very recently, a few works [19, 14] have successfully avoided the  $\mathcal{O}(e^B)$  dependency. However, their methods are limited to pairwise comparisons. To the best of our knowledge, this is the first work in PbRL with ranking feedback involving more than two options that avoids  $\mathcal{O}(e^B)$  dependence.

**Lower bound:** We establish a near-matching lower bound of  $\Omega\left(\frac{d}{K\sqrt{T}}\right)$  under PL model.

## 2 Problem Setting and Preliminaries

We have a set of contexts (or prompts), denoted by  $\mathcal{X}$ , and a set of possible actions (or answers), denoted by  $\mathcal{A} := \{a_1, \dots, a_N\}$ .<sup>1</sup> We consider preference feedback in the form of partial rankings over subsets of  $\mathcal{A}$ , and model this feedback using the Plackett-Luce (PL) distribution:

<sup>1</sup>For simplicity, while we consider the action space  $\mathcal{A}$  to be stationary in this paper, it is important to note that  $\mathcal{A}$  can vary depending on the context  $x \in \mathcal{X}$  or even across rounds  $t \in [T]$ .

---

**Algorithm 1** M-AUP0: Maximizing Average Uncertainty for Preference Optimization
 

---

1: **Inputs:** maximum assortment size  $K$ , regularization parameter  $\lambda$ , step size  $\eta$   
 2: **for** round  $t = 1$  to  $T$  **do**  
 3:   Observe  $x_t$  and select  $(\bar{a}_t, S_t)$  via (4)  
 4:   Receive ranking feedback  $\sigma_t$  for  $S_t$   
 5:   Compute  $\hat{\theta}_{t+1}$  by OMD and update  $H_{t+1} \leftarrow H_t + \sum_{j=1}^{|S_t|} \nabla^2 \ell_t^{(j)}(\hat{\theta}_t^{(j+1)})$   
 6: **end for**  
 7: **Return:**  $\hat{\pi}_T(x) \leftarrow \operatorname{argmax}_{a \in \mathcal{A}} \phi(x, a)^\top \hat{\theta}_{T+1}$

---

**Definition 1** (PL model). Let  $\mathcal{S} := \{S \subseteq \mathcal{A} \mid 2 \leq |S| \leq K\}$  be the collection of all action subsets whose sizes range from 2 to  $K$ . For any  $S \in \mathcal{S}$ , let  $\sigma$  denote the labeler’s ranking feedback—that is, a permutation of the elements in  $S$ . We write  $\sigma_j$  for the  $j$ -th most preferred action under  $\sigma$ . We model the distribution of such rankings using the Plackett-Luce (PL) model [59, 45], defined as:

$$\mathbb{P}(\sigma|x, S; \theta^*) = \prod_{j=1}^{|S|} \frac{\exp(r_{\theta^*}(x, \sigma_j))}{\sum_{k=j}^{|S|} \exp(r_{\theta^*}(x, \sigma_k))}, \quad \text{where } (x, S) \in \mathcal{X} \times \mathcal{S}. \quad (1)$$

Here,  $r_{\theta^*}$  represents a reward model parameterized by the unknown parameter  $\theta^*$ .

**Assumption 1** (Linear reward). Let  $\phi : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^d$  be a known feature map satisfying  $\max_{x,a} \|\phi(x, a)\|_2 \leq 1$ , and let  $\theta^* \in \mathbb{R}^d$  denote the true but unknown parameter. The reward is assumed to follow a linear structure given by  $r_{\theta^*}(x, a) = \phi(x, a)^\top \theta^*$ . To ensure identifiability of  $\theta^*$ , we assume that  $\theta^* \in \Theta$ , where  $\Theta := \{\theta \in \mathbb{R}^d \mid \|\theta\|_2 \leq B\}$ .

At each round  $t \in [T]$ , a context  $x_t \in \mathcal{X}$  is drawn from a fixed but unknown distribution  $\rho$ . Given the context  $x_t$ , the learning agent selects a subset of actions  $S_t \in \mathcal{S}$ —referred to as an *assortment* throughout the paper—and receives a ranking over  $S_t$  as feedback, generated according to the PL model. After  $T$  rounds of interaction with the labeler, the goal is to output a policy  $\hat{\pi}_T : \mathcal{X} \rightarrow \mathcal{A}$  that minimizes the *suboptimality gap*, defined as:

$$\text{SubOpt}(T) := \mathbb{E}_{x \sim \rho} [r_{\theta^*}(x, \pi^*(x)) - r_{\theta^*}(x, \hat{\pi}_T(x))], \quad \text{where } \pi^*(x) = \operatorname{argmax}_{a \in \mathcal{A}} r_{\theta^*}(x, a).$$

## 2.1 Plackett-Luce (PL) Loss And Online Parameter Estimation

The PL loss function for round  $t$  is defined as follows:

$$\ell_t(\theta) := \sum_{j=1}^{|S_t|} \ell_t^{(j)}(\theta), \quad \text{where } \ell_t^{(j)}(\theta) := -\log \left( \frac{\exp(\phi(x_t, \sigma_{tj})^\top \theta)}{\sum_{k=j}^{|S_t|} \exp(\phi(x_t, \sigma_{tk})^\top \theta)} \right). \quad (2)$$

Here,  $\ell_t^{(j)}(\theta)$  denotes the negative log-likelihood loss under the Multinomial Logit (MNL) model [46], conditioned on the assortment being the remaining actions in  $S_t$  after removing the previously selected actions  $\sigma_{t1}, \dots, \sigma_{t(j-1)}$ —that is, over the set  $S_t \setminus \{\sigma_{t1}, \dots, \sigma_{t(j-1)}\}$ .

Motivated by recent advances in Multinomial Logit (MNL) bandits [93, 39, 41], we adopt an online mirror descent (OMD) algorithm to estimate the underlying parameter  $\theta^*$ :

$$\hat{\theta}_t^{(j+1)} = \operatorname{argmin}_{\theta \in \Theta} \langle \nabla \ell_t^{(j)}(\hat{\theta}_t^{(j)}), \theta \rangle + \frac{1}{2\eta} \|\theta - \hat{\theta}_t^{(j)}\|_{\tilde{H}_t^{(j)}}^2, \quad j = 1, \dots, |S_t|, \quad (3)$$

where we write  $\hat{\theta}_t^{(|S_t|+1)} = \hat{\theta}_{t+1}^{(1)}$ , and  $\eta$  is the step-size parameter to be specified later. The matrix  $\tilde{H}_t^{(j)}$  is given by  $\tilde{H}_t^{(j)} := H_t + \eta \sum_{j'=1}^j \nabla^2 \ell_t^{(j')}(\hat{\theta}_t^{(j')})$ , where  $H_t := \sum_{s=1}^{t-1} \sum_{j=1}^{|S_s|} \nabla^2 \ell_s^{(j)}(\hat{\theta}_s^{(j+1)}) + \lambda \mathbf{I}_d$ ,  $\lambda > 0$ . The optimization problem (3) can be solved using a single projected gradient step [55], which enjoys a computational cost of only  $\mathcal{O}(Kd^3)$ —independent of  $t$  [48], unlike MLE—and requires only  $\mathcal{O}(d^2)$  storage, thanks to the incremental updates of  $\tilde{H}_t^{(j)}$  and  $H_t$ .

## 3 M-AUP0: Maximizing Average Uncertainty

In this section, we propose a new algorithm, M-AUP0, designed to select an assortment that maximizes *average uncertainty* of  $S_t$ , thereby leveraging the potential benefits of a large  $K$ . At each round  $t$ , a

context  $x_t$  is drawn from a fixed but unknown distribution  $\rho$ . The algorithm then selects a *reference action-assortment pair*  $(\bar{a}_t, S_t)$  by maximizing the average feature uncertainty—measured in the  $H_t^{-1}$ -norm—relative to a candidate reference action  $\bar{a}$  (Line 3):

$$(\bar{a}_t, S_t) = \operatorname{argmax}_{\bar{a} \in \mathcal{A}} \operatorname{argmax}_{\substack{S \subseteq \mathcal{S} \\ \bar{a} \in S}} \frac{1}{|S|} \sum_{a \in S} \|\phi(x_t, a) - \phi(x_t, \bar{a})\|_{H_t^{-1}}. \quad (4)$$

By construction, the reference action  $\bar{a}_t$  is always included in the selected assortment  $S_t$ . This selection strategy plays a key role in our algorithm, as it promotes rapid reduction in reward uncertainty—particularly when the assortment size  $|S_t|$  is large—by encouraging informative comparisons centered around the reference action. Importantly, the assortment selection rule in Equation (4) can be computed efficiently, without enumerating all  $\binom{N}{K}$  possible subsets. Then, we observe the ranking feedback  $\sigma_t$  from a labeler and update the parameter by OMD. After  $T$  rounds, the algorithm returns the final policy  $\hat{\pi}_T$ , which selects actions by maximizing the estimated reward under the final parameter estimate, i.e.,  $\hat{\pi}_T(x) := \operatorname{argmax}_a \phi(x, a)^\top \hat{\theta}_{T+1}$  (Line 7).

## 4 Main Results

**Theorem 1.** *Let  $\delta \in (0, 1]$ . We set  $\lambda = \Omega(d \log(KT/\delta) + \eta(B + d))$  and  $\eta = \frac{1}{2}(1 + 3\sqrt{2}B)$ . Define  $\kappa := e^{-4B}$ . If Assumption 1 holds, then, with probability at least  $1 - \delta$ , M-AUPD (Algorithm 1) achieves the following suboptimality gap:*

$$\text{SubOpt}(T) = \tilde{\mathcal{O}} \left( \frac{d}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} + \frac{d^2 K^2}{\kappa T} \right).$$

**Discussion of Theorem 1.** The proof is provided in Appendix D. For sufficiently large  $T$ , the second (non-leading) term becomes negligible, and Theorem 1 shows that the suboptimality gap of M-AUPD decreases as the assortment size  $|S_t|$  increases. This establishes a strict advantage of receiving ranking feedback over larger assortments. Moreover, our result does not involve any  $\mathcal{O}(e^B)$  dependency in the leading term, a harmful dependency that commonly appears in prior works [65, 70, 94, 87, 92, 18, 77, 37]. Although very recent works [19, 14] achieve  $\mathcal{O}(e^B)$ -free performance, they are limited to pairwise preference feedback. Extending their methods to accommodate richer ranking feedback is non-trivial. To the best of our knowledge, this is the first theoretical work to explicitly demonstrate both the performance improvements enabled by ranking feedback over larger assortments and the elimination of the  $\mathcal{O}(e^B)$  dependency in PbRL framework when handling multiple options—specifically, more than two.

**Theorem 2** (Lower bound). *Suppose  $T \geq d^2/(8K^2)$ . Define the feature space as  $\Phi := \mathcal{S}^{d-1}$ , the unit sphere in  $\mathbb{R}^d$ , and let the parameter space be  $\Theta = \{-\mu, \mu\}^d$ , where  $\mu = \sqrt{d/(8K^2T)}$ . Then, for any policy  $\hat{\pi}_T \in \Delta_\Phi$  returned after collecting  $T$  samples (using any sampling policy), the expected suboptimality gap is lower bounded as:*

$$\text{SubOpt}(T) = \Omega \left( \frac{d}{K\sqrt{T}} \right).$$

**Discussion of Theorem 2.** The proof is deferred to Appendix G. Theorem 2 provides theoretical support for our upper bounds, particularly with respect to the dependency on  $K$ . Compared to the upper bound in Theorems 1, the remaining gap is only a factor of  $\frac{1}{\sqrt{K}}$ . To the best of our knowledge, this is the first lower bound on the suboptimality gap that incorporates PL ranking feedback in PbRL and formally shows that the suboptimality gap can diminish as  $K$  grows, highlighting the advantage of utilizing ranking feedback over simple pairwise comparisons.

## 5 Experiments and Conclusion

Due to space constraints, we present our experimental results in Appendix I. In this work, to the best of our knowledge, we present the first theoretical result in online PbRL showing that the suboptimality gap decreases as more options are revealed to the labeler for ranking feedback. This finding highlights the power of richer feedback in accelerating learning and opens new avenues for the development of more data-efficient AI systems guided by human preferences.

## References

- [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24:2312–2320, 2011.
- [2] Marc Abeille, Louis Faury, and Clément Calauzènes. Instance-wise minimax-optimal algorithms for logistic bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 3691–3699. PMLR, 2021.
- [3] Priyank Agrawal, Theja Tulabandhula, and Vashist Avadhanula. A tractable online learning algorithm for the multinomial logit contextual bandit. *European Journal of Operational Research*, 310(2):737–750, 2023.
- [4] Shipra Agrawal, Vashist Avadhanula, Vineet Goyal, and Assaf Zeevi. Thompson sampling for the mnl-bandit. In *Conference on learning theory*, pages 76–78. PMLR, 2017.
- [5] Shipra Agrawal, Vashist Avadhanula, Vineet Goyal, and Assaf Zeevi. Mnl-bandit: A dynamic learning approach to assortment selection. *Operations Research*, 67(5):1453–1485, 2019.
- [6] Hossein Azari Soufiani, William Chen, David C Parkes, and Lirong Xia. Generalized method-of-moments for rank aggregation. *Advances in Neural Information Processing Systems*, 26, 2013.
- [7] Akshay Balsubramani, Zohar Karnin, Robert E Schapire, and Masrour Zoghi. Instance-dependent regret bounds for dueling bandits. In *Conference on Learning Theory*, pages 336–360. PMLR, 2016.
- [8] Viktor Bengs, Róbert Busa-Fekete, Adil El Mesaoudi-Paul, and Eyke Hüllermeier. Preference-based online learning with dueling bandits: A survey. *Journal of Machine Learning Research*, 22(7):1–108, 2021.
- [9] Viktor Bengs, Aadirupa Saha, and Eyke Hüllermeier. Stochastic contextual dueling bandits under linear stochastic transitivity models. In *International Conference on Machine Learning*, pages 1764–1786. PMLR, 2022.
- [10] Alina Beygelzimer, John Langford, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandit algorithms with supervised learning guarantees. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 19–26. JMLR Workshop and Conference Proceedings, 2011.
- [11] Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- [12] Róbert Busa-Fekete, Balázs Szörényi, Paul Weng, Weiwei Cheng, and Eyke Hüllermeier. Preference-based reinforcement learning: evolutionary direct policy search using a preference-based racing algorithm. *Machine learning*, 97:327–351, 2014.
- [13] Shicong Cen, Jincheng Mei, Katayoon Goshvadi, Hanjun Dai, Tong Yang, Sherry Yang, Dale Schuurmans, Yuejie Chi, and Bo Dai. Value-incentivized preference optimization: A unified approach to online and offline RLHF. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=SQnitDuow6>.
- [14] Mingyu Chen, Yiding Chen, Wen Sun, and Xuezhou Zhang. Avoiding  $\exp(r_{\max})$  dependency in preference-based reinforcement learning. *arXiv preprint arXiv:2502.00666*, 2025.
- [15] Xi Chen, Yining Wang, and Yuan Zhou. Dynamic assortment optimization with changing contextual information. *The Journal of Machine Learning Research*, 21(1):8918–8961, 2020.
- [16] Xiaoyu Chen, Han Zhong, Zhuoran Yang, Zhaoran Wang, and Liwei Wang. Human-in-the-loop: Provably efficient preference-based reinforcement learning with general function approximation. In *International Conference on Machine Learning*, pages 3773–3793. PMLR, 2022.
- [17] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

- [18] Nirjhar Das, Souradip Chakraborty, Aldo Pacchiano, and Sayak Ray Chowdhury. Active preference optimization for sample efficient rlhf. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 96–112. Springer, 2025.
- [19] Qiwei Di, Jiafan He, and Quanquan Gu. Nearly optimal algorithms for contextual dueling bandits from adversarial feedback. In *Forty-second International Conference on Machine Learning*, 2024.
- [20] Qiwei Di, Tao Jin, Yue Wu, Heyang Zhao, Farzad Farnoud, and Quanquan Gu. Variance-aware regret bounds for stochastic contextual dueling bandits. In *The Twelfth International Conference on Learning Representations*, 2024.
- [21] Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf. *Transactions on Machine Learning Research*, 2024.
- [22] Shi Dong, Tengyu Ma, and Benjamin Van Roy. On the performance of thompson sampling on logistic bandits. In *Conference on Learning Theory*, pages 1158–1160. PMLR, 2019.
- [23] Miroslav Dudík, Katja Hofmann, Robert E Schapire, Aleksandrs Slivkins, and Masrour Zoghi. Contextual dueling bandits. In *Conference on Learning Theory*, pages 563–587. PMLR, 2015.
- [24] Moein Falahatgar, Yi Hao, Alon Orlitsky, Venkatadheeraj Pichapati, and Vaishakh Ravindrakumar. Maxing and ranking with few assumptions. *Advances in Neural Information Processing Systems*, 30, 2017.
- [25] Louis Faury, Marc Abeille, Clément Calauzènes, and Olivier Fercoq. Improved optimistic algorithms for logistic bandits. In *International Conference on Machine Learning*, pages 3052–3060. PMLR, 2020.
- [26] Louis Faury, Marc Abeille, Kwang-Sung Jun, and Clément Calauzènes. Jointly efficient and optimal algorithms for logistic bandits. In *International Conference on Artificial Intelligence and Statistics*, pages 546–580. PMLR, 2022.
- [27] Reinhard Heckel, Max Simchowitz, Kannan Ramchandran, and Martin Wainwright. Approximate ranking from pairwise comparisons. In *International Conference on Artificial Intelligence and Statistics*, pages 1057–1066. PMLR, 2018.
- [28] Ashesh Jain, Brian Wojcik, Thorsten Joachims, and Ashutosh Saxena. Learning trajectory preferences for manipulators via iterative improvement. *Advances in neural information processing systems*, 26, 2013.
- [29] Kevin Jamieson, Sumeet Katariya, Atul Deshpande, and Robert Nowak. Sparse dueling bandits. In *Artificial Intelligence and Statistics*, pages 416–424. PMLR, 2015.
- [30] Minje Jang, Sunghyun Kim, Changho Suh, and Sewoong Oh. Optimal sample complexity of m-wise data for top-k ranking. *Advances in Neural Information Processing Systems*, 30, 2017.
- [31] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL <https://arxiv.org/abs/2310.06825>.
- [32] Ashish Khetan and Sewoong Oh. Data-driven rank breaking for efficient rank aggregation. *Journal of Machine Learning Research*, 17(193):1–54, 2016.
- [33] Yeoneung Kim, Insoon Yang, and Kwang-Sung Jun. Improved regret analysis for variance-adaptive linear bandits and horizon-free linear mixture mdps. *Advances in Neural Information Processing Systems*, 35:1060–1072, 2022.
- [34] Junpei Komiyama, Junya Honda, Hisashi Kashima, and Hiroshi Nakagawa. Regret lower bound and optimal algorithm in dueling bandit problem. In *Conference on learning theory*, pages 1141–1154. PMLR, 2015.

- [35] Junpei Komiyama, Junya Honda, and Hiroshi Nakagawa. Copeland dueling bandit problem: Regret lower bound, optimal algorithm, and computationally efficient algorithm. In *International Conference on Machine Learning*, pages 1235–1244. PMLR, 2016.
- [36] Andras Kupcsik, David Hsu, and Wee Sun Lee. Learning dynamic robot-to-human object handover from human feedback. In *Robotics Research: Volume 1*, pages 161–176. Springer, 2017.
- [37] Branislav Kveton, Xintong Li, Julian McAuley, Ryan Rossi, Jingbo Shang, Junda Wu, and Tong Yu. Active learning for direct preference optimization. *arXiv preprint arXiv:2503.01076*, 2025.
- [38] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- [39] Joongkyu Lee and Min-hwan Oh. Nearly minimax optimal regret for multinomial logistic bandit. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [40] Joongkyu Lee and Min-hwan Oh. Combinatorial reinforcement learning with preference feedback. In *Forty-second International Conference on Machine Learning*, 2025.
- [41] Joongkyu Lee and Min-hwan Oh. Improved online confidence bounds for multinomial logistic bandits. In *Forty-second International Conference on Machine Learning*, 2025.
- [42] Junghyun Lee, Se-Young Yun, and Kwang-Sung Jun. Improved regret bounds of (multinomial) logistic bandits via regret-to-confidence-set conversion. In *International Conference on Artificial Intelligence and Statistics*, pages 4474–4482. PMLR, 2024.
- [43] Junghyun Lee, Se-Young Yun, and Kwang-Sung Jun. A unified confidence sequence for generalized linear models, with applications to bandits. *Advances in Neural Information Processing Systems*, 37:124640–124685, 2024.
- [44] Zhihan Liu, Miao Lu, Shenao Zhang, Boyi Liu, Hongyi Guo, Yingxiang Yang, Jose Blanchet, and Zhaoran Wang. Provably mitigating overoptimization in RLHF: Your SFT loss is implicitly an adversarial regularizer. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=2cQ3lPhke0>.
- [45] R Duncan Luce et al. *Individual choice behavior*, volume 4. Wiley New York, 1959.
- [46] Daniel McFadden. Modelling the choice of residential location. 1977.
- [47] Viraj Mehta, Vikramjeet Das, Ojash Neopane, Yijia Dai, Ilija Bogunovic, Jeff Schneider, and Willie Neiswanger. Sample efficient reinforcement learning from human feedback via active exploration. *arXiv preprint arXiv:2312.00267*, 2023.
- [48] Zakaria Mhammedi, Wouter M Koolen, and Tim Van Erven. Lipschitz adaptivity with multiple learning rates in online learning. In *Conference on Learning Theory*, pages 2490–2511. PMLR, 2019.
- [49] Subhojyoti Mukherjee, Anusha Lalitha, Kousha Kalantari, Aniket Anand Deshmukh, Ge Liu, Yifei Ma, and Branislav Kveton. Optimal design for human preference elicitation. *Advances in Neural Information Processing Systems*, 37:90132–90159, 2024.
- [50] William Muldrew, Peter Hayes, Mingtian Zhang, and David Barber. Active preference learning for large language models. In *Forty-first International Conference on Machine Learning*, 2024.
- [51] Rémi Munos, Michal Valko, Daniele Calandriello, Mohammad Gheshlaghi Azar, Mark Rowland, Zhaohan Daniel Guo, Yunhao Tang, Matthieu Geist, Thomas Mesnard, Côme Fiegl, et al. Nash learning from human feedback. In *Forty-first International Conference on Machine Learning*, 2023.
- [52] Ellen Novoseller, Yibing Wei, Yanan Sui, Yisong Yue, and Joel Burdick. Dueling posterior sampling for preference-based reinforcement learning. In *Conference on Uncertainty in Artificial Intelligence*, pages 1029–1038. PMLR, 2020.

- [53] Min-hwan Oh and Garud Iyengar. Thompson sampling for multinomial logit contextual bandits. *Advances in Neural Information Processing Systems*, 32, 2019.
- [54] Min-hwan Oh and Garud Iyengar. Multinomial logit contextual bandits: Provable optimality and practicality. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 9205–9213, 2021.
- [55] Francesco Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*, 2019.
- [56] Mingdong Ou, Nan Li, Shenghuo Zhu, and Rong Jin. Multinomial logit bandit with linear utility functions. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*, pages 2602–2608. International Joint Conferences on Artificial Intelligence Organization, 2018.
- [57] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- [58] Noemie Perivier and Vineet Goyal. Dynamic pricing and assortment under a contextual mnl demand. *Advances in Neural Information Processing Systems*, 35:3461–3474, 2022.
- [59] Robin L Plackett. The analysis of permutations. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 24(2):193–202, 1975.
- [60] Filip Radlinski, Robert Kleinberg, and Thorsten Joachims. Learning diverse rankings with multi-armed bandits. In *Proceedings of the 25th international conference on Machine learning*, pages 784–791, 2008.
- [61] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741, 2023.
- [62] Siddhartha Y Ramamohan, Arun Rajkumar, and Shivani Agarwal. Dueling bandits: Beyond condorcet winners to general tournament solutions. *Advances in Neural Information Processing Systems*, 29, 2016.
- [63] Sayak Ray Chowdhury, Anush Kini, and Nagarajan Natarajan. Provably robust DPO: Aligning language models with noisy feedback. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 42258–42274. PMLR, 21–27 Jul 2024.
- [64] Dorsa Sadigh, Anca Dragan, Shankar Sastry, and Sanjit Seshia. *Active preference-based learning of reward functions*. 2017.
- [65] Aadirupa Saha. Optimal algorithms for stochastic contextual preference bandits. *Advances in Neural Information Processing Systems*, 34:30050–30062, 2021.
- [66] Aadirupa Saha and Pierre Gaillard. Finally rank-breaking conquers mnl bandits: Optimal and efficient algorithms for mnl assortment. In *The Thirteenth International Conference on Learning Representations*, 2025.
- [67] Aadirupa Saha and Aditya Gopalan. Battle of bandits. In *UAI*, pages 805–814, 2018.
- [68] Aadirupa Saha and Akshay Krishnamurthy. Efficient and optimal algorithms for contextual dueling bandits under realizability. In *International Conference on Algorithmic Learning Theory*, pages 968–994. PMLR, 2022.
- [69] Aadirupa Saha, Tomer Koren, and Yishay Mansour. Adversarial dueling bandits. In *International Conference on Machine Learning*, pages 9235–9244. PMLR, 2021.
- [70] Aadirupa Saha, Aldo Pacchiano, and Jonathan Lee. Dueling rl: Reinforcement learning with trajectory preferences. In *International conference on artificial intelligence and statistics*, pages 6263–6289. PMLR, 2023.



- [71] Antoine Scheid, Etienne Boursier, Alain Durmus, Michael I Jordan, Pierre Ménard, Eric Moulines, and Michal Valko. Optimal design for reward modeling in rlhf. *arXiv preprint arXiv:2410.17055*, 2024.
- [72] Ayush Sekhari, Karthik Sridharan, Wen Sun, and Runzhe Wu. Contextual bandits and imitation learning with preference-based active queries. *Advances in Neural Information Processing Systems*, 36:11261–11295, 2023.
- [73] Ohad Shamir. On the complexity of bandit and derivative-free stochastic convex optimization. In *Conference on learning theory*, pages 3–24. PMLR, 2013.
- [74] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in neural information processing systems*, 33:3008–3021, 2020.
- [75] Gokul Swamy, Christoph Dann, Rahul Kidambi, Zhiwei Steven Wu, and Alekh Agarwal. A minimaximalist approach to reinforcement learning from human feedback. In *Proceedings of the 41st International Conference on Machine Learning*, pages 47345–47377, 2024.
- [76] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- [77] Kiran Koshy Thekumparampil, Gaurush Hiranandani, Kousha Kalantari, Shoham Sabach, and Branislav Kveton. Comparing few to rank many: Active human preference learning using randomized frank-wolfe method. In *Forty-second International Conference on Machine Learning*, 2024.
- [78] Joel A Tropp. User-friendly tail bounds for matrix martingales. *ACM Report*, 1, 2011.
- [79] Andrew J Wagenmaker, Yifang Chen, Max Simchowitz, Simon Du, and Kevin Jamieson. Reward-free rl is no harder than reward-aware rl in linear markov decision processes. In *International Conference on Machine Learning*, pages 22430–22456. PMLR, 2022.
- [80] Christian Wirth and Johannes Fürnkranz. Preference-based reinforcement learning: A preliminary survey. In *Proceedings of the ECML/PKDD-13 Workshop on Reinforcement Learning from Generalized Feedback: Beyond Numeric Rewards*, 2013.
- [81] Christian Wirth, Johannes Fürnkranz, and Gerhard Neumann. Model-free preference-based reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
- [82] Christian Wirth, Riad Akrou, Gerhard Neumann, and Johannes Fürnkranz. A survey of preference-based reinforcement learning methods. *Journal of Machine Learning Research*, 18 (136):1–46, 2017.
- [83] Huasen Wu and Xin Liu. Double thompson sampling for dueling bandits. *Advances in neural information processing systems*, 29, 2016.
- [84] Runzhe Wu and Wen Sun. Making rl with preference-based feedback efficient via randomization. In *The Twelfth International Conference on Learning Representations*, 2023.
- [85] Yue Wu, Tao Jin, Qiwei Di, Hao Lou, Farzad Farnoud, and Quanquan Gu. Borda regret minimization for generalized linear dueling bandits. In *International Conference on Machine Learning*, pages 53571–53596. PMLR, 2024.
- [86] Tengyang Xie, Dylan J Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Hassan Awadallah, and Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit  $q^*$ -approximation for sample-efficient rlhf. In *The Thirteenth International Conference on Learning Representations*, 2024.
- [87] Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint. In *Forty-first International Conference on Machine Learning*, 2024.

- [88] Yichong Xu, Ruosong Wang, Lin Yang, Aarti Singh, and Artur Dubrawski. Preference-based reinforcement learning with finite-time guarantees. *Advances in Neural Information Processing Systems*, 33:18784–18794, 2020.
- [89] Yisong Yue, Josef Broder, Robert Kleinberg, and Thorsten Joachims. The k-armed dueling bandits problem. *Journal of Computer and System Sciences*, 78(5):1538–1556, 2012.
- [90] Andrea Zanette, Ching-An Cheng, and Alekh Agarwal. Cautiously optimistic policy optimization and exploration with linear function approximation. In *Conference on Learning Theory*, pages 4473–4525. PMLR, 2021.
- [91] Wenhao Zhan, Masatoshi Uehara, Nathan Kallus, Jason D Lee, and Wen Sun. Provable offline preference-based reinforcement learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- [92] Wenhao Zhan, Masatoshi Uehara, Wen Sun, and Jason D Lee. Provable reward-agnostic preference-based reinforcement learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- [93] Yu-Jie Zhang and Masashi Sugiyama. Online (multinomial) logistic bandit: Improved regret and constant computation cost. *Advances in Neural Information Processing Systems*, 36, 2024.
- [94] Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human feedback from pairwise or k-wise comparisons. In *International Conference on Machine Learning*, pages 43037–43067. PMLR, 2023.
- [95] Masrour Zoghi, Shimon Whiteson, Remi Munos, and Maarten Rijke. Relative upper confidence bound for the k-armed dueling bandit problem. In *International conference on machine learning*, pages 10–18. PMLR, 2014.
- [96] Masrour Zoghi, Zohar S Karnin, Shimon Whiteson, and Maarten De Rijke. Copeland dueling bandits. *Advances in neural information processing systems*, 28, 2015.

# Appendix

## Table of Contents

---

<b>A Related Works</b>	<b>11</b>
<b>B Efficient Assortment Selection</b>	<b>12</b>
<b>C Notation</b>	<b>13</b>
<b>D Proof for PL Loss</b>	<b>15</b>
D.1 Main Proof of Theorem 1 . . . . .	15
D.2 Proofs of Lemmas for Theorem 1 . . . . .	19
<b>E Proof for RB Loss</b>	<b>24</b>
E.1 Rank-Breaking (RB) Loss. . . . .	25
E.2 Main Proof of Theorem E.1 . . . . .	25
E.3 Proofs of Lemmas for Theorem E.1 . . . . .	27
<b>F Technical Lemmas</b>	<b>30</b>
<b>G Proof of Theorem 2</b>	<b>32</b>
G.1 Main Proof of Theorem 2 . . . . .	32
G.2 Proof of Lemmas for Theorem 2 . . . . .	35
<b>H Additional Discussions</b>	<b>36</b>
H.1 Arbitrary Reference Action for More Efficient Assortment Selection . . . . .	36
H.2 Suboptimality Gap Under Sufficient Diversity Condition . . . . .	37
H.3 Extension to Active Learning Setting . . . . .	38
<b>I Numerical Experiments</b>	<b>41</b>
I.1 Synthetic Data . . . . .	41
I.2 Real-World Dataset . . . . .	44
<b>J Limitations</b>	<b>46</b>

---

## A Related Works

**PbRL.** Fueled by the remarkable success of LLMs [17, 57, 61], the theoretical study of PbRL has rapidly emerged as a central focus within the research community. Early work in this area traces back to the dueling bandits literature [89, 96, 67, 8]. Subsequent studies have extended this line of research to the RL framework, considering both online settings [88, 52, 16, 70, 84] and offline settings [94, 92, 44]. More recently, under the active learning framework—where the full set of contexts  $\mathcal{X}$  is accessible—many studies aim to improve sample efficiency by selecting prompts either based on the differences in estimated rewards for their responses [50] or through D-optimal design methods [47, 71, 18, 49, 77, 37]. However, most of these works focus exclusively on pairwise preference feedback and cannot be extended to more general ranking feedback cases. Mukherjee et al. [49] study the online learning-to-rank problem when prompts are given along with  $K$  candidate answers, while Thekumparampil et al. [77] investigate learning to rank  $N \geq K$  answers from partial rankings over  $K$  answers, but under a context-free setting. In this paper, we consider a stochastic contextual setting (more general than Thekumparampil et al. [77]), where contexts are sampled from an unknown but fixed distribution, and aim to minimize the suboptimality gap using ranking feedback of up to length  $K$ .

**Dueling bandits.** The dueling bandit framework, introduced by Yue et al. [89], departs from the classical multi-armed bandit setting by requiring the learner to select two arms and observe only their pairwise preference. For general preferences, a single best arm that is globally dominant may not exist. To address this, various alternative winners have been proposed, including the Condorcet winner [95, 34], Copeland winner [96, 83, 35], Borda winner [29, 24, 27, 69, 85], and von Neumann winner [62, 23, 7], each with its own corresponding performance metric.

To address scalability and contextual information, Saha [65] proposed a structured contextual dueling bandit setting in which preferences are modeled using a Bradley–Terry–Luce (BTL) model [11] based on the unknown intrinsic rewards of each arm. In a similar setting, Bengs et al. [9] studied a contextual linear stochastic transitivity model, and Di et al. [20] proposed a layered algorithm that achieves variance-aware regret bounds. More recently, Di et al. [19] studied the linear contextual dueling bandit problem under adversarial feedback and developed an algorithm that achieves regret bounds independent of  $\mathcal{O}(e^B)$  under the BTL model.

Another line of work moves beyond reward-based assumptions and instead models preference feedback using a general function class. For example, Saha and Krishnamurthy [68] proposed an algorithm that achieves optimal regret for the  $K$ -armed contextual dueling bandit setting. Sekhari et al. [72] further generalized this framework and introduced an algorithm with theoretical guarantees for both regret and query complexity.

**Logistic and MNL bandits.** Our work is also closely related to logistic bandits and multinomial logit (MNL) bandits. The logistic bandit problem [22, 25, 2, 26, 42, 43] is a special case of the MNL bandit model in which the agent offers only a single item (i.e.,  $K = 1$ ) at each round and receives binary feedback indicating whether the item was selected (1) or not (0). Fauray et al. [25] examined how the regret in logistic bandits depends on the non-linearity parameter  $\kappa$  of the logistic link function and proposed the first algorithm whose regret bound eliminates explicit dependence on  $1/\kappa = \mathcal{O}(e^B)$ . Abeille et al. [2] further improved the theoretical dependency on  $1/\kappa$  and established a matching, problem-dependent lower bound. Building on this, Fauray et al. [26] developed a more computationally efficient algorithm whose regret still matches the lower bound proved by Abeille et al. [2].

Multinomial logit (MNL) bandits tackle a more sophisticated problem than logistic bandits. Instead of offering a single item and observing binary feedback, the learner chooses a subset of items—underscoring the combinatorial nature of the task—and receives non-uniform rewards driven by an MNL choice model [5, 4, 56, 15, 53, 54, 58, 3, 39, 41]. A recent breakthrough by Lee and Oh [39] closed a long-standing gap by providing a computationally efficient algorithm that attains the minimax-optimal regret for this setting. Building on this result, Lee and Oh [41] further reduced the regret bound by a factor polynomial in  $B$  and logarithmic in  $K$ , and established the first variance-dependent regret bounds for MNL bandits. Extending this line of work to the RL setting, Lee and Oh [40] proposed a new framework, *combinatorial RL with preference feedback*, in which the agent selects a subset of items to maximize long-term cumulative rewards.

Our work extends the online confidence bound analysis of Lee and Oh [41] to the Plackett–Luce (PL) model. This extension is natural because the PL probability distribution decomposes into a sequence of MNL probabilities over successive choices. Crucially, we leverage their key insight—that the MNL loss exhibits an  $\ell_\infty$ -self-concordant property—to eliminate the harmful  $\mathcal{O}(e^B)$  dependence. This is one of the main contributions of our work (see Lemma D.2).

## B Efficient Assortment Selection

In this section, we describe how the assortment selection rule in Equation (4) can be solved efficiently.

Given  $x_t$ , the reference action–assortment pair  $(\bar{a}_t, S_t)$  is selected by evaluating each candidate reference action  $\bar{a} \in \mathcal{A}$ . For each  $\bar{a}$ , we construct an assortment  $S$  beginning with the singleton  $\bar{a}$ , and iteratively add actions  $a \in \mathcal{A} \setminus \{\bar{a}\}$  in decreasing order of their uncertainty relative to  $\bar{a}$ , measured by

$$\|\phi(x_t, a) - \phi(x_t, \bar{a})\|_{H_t^{-1}}.$$

Let  $a_{tk}(\bar{a})$  denote the action with the  $k$ -th highest uncertainty with respect to  $\bar{a}$  at round  $t$ . For example,  $a_{t1}(\bar{a}) = \operatorname{argmax}_{a \in \mathcal{A} \setminus \{\bar{a}\}} \|\phi(x_t, a) - \phi(x_t, \bar{a})\|_{H_t^{-1}}$ . We add actions greedily to the set  $S$ ,

---

**Procedure B.1** Greedy Selection of Reference Action and Assortment
 

---

```

1: Input:  $x_t, H_t^{-1}, \mathcal{A}, K$ 
2: Initialize:  $(\bar{a}_t^*, S_t^*, \text{max\_avg}) \leftarrow (\text{None}, \text{None}, -\infty)$ 
3: for all  $\bar{a} \in \mathcal{A}$  do
4:   Initialize  $S \leftarrow \{\bar{a}\}$ ,  $\text{prev\_avg} \leftarrow 0$ 
5:   while  $|S| < K$  do
6:     Find
      
$$a^* \leftarrow \operatorname{argmax}_{a \in \mathcal{A} \setminus S} \|\phi(x_t, a) - \phi(x_t, \bar{a})\|_{H_t^{-1}}$$

7:     Tentatively update  $S' \leftarrow S \cup \{a^*\}$ 
8:     Compute
      
$$\text{cur\_avg} \leftarrow \frac{1}{|S'|} \sum_{a \in S'} \|\phi(x_t, a) - \phi(x_t, \bar{a})\|_{H_t^{-1}}$$

9:     if  $\text{cur\_avg} < \text{prev\_avg}$  then
10:      break
11:     else
12:       $S \leftarrow S'$ 
13:       $\text{prev\_avg} \leftarrow \text{cur\_avg}$ 
14:     end if
15:   end while
16:   if  $\text{prev\_avg} > \text{max\_avg}$  then
17:      $(\bar{a}_t^*, S_t^*, \text{max\_avg}) \leftarrow (\bar{a}, S, \text{prev\_avg})$ 
18:   end if
19: end for
20: return  $(\bar{a}_t^*, S_t^*)$ 

```

---

as long as the average uncertainty continues to increase:

$$\frac{1}{|S|} \sum_{a \in S} \|\phi(x_t, a) - \phi(x_t, \bar{a})\|_{H_t^{-1}}, \quad \text{where } \bar{a} \in S.$$

Among all candidates  $\bar{a} \in \mathcal{A}$ , we select the pair  $(\bar{a}_t, S_t)$  that achieves the highest average uncertainty. The pseudocode is given in Procedure B.1.

For each candidate reference action, the algorithm incrementally constructs a subset of actions by greedily adding those with the highest uncertainty relative to the reference—stopping once the average uncertainty no longer increases. This greedy strategy guarantees that, for each reference, the selected subset maximizes the average uncertainty. By applying this procedure across all possible reference actions and selecting the pair that achieves the highest score, the algorithm obtains the global optimum over all reference–assortment combinations.

As for the computational cost, each greedy addition step involves searching over  $\mathcal{O}(N)$  candidate actions, resulting in a total of  $\tilde{\mathcal{O}}(NK)$  operations per each reference action  $\bar{a}$ . Repeating this process for all  $N$  candidate references yields a total cost of  $\tilde{\mathcal{O}}(N^2K)$ .

## C Notation

Let  $T$  denote the total number of rounds, with  $t \in [T]$  representing the current round. We use  $N$  for the total number of items,  $K$  for the maximum assortment size,  $d$  for the feature vector dimension, and  $B$  as an upper bound on the norm of the unknown parameter. For notational convenience, we provide Table C.1.

For clarity, we derive the first- and second-order derivatives (i.e., gradients and Hessians) of the loss functions. For the PL loss at round  $t$  for the  $j$ 'th ranking, let  $y_{ti}^{(j)} = 1$  if  $i = j$ , and  $y_{ti}^{(j)} = 0$  for

Table C.1: Symbols

$\mathcal{X}, \mathcal{A}, \mathcal{S}$	context (prompt) space, action (answer) space, assortment space
$\phi(x, a) \in \mathbb{R}^d$	feature representation of context-action pair $(x, a)$
$z_{tjk}$	$:= \phi(x_t, \sigma_{tj}) - \phi(x_t, \sigma_{tk})$ , feature difference between $\sigma_{tj}$ and $\sigma_{tk}$ under context $x_t$
$S_t$	assortment chosen by an algorithm at round $t$
$\ell_t^{(j)}(\theta)$	$:= -\log \left( \frac{\exp(\phi(x_t, \sigma_{tj})^\top \theta)}{\sum_{k=j}^{ S_t } \exp(\phi(x_t, \sigma_{tk})^\top \theta)} \right)$ , PL loss at round $t$ for $j$ 'th ranking
$\ell_t^{(j,k)}(\theta)$	$:= -\log \left( \frac{\exp(\phi(x_t, \sigma_{tj})^\top \theta)}{\sum_{m \in \{j,k\}} \exp(\phi(x_t, \sigma_{tm})^\top \theta)} \right)$ , RB loss at round $t$ for comparison $\sigma_{tj}$ vs $\sigma_{tk}$
$\nabla^2 \ell_t^{(j)}(\theta)$	$= \sum_{k=j}^{ S_t } \sum_{k'=j}^{ S_t } \frac{\exp((\phi(x_t, \sigma_{tk}) + \phi(x_t, \sigma_{tk'}))^\top \theta)}{2 \left( \sum_{k'=j}^{ S_t } \exp(\phi(x_t, \sigma_{tk'})^\top \theta) \right)^2} \cdot z_{tkk'} z_{tkk'}^\top$
$\nabla^2 \ell_t^{(j,k)}(\theta)$	$= \mu \left( z_{tjk}^\top \theta \right) z_{tjk} z_{tjk}^\top$ , where $\mu(w) = \frac{1}{1+e^{-w}}$ is sigmoid function
$\hat{\theta}_t^{(j+1)}$	online parameter estimate using PL loss at round $t$ , after $j$ 'th update
$\hat{\theta}_t^{(j,k+1)}$	online parameter estimate using RB loss at round $t$ , after $(j, k)$ 'th comparison update
$\eta$	$:= \frac{1}{2}(1 + 3\sqrt{2}B)$ , step-size parameter
$\lambda$	$:= \Omega(d \log(KT/\delta) + \eta(B + d))$ , regularization parameter
$H_t$	$:= \sum_{s=1}^{t-1} \sum_{j=1}^{ S_s } \nabla^2 \ell_s^{(j)}(\hat{\theta}_s^{(j+1)}) + \lambda \mathbf{I}_d$ (or $\sum_{s=1}^{t-1} \sum_{j=1}^{ S_s -1} \sum_{k=j+1}^{ S_s } \nabla^2 \ell_s^{(j,k)}(\hat{\theta}_s^{(j,k+1)}) + \lambda \mathbf{I}_d$ )
$\tilde{H}_t^{(j)}$	$:= H_t + \eta \sum_{j'=1}^j \nabla^2 \ell_t^{(j')}(\hat{\theta}_t^{(j')})$ (for PL loss)
$\tilde{H}_t^{(j,k)}$	$:= H_t + \eta \sum_{(j',k') \leq (j,k)} \nabla^2 \ell_t^{(j',k')}(\hat{\theta}_t^{(j',k')})$ (for RB loss)
$\beta_t(\delta)$	$:= \mathcal{O} \left( B\sqrt{d \log(tK/\delta)} + B\sqrt{\lambda} \right)$ , confidence radius for $\theta_t$ at round $t$
$\mathcal{T}^w$	$:= \left\{ t \in [T] : \max_{a \in \mathcal{A}} \ \phi(x_t, a) - \phi(x_t, \bar{a}_t)\ _{H_t^{-1}} \geq \frac{1}{3\sqrt{2}\beta_{T+1}(\delta)} \right\}$ , warm-up rounds
$\Lambda_t$	$:= \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{a \in S_s} (\phi(x_s, a) - \phi(x_s, \bar{a}_s)) (\phi(x_s, a) - \phi(x_s, \bar{a}_s))^\top + \lambda \mathbf{I}_d$
$\mathcal{T}_0$	$:= \left\{ t \in [T] : \sum_{a \in S_t} \ \phi(x_t, a) - \phi(x_t, \bar{a}_t)\ _{\Lambda_t^{-1}} \geq 1 \right\}$ , large EP rounds

otherwise. Then, we have

$$\begin{aligned}
\ell_t^{(j)}(\theta) &= -\log \left( \frac{\exp(\phi(x_t, \sigma_{tj})^\top \theta)}{\sum_{k=j}^{|S_t|} \exp(\phi(x_t, \sigma_{tk})^\top \theta)} \right) = -\sum_{i=j}^{|S_t|} y_{ti}^{(j)} \log \left( \underbrace{\frac{\exp(\phi(x_t, \sigma_{ti})^\top \theta)}{\sum_{k=j}^{|S_t|} \exp(\phi(x_t, \sigma_{tk})^\top \theta)}}_{=: P_{t,\theta}^{(j)}(\sigma_{ti})} \right) \\
&= -\sum_{i=j}^{|S_t|} y_{ti}^{(j)} \log P_{t,\theta}^{(j)}(\sigma_{ti}), \\
\nabla \ell_t^{(j)}(\theta) &= \sum_{i=j}^{|S_t|} \left( P_{t,\theta}^{(j)}(\sigma_{ti}) - y_{ti}^{(j)} \right) \phi(x_t, \sigma_{ti}), \\
\nabla^2 \ell_t^{(j)}(\theta) &= \sum_{i=j}^{|S_t|} P_{t,\theta}^{(j)}(\sigma_{ti}) \phi(x_t, \sigma_{ti}) \phi(x_t, \sigma_{ti})^\top - \sum_{i=j}^{|S_t|} \sum_{k=j}^{|S_t|} P_{t,\theta}^{(j)}(\sigma_{ti}) P_{t,\theta}^{(j)}(\sigma_{tk}) \phi(x_t, \sigma_{ti}) \phi(x_t, \sigma_{tk})^\top \\
&= \frac{1}{2} \sum_{i=j}^{|S_t|} \sum_{k=j}^{|S_t|} P_{t,\theta}^{(j)}(\sigma_{ti}) P_{t,\theta}^{(j)}(\sigma_{tk}) (\phi(x_t, \sigma_{ti}) - \phi(x_t, \sigma_{tk})) (\phi(x_t, \sigma_{ti}) - \phi(x_t, \sigma_{tk}))^\top.
\end{aligned}$$

For the RB loss at round  $t$  for the pairwise comparison between  $\sigma_{tj}$  and  $\sigma_{tk}$ , let  $y_{ti}^{(j,k)} = 1$  if  $i = j$ , and  $y_{ti}^{(j,k)} = 0$  for otherwise (i.e., when  $i = k$ ). Then, we have

$$\begin{aligned}\ell_t^{(j,k)}(\boldsymbol{\theta}) &= -\log \left( \frac{\exp(\phi(x_t, \sigma_{tj})^\top \boldsymbol{\theta})}{\exp(\phi(x_t, \sigma_{tj})^\top \boldsymbol{\theta}) + \exp(\phi(x_t, \sigma_{tk})^\top \boldsymbol{\theta})} \right) \\ &= -\log \mu \left( (\phi(x_t, \sigma_{tj}) - \phi(x_t, \sigma_{tk}))^\top \boldsymbol{\theta} \right), \quad \text{where } \mu(w) = \frac{1}{1 + e^{-w}}, \\ \nabla \ell_t^{(j,k)}(\boldsymbol{\theta}) &= \left( \mu \left( (\phi(x_t, \sigma_{tj}) - \phi(x_t, \sigma_{tk}))^\top \boldsymbol{\theta} \right) - 1 \right) (\phi(x_t, \sigma_{tj}) - \phi(x_t, \sigma_{tk})), \\ \nabla^2 \ell_t^{(j,k)}(\boldsymbol{\theta}) &= \dot{\mu} \left( (\phi(x_t, \sigma_{tj}) - \phi(x_t, \sigma_{tk}))^\top \boldsymbol{\theta} \right) (\phi(x_t, \sigma_{tj}) - \phi(x_t, \sigma_{tk})) (\phi(x_t, \sigma_{tj}) - \phi(x_t, \sigma_{tk}))^\top.\end{aligned}$$

## D Proof for PL Loss

In this section, we present the proof of Theorem 1.

### D.1 Main Proof of Theorem 1

**PL loss and OMD.** We begin by recalling the loss function and the parameter update rule. Specifically, we use the PL loss defined in Equation (2) and update the parameter according to Equation (3).

$$\ell_t(\boldsymbol{\theta}) := \sum_{j=1}^{|S_t|} -\log \underbrace{\left( \frac{\exp(\phi(x_t, \sigma_{tj})^\top \boldsymbol{\theta})}{\sum_{k=j}^{|S_t|} \exp(\phi(x_t, \sigma_{tk})^\top \boldsymbol{\theta})} \right)}_{=: \ell_t^{(j)}(\boldsymbol{\theta})} = \sum_{j=1}^{|S_t|} \ell_t^{(j)}(\boldsymbol{\theta}).$$

and

$$\hat{\boldsymbol{\theta}}_t^{(j+1)} = \operatorname{argmin}_{\boldsymbol{\theta} \in \Theta} \langle \nabla \ell_t^{(j)}(\hat{\boldsymbol{\theta}}_t^{(j)}), \boldsymbol{\theta} \rangle + \frac{1}{2\eta} \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_t^{(j)}\|_{\tilde{H}_t^{(j)}}^2, \quad j = 1, \dots, |S_t|,$$

where  $\hat{\boldsymbol{\theta}}_t^{(|S_t|+1)} = \hat{\boldsymbol{\theta}}_{t+1}^{(1)}$ , and  $\eta := \frac{1}{2}(1 + 3\sqrt{2}B)$  is the step-size parameter. The matrix  $\tilde{H}_t^{(j)}$  is given by  $\tilde{H}_t^{(j)} := H_t + \eta \sum_{j'=1}^j \nabla^2 \ell_t^{(j')}(\hat{\boldsymbol{\theta}}_t^{(j')})$ , where

$$H_t := \sum_{s=1}^{t-1} \sum_{j=1}^{|S_s|} \nabla^2 \ell_s^{(j)}(\hat{\boldsymbol{\theta}}_s^{(j+1)}) + \lambda \mathbf{I}_d, \quad \lambda > 0.$$

**Online confidence bound for PL loss.** Now, we present the confidence bound for online parameter estimation in MNL models, as recently proposed by Lee and Oh [41].

**Lemma D.1** (Online confidence bound, Theorem 4.2 of Lee and Oh 41). *Let  $\delta \in (0, 1]$ . We set  $\eta = (1 + 3\sqrt{2}B)/2$  and  $\lambda = \max\{12\sqrt{2}B\eta, 144\eta d, 2\}$ . Then, under Assumption 1, with probability at least  $1 - \delta$ , we have*

$$\|\hat{\boldsymbol{\theta}}_t - \boldsymbol{\theta}^*\|_{H_t} \leq \beta_t(\delta) = \mathcal{O} \left( B\sqrt{d \log(t/\delta)} + B\sqrt{\lambda} \right), \quad \forall t \geq 1.$$

We now extend this result to our setting. Since the total number of updates up to round  $t$  is  $\sum_{s=1}^t |S_s|$ , the corresponding confidence bound can be expressed as follows:

**Corollary D.1** (Online confidence bound for PL loss). *Let  $\delta \in (0, 1]$ . We set  $\eta = (1 + 3\sqrt{2}B)/2$  and  $\lambda = \max\{12\sqrt{2}B\eta, 144\eta d, 2\}$ . Then, under Assumption 1, with probability at least  $1 - \delta$ , we have*

$$\|\hat{\boldsymbol{\theta}}_t^{(j)} - \boldsymbol{\theta}^*\|_{H_t^{(j)}} \leq \beta_t(\delta) = \mathcal{O} \left( B\sqrt{d \log(tK/\delta)} + B\sqrt{\lambda} \right), \quad \forall t \geq 1, j \leq |S_t|,$$

where  $H_t^{(j)} := H_t + \sum_{j'=1}^{j-1} \nabla^2 \ell_s^{(j')}(\hat{\boldsymbol{\theta}}_s^{(j'+1)}) + \lambda \mathbf{I}_d$  and  $\hat{\boldsymbol{\theta}}_t^{(1)} = \hat{\boldsymbol{\theta}}_t$ .

**Useful definitions.** We define the set of *warm-up rounds*, denoted by  $\mathcal{T}^w$ , which consists of rounds with large uncertainty, as follows:

$$\mathcal{T}^w := \left\{ t \in [T] : \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} \geq \frac{1}{3\sqrt{2}\beta_{T+1}(\delta)} \right\}, \quad (\text{D.1})$$

where  $\beta_{T+1}(\delta)$  denotes the confidence radius as defined in Corollary D.1. Furthermore, we define the regularized sample covariance matrix of feature differences (with respect to  $\phi(x_s, \bar{a}_s)$ ) over the non-warm-up rounds as:

$$\Lambda_t := \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{a \in S_s} (\phi(x_s, a) - \phi(x_s, \bar{a}_s)) (\phi(x_s, a) - \phi(x_s, \bar{a}_s))^\top + \lambda \mathbf{I}_d. \quad (\text{D.2})$$

To control the elliptical potentials, we also define the set of *large elliptical potential (EP) rounds*, denoted by  $\mathcal{T}_0$ , as follows:

$$\mathcal{T}_0 := \left\{ t \in [T] : \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{\Lambda_t^{-1}} \geq 1 \right\}, \quad (\text{D.3})$$

**Key lemmas.** We now present key lemmas needed to prove Theorem 1. The following lemma, one of our main contributions, is crucial for avoiding the  $1/\kappa = \mathcal{O}(e^B)$  dependency in the leading term.

**Lemma D.2.** *Let  $\Lambda_t$  be defined as in Equation (D.2). Set  $\lambda = \Omega(d \log(KT/\delta))$ . Then, for all  $t \in [T]$ , with probability at least  $1 - \delta$ , we have*

$$H_t \geq \frac{1}{50} \Lambda_t.$$

The proof is deferred to Appendix D.2.1.

The following lemma is a variant of the elliptical potential lemma [1], adapted specifically to the assortment offering setting. For completeness, we provide the proof tailored to our setting.

**Lemma D.3** (Elliptical potential lemma for  $S_t$ ). *Let  $\{z_{ta}\}_{t \geq 1, a \in S_t}$  be a bounded sequence in  $\mathbb{R}^d$  satisfying  $\max_{t \geq 1} \|z_{ta}\|_2 \leq X$ . For any  $t \geq 1$ , we define  $\Lambda_t := \sum_{s=1}^{t-1} \sum_{a \in S_s} z_{sa} z_{sa}^\top + \lambda \mathbf{I}_d$  with  $\lambda > 0$ . Then, we have*

$$\sum_{t=1}^T \min \left\{ 1, \sum_{a \in S_t} \|z_{ta}\|_{\Lambda_t^{-1}}^2 \right\} \leq 2d \log \left( 1 + \frac{X^2 KT}{d\lambda} \right).$$

The proof is deferred to Appendix D.2.2.

The cardinality of the set  $\mathcal{T}_0$  can be bounded by a variant of the elliptical potential counting lemma [38, 33].

**Lemma D.4** (Elliptical potential count lemma for  $S_t$ ). *Let  $\{z_{ta}\}_{t \geq 1, a \in S_t}$  be a bounded sequence in  $\mathbb{R}^d$  satisfying  $\max_{t \geq 1} \|z_{ta}\|_2 \leq X$ . For any  $t \geq 1$ , we define  $\Lambda_t := \sum_{s=1}^{t-1} \sum_{a \in S_s} z_{sa} z_{sa}^\top + \lambda \mathbf{I}_d$  with  $\lambda > 0$ . Let  $\mathcal{T}_0 \subseteq [T]$  be the set of indices where  $\sum_{a \in S_t} \|z_{ta}\|_{\Lambda_t^{-1}}^2 \geq L$ . Then,*

$$|\mathcal{T}_0| \leq \frac{2d}{\log(1+L)} \log \left( 1 + \frac{X^2 K}{\log(1+L)\lambda} \right).$$

The proof is deferred to Appendix D.2.3.

The size of the set  $\mathcal{T}^w \cap (\mathcal{T}_0)^c$  is bounded as described in the following lemma:

**Lemma D.5.** *Let  $\mathcal{T}_0 := \{t \in [T] : \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{\Lambda_t^{-1}} \geq 1\}$  and  $\mathcal{T}^w = \{t \in [T] : \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} \geq \frac{1}{3\sqrt{2}\beta_{T+1}(\delta)}\}$ . Define  $\kappa := e^{-4B}$ . Then, the size of the set  $\mathcal{T}^w \cap (\mathcal{T}_0)^c$  is bounded as follows:*

$$|\mathcal{T}^w \cap (\mathcal{T}_0)^c| \leq \frac{12\sqrt{2}K^2}{\kappa} \beta_{T+1}(\delta)^2 d \log \left( 1 + \frac{2KT}{d\lambda} \right).$$



The proof is deferred to Appendix D.2.4.

We are now ready to provide the proof of Theorem 1.

*Proof of Theorem 1.* To begin, we define a martingale difference sequence (MDS)  $\zeta_t$  as follows:

$$\zeta_t := \mathbb{E}_{x \sim \rho} \left[ \left( \phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)) \right)^\top \boldsymbol{\theta}^* \right] - \left( \phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)) \right)^\top \boldsymbol{\theta}^*,$$

which satisfies  $|\zeta_t| \leq 2B$ . Then, by the definition of the suboptimality gap, we have

$$\begin{aligned} \text{SubOpt}(T) &= \mathbb{E}_{x \sim \rho} \left[ \left( \phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)) \right)^\top \boldsymbol{\theta}^* \right] \\ &= \frac{1}{T} \sum_{t=1}^T \left( \phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)) \right)^\top \boldsymbol{\theta}^* + \frac{1}{T} \sum_{t=1}^T \zeta_t \quad (\text{Def. of } \zeta_t) \\ &\leq \frac{1}{T} \sum_{t=1}^T \left( \phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)) \right)^\top \left( \boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1} \right) + \frac{1}{T} \sum_{t=1}^T \zeta_t \\ &\quad (\hat{\pi}_T(x_t) = \arg\max_{a \in \mathcal{A}} \phi(x_t, a)^\top \hat{\boldsymbol{\theta}}_{T+1}) \\ &\leq \frac{1}{T} \sum_{t=1}^T \left( \phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)) \right)^\top \left( \boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1} \right) + \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{T}} \right), \quad (\text{D.4}) \end{aligned}$$

where the last inequality follows from the Azuma–Hoeffding inequality. Specifically, for any  $T \geq 1$ , with probability at least  $1 - \delta$ , we have

$$\frac{1}{T} \sum_{t=1}^T \zeta_t \leq \frac{1}{T} \sqrt{8B^2 T \log(1/\delta)} = \tilde{\mathcal{O}} \left( \frac{1}{\sqrt{T}} \right).$$

To complete the proof, it remains to bound the first term in Equation (D.4).

Recall the definitions of the set of *large elliptical potential (EP) rounds* (Equation (D.3)), denoted by  $\mathcal{T}_0$ , and the set of *warm-up rounds* (Equation (D.1)), denoted by  $\mathcal{T}^w$ :

$$\begin{aligned} \mathcal{T}_0 &:= \left\{ t \in [T] : \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{\Lambda_t^{-1}} \geq 1 \right\}, \quad (\text{large EP rounds}) \\ \mathcal{T}^w &:= \left\{ t \in [T] : \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} \geq \frac{1}{3\sqrt{2}\beta_{T+1}(\delta)} \right\}, \quad (\text{warm-up rounds}) \end{aligned}$$

where  $\Lambda_t$  is defined in Equation (D.2). Then, by applying the elliptical potential count lemma (Lemma D.4) and the bound on the cardinality of the set  $|\mathcal{T}^w \cap (\mathcal{T}_0)^c|$  lemma (Lemma D.5), we

obtain

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T (\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)))^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1}) \\
&= \frac{1}{T} \sum_{t \in \mathcal{T}_0} (\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)))^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1}) \\
&+ \frac{1}{T} \sum_{t \in \mathcal{T}^w \cap (\mathcal{T}_0)^c} (\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)))^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1}) \\
&+ \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} (\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)))^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1}) \\
&\leq \frac{4B}{T} |\mathcal{T}_0| + \frac{4B}{T} |\mathcal{T}^w \cap (\mathcal{T}_0)^c| + \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} (\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)))^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1}) \\
&\hspace{15em} (\text{Assumption 1}) \\
&\leq \frac{8B}{\log(2)T} d \log \left( 1 + \frac{2K}{\log(2)\lambda} \right) + \frac{48\sqrt{2}BK^2}{\kappa T} \beta_{T+1}(\delta)^2 d \log \left( 1 + \frac{2KT}{d\lambda} \right) \\
&\hspace{15em} (\text{Lemma D.4 and D.5}) \\
&+ \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} (\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)))^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1}). \tag{D.5}
\end{aligned}$$

To further bound the last term of Equation (D.5), we get

$$\begin{aligned}
& \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} (\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)))^\top (\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1}) \\
&\leq \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \|\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t))\|_{H_{T+1}^{-1}} \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1}\|_{H_{T+1}} \\
&\hspace{15em} (\text{H\"older's ineq.}) \\
&\leq \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \|\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t))\|_{H_t^{-1}} \|\boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1}\|_{H_{T+1}} \\
&\hspace{15em} (H_{T+1} \geq H_t) \\
&\leq \frac{\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \|\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t))\|_{H_t^{-1}}. \\
&\hspace{15em} (\text{Corollary D.1, with prob. } 1 - \delta)
\end{aligned}$$

We denote  $S_t^* = \{\pi^*(x_t), \hat{\pi}_T(x_t)\}$ . Then, we have

$$\begin{aligned}
& \frac{\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \|\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t))\|_{H_t^{-1}} \\
&= \frac{\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \sum_{a \in S_t^*} \|\phi(x_t, a) - \phi(x_t, \hat{\pi}_T(x_t))\|_{H_t^{-1}} \\
&= \frac{\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{|S_t^*|}{|S_t^*|} \sum_{a \in S_t^*} \|\phi(x_t, a) - \phi(x_t, \hat{\pi}_T(x_t))\|_{H_t^{-1}} \\
&= \frac{2\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{1}{|S_t^*|} \sum_{a \in S_t^*} \|\phi(x_t, a) - \phi(x_t, \hat{\pi}_T(x_t))\|_{H_t^{-1}}, \tag{D.6}
\end{aligned}$$

where the last equality holds due to the fact that  $|S_t^*| = 2$ . To proceed, by our efficient assortment selection rule in Equation (4), we obtain

$$\begin{aligned}
& \frac{2\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{1}{|S_t^*|} \sum_{a \in S_t^*} \|\phi(x_t, a) - \phi(x_t, \hat{\pi}_T(x_t))\|_{H_t^{-1}} \\
& \leq \frac{2\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{1}{|S_t|} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} \\
& \hspace{15em} (S_t \text{ selection rule, Eqn. (4)}) \\
& \leq \frac{2\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \left(\frac{1}{|S_t|}\right)^2 |S_t|} \sqrt{\sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}^2} \\
& \hspace{15em} (\text{Cauchy-Schwartz ineq.}) \\
& = \frac{2\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \left(\frac{1}{|S_t|}\right)^2 |S_t|} \sqrt{50 \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{\Lambda_t^{-1}}^2} \\
& \hspace{15em} (\text{Lemma D.2, with prob. } 1 - \delta) \\
& \leq \frac{15\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} \sqrt{\sum_{t=1}^T \min \left\{ 1, \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{\Lambda_t^{-1}}^2 \right\}} \\
& \leq \frac{15\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} \sqrt{2d \log \left( 1 + \frac{2KT}{d\lambda} \right)} \\
& \hspace{15em} (\text{Lemma D.3}) \\
& = \mathcal{O} \left( \frac{\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} \cdot \sqrt{d \log(KT)} \right). \tag{D.7}
\end{aligned}$$

By combining Equations (D.4), (D.5), and (D.7), and setting  $\beta_{T+1}(\delta) = \mathcal{O}(B\sqrt{d \log(KT)} + B\sqrt{\lambda})$ , we derive that, with probability at least  $1 - 3\delta$  (omitting logarithmic terms and polynomial dependencies on  $B$  for brevity),

$$\text{SubOpt}(T) = \tilde{\mathcal{O}} \left( \frac{d}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} + \frac{d^2 K^2}{\kappa T} \right).$$

Substituting  $\delta \leftarrow \frac{\delta}{3}$ , we conclude the proof of Theorem 1.  $\square$

## D.2 Proofs of Lemmas for Theorem 1

### D.2.1 Proof of Lemma D.2

*Proof of Lemma D.2.* Recall the definition of  $H_t$ .

$$H_t = \sum_{s=1}^{t-1} \sum_{j=1}^{|S_s|} \nabla^2 \ell_s^{(j)}(\hat{\theta}_s^{(j+1)}) + \lambda \mathbf{I}_d \geq \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{|S_s|} \nabla^2 \ell_s^{(j)}(\hat{\theta}_s^{(j+1)}) + \lambda \mathbf{I}_d$$

Here, we can equivalently express the MNL loss at step  $j$  and round  $s$ , denoted by  $\nabla^2 \ell_s^{(j)}(\hat{\theta}_s^{(j+1)})$ , as follows:

$$\begin{aligned}
\ell_s^{(j)}(\hat{\theta}_s^{(j+1)}) &= -\log \left( \frac{\exp(\phi(x_s, \sigma_{sj})^\top \hat{\theta}_s^{(j+1)})}{\sum_{k=j}^{|S_s|} \exp(\phi(x_s, \sigma_{sk})^\top \hat{\theta}_s^{(j+1)})} \right) = -\log \left( \frac{\exp(a_{sj})}{\sum_{k=j}^{|S_s|} \exp(a_{sk})} \right) \\
&=: \bar{\ell}_s^{(j)}(\mathbf{a}_s^{(j)}), \tag{D.8}
\end{aligned}$$

where  $a_{sj} = \phi(x_s, \sigma_{sj})^\top \hat{\theta}_s^{(j+1)}$ ,  $\mathbf{a}_s^{(j)} = (a_{sk})_{k=j}^{|S_s|} \in \mathbb{R}^{|S_s|-j+1}$ . Define the matrix

$$\Phi_s^{(j)} = \begin{pmatrix} \phi(x_s, \sigma_{sj})^\top \\ \vdots \\ \phi(x_s, \sigma_{s|S_s|})^\top \end{pmatrix} \in \mathbb{R}^{(|S_s|-j+1) \times d},$$

where each row corresponds to the feature vector of an action ranked from position  $j$  to  $|S_s|$  in the ranking  $\sigma_s$ . Moreover, we define  $\mathbf{a}_{sj}^* = \phi(x_s, \sigma_{sj})^\top \theta^*$ ,  $\mathbf{a}_s^{*,(j)} = (a_{sk}^*)_{k=j}^{|S_s|} \in \mathbb{R}^{|S_s|-j+1}$ .

Then, using the  $\ell_\infty$ -norm self-concordant property of the MNL loss [41], for any  $s \in [t-1] \setminus \mathcal{T}^w$ , we obtain

$$\begin{aligned} \nabla^2 \ell_s^{(j)}(\hat{\theta}_s^{(j+1)}) &= \left( \Phi_s^{(j)} \right)^\top \nabla_{\mathbf{a}}^2 \bar{\ell}_s^{(j)}(\mathbf{a}_s^{(j)}) \Phi_s^{(j)} & (\text{Eqn. (D.8)}) \\ &\geq e^{-3\sqrt{2}\|\mathbf{a}_s^{(j)} - \mathbf{a}_s^{*,(j)}\|_\infty} \left( \Phi_s^{(j)} \right)^\top \nabla_{\mathbf{a}}^2 \bar{\ell}_s^{(j)}(\mathbf{a}_s^{*,(j)}) \Phi_s^{(j)} & (\text{Lemma F.1}) \\ &\geq \frac{1}{e} \left( \Phi_s^{(j)} \right)^\top \nabla_{\mathbf{a}}^2 \bar{\ell}_s^{(j)}(\mathbf{a}_s^{*,(j)}) \Phi_s^{(j)} & (\|\mathbf{a}_s^{(j)} - \mathbf{a}_s^{*,(j)}\|_\infty \leq \frac{1}{3\sqrt{2}}) \\ &= \frac{1}{e} \nabla^2 \ell_s^{(j)}(\theta^*), & (\text{Eqn. (D.8)}) \end{aligned}$$

where the last inequality holds because, for any  $s \in [t-1] \setminus \mathcal{T}^w$  and  $j \leq |S_s|$ , the following holds:

$$\begin{aligned} \|\mathbf{a}_s^{(j)} - \mathbf{a}_s^{*,(j)}\|_\infty &= \max_{k=j, \dots, |S_s|} \left| \phi(x_k, \sigma_{sk})^\top (\hat{\theta}_s^{(k+1)} - \theta^*) \right| \\ &\leq \max_{k=j, \dots, |S_s|} \|\phi(x_k, \sigma_{sk})\|_{H_s^{-1}} \left\| \hat{\theta}_s^{(k+1)} - \theta^* \right\|_{H_s} & (\text{Hölder's inequality}) \\ &\leq \frac{1}{3\sqrt{2}\beta_{T+1}(\delta)} \max_{k=j, \dots, |S_s|} \left\| \hat{\theta}_s^{(k+1)} - \theta^* \right\|_{H_s^{(k+1)}} & (s \notin \mathcal{T}^w, H_s \leq H_s^{(k+1)}) \\ &\leq \frac{\beta_{T+1}(\delta)}{3\sqrt{2}\beta_{T+1}(\delta)} & (\text{Corollary D.1, } \beta_t(\delta) \text{ is non-decreasing}) \\ &= \frac{1}{3\sqrt{2}}. \end{aligned}$$

Therefore, we get

$$H_t \geq \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{|S_s|} \nabla^2 \ell_s^{(j)}(\hat{\theta}_s^{(j+1)}) + \lambda \mathbf{I}_d \geq \frac{1}{e} \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{|S_s|} \nabla^2 \ell_s^{(j)}(\theta^*) + \lambda \mathbf{I}_d. \quad (\text{D.9})$$

Now, for better presentation, we define the Multinomial Logit (MNL) choice probability [46] for a given assortment  $S$  at round  $s$  as follows:

$$P_s(a|S; \theta) := \frac{\exp(\phi(x_s, a)^\top \theta)}{\sum_{a' \in S} \exp(\phi(x_s, a')^\top \theta)}, \quad \forall a \in S.$$

Let  $S_s = \{\sigma_{s1}, \dots, \sigma_{s|S_s|}\}$ . Thus, the PL model in Equation (1) can be rewritten as follows:

$$\begin{aligned} \mathbb{P}(\sigma_s | x_s, S_s; \theta) &= P_s(\sigma_{s1} | S_s; \theta) \cdot P_s(\sigma_{s2} | S_s \setminus \{\sigma_{s1}\}; \theta) \cdot \dots \cdot P_s(\sigma_{s|S_s|} | \{\sigma_{s|S_s|}\}; \theta) \\ &= \prod_{j=1}^{|S_s|} P_s(\sigma_{sj} | \{\sigma_{sj}, \dots, \sigma_{s|S_s|}\}; \theta). \end{aligned}$$

For simplicity, we define  $S_s^{(j)} := \{\sigma_{sj}, \dots, \sigma_{s|S_s|}\}$ , and let  $P_s^{(j)}$  denote the (true) MNL distribution over the remaining actions in  $S_s$  after removing the first  $j-1$  selected actions, i.e.,  $P_s^{(j)} = P_s(\cdot | S_s^{(j)}; \theta^*)$ . Then, to further lower bound the right-hand side of Equation (D.9), we

proceed as follows:

$$\begin{aligned}
\sum_{j=1}^{|S_s|} \nabla^2 \ell_s^{(j)}(\theta^*) &= \sum_{j=1}^{|S_s|} \sum_{k=j}^{|S_s|} \sum_{k'=j}^{|S_s|} \frac{\exp\left((\phi(x_s, \sigma_{sk}) + \phi(x_s, \sigma_{sk'}))^\top \theta^*\right)}{2 \left(\sum_{k'=j}^{|S_s|} \exp(\phi(x_s, \sigma_{sk'})^\top \theta^*)\right)^2} \cdot z_{skk'} z_{skk'}^\top \\
&= \frac{1}{2} \sum_{j=1}^{|S_s|} \sum_{k=j}^{|S_s|} \sum_{k'=j}^{|S_s|} P_s(\sigma_{sk} | S_s^{(j)}; \theta^*) P_s(\sigma_{sk'} | S_s^{(j)}; \theta^*) z_{skk'} z_{skk'}^\top \\
&= \frac{1}{2} \sum_{j=1}^{|S_s|} \mathbb{E}_{(a, a') \sim P_s^{(j)} \times P_s^{(j)}} \left[ (\phi(x_s, a) - \phi(x_s, a')) (\phi(x_s, a) - \phi(x_s, a'))^\top \right],
\end{aligned} \tag{D.10}$$

where  $z_{skk'} = \phi(x_s, \sigma_{sk}) - \phi(x_s, \sigma_{sk'})$ . Let the action  $\bar{a}_s \in S_s$  be ranked at position  $\bar{k}_s$  in the ranking  $\sigma_s$ . That is,

$$\sigma_s = (\underbrace{\sigma_{s1}, \dots, \sigma_{s\bar{k}_s-1}}_{\bar{k}_s-1 \text{ actions}}, \bar{a}_s, \sigma_{s\bar{k}_s+1}, \dots, \sigma_{s|S_s|-1}).$$

Note that  $\bar{a}_s \in S_s^{(j)}$  for  $j \leq \bar{k}_s$ . We also note that  $P_s^{(j)}$  is measurable with respect to the filtration  $\mathcal{F}'_{s-1, j-1} = \sigma(S_1, \sigma_{11}, \sigma_{12}, \dots, S_s, \sigma_{s1}, \dots, \sigma_{sj-1})$ . Then, by plugging Equation (D.10) into Equation (D.9) and applying the covariance matrix concentration result (Corollary F.1), since  $\lambda = \Omega(d \log(KT/\delta))$ , we have, with probability at least  $1 - \delta$ ,

$$\begin{aligned}
H_t &\geq \frac{1}{2e} \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{|S_s|} \mathbb{E}_{(a, a') \sim P_s^{(j)} \times P_s^{(j)}} \left[ (\phi(x_s, a) - \phi(x_s, a')) (\phi(x_s, a) - \phi(x_s, a'))^\top \right] + \lambda \mathbf{I}_d \\
&\geq \frac{1}{2e} \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{\bar{k}_s} \mathbb{E}_{(a, a') \sim P_s^{(j)} \times P_s^{(j)}} \left[ (\phi(x_s, a) - \phi(x_s, a')) (\phi(x_s, a) - \phi(x_s, a'))^\top \right] + \lambda \mathbf{I}_d \\
&\hspace{25em} (\bar{k}_s \leq |S_s|) \\
&\geq \frac{3}{10e} \left( \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{\bar{k}_s} (\phi(x_s, \sigma_{sj}) - \phi(x_s, \bar{a}_s)) (\phi(x_s, \sigma_{sj}) - \phi(x_s, \bar{a}_s))^\top + \lambda \mathbf{I}_d \right) \\
&\hspace{15em} (\text{Corollary F.1, } \bar{a}_s \in S_s^{(j)} \text{ for } j \leq \bar{k}_s) \\
&= \frac{3K}{10e} \left( \sum_{s \in [t-1] \setminus \mathcal{T}^w} \underbrace{\frac{1}{K} \sum_{j=1}^{\bar{k}_s} (\phi(x_s, \sigma_{sj}) - \phi(x_s, \bar{a}_s)) (\phi(x_s, \sigma_{sj}) - \phi(x_s, \bar{a}_s))^\top}_{=: X(\sigma_s)} + \lambda \mathbf{I}_d \right).
\end{aligned}$$

Here,  $\{X(\sigma_s)\}_{s \in [t-1] \setminus \mathcal{T}^w}$  is a sequence of positive semi-definite (PSD) random matrices, where each matrix  $X(\sigma_s)$  depends on the sampled ranking  $\sigma_s$ , and satisfies  $\lambda_{\max}(X(\sigma_s)) \leq 1$ .

Note that the ranking  $\sigma_s$  is drawn from the PL distribution  $\mathbb{P}(\cdot \mid x_s, S_s; \theta^*)$ , which is measurable with respect to the filtration  $\mathcal{F}_{s-1} = \sigma(S_1, \sigma_1, \dots, S_s)$ . Furthermore,  $X(\sigma_s)$  is measurable with respect to  $\sigma(\mathcal{F}_{s-1}, \sigma_s)$ . Then, by applying the concentration lemma for PSD matrices (Lemma F.4) two times, with probability at least  $1 - 2\delta$ , we get

$$\begin{aligned}
H_t &\geq \frac{3K}{10e} \left( \sum_{s \in [t-1] \setminus \mathcal{T}^w} X(\sigma_s) + \lambda \mathbf{I}_d \right) \\
&\geq \frac{K}{10e} \left( \sum_{s \in [t-1] \setminus \mathcal{T}^w} \mathbb{E}_{\sigma \sim \mathbb{P}(\cdot \mid x_s, S_s; \theta^*)} [X(\sigma)] + \lambda \mathbf{I}_d \right) \tag{Lemma F.4} \\
&\geq \frac{3K}{50e} \left( \sum_{s \in [t-1] \setminus \mathcal{T}^w} X(\tilde{\sigma}_s) + \lambda \mathbf{I}_d \right), \tag{Lemma F.4}
\end{aligned}$$

where  $\tilde{\sigma}_s$  denotes an arbitrary ranking in which  $\bar{a}_s$  is placed last. For example,  $\tilde{\sigma}_s = (\sigma_{s1}, \dots, \sigma_{s\bar{k}_s-1}, \sigma_{s\bar{k}_s+1}, \sigma_{s|S_s|-1}, \bar{a}_s)$ . Note that  $\tilde{\sigma}_s$  is a possible *virtual* ranking feedback for the assortment  $S_s$ , whereas  $\sigma_s$  denotes the *actual* ranking feedback observed at round  $s$ . Hence, since  $\bar{a}_s$  occupies the final position in the virtual sequence  $\tilde{\sigma}_s$ , it follows that:

$$\begin{aligned}
H_t &\geq \frac{3K}{50e} \left( \sum_{s \in [t-1] \setminus \mathcal{T}^w} X(\tilde{\sigma}_s) + \lambda \mathbf{I}_d \right) \\
&= \frac{3}{50e} \left( \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{|S_s|} (\phi(x_s, \tilde{\sigma}_{sj}) - \phi(x_s, \bar{a}_s)) (\phi(x_s, \tilde{\sigma}_{sj}) - \phi(x_s, \bar{a}_s))^\top + \lambda \mathbf{I}_d \right) \\
&= \frac{3}{50e} \left( \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{a \in S_s} (\phi(x_s, a) - \phi(x_s, \bar{a}_s)) (\phi(x_s, a) - \phi(x_s, \bar{a}_s))^\top + \lambda \mathbf{I}_d \right) \\
&= \frac{3}{50e} \Lambda_t \geq \frac{1}{50} \Lambda_t. \tag{Def. of  $\Lambda_t$ , Eqn. (D.2)}
\end{aligned}$$

By substituting  $\delta \leftarrow \frac{\delta}{3}$ , we conclude the proof of Lemma D.2.  $\square$

### D.2.2 Proof of Lemma D.3

*Proof of Lemma D.3.* By the definition of  $\Lambda_t$ , we have

$$\begin{aligned}
\det(\Lambda_{t+1}) &= \det \left( \Lambda_t + \sum_{a \in S_t} z_{ta} z_{ta}^\top \right) \\
&\geq \det(\Lambda_t) \left( 1 + \sum_{a \in S_t} \|z_{ta}\|_{\Lambda_t^{-1}}^2 \right) \\
&\geq \det(\lambda \mathbf{I}_d) \prod_{s=1}^t \left( 1 + \sum_{a \in S_s} \|z_{sa}\|_{\Lambda_s^{-1}}^2 \right) \\
&\geq \det(\lambda \mathbf{I}_d) \prod_{s=1}^t \left( 1 + \min \left\{ 1, \sum_{a \in S_s} \|z_{sa}\|_{\Lambda_s^{-1}}^2 \right\} \right). \tag{D.11}
\end{aligned}$$

Then, using the fact that  $a \leq 2 \log(1+a)$  for any  $a \in [0, 1]$ , we get

$$\begin{aligned}
\sum_{t=1}^T \min \left\{ 1, \sum_{a \in S_t} \|z_{ta}\|_{\Lambda_t^{-1}}^2 \right\} &\leq 2 \sum_{t=1}^T \log \left( 1 + \min \left\{ 1, \sum_{a \in S_t} \|z_{ta}\|_{\Lambda_t^{-1}}^2 \right\} \right) \\
&\leq 2 \log \left( \frac{\det(\Lambda_{T+1})}{\det(\lambda \mathbf{I}_d)} \right) \tag{Eqn. (D.11)} \\
&\leq 2d \log \left( 1 + \frac{X^2 K T}{d\lambda} \right),
\end{aligned}$$

where the last inequality holds because

$$\begin{aligned}
\det(\Lambda_{T+1}) &\leq \left( \frac{\lambda_1 + \dots + \lambda_d}{d} \right)^d \quad (\lambda_1, \dots, \lambda_d \text{ are eigenvalues of } \Lambda_{T+1}, \text{ AM-GM ineq.}) \\
&= \left( \frac{\text{trace}(\Lambda_{T+1})}{d} \right)^d \\
&= \left( \frac{\lambda d + \sum_{t=1}^T \sum_{a \in S_t} \|z_{ta}\|_2^2}{d} \right)^d \leq \left( \lambda + \frac{X^2 K T}{d} \right)^d.
\end{aligned}$$

This concludes the proof of Lemma D.3.  $\square$

### D.2.3 Proof of Lemma D.4

*Proof of Lemma D.4.* Let  $W_t := \lambda \mathbf{I}_d + \sum_{s \in \mathcal{T}_0, s < t} \sum_{a \in S_s} z_{sa} z_{sa}^\top + \lambda \mathbf{I}_d$ . Then, we have

$$\begin{aligned}
\left( \lambda + \frac{X^2 |\mathcal{T}_0| K}{d} \right)^d &\geq \left( \frac{\lambda d + \sum_{t \in \mathcal{T}_0} \sum_{a \in S_t} \|z_{ta}\|_2^2}{d} \right)^d \\
&= \left( \frac{\text{trace}(W_{T+1})}{d} \right)^d \\
&\geq \det(W_{T+1}) && \text{(AM-GM ineq.)} \\
&= \det(\lambda \mathbf{I}_d) \prod_{t \in \mathcal{T}_0} \left( 1 + \sum_{a \in S_t} \|z_{ta}\|_{W_t^{-1}}^2 \right) && \text{(update equality for det.)} \\
&\geq \det(\lambda \mathbf{I}_d) \prod_{t \in \mathcal{T}_0} \left( 1 + \sum_{a \in S_t} \|z_{ta}\|_{\Lambda_t^{-1}}^2 \right) && (W_t \leq \Lambda_t) \\
&\geq \lambda^d (1+L)^{|\mathcal{T}_0|}. && (\sum_{a \in S_t} \|z_{ta}\|_{\Lambda_t^{-1}}^2 \geq L \text{ for } t \in \mathcal{T}_0)
\end{aligned}$$

Hence, we get

$$\begin{aligned}
|\mathcal{T}_0| &\leq \frac{d}{\log(1+L)} \log \left( 1 + \frac{X^2 |\mathcal{T}_0| K}{d\lambda} \right) && \text{(D.12)} \\
&= \frac{d}{\log(1+L)} \left( \log \left( \frac{|\mathcal{T}_0|}{2d/\log(1+L)} \right) + \log \left( \frac{2d}{\log(1+L)} \left( \frac{1}{|\mathcal{T}_0|} + \frac{X^2 K}{d\lambda} \right) \right) \right) \\
&\leq \frac{|\mathcal{T}_0|}{2} + \frac{d}{\log(1+L)} \log \left( \frac{2d}{e \log(1+L)} \left( \frac{1}{|\mathcal{T}_0|} + \frac{X^2 K}{d\lambda} \right) \right),
\end{aligned}$$

which implies that

$$|\mathcal{T}_0| \leq \frac{2d}{\log(1+L)} \log \left( \frac{2d}{e \log(1+L)} \left( \frac{1}{|\mathcal{T}_0|} + \frac{X^2 K}{d\lambda} \right) \right). \quad \text{(D.13)}$$

Now, we fix  $c > 0$  and consider two cases:

- Case 1:  $|\mathcal{T}_0| < cd$

In this case, from Equation (D.12), we have  $|\mathcal{T}_0| \leq \frac{d}{\log(1+L)} \log \left( 1 + \frac{X^2 cK}{\lambda} \right)$ .

- Case 2:  $|\mathcal{T}_0| \geq cd$

In this case, from Equation (D.13), we have  $|\mathcal{T}_0| \leq \frac{2d}{\log(1+L)} \log \left( \frac{2}{e \log(1+L)} \left( \frac{1}{c} + \frac{X^2 K}{\lambda} \right) \right)$ .

By setting  $c = \frac{2}{e \log(1+L)}$ , we obtain

$$|\mathcal{T}_0| \leq \frac{2d}{\log(1+L)} \log \left( 1 + \frac{X^2 K}{\log(1+L)\lambda} \right),$$

which concludes the proof of Lemma D.4.  $\square$

### D.2.4 Proof of Lemma D.5

*Proof of Lemma D.5.* For simplicity, let  $\tilde{\mathcal{T}}_t^w = \{s \in [t-1] \mid s \in \mathcal{T}^w \cap (\mathcal{T}_0)^c\}$ . Clearly,  $\tilde{\mathcal{T}}_{T+1}^w = \mathcal{T}^w \cap (\mathcal{T}_0)^c$ . Recall that by the definition of  $H_t$ , we have

$$\begin{aligned}
H_t &= \sum_{s=1}^{t-1} \sum_{j=1}^{|S_s|} \nabla^2 \ell_s^{(j)}(\hat{\theta}_s^{(j+1)}) + \lambda \mathbf{I}_d \\
&= \sum_{s=1}^{t-1} \sum_{j=1}^{|S_s|} \sum_{k=j}^{|S_s|} \sum_{k'=j}^{|S_s|} \frac{\exp\left((\phi(x_s, \sigma_{sk}) + \phi(x_s, \sigma_{sk'}))^{\top} \hat{\theta}_s^{(j+1)}\right)}{2 \left(\sum_{k'=j}^{|S_s|} \exp\left(\phi(x_s, \sigma_{sk'})^{\top} \hat{\theta}_s^{(j+1)}\right)\right)^2} \cdot z_{skk'} z_{skk'}^{\top} + \lambda \mathbf{I}_d \\
&\geq \frac{\kappa}{2K^2} \sum_{s=1}^{t-1} \sum_{k=j}^{|S_s|} \sum_{k'=j}^{|S_s|} z_{skk'} z_{skk'}^{\top} + \lambda \mathbf{I}_d \quad (\kappa = e^{-4B}) \\
&\geq \frac{\kappa}{2K^2} \sum_{s=1}^{t-1} \sum_{a \in S_s} (\phi(x_s, a) - \phi(x_s, \bar{a}_s)) (\phi(x_s, a) - \phi(x_s, \bar{a}_s))^{\top} + \lambda \mathbf{I}_d \quad (\bar{a}_s \in S_s \text{ by Eqn. (4)}) \\
&\geq \frac{\kappa}{2K^2} \underbrace{\left( \sum_{s \in \tilde{\mathcal{T}}_t^w} \sum_{a \in S_s} (\phi(x_s, a) - \phi(x_s, \bar{a}_s)) (\phi(x_s, a) - \phi(x_s, \bar{a}_s))^{\top} + \lambda \mathbf{I}_d \right)}_{=: \Lambda_t^w}, \quad (\text{D.14})
\end{aligned}$$

where  $z_{skk'} = \phi(x_s, \sigma_{sk}) - \phi(x_s, \sigma_{sk'})$ .

Let  $\tilde{a}_t = \arg\max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}$ . Then, we get

$$\begin{aligned}
&\sum_{t \in \tilde{\mathcal{T}}_{T+1}^w} \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}^2 \\
&\leq \sum_{t \in \tilde{\mathcal{T}}_{T+1}^w} \|\phi(x_t, \tilde{a}_t) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}^2 \\
&\leq \sum_{t \in \tilde{\mathcal{T}}_{T+1}^w} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}^2 \quad (\bar{a}_t, \tilde{a}_t \in S_t \text{ by Eqn. (4)}) \\
&\leq \frac{2K^2}{\kappa} \sum_{t \in \tilde{\mathcal{T}}_{T+1}^w} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{(\Lambda_t^w)^{-1}}^2 \quad (\text{Eqn. (D.14)}) \\
&\leq \frac{2K^2}{\kappa} \sum_{t=1}^T \min \left\{ 1, \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{(\Lambda_t^w)^{-1}}^2 \right\} \quad (t \notin \mathcal{T}_0 \text{ and } \tilde{\mathcal{T}}_{T+1}^w \subseteq [T]) \\
&\leq \frac{4K^2}{\kappa} d \log \left( 1 + \frac{2KT}{d\lambda} \right). \quad (\text{Lemma D.3})
\end{aligned}$$

On the other hand, for  $t \in \tilde{\mathcal{T}}_{T+1}^w = \mathcal{T}^w \cap (\mathcal{T}_0)^c$ , we know that

$$\sum_{t \in \tilde{\mathcal{T}}_{T+1}^w} \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}^2 \geq \frac{|\tilde{\mathcal{T}}_{T+1}^w|}{3\sqrt{2}\beta_{T+1}(\delta)^2}.$$

By combining the two results above, we obtain

$$|\tilde{\mathcal{T}}_{T+1}^w| = |\mathcal{T}^w \cap (\mathcal{T}_0)^c| \leq \frac{12\sqrt{2}K^2}{\kappa} \beta_{T+1}(\delta)^2 d \log \left( 1 + \frac{2KT}{d\lambda} \right),$$

which concludes the proof.  $\square$

## E Proof for RB Loss

In this paper, we also consider the rank-breaking (RB) loss, obtained by decomposing ranking feedback into pairwise comparisons.



### E.1 Rank-Breaking (RB) Loss.

In addition to this standard approach, one can replace the full  $|S_t|$ -action ranking with its  $\binom{|S_t|}{2}$  pairwise comparisons. This technique, referred to as *rank breaking* (RB), decomposes (partial) ranking data into individual pairwise comparisons, treating each comparison as independent [6, 32, 30, 66]. Thus, the RB loss is defined as:

$$\ell_t(\theta) := \sum_{j=1}^{|S_t|-1} \sum_{k=j+1}^{|S_t|} \ell_t^{(j,k)}(\theta), \quad \text{where } \ell_t^{(j,k)}(\theta) := -\log \left( \frac{\exp(\phi(x_t, \sigma_{tj})^\top \theta)}{\sum_{m \in \{j,k\}} \exp(\phi(x_t, \sigma_{tm})^\top \theta)} \right). \quad (\text{E.1})$$

This approach is applied in the current RLHF for LLM (e.g., Ouyang et al. [57]) and is also studied in the theoretical RLHF paper [94] under the offline setting.

**OMD update for RB loss.** Similarly, for the RB loss (E.1), we estimate the underlying parameter as:

$$\hat{\theta}_t^{(j,k+1)} = \underset{\theta \in \Theta}{\operatorname{argmin}} \langle \nabla \ell_t^{(j,k)}(\hat{\theta}_t^{(j,k)}), \theta \rangle + \frac{1}{2\eta} \|\theta - \hat{\theta}_t^{(j,k)}\|_{\tilde{H}_t^{(j,k)}}^2, \quad 1 \leq j < k \leq |S_t|, \quad (\text{E.2})$$

where we set  $\hat{\theta}_t^{(j,|S_t|+1)} = \hat{\theta}_t^{(j+1,j+2)}$  for all  $j < |S_t| - 1$  and for the final pair, let  $\hat{\theta}_t^{(|S_t|-1,|S_t|+1)} = \hat{\theta}_{t+1}^{(1,2)}$ . Also, the matrix  $\tilde{H}_t^{(j,k)}$  is defined as  $\tilde{H}_t^{(j,k)} := H_t + \eta \sum_{(j',k') \leq (j,k)} \nabla^2 \ell_t^{(j',k')}(\hat{\theta}_t^{(j',k')})^2$ , where

$$H_t := \sum_{s=1}^{t-1} \sum_{j=1}^{|S_s|-1} \sum_{k=j+1}^{|S_s|} \nabla^2 \ell_s^{(j,k)}(\hat{\theta}_s^{(j,k+1)}) + \lambda \mathbf{I}_d, \quad \lambda > 0. \quad (\text{E.3})$$

The cost for the RB parameter update is  $\mathcal{O}(K^3 d^3)$ , as the parameter is updated  $\binom{|S_t|}{2}$  times per round.

**Theorem E.1.** *Under the same setting as Theorem 1, let  $\kappa := \frac{e^{-4B}}{4}$ . Then, with probability at least  $1 - \delta$ , M-AUPO (Algorithm 1) achieves the following suboptimality gap:*

$$\text{SubOpt}(T) = \tilde{\mathcal{O}} \left( \frac{d}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} + \frac{d^2}{\kappa T} \right).$$

**Discussion of Theorem E.1.** For sufficiently large  $T$ , the suboptimality gap in Theorem E.1 matches the leading-order term of Theorem 1, while its second (non-leading) term is tighter by a factor of  $\mathcal{O}(K^2)$ . However, the per-round computational cost of the RB parameter update is  $K$  times higher than that of the PL parameter update. Despite this, the result is particularly notable as it offers a rigorous theoretical explanation for the empirical success of RLHF in LLMs (e.g., Ouyang et al. [57]), where ranking feedback is decomposed into pairwise comparisons for parameter estimation.

### E.2 Main Proof of Theorem E.1

In this section, we present the proof of Theorem E.1, which is obtained by using the RB loss (E.1) instead of the PL loss (2). Note that this approach is based on the concept of *rank breaking* (RB), which decomposes (partial) ranking data into individual pairwise comparisons, treats each comparison as independent, and has been extensively studied in previous works [6, 32, 30, 66]. Moreover, this RB approach is applied in the current RLHF for LLM (e.g., Ouyang et al. [57]) and is also studied theoretically in Zhu et al. [94] under the offline setting.

**Online confidence bound for RB loss.** Now, we introduce the online confidence bound for RB loss. Since the total number of updates up to round  $t$  is  $\sum_{s=1}^t \binom{|S_s|}{2}$ , a modification of Lemma D.1 yields the following result:

**Corollary E.1** (Online confidence bound for RB loss). *Let  $\delta \in (0, 1]$ . We set  $\eta = (1 + 3\sqrt{2}B)/2$  and  $\lambda = \max\{12\sqrt{2}B\eta, 144\eta d, 2\}$ . Then, under Assumption 1, with probability at least  $1 - \delta$ , we have*

$$\|\hat{\theta}_t^{(j,k)} - \theta^*\|_{H_t^{(j,k)}} \leq \beta_t(\delta) = \mathcal{O} \left( B\sqrt{d \log(tK/\delta)} + B\sqrt{\lambda} \right), \quad \forall t \geq 1, 1 \leq j < k \leq |S_t|,$$

where  $H_t^{(j,k)} := H_t + \sum_{(j',k') \leq (j,k)} \nabla^2 \ell_t^{(j',k')}(\hat{\theta}_t^{(j',k'+1)}) + \lambda \mathbf{I}_d$  and  $\hat{\theta}_t^{(1,2)} = \hat{\theta}_t$ .

<sup>2</sup>We write  $(j', k') \leq (j, k)$  to indicate lexicographic order, i.e.,  $j' < j$  or  $j' = j$  and  $k' \leq k$ .

**Useful definitions.** We use the same or similar definitions for the set of *warm-up rounds*  $\mathcal{T}^w$  (given in Equation (D.1)), the set of *large elliptical potential (EP) rounds*  $\mathcal{T}_0$  (given in Equation (D.3)), and the regularized covariance matrix  $\Lambda_t$  (given in Equation (D.2)).

$$\begin{aligned}\mathcal{T}^w &:= \left\{ t \in [T] : \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} \geq \frac{1}{\beta_{T+1}(\delta)} \right\}, & (\text{warm-up rounds}) \\ \mathcal{T}_0 &:= \left\{ t \in [T] : \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{\Lambda_t^{-1}} \geq 1 \right\}, & (\text{large EP rounds}) \\ \Lambda_t &:= \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{a \in S_s} (\phi(x_s, a) - \phi(x_s, \bar{a}_s)) (\phi(x_s, a) - \phi(x_s, \bar{a}_s))^\top + \lambda \mathbf{I}_d.\end{aligned}$$

**Key Lemmas.** We can avoid the  $1/\kappa = \mathcal{O}(e^B)$  dependency in the leading term, thanks to the following lemma.

**Lemma E.1.** *Let  $\Lambda_t$  be defined as in Equation (D.2). Set  $\lambda = \Omega(d \log(KT/\delta))$ . Then, for all  $t \in [T]$ , with probability at least  $1 - \delta$ , we have*

$$H_t \geq \frac{1}{10} \Lambda_t.$$

The proof is deferred to Appendix E.3.1.

**Lemma E.2.** *Let  $\mathcal{T}_0 := \{t \in [T] : \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{\Lambda_t^{-1}} \geq 1\}$  and  $\mathcal{T}^w = \{t \in [T] : \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} \geq \frac{1}{\beta_{T+1}(\delta)}\}$ . Define  $\kappa := \frac{e^{-4B}}{4}$ . Then, the size of the set  $\mathcal{T}^w \cap (\mathcal{T}_0)^c$  is bounded as follows:*

$$|\mathcal{T}^w \cap (\mathcal{T}_0)^c| \leq \frac{2}{\kappa} \beta_{T+1}(\delta)^2 d \log \left( 1 + \frac{2KT}{d\lambda} \right).$$

The proof is deferred to Appendix E.3.2.

We are now ready to provide the proof of Theorem E.1.

*Proof of Theorem E.1.* The overall proof structure is similar to that of Theorem 1. We begin with Equation (D.5), but apply Lemma E.2 instead of Lemma D.5. With probability at least  $1 - \delta$ , we have

$$\begin{aligned} \mathbf{SubOpt}(T) &= \mathbb{E}_{x \sim \rho} \left[ \left( \phi(x, \pi^\star(x)) - \phi(x, \hat{\pi}_T(x)) \right)^\top \boldsymbol{\theta}^\star \right] \\ &\leq \tilde{O} \left( \frac{1}{\sqrt{T}} \right) + \frac{8B}{\log(2)T} d \log \left( 1 + \frac{2K}{\log(2)\lambda} \right) + \frac{8B}{\kappa T} \beta_{T+1}(\delta)^2 d \log \left( 1 + \frac{2KT}{d\lambda} \right) \\ &\quad \text{(Lemma E.2)} \\ &+ \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \left( \phi(x_t, \pi^\star(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)) \right)^\top \left( \boldsymbol{\theta}^\star - \hat{\boldsymbol{\theta}}_{T+1} \right). \end{aligned} \quad (\text{E.4})$$

To further bound the last term of Equation (E.4), by following the same logic from Equation (D.5) to Equation (D.6), with probability at least  $1 - \delta$ , we obtain

$$\begin{aligned} & \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \left( \phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)) \right)^\top \left( \theta^* - \hat{\theta}_{T+1} \right) \\ & \leq \frac{2\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{1}{|S_t^*|} \sum_{a \in S_t^*} \|\phi(x_t, a) - \phi(x_t, \hat{\pi}_T(x_t))\|_{H_t^{-1}} \\ & \hspace{20em} (S_t^* := \{\pi^*(x_t), \hat{\pi}_T(x_t)\}) \\ & \leq \frac{2\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{1}{|S_t|} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} . \\ & \hspace{20em} (S_t \text{ selection rule, Eqn. (4)}) \end{aligned}$$

To further bound the right-hand side, by applying the Cauchy-Schwartz inequality, we get

$$\begin{aligned}
& \frac{2\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{1}{|S_t|} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} \\
& \leq \frac{2\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \left(\frac{1}{|S_t|}\right)^2 |S_t|} \sqrt{\sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}^2} \\
& \quad \text{(Cauchy-Schwartz ineq.)} \\
& = \frac{2\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \left(\frac{1}{|S_t|}\right)^2 |S_t|} \sqrt{10 \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{\Lambda_t^{-1}}^2} \\
& \quad \text{(Lemma E.1, with probability at least } 1 - \delta) \\
& \leq \frac{2\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} \sqrt{10 \sum_{t \notin \mathcal{T}^w} \min \left\{ 1, \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{\Lambda_t^{-1}}^2 \right\}} \\
& \quad (t \notin \mathcal{T}_0 \text{ and } \mathcal{T}_0 \cup \mathcal{T}^w \subseteq [T]) \\
& \leq \frac{2\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} \sqrt{20d \log \left( 1 + \frac{2K(T - |\mathcal{T}^w|)}{d\lambda} \right)} \quad \text{(Lemma D.3)} \\
& = \mathcal{O} \left( \frac{\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} \cdot \sqrt{d \log(KT)} \right). \quad \text{(E.5)}
\end{aligned}$$

By plugging Equation (E.5) into Equation (E.4) and setting  $\beta_{T+1}(\delta) = \mathcal{O}(B\sqrt{d \log(KT)} + B\sqrt{\lambda})$ , then with probability at least  $1 - 3\delta$ , we derive that

$$\mathbf{SubOpt}(T) = \tilde{\mathcal{O}} \left( \frac{d}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} + \frac{d^2}{\kappa T} \right).$$

Substituting  $\delta \leftarrow \frac{\delta}{3}$ , we conclude the proof of Theorem E.1.  $\square$

### E.3 Proofs of Lemmas for Theorem E.1

#### E.3.1 Proof of Lemma E.1

*Proof of Lemma E.1.* Recall that, under the Bradley–Terry–Luce (BTL) model, the probability that action  $a$  is preferred over action  $a'$  is given by:

$$\mathbb{P}(a > a' | x_t, ; \boldsymbol{\theta}) = \frac{\exp(\phi(x_t, a)^\top \boldsymbol{\theta})}{\exp(\phi(x_t, a)^\top \boldsymbol{\theta}) + \exp(\phi(x_t, a')^\top \boldsymbol{\theta})} = \mu \left( (\phi(x_t, a) - \phi(x_t, a'))^\top \boldsymbol{\theta} \right).$$

Then, we can derive a lower bound on the matrix  $H_t$  as follows:

$$\begin{aligned}
H_t &= \sum_{s=1}^{t-1} \sum_{j=1}^{|S_s|-1} \sum_{k=j+1}^{|S_s|} \nabla^2 \ell_s^{(j,k)}(\hat{\theta}_s^{(j,k+1)}) + \lambda \mathbf{I}_d \\
&\geq \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{|S_s|-1} \sum_{k=j+1}^{|S_s|} \nabla^2 \ell_s^{(j,k)}(\hat{\theta}_s^{(j,k+1)}) + \lambda \mathbf{I}_d \\
&= \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{|S_s|-1} \sum_{k=j+1}^{|S_s|} \dot{\mu} \left( z_{sjk}^\top \hat{\theta}_s^{(j,k+1)} \right) z_{sjk} z_{sjk}^\top + \lambda \mathbf{I}_d \\
&\geq \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{|S_s|-1} \sum_{k=j+1}^{|S_s|} \dot{\mu} \left( z_{sjk}^\top \theta^* \right) e^{-\left| z_{sjk}^\top (\hat{\theta}_s^{(j,k+1)} - \theta^*) \right|} z_{sjk} z_{sjk}^\top + \lambda \mathbf{I}_d \quad (\text{Lemma F.2}) \\
&\geq e^{-1} \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{j=1}^{|S_s|-1} \sum_{k=j+1}^{|S_s|} \dot{\mu} \left( z_{sjk}^\top \theta^* \right) z_{sjk} z_{sjk}^\top + \lambda \mathbf{I}_d, \tag{E.6}
\end{aligned}$$

where the last inequality holds because, for any  $s \notin \mathcal{T}^w$ , the following property is satisfied:

$$\begin{aligned}
\left| z_{sjk}^\top (\hat{\theta}_s^{(j,k+1)} - \theta^*) \right| &= \left| (\phi(x_s, \sigma_{sj}) - \phi(x_s, \sigma_{sk}))^\top (\hat{\theta}_s^{(j,k+1)} - \theta^*) \right| \\
&\leq \|\phi(x_s, \sigma_{sj}) - \phi(x_s, \sigma_{sk})\|_{H_s^{-1}} \|\hat{\theta}_s^{(j,k+1)} - \theta^*\|_{H_s} \quad (\text{H\"older's inequality}) \\
&\leq \frac{1}{\beta_{T+1}(\delta)} \|\hat{\theta}_s^{(j,k+1)} - \theta^*\|_{H_s^{(j,k+1)}} \quad (s \neq \mathcal{T}^w, H_s \leq H_s^{(j,k+1)}) \\
&\leq \frac{\beta_t(\delta)}{\beta_{T+1}(\delta)} \quad (\text{Corollary E.1}) \\
&\leq 1. \quad (\beta_t(\delta) \text{ is non-decreasing})
\end{aligned}$$

For simplicity, we write  $\mathbb{P}_s(a > a') = \mathbb{P}(a > a' | x_s; \theta^*)$ . Let  $P_{s, \{a, a'\}}$  denote the Bernoulli distribution over the support  $\{a, a'\}$ , where  $a$  occurs with probability  $\mu((\phi(x_s, a) - \phi(x_s, a'))^\top \theta^*)$ . Then, to further lower bound the right-hand side of Equation (E.6), we proceed as follows:

$$\begin{aligned}
&\sum_{j=1}^{|S_s|-1} \sum_{k=j+1}^{|S_s|} \dot{\mu} \left( z_{sjk}^\top \theta^* \right) z_{sjk} z_{sjk}^\top \\
&= \sum_{j=1}^{|S_s|-1} \sum_{k=j+1}^{|S_s|} \mu \left( z_{sjk}^\top \theta^* \right) \mu \left( z_{skj}^\top \theta^* \right) z_{sjk} z_{sjk}^\top \\
&= \frac{1}{2} \sum_{a \in S_s} \sum_{a' \in S_s} \mathbb{P}_s(a > a') \mathbb{P}_s(a' > a) (\phi(x_s, a) - \phi(x_s, a')) (\phi(x_s, a) - \phi(x_s, a'))^\top \\
&\geq \frac{1}{2} \sum_{a \in S_s} 2\mathbb{P}_s(a > \bar{a}_s) \mathbb{P}_s(\bar{a}_s > a) (\phi(x_s, a) - \phi(x_s, \bar{a}_s)) (\phi(x_s, a) - \phi(x_s, \bar{a}_s))^\top \\
&\quad (\bar{a}_s \in S_s \text{ by Eqn. (4)}) \\
&= \frac{1}{2} \sum_{a \in S_s} \mathbb{E}_{(a', a'') \sim P_{s, \{a, \bar{a}_s\}}^{\otimes 2}} \left[ (\phi(x_s, a') - \phi(x_s, a'')) (\phi(x_s, a') - \phi(x_s, a''))^\top \right], \tag{E.7}
\end{aligned}$$

where  $P_{s, \{a, \bar{a}_s\}}^{\otimes 2} = P_{s, \{a, \bar{a}_s\}} \times P_{s, \{a, \bar{a}_s\}}$  denotes the the product distribution over two independent samples from  $P_{s, \{a, \bar{a}_s\}}$ . Note that the Bernoulli distribution  $P_{s, \{a, \bar{a}_s\}}$ , where  $a \in S_s$ , is measurable with respect to the filtration  $\mathcal{F}_{s-1} = \sigma(S_1, \sigma_1, \dots, S_{s-1}, \sigma_{s-1}, S_s)$ . Then, plugging Equation (E.7)

into Equation (E.6), we get

$$\begin{aligned}
H_t &\geq \frac{e^{-1}}{2} \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{a \in S_s} \mathbb{E}_{(a', a'') \sim P_{s, \{a, \bar{a}_s\}}^{\otimes 2}} \left[ (\phi(x_s, a') - \phi(x_s, a'')) (\phi(x_s, a') - \phi(x_s, a''))^\top \right] + \lambda \mathbf{I}_d \\
&\geq \frac{3e^{-1}}{10} \underbrace{\left( \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{a \in S_s} (\phi(x_s, a) - \phi(x_s, \bar{a}_s)) (\phi(x_s, a) - \phi(x_s, \bar{a}_s))^\top + \lambda \mathbf{I}_d \right)}_{=\Lambda_t} \\
&\quad \text{(covariance matrix concentration lemma (Corollary F.1))} \\
&\geq \frac{1}{10} \Lambda_t,
\end{aligned}$$

which conclude the proof of Lemma E.1.  $\square$

### E.3.2 Proof of Lemma E.2

*Proof of Lemma E.2.* For simplicity, let  $\tilde{\mathcal{T}}_t^w = \{s \in [t-1] \mid s \in \mathcal{T}^w \cap (\mathcal{T}_0)^c\}$ . Clearly,  $\tilde{\mathcal{T}}_{T+1}^w = \mathcal{T}^w \cap (\mathcal{T}_0)^c$ . Recall that by the definition of  $H_t$ , we have

$$\begin{aligned}
H_t &= \sum_{s=1}^{t-1} \sum_{j=1}^{|S_s|-1} \sum_{k=j+1}^{|S_s|} \nabla^2 \ell_s^{(j,k)}(\hat{\theta}_s^{(j,k+1)}) + \lambda \mathbf{I}_d \\
&\geq \kappa \sum_{s=1}^{t-1} \sum_{j=1}^{|S_s|-1} \sum_{k=j+1}^{|S_s|} z_{sjk} z_{sjk}^\top + \lambda \mathbf{I}_d \quad (\kappa = e^{-4B}/4) \\
&\geq \kappa \sum_{s=1}^{t-1} \sum_{a \in S_s} (\phi(x_s, a) - \phi(x_s, \bar{a}_s)) (\phi(x_s, a) - \phi(x_s, \bar{a}_s))^\top + \lambda \mathbf{I}_d \quad (\bar{a}_s \in S_s) \\
&\geq \kappa \underbrace{\left( \sum_{s \in \tilde{\mathcal{T}}_t^w} \sum_{a \in S_s} (\phi(x_s, a) - \phi(x_s, \bar{a}_s)) (\phi(x_s, a) - \phi(x_s, \bar{a}_s))^\top + \lambda \mathbf{I}_d \right)}_{=: \Lambda_t^w}. \quad (\text{E.8})
\end{aligned}$$

Let  $\tilde{a}_t = \arg\max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}$ . Then, we get

$$\begin{aligned}
&\sum_{t \in \tilde{\mathcal{T}}_{T+1}^w} \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}^2 \\
&\leq \sum_{t \in \tilde{\mathcal{T}}_{T+1}^w} \|\phi(x_t, \tilde{a}_t) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}^2 \\
&\leq \sum_{t \in \tilde{\mathcal{T}}_{T+1}^w} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}^2 \quad (\tilde{a}_t, \bar{a}_t \in S_t \text{ by Eqn. (4)}) \\
&\leq \frac{1}{\kappa} \sum_{t \in \tilde{\mathcal{T}}_{T+1}^w} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{(\Lambda_t^w)^{-1}}^2 \quad (\text{Eqn. (E.8)}) \\
&\leq \frac{1}{\kappa} \sum_{t=1}^T \min \left\{ 1, \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{(\Lambda_t^w)^{-1}}^2 \right\} \quad (t \neq \mathcal{T}_0 \text{ and } \tilde{\mathcal{T}}_{T+1}^w \subseteq [T]) \\
&\leq \frac{2}{\kappa} d \log \left( 1 + \frac{2KT}{d\lambda} \right). \quad (\text{Lemma D.3})
\end{aligned}$$

On the other hand, for  $t \in \tilde{\mathcal{T}}_{T+1}^w = \mathcal{T}^w \cap (\mathcal{T}_0)^c$ , we know that

$$\sum_{t \in \tilde{\mathcal{T}}_{T+1}^w} \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}^2 \geq \frac{|\tilde{\mathcal{T}}_{T+1}^w|}{\beta_{T+1}(\delta)^2}.$$

By combining the two results above, we get

$$|\tilde{\mathcal{T}}_{T+1}^w| \leq \frac{2}{\kappa} \beta_{T+1}(\delta)^2 d \log \left( 1 + \frac{2KT}{d\lambda} \right),$$

which concludes the proof.  $\square$

## F Technical Lemmas

**Lemma F.1** (Proposition B.5 of Lee and Oh 41). *The Hessian of the multinomial logistic loss  $\bar{\ell} : \mathbb{R}^M \rightarrow \mathbb{R}$  satisfies that, for any  $\mathbf{a}_1, \mathbf{a}_2 \in \mathbb{R}^M$ , we have:*

$$e^{-3\sqrt{2}\|\mathbf{a}_1 - \mathbf{a}_2\|_\infty} \nabla^2 \bar{\ell}(\mathbf{a}_1) \leq \nabla^2 \bar{\ell}(\mathbf{a}_2) \leq e^{3\sqrt{2}\|\mathbf{a}_1 - \mathbf{a}_2\|_\infty} \nabla^2 \bar{\ell}(\mathbf{a}_1).$$

**Lemma F.2** (Lemma 9 of Abeille et al. 2). *Let  $f$  be a strictly increasing function such that  $|\ddot{f}| \leq \dot{f}$ , and let  $\mathcal{Z}$  be any bounded interval of  $\mathbb{R}$ . Then, for all  $z_1, z_2 \in \mathcal{Z}$ , we have*

$$\dot{f}(z_2) \exp(-|z_2 - z_1|) \leq \dot{f}(z_1) \leq \dot{f}(z_2) \exp(|z_2 - z_1|).$$

**Lemma F.3** (Concentration of covariances, Lemma 39 of Zanette et al. 90). *Let  $\mu_i$  be the conditional distribution of  $\phi \in \mathbb{R}^d$  given the sampled  $\phi_1, \dots, \phi_{i-1}$ . Assume  $\|\phi\|_2 \leq 1$ . Define  $\Sigma = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\phi \sim \mu_i} \phi \phi^\top$ . If  $\lambda = \Omega(d \log(n/\delta))$ , then, with probability at least  $1 - \delta$ , for any  $n \geq 1$ , we have*

$$\frac{1}{3} (n\Sigma + \lambda \mathbf{I}_d) \leq \sum_{i=1}^n \phi_i \phi_i^\top + \lambda \mathbf{I}_d \leq \frac{5}{3} (n\Sigma + \lambda \mathbf{I}_d).$$

We now extend the previous sequence  $\phi_1, \dots, \phi_{i-1}$  to a general setting where conditioning is performed with respect to an arbitrary filtration  $\mathcal{F}_{i-1}$ . For instance,  $\mathcal{F}_{i-1}$  may be the  $\sigma$ -algebra generated by previous samples  $\tilde{\phi}_1, \dots, \tilde{\phi}_{i-1}$ , where each  $\tilde{\phi}_j$  is drawn from a potentially different distribution  $\tilde{\mu}_j$ . This generalization is valid because a martingale difference sequence can be defined with respect to any filtration provided that the  $\sigma$ -algebras satisfy the usual properties (e.g., nestedness and making  $\{\phi_i\}$  adapted).

**Corollary F.1** (Generalized version of covariance concentration). *Let  $\mu_i$  denote the conditional distribution of  $\phi \in \mathbb{R}^d$  conditioned on the filtration  $\mathcal{F}_{i-1}$ . Assume  $\|\phi\|_2 \leq 1$ . Define  $\Sigma = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\phi \sim \mu_i} \phi \phi^\top$ . If  $\lambda = \Omega(d \log(n/\delta))$ , then with probability at least  $1 - \delta$ , for any  $n \geq 1$ , we have*

$$\frac{1}{3} (n\Sigma + \lambda \mathbf{I}_d) \leq \sum_{i=1}^n \phi_i \phi_i^\top + \lambda \mathbf{I}_d \leq \frac{5}{3} (n\Sigma + \lambda \mathbf{I}_d).$$

We also provide a concentration lemma for the more general case of positive semi-definite (PSD) random matrices.

**Lemma F.4** (Concentration of PSD matrices). *Let  $\mu_i$  denote the conditional distribution of a positive semi-definite  $M \in \mathbb{R}^{d \times d}$  conditioned on the filtration  $\mathcal{F}_{i-1}$ . Assume  $\lambda_{\max}(M) \leq 1$ . Define  $\bar{M} := \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{M \sim \mu_i} M$ . If  $\lambda = \Omega(d \log(n/\delta))$ , then with probability at least  $1 - \delta$ , for any  $n \geq 1$ , we have*

$$\frac{1}{3} (n\bar{M} + \lambda \mathbf{I}_d) \leq \sum_{i=1}^n M_i + \lambda \mathbf{I}_d \leq \frac{5}{3} (n\bar{M} + \lambda \mathbf{I}_d).$$

*Proof of Lemma F.4.* The overall structure of the proof closely follows that of Lemma 39 in Zanette et al. 90. For completeness, we provide the full proof below.

Fix  $x \in \mathbb{R}^d$  such that  $\|x\|_2 = 1$ . Let  $\bar{M}_i = \mathbb{E}_{M \sim \mu_i} M$  and  $\bar{M} = \frac{1}{n} \sum_{i=1}^n \bar{M}_i$ . Then, we have

$$\mathbb{E}_{M \sim \mu_i} x^\top M x = x^\top \mathbb{E}_{M \sim \mu_i} M x = x^\top \bar{M}_i x.$$

Since  $M$  is a positive semi-definite matrix, the random variable  $x^\top Mx$  is non-negative with the maximum value  $x^\top Mx \leq \lambda_{\max}(M)\|x\|_2^2 \leq 1$ . Thus, the conditional variance is at most  $x^\top \bar{M}_i x$  because

$$\text{Var}_{M \sim \mu_i}(x^\top Mx) \leq \mathbb{E}_{M \sim \mu_i}(x^\top Mx)^2 \leq \mathbb{E}_{M \sim \mu_i} x^\top Mx = x^\top \bar{M}_i x.$$

Then, by Lemma F.5, with probability at least  $1 - \delta$ , for some constant  $c$ , we have

$$\left| \frac{1}{n} \sum_{i=1}^n (x^\top M_i x - x^\top \bar{M}_i x) \right| = \left| \frac{1}{n} \sum_{i=1}^n x^\top M_i x - x^\top \bar{M} x \right| \leq c \left( \sqrt{2 \frac{x^\top \bar{M} x}{n} \log(2/\delta)} + \frac{\log(2/\delta)}{3n} \right).$$

Now, we will show that if  $\lambda = \Omega(\log(1/\delta))$ , we can derive

$$c \left( \sqrt{2 \frac{x^\top \bar{M} x}{n} \log(2/\delta)} + \frac{\log(2/\delta)}{3n} \right) \leq \frac{1}{2} \left( x^\top \bar{M} x + \frac{\lambda}{n} \right). \quad (\text{F.1})$$

**Case 1.**  $x^\top \bar{M} x \leq \frac{\lambda}{n}$ .

In this case, it is sufficient to satisfy for some constants  $c', c''$

$$\begin{aligned} \sqrt{2 \frac{\log(2/\delta)}{n}} &\leq c' \sqrt{\frac{\lambda}{n}} \iff \Omega(\log(1/\delta)) \leq \lambda \\ \frac{\log(2/\delta)}{3n} &\leq c'' \left( \frac{\lambda}{n} \right) \iff \Omega(\log(1/\delta)) \leq \lambda. \end{aligned}$$

**Case 2.**  $x^\top \bar{M} x > \frac{\lambda}{n}$ .

In this case, it is sufficient to satisfy for some constants  $c''', c'''$

$$\begin{aligned} \sqrt{2 \frac{x^\top \bar{M} x}{n} \log(2/\delta)} &\leq c''' \left( \frac{\lambda}{n} \right) \iff \Omega(\log(1/\delta)) \leq \lambda \\ \frac{\log(2/\delta)}{3n} &\leq c''' \left( \frac{\lambda}{n} \right) \iff \Omega(\log(1/\delta)) \leq \lambda. \end{aligned}$$

Therefore, Equation (F.1) is satisfied. Since  $\|x\|_2 \leq 1$ , this implies

$$\left| x^\top \left( \frac{1}{n} \sum_{i=1}^n M_i - \bar{M} \right) x \right| \leq \frac{1}{2} x^\top \left( \bar{M} + \frac{\lambda}{n} \mathbf{I}_d \right) x. \quad (\text{F.2})$$

We denote the boundary of the unit ball by  $\partial\mathcal{B} = \{\|x\|_2 = 1\}$ . Then, for any  $x \in \partial\mathcal{B}$ , we know there exists a  $x'$  in the  $\epsilon$ -covering such that  $\|x - x'\|_2 \leq \epsilon$ . Let  $\mathcal{N}_\epsilon$  be the  $\epsilon$ -covering number of  $\partial\mathcal{B}$ . Then, by the covering number of Euclidean ball lemma (Lemma F.6), we get

$$\mathcal{N}_\epsilon \leq \left( \frac{3}{\epsilon} \right)^d. \quad (\text{F.3})$$

Taking a union bound over  $x'$  and the number of samples  $n$ , with probability at least  $1 - n\mathcal{N}_\epsilon\delta$ , we obtain

$$\begin{aligned} \left| x^\top \left( \frac{1}{n} \sum_{i=1}^n M_i - \bar{M} \right) x \right| &\leq \left| (x')^\top \left( \frac{1}{n} \sum_{i=1}^n M_i - \bar{M} \right) x' \right| + \left| (x - x')^\top \left( \frac{1}{n} \sum_{i=1}^n M_i - \bar{M} \right) x' \right| \\ &\quad + \left| (x')^\top \left( \frac{1}{n} \sum_{i=1}^n M_i - \bar{M} \right) (x - x') \right| \\ &\leq \left| (x')^\top \left( \frac{1}{n} \sum_{i=1}^n M_i - \bar{M} \right) x' \right| + 4\epsilon. \\ &\quad (\|x - x'\|_2 \leq \epsilon \text{ and } M_i, \|\bar{M}\|_2 \leq 1) \\ &\leq \frac{1}{2} (x')^\top \left( \bar{M} + \frac{\lambda}{n} \mathbf{I}_d \right) x' + 4\epsilon \quad (\text{Eqn. (F.2)}) \\ &\leq \frac{1}{2} x^\top \left( \bar{M} + \frac{\lambda}{n} \mathbf{I}_d \right) x + \frac{9}{2} \epsilon \quad (\|x - x'\|_2 \leq \epsilon \text{ and } \|\bar{M}\|_2 \leq 1) \\ &\leq \frac{2}{3} x^\top \left( \bar{M} + \frac{\lambda}{n} \mathbf{I}_d \right) x, \quad (\text{set } \epsilon = \mathcal{O}(\frac{1}{n})) \end{aligned}$$

where  $\lambda = \Omega\left(\log\left(\frac{2n\mathcal{N}_\epsilon}{\delta}\right)\right)$ . By substituting  $\delta \leftarrow \delta/(n\mathcal{N}_\epsilon + 1)$  and combining this with Equation (F.3), we obtain:

$$\frac{1}{3} \left( \bar{M} + \frac{\lambda}{n} \mathbf{I}_d \right) \leq \frac{1}{n} \sum_{i=1}^n M_i + \frac{\lambda}{n} \mathbf{I}_d \leq \frac{5}{3} \left( \bar{M} + \frac{\lambda}{n} \mathbf{I}_d \right),$$

which concludes the proof.  $\square$

**Lemma F.5** (Bernstein for martingales, Theorem 1 of Beygelzimer et al. 10 and Lemma 45 of Zanette et al. 90). *Consider the stochastic process  $\{X_n\}$  adapted to the filtration  $\{\mathcal{F}_n\}$ . Assume  $\mathbb{E}X_n = 0$  and  $cX_n \leq 1$  for every  $n$ ; then for every constant  $z \neq 0$  it holds that*

$$\Pr \left( \sum_{n=1}^N X_n \leq z \sum_{n=1}^N \mathbb{E}(X_n^2 | \mathcal{F}_n) + \frac{1}{z} \log \frac{1}{\delta} \right) \geq 1 - \delta.$$

By optimizing the bound as a function of  $z$ , we also have

$$\Pr \left( \sum_{n=1}^N X_n \leq c \sqrt{\sum_{n=1}^N \mathbb{E}(X_n^2 | \mathcal{F}_n) \log \frac{1}{\delta}} + \log \frac{1}{\delta} \right) \geq 1 - \delta.$$

**Lemma F.6** (Covering number of Euclidean ball). *For any  $\epsilon > 0$ , the  $\epsilon$ -covering number of the Euclidean ball in  $\mathbb{R}^d$  with radius  $R > 0$  is upper bounded by  $(1 + 2R/\epsilon)^2$ .*

## G Proof of Theorem 2

### G.1 Main Proof of Theorem 2

Throughout the proof, we consider the setting where the context space is a singleton, i.e.,  $\mathcal{X} = \{x\}$ . As a result, the problem reduces to a context-free setting, and we focus solely on the action space  $\mathcal{A}$ . Note that this is equivalent to assuming that  $\rho$  is a Dirac distribution.

We first present the following theorem, which serves as the foundation for our analysis.

**Theorem G.1** (Lower bound on adaptive PL model parameter estimation). *Let  $\Phi = \mathcal{S}^{d-1}$  be the unit sphere in  $\mathbb{R}^d$ , and let  $\Theta = \{-\mu, \mu\}^d$  for some  $\mu \in (0, 1/\sqrt{d}]$ . We consider a query model where, at each round  $t = 1, \dots, T$ , the learner selects a subset  $S_t \subseteq \Phi$  of feature vectors, with cardinality satisfying  $2 \leq |S_t| \leq K$ , and then receives a ranking feedback  $\sigma_t$  drawn from the Plackett–Luce (PL) model defined as:*

$$\mathbb{P}(\sigma_t | S_t; \boldsymbol{\theta}) = \prod_{j=1}^{|S_t|} \frac{\exp(\phi_{\sigma_{tj}}^\top \boldsymbol{\theta})}{\sum_{k=j}^{|S_t|} \exp(\phi_{\sigma_{tk}}^\top \boldsymbol{\theta})},$$

where  $\sigma_t = (\sigma_{t1}, \dots, \sigma_{t|S_t|})$  is a permutation of the actions in  $S_t$ ,  $\phi_a \in \Phi$  denotes the feature vector associated with action  $a \in \mathcal{A}$  in the selected subset at round  $t$ , and  $\boldsymbol{\theta} \in \Theta$ . Then, we have

$$\inf_{\hat{\boldsymbol{\theta}}, \pi} \max_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_{\boldsymbol{\theta}} \left[ \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|_2^2 \right] \geq \frac{d\mu^2}{2} \left( 1 - \sqrt{\frac{2K^2T\mu^2}{d}} \right),$$

where the infimum is over all measurable estimators  $\hat{\boldsymbol{\theta}}$  and measurable (but possibly adaptive) query rules  $\pi$ , and  $\mathbb{E}_{\boldsymbol{\theta}}[\cdot]$  denotes the expectation over the randomness in the observations and decision rules if  $\boldsymbol{\theta}$  is the true instance. In particular, if  $T \geq \frac{d^2}{8K^2}$ , by choosing  $\mu = \sqrt{d/(8K^2T)}$ , we obtain

$$\inf_{\hat{\boldsymbol{\theta}}, \pi} \max_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_{\boldsymbol{\theta}} \left[ \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|_2^2 \right] \geq \frac{d^2}{32K^2T}.$$

*Proof of Theorem G.1.* The analysis of this result closely follows the proof of Theorem 3 in Shamir [73]. The key distinction lies in the input structure: our setting involves a set of feature vectors, while theirs is restricted to a single feature vector.



To begin with, since the worst-case expected regret with respect to  $\theta$  can be lower bounded by the average regret under the uniform prior over  $\Theta$ , we have:

$$\begin{aligned} \max_{\theta \in \Theta} \mathbb{E}_{\theta} \left[ \|\theta - \hat{\theta}\|_2^2 \right] &\geq \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\theta} \left[ \|\theta - \hat{\theta}\|_2^2 \right] \\ &= \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\theta} \left[ \sum_{i=1}^d \left( \theta_i - \hat{\theta}_i \right)^2 \right] \\ &\geq \mu^2 \cdot \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\theta} \left[ \sum_{i=1}^d \mathbb{I} \left\{ \theta_i \hat{\theta}_i < 0 \right\} \right]. \end{aligned} \quad (\text{G.1})$$

As in Shamir [73], we assume that the query strategy is deterministic conditioned on the past: that is,  $S_t$  is a deterministic function of the previous queries and observations, i.e.,  $S_1, \sigma_1, \dots, S_{t-1}, \sigma_{t-1}$ . This assumption is made without loss of generality, since any randomized querying strategy can be viewed as a distribution over deterministic strategies. Therefore, a lower bound that holds uniformly for all deterministic strategies also applies to any randomized strategy. Then, we use the following lemma.

**Lemma G.1** (Lemma 4 of Shamir 73). *Let  $\theta$  be a random vector, none of whose coordinates is supported on 0, and let  $y_1, y_2, \dots, y_T$  be a sequence of queries obtained by a deterministic strategy returning a point  $\hat{\theta}$  (that is,  $\psi_t$  is a deterministic function of  $\psi_1, y_1, \dots, \psi_{t-1}, y_{t-1}$ , and  $\hat{\theta}$  is a deterministic function of  $y_1, \dots, y_T$ ). Then, we have*

$$\mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\theta} \left[ \sum_{i=1}^d \mathbb{I} \left\{ \theta_i \hat{\theta}_i < 0 \right\} \right] \geq \frac{d}{2} \left( 1 - \sqrt{\frac{1}{d} \sum_{i=1}^d \sum_{t=1}^T U_{ti}} \right),$$

where

$$U_{ti} := \sup_{\theta_j, j \neq i} D_{\text{KL}} \left( P(y_t | \theta_i > 0, \{\theta_j\}_{j \neq i}, \{y_s\}_{s=1}^{t-1}) \parallel P(y_t | \theta_i < 0, \{\theta_j\}_{j \neq i}, \{y_s\}_{s=1}^{t-1}) \right).$$

In our setting, we interpret  $y_t = \sigma_t$ , and  $\psi_t = \{\phi_a\}_{a \in S_t} \subseteq \Phi$ . Then, we can write  $U_{ti}$  as follows:

$$U_{ti} = \sup_{\theta_j, j \neq i} D_{\text{KL}} \left( \mathbb{P}(\sigma_t | S_t; \theta_i > 0, \{\theta_j\}_{j \neq i},) \parallel \mathbb{P}(\sigma_t | S_t; \theta_i < 0, \{\theta_j\}_{j \neq i},) \right).$$

For simplicity, let  $\mathbb{P}_{\theta}(\sigma | S) = \mathbb{P}(\sigma | S; \theta)$ . Then, we can upper bound  $U_{ti}$  using the following lemma.

**Lemma G.2.** *For any  $\theta, \theta' \in \mathbb{R}^d$ , let  $\mathbb{P}_{\theta}(\cdot | S)$  denote the PL distribution over rankings induced by the action set  $S$  and parameter vector  $\theta$ . Then, we have*

$$D_{\text{KL}}(\mathbb{P}_{\theta}(\cdot | S) \parallel \mathbb{P}_{\theta'}(\cdot | S)) \leq \frac{K}{2} \sum_{a \in S} (\phi_a^{\top} (\theta' - \theta))^2.$$

The proof is deferred to Appendix G.2.1.

By applying Lemma G.2, we have

$$\begin{aligned} \sum_{i=1}^d U_{ti} &\leq \frac{K}{2} \sum_{i=1}^d \sum_{a \in S_t} (2\mu \cdot [\phi_a]_i)^2 = 2K\mu^2 \sum_{a \in S_t} \underbrace{\sum_{i=1}^d ([\phi_a]_i)^2}_{=1} \\ &= 2K\mu^2 \cdot |S_t| \quad (\phi_a \in \mathcal{S}^{d-1}) \\ &\leq 2K^2\mu^2. \quad (|S_t| \leq K) \end{aligned}$$

Hence, by Lemma G.1, we get

$$\begin{aligned} \mathbb{E}_{\theta \sim \text{Unif}(\Theta)} \mathbb{E}_{\theta} \left[ \sum_{i=1}^d \mathbb{I} \left\{ \theta_i \hat{\theta}_i < 0 \right\} \right] &\geq \frac{d}{2} \left( 1 - \sqrt{\frac{1}{d} \sum_{i=1}^d \sum_{t=1}^T U_{ti}} \right) \\ &\geq \frac{d}{2} \left( 1 - \sqrt{\frac{2K^2 T \mu^2}{d}} \right). \end{aligned} \quad (\text{G.2})$$

Combining Equation (G.1) and (G.2), we prove the first inequality of Theorem G.1. The second inequality directly follows by choosing  $\mu = \sqrt{d/(8K^2 T)}$ .  $\square$

We are now ready to present the proof of Theorem 2.

*Proof of Theorem 2.* The structure of our proof is similar to that of Theorem 2 in Wagenmaker et al. [79]. However, while they consider the linear bandit setting, we focus on the Plackett–Luce (PL) bandit setting.

We adopt the same instance construction as in Theorem G.1, where  $\Phi = \mathcal{S}^{d-1}$  and  $\Theta = \{-\mu, \mu\}^d$ . Define  $\phi^*(\theta) = \operatorname{argmax}_{a \in \mathcal{A}} \phi_a^\top \theta$ . Then, since  $\phi^*(\theta) \in \Phi$  and  $\theta \in \Theta$ , it is clear that

$$\phi^*(\theta) = \theta / \|\theta\|_2 = \theta / (\sqrt{d}\mu), \quad \text{and} \quad \phi^*(\theta)^\top \theta = \sqrt{d}\mu. \quad (\text{G.3})$$

Fix the suboptimality gap  $\epsilon > 0$ . By definition, a policy  $\pi \in \Delta_\Phi$  is said to be  $\epsilon$ -optimal if it satisfies

$$\mathbb{E}_{\phi \sim \pi} [\phi^\top \theta] = \underbrace{(\mathbb{E}_{\phi \sim \pi} [\phi])^\top}_{=: \phi_\pi} \theta \geq \phi^*(\theta)^\top \theta - \epsilon = \sqrt{d}\mu - \epsilon. \quad (\text{G.4})$$

Moreover, by Jensen's inequality, we have

$$\|\phi_\pi\|_2^2 \leq \mathbb{E}_{\phi \sim \pi} [\|\phi\|_2^2] = 1.$$

Let  $\Delta = \phi_\pi - \phi^*(\theta)$ . Then, we get

$$\begin{aligned} 1 &\geq \|\phi_\pi\|_2^2 = \|\phi^*(\theta) + \Delta\|_2^2 = 1 + \|\Delta\|_2^2 + 2\phi^*(\theta)^\top \Delta \\ &\iff \phi^*(\theta)^\top \Delta \leq -\frac{1}{2}\|\Delta\|_2^2 \\ &\iff \theta^\top \Delta \leq -\frac{\sqrt{d}\mu}{2}\|\Delta\|_2^2. \end{aligned} \quad (\text{Eqn. (G.3)})$$

Hence, if a policy  $\pi$  is  $\epsilon$ -optimal for a parameter  $\theta$ , then the following bound holds:

$$\begin{aligned} -\epsilon &\leq -\frac{\sqrt{d}\mu}{2}\|\Delta\|_2^2. \\ &\iff \|\Delta\|_2^2 \leq \frac{2\epsilon}{\sqrt{d}\mu}, \quad \text{where } \theta = \sqrt{d}\mu(\phi_\pi - \Delta). \end{aligned} \quad (\text{Eqn. (G.4)})$$

We now assume that we are given an  $\epsilon$ -optimal policy  $\hat{\pi}$ . Define  $\hat{\phi} := \phi_{\hat{\pi}}$  and the following estimator

$$\hat{\theta} = \begin{cases} \theta' & \text{if } \exists \theta' \in \Theta \text{ with } \theta' = \sqrt{d}\mu(\hat{\phi} - \Delta') \text{ for some } \Delta' \in \mathbb{R}^d, \|\Delta'\|_2^2 \leq \frac{2\epsilon}{\sqrt{d}\mu}; \\ \text{any } \theta' \in \Theta & \text{otherwise.} \end{cases}$$

If  $\hat{\pi}$  is indeed  $\epsilon$ -optimal for some  $\theta \in \Theta$ , then the first condition is satisfied, and we have:

$$\|\hat{\theta} - \theta\|_2 = \|\sqrt{d}\mu(\hat{\phi} - \Delta') - \sqrt{d}\mu(\hat{\phi} - \Delta)\|_2 \leq 2\sqrt{d}\mu \sqrt{\frac{2\epsilon}{\sqrt{d}\mu}} = \sqrt{8\sqrt{d}\mu\epsilon}. \quad (\text{G.5})$$

We denote  $\mathcal{E}$  as the event that  $\hat{\pi}$  is  $\epsilon$ -optimal for  $\theta \in \Theta$ . Then, we get

$$\begin{aligned} \mathbb{E}_\theta [\|\hat{\theta} - \theta\|_2^2] &= \mathbb{E}_\theta [\|\hat{\theta} - \theta\|_2^2 \cdot \mathbb{I}\{\mathcal{E}\} + \|\hat{\theta} - \theta\|_2^2 \cdot \mathbb{I}\{\mathcal{E}^c\}] \\ &\leq 8\sqrt{d}\mu\epsilon + \mathbb{E}_\theta [\|\hat{\theta} - \theta\|_2^2 \cdot \mathbb{I}\{\mathcal{E}^c\}] \\ &\leq 8\sqrt{d}\mu\epsilon + 2d\mu^2 \cdot P_\theta[\mathcal{E}^c]. \end{aligned} \quad (\text{Eqn. (G.5)})$$

(max $\{\|\hat{\theta}\|_2^2, \|\theta\|_2^2\} \leq d\mu^2$ )

On the other hand, by Theorem G.1, there exists a parameter  $\theta \in \Theta$  such that, if we collect  $T$  samples and set  $\mu = \sqrt{d/(8K^2T)}$ , then the following lower bound holds:

$$\mathbb{E}_\theta [\|\hat{\theta} - \theta\|_2^2] \geq \frac{d^2}{32K^2T}.$$

To satisfy both inequalities, we require:

$$\begin{aligned} \frac{2\sqrt{2}d\epsilon}{\sqrt{K^2T}} + \frac{d^2}{4K^2T} \cdot P_\theta[\mathcal{E}^c] &\geq \frac{d^2}{32K^2T} \\ &\iff P_\theta[\mathcal{E}^c] \geq \frac{1}{8} - \frac{4\sqrt{2}K\sqrt{T}\epsilon}{d}. \end{aligned}$$

It follows that if

$$\frac{1}{8} - \frac{4\sqrt{2}K\sqrt{T}\epsilon}{d} \geq 0.1 \iff \frac{0.025^2}{32} \cdot \frac{d^2}{K^2\epsilon^2} \geq T,$$

then we have that  $P_\theta[\mathcal{E}^c] \geq 0.1$ . In words, this means that with constant probability, any algorithm must either collect more than  $c \cdot \frac{d^2}{K^2\epsilon^2}$  samples, or output a policy that is not  $\epsilon$ -optimal. This implies that  $T = \Omega(\frac{d^2}{K^2\epsilon^2})$  samples are necessary to guarantee an  $\epsilon$ -optimal policy. Equivalently, after  $T$  rounds, the suboptimality gap  $\epsilon$  is lower bounded as

$$\text{SubOpt}(T) = \Omega\left(\frac{d}{K\sqrt{T}}\right).$$

This concludes the proof of Theorem 2.  $\square$

## G.2 Proof of Lemmas for Theorem 2

### G.2.1 Proof of Lemma G.2

*Proof of Lemma G.2.* By the definition of KL divergence, we have

$$D_{\text{KL}}(\mathbb{P}_\theta(\cdot|S) \parallel \mathbb{P}_{\theta'}(\cdot|S)) = \mathbb{E}_{\sigma \sim \mathbb{P}_\theta(\cdot|S)} \left[ \sum_{j=1}^{|S|} \left( \phi_{\sigma_j}^\top (\theta - \theta') - \log \frac{\sum_{k=j}^{|S|} e^{\phi_{\sigma_k}^\top \theta}}{\sum_{k=j}^{|S|} e^{\phi_{\sigma_k}^\top \theta'}} \right) \right]. \quad (\text{G.6})$$

Fix a stage  $j$  and a ranking  $\sigma$ . We define

$$p_{k'}(\theta) := \frac{\exp(\phi_{\sigma_{k'}}^\top \theta)}{\sum_{k=j}^{|S|} \exp(\phi_{\sigma_k}^\top \theta)}, \quad \text{where } k' \in \{j, \dots, |S|\},$$

which corresponds to the Multinomial Logit (MNL) probability of selecting action  $\sigma_{k'}$  at position  $j$ , given the parameter  $\theta$  and the choice set  $S$ . Moreover, we define

$$f(\theta) := \log \left( \sum_{k=j}^{|S|} e^{\phi_{\sigma_k}^\top \theta} \right).$$

Then, by applying the mean value form of Taylor's theorem, there exists  $\bar{\theta} = (1-c)\theta + c\theta'$  for some  $c \in (0, 1)$  such that

$$\begin{aligned} -\log \frac{\sum_{k=j}^{|S|} e^{\phi_{\sigma_k}^\top \theta}}{\sum_{k=j}^{|S|} e^{\phi_{\sigma_k}^\top \theta'}} &= f(\theta') - f(\theta) \\ &= \nabla_\theta f(\theta)^\top (\theta' - \theta) + \frac{1}{2} (\theta' - \theta)^\top \nabla_\theta^2 f(\bar{\theta}) (\theta' - \theta) \quad (\text{Taylor's theorem}) \\ &\leq \nabla_\theta f(\theta)^\top (\theta' - \theta) + \frac{1}{2} \sum_{k=j}^{|S|} p_k(\bar{\theta}) (\phi_{\sigma_k}^\top (\theta' - \theta))^2 \\ &\leq \sum_{k=j}^{|S|} p_k(\theta) \phi_{\sigma_k}^\top (\theta' - \theta) + \frac{1}{2} \sum_{a \in S} (\phi_a^\top (\theta' - \theta))^2, \end{aligned} \quad (\text{G.7})$$

where the first inequality holds because

$$\nabla_\theta^2 f(\bar{\theta}) = \sum_{k=j}^{|S|} p_k(\bar{\theta}) \phi_{\sigma_k} \phi_{\sigma_k}^\top - \left( \sum_{k=j}^{|S|} p_k(\bar{\theta}) \phi_{\sigma_k} \right) \left( \sum_{k=j}^{|S|} p_k(\bar{\theta}) \phi_{\sigma_k} \right)^\top \leq \sum_{k=j}^{|S|} p_k(\bar{\theta}) \phi_{\sigma_k} \phi_{\sigma_k}^\top.$$

Plugging Equation (G.7) into Equation (G.6), we get

$$\begin{aligned}
& D_{\text{KL}}(\mathbb{P}_{\boldsymbol{\theta}}(\cdot|S) \|\mathbb{P}_{\boldsymbol{\theta}'}(\cdot|S)) \\
& \leq \mathbb{E}_{\boldsymbol{\sigma} \sim \mathbb{P}_{\boldsymbol{\theta}}(\cdot|S)} \left[ \sum_{j=1}^{|S|} \left( \phi_{\sigma_j}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}') - \sum_{k=j}^{|S|} p_k(\boldsymbol{\theta}) \phi_{\sigma_k}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}') + \frac{1}{2} \sum_{a \in S} (\phi_a^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}))^2 \right) \right] \\
& = \mathbb{E}_{\boldsymbol{\sigma} \sim \mathbb{P}_{\boldsymbol{\theta}}(\cdot|S)} \left[ \underbrace{\sum_{j=1}^{|S|} \mathbb{E}_{\sigma_j} \left[ \phi_{\sigma_j}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}') - \sum_{k=j}^{|S|} p_k(\boldsymbol{\theta}) \phi_{\sigma_k}^\top (\boldsymbol{\theta} - \boldsymbol{\theta}') \mid \sigma_1, \dots, \sigma_{j-1} \right]}_{=0} \right] \quad (\text{Tower rule}) \\
& \quad + \frac{|S|}{2} \sum_{a \in S} (\phi_a^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}))^2 \\
& \leq \frac{K}{2} \sum_{a \in S} (\phi_a^\top (\boldsymbol{\theta}' - \boldsymbol{\theta}))^2, \quad (|S| \leq K)
\end{aligned}$$

which concludes the proof.  $\square$

## H Additional Discussions

In this section, we provide additional discussion of our approach. In Subsection H.1, we propose a more efficient assortment selection rule than Equation (4), by using an arbitrary reference action  $\bar{a}_t \in \mathcal{A}$  instead of selecting the one that maximizes average uncertainty. In Subsection H.2, we show that under a sufficient feature diversity condition, selecting  $S_t$  uniformly at random can still achieve a comparable suboptimality gap. Finally, in Subsection H.3, we extend our approach to the active learning setting, as studied in [18].

### H.1 Arbitrary Reference Action for More Efficient Assortment Selection

As described in the main paper, the reference action  $\bar{a}_t$  is selected to maximize the average uncertainty across the subset  $S_t$ , according to Equation (4). This selection incurs a computational cost of  $\tilde{O}(N^2 K)$ .

However, in this subsection, we show that  $\bar{a}_t$  can, in fact, be selected arbitrarily—i.e., any  $\bar{a}_t \in \mathcal{A}$  is valid. Specifically, we modify our assortment selection rule as follows:

$$S_t = \underset{\substack{\bar{S} \in \mathcal{S} \\ \bar{a}_t \in \bar{S}}}{\text{argmax}} \frac{1}{|\bar{S}|} \sum_{a \in \bar{S}} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}, \quad \text{for any } \bar{a}_t \in \mathcal{A}. \quad (\text{H.1})$$

This results in only a constant-factor increase (specifically, by a factor of 2) in the suboptimality gap, while reducing the computational cost to  $\tilde{O}(NK)$ , as it removes the need to enumerate over all possible reference actions.

To show this explicitly, we return to Equation (D.6). Let  $\bar{a}_t$  be an arbitrary action in  $\mathcal{A}$  (e.g., selected uniformly at random). Then, we have

$$\begin{aligned}
& \frac{\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \|\phi(x_t, \pi^\star(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)) \pm \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} \\
& \leq \frac{\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \left( \|\phi(x_t, \pi^\star(x_t)) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} + \|\phi(x_t, \hat{\pi}_T(x_t)) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} \right) \\
& = \frac{\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \left( \sum_{a \in S_t^\star} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} + \sum_{a' \in \hat{S}_t} \|\phi(x_t, a') - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}} \right) \\
& \quad (\text{Let } S_t^\star := \{\pi^\star(x_t), \bar{a}_t\} \text{ and } \hat{S}_t := \{\hat{\pi}_T(x_t), \bar{a}_t\}) \\
& \leq \frac{4\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{1}{|S_t|} \sum_{a \in S_t} \|\phi(x_t, a) - \phi(x_t, \bar{a}_t)\|_{H_t^{-1}}. \quad (S_t \text{ selection rule, Eqn. (H.1)})
\end{aligned}$$

The remaining steps of the proof follow exactly as in the proof of Theorem 1 (or Theorem E.1).

## H.2 Suboptimality Gap Under Sufficient Diversity Condition

So far, we have considered the general case where the feature vectors  $\phi$  are not required to be diverse, and as a result, the induced matrix  $H_t - \lambda \mathbf{I}_d$  (or  $\Lambda_t - \lambda \mathbf{I}_d$ ) may be singular. In this subsection, we discuss the case where the following diversity assumption holds:

**Assumption H.1** (Diverse features). *For any  $S \in \mathcal{S}$  and  $a' \in S$ , there exists a constant  $\lambda_0 > 0$  such that  $\lambda_{\min} \left( \mathbb{E}_{x \sim \rho} \left[ \frac{1}{|S|} \sum_{a \in S} (\phi(x, a) - \phi(x, a'))(\phi(x, a) - \phi(x, a'))^\top \right] \right) \geq \lambda_0$ .*

Under this condition, it is sufficient to randomly select *exactly*  $K$  actions, rather than solving the optimization problem in Equation (4) to construct the assortment. Specifically, we can select  $S_t$  as:

$$S_t \sim \text{Unif}(\{S \subseteq \mathcal{A} : |S| = K\}), \quad \forall t \in [T]. \quad (\text{H.2})$$

**Theorem H.1** (Suboptimality Gap of Random Assortment Selection Under Diversity). *Let  $\mathcal{T}^w := \{t \in [T] : \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, a'_t)\|_{H_t^{-1}} \geq \frac{1}{\beta_{T+1}(\delta)}\}$ , where  $a'_t \in S_t$  is an arbitrary action selected from the assortment  $S_t$ . Suppose  $T = \Omega(\log(dT)/\lambda_0)$  and  $T > |\mathcal{T}^w|$ . Then, under the same setting as Theorem 1 and Assumption H.1, if  $S_t$  is randomly selected according to Equation (H.2), then with probability at least  $1 - \delta$ , we have:*

$$\text{SubOpt}(T) = \tilde{O} \left( \sqrt{\frac{d}{\lambda_0(T - |\mathcal{T}^w|)K}} \right) = \tilde{O} \left( \sqrt{\frac{d}{\lambda_0(T - \min\{(dK)^2/\kappa, T - 1\})K}} \right).$$

**Discussion of Theorem H.1.** Theorem H.1 shows that for sufficiently large  $T$  (i.e.,  $T = \Omega((dK)^2/\kappa + \log(dT)/\lambda_0)$ ), the suboptimality gap under the uniform random assortment selection strategy achieves  $\tilde{O}(\sqrt{\frac{d}{\lambda_0 T K}})$ . This result suggests that when the feature space is sufficiently diverse, uniform random selection is effective for learning. It also provides a theoretical explanation for the empirical success of many RLHF implementations [74, 57], where the feature space is sufficiently diverse and prompt-action (sub)set pairs are often selected uniformly at random.

Note that the lower bound we establish in Theorem 2 does not rely on the diversity assumption (Assumption H.1). As a result, deriving a lower bound under the diversity assumption remains an open question, which we leave for future work.

*Proof of Theorem H.1.* To provide the proof of Theorem H.1, we first introduce useful concentration inequalities.

**Lemma H.1** (Matrix Chernoff, Adapted Sequence from Tropp 78). *Consider a finite adapted sequence  $\{X_k\}$  with filtration  $\{\mathcal{F}_t\}_{t \geq 0}$  of positive-semi definite matrices with dimension  $d$ , and suppose that  $\lambda_{\max}(X_k) \leq R$  almost surely. Define the finite series*

$$Y := \sum_k X_k, \quad \text{and} \quad W := \sum_k \mathbb{E}_{k-1} X_k.$$

*Then, for all  $\mu \geq 0$ , we have*

$$\mathbb{P} \{ \lambda_{\min}(Y) \leq (1 - \delta)\mu \text{ and } \lambda_{\min}(W) \geq \mu \} \leq d \left( \frac{e^{-\delta}}{(1 - \delta)^{1-\delta}} \right)^{\mu/R}, \quad \delta \in [0, 1).$$

By setting  $\delta = \frac{3}{4}$ ,  $\mu = \lambda_0 t$ ,  $R = x_{\max}$  in Lemma H.1, we obtain the following result.

**Corollary H.1** (Eigenvalue Growth of Adaptive Gram Matrix). *If  $\|X_i\|_2 \leq x_{\max}$  and  $\lambda_{\min}(\mathbb{E}[X_i X_i^\top | \mathcal{F}_{i-1}]) \geq \lambda_0$ , then, with probability at least  $1 - d \exp(-c_1 \frac{\lambda_0 t}{x_{\max}})$ ,*

$$\lambda_{\min} \left( \sum_{i=1}^t X_i X_i^\top \right) \geq \frac{\lambda_0}{4} t$$

*holds for some absolute constant  $c_1$ .*

Now we are ready to provide the proof of Theorem H.1. For simplicity, we present only the case of the PL loss, since the extension to the RB loss directly follows from similar arguments in the proof of Theorem E.1. By the definition of the suboptimality gap, we have

$$\begin{aligned}
\mathbf{SubOpt}(T) &= \mathbb{E}_{x \sim \rho} \left[ \left( \phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)) \right)^\top \boldsymbol{\theta}^* \right] \\
&\leq \mathbb{E}_{x \sim \rho} \left[ \left( \phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)) \right)^\top \left( \boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1} \right) \right] \\
&\quad (\hat{\pi}_T(x) = \operatorname{argmax}_{a \in \mathcal{A}} \phi(x, a)^\top \hat{\boldsymbol{\theta}}_{T+1}) \\
&\leq \mathbb{E}_{x \sim \rho} [\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x))]_{H_{T+1}^{-1}} \left\| \boldsymbol{\theta}^* - \hat{\boldsymbol{\theta}}_{T+1} \right\|_{H_{T+1}} \quad (\text{H\"older's ineq.}) \\
&\leq \beta_{T+1}(\delta) \mathbb{E}_{x \sim \rho} [\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x))]_{H_{T+1}^{-1}}. \\
&\quad (\text{Corollary D.1, with prob. } 1 - \delta)
\end{aligned}$$

To proceed, we slightly modify the definition of  $\Lambda_t$ , as we no longer compute the reference action explicitly. Let  $a'_s$  be an arbitrary action selected from  $S_s$ , which can simply be chosen by sampling uniformly from  $S_s$ . Additionally, the regularization term  $\lambda$  is no longer required. Then, we redefine  $\Lambda_t$  as follows:

$$\Lambda_t := \sum_{s \in [t-1] \setminus \mathcal{T}^w} \sum_{a \in S_s} (\phi(x_s, a) - \phi(x_s, a'_s)) (\phi(x_s, a) - \phi(x_s, a'_s))^\top, \quad a'_s \in S$$

where

$$\mathcal{T}^w := \left\{ t \in [T] : \max_{a \in \mathcal{A}} \|\phi(x_t, a) - \phi(x_t, a'_t)\|_{H_t^{-1}} \geq \frac{1}{\beta_{T+1}(\delta)} \right\}.$$

Then, by Lemma D.2, we obtain

$$\begin{aligned}
\mathbb{E}_{x \sim \rho} [\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x))]_{H_{T+1}^{-1}} &\leq \sqrt{50} \mathbb{E}_{x \sim \rho} [\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x))]_{\Lambda_{T+1}^{-1}} \\
&\quad (\text{Lemma D.2, with prob. } 1 - \delta) \\
&\leq \frac{10\sqrt{2}}{\lambda_{\min}(\Lambda_{T+1})} \quad (\|\phi(x, a)\|_2 \leq 1)
\end{aligned}$$

Under the diversity assumption (Assumption H.1), for  $T \geq \frac{c}{\lambda_0} \log \frac{d}{\delta}$  with some constant  $c > 0$ , by Corollary H.1, we have, with probability at least  $1 - \delta$ ,

$$\lambda_{\min}(\Lambda_{T+1}) \geq \frac{\lambda_0}{4} (T - |\mathcal{T}^w|) K.$$

Suppose  $T - |\mathcal{T}^w| > 0$ . Then, combining the above results, we get

$$\begin{aligned}
\mathbf{SubOpt}(T) &\leq \beta_{T+1}(\delta) \frac{20\sqrt{2}}{\sqrt{\lambda_0(T - |\mathcal{T}^w|)K}} = \tilde{O} \left( \sqrt{\frac{d}{\lambda_0(T - |\mathcal{T}^w|)K}} \right) \\
&= \tilde{O} \left( \sqrt{\frac{d}{\lambda_0(T - \min\{(dK)^2/\kappa, T-1\})K}} \right),
\end{aligned}$$

where in the last inequality we use the fact that it holds that  $|\mathcal{T}^w| \leq |\mathcal{T}^w \cap (\mathcal{T}_0)^c| + |\mathcal{T}_0| = \tilde{O} \left( \frac{d^2 K^2}{\kappa} \right)$ , which follows from Lemmas D.4 and D.5. This concludes the proof of Theorem H.1.  $\square$

### H.3 Extension to Active Learning Setting

In this subsection, we consider a different setting—referred to as the *active learning setting*—where the learner has access to the entire context set  $\mathcal{X}$ , and the objective is to minimize the following *worst-case suboptimality gap*, defined as:

$$\mathbf{WorstSubOpt}(T) := \max_{x \in \mathcal{X}} [r_{\boldsymbol{\theta}^*}(x, \pi^*(x)) - r_{\boldsymbol{\theta}^*}(x, \hat{\pi}(x))].$$

This setting has received increasing attention in recent work [50, 47, 71, 18, 49, 77, 37]. However, most existing approaches focus exclusively on pairwise preference feedback. Mukherjee et al. [49]

study an online learning-to-rank problem where, for each context, a fixed set of  $K$  actions is provided, and the goal is to recover the true ranking based on feedback over these  $K$  actions. In contrast, we consider a more general setting in which, for each context, a set of  $N$  actions is available. The learner selects at most  $K$  actions from this set and receives ranking feedback over the selected subset. Thekumparampil et al. [77] investigate the problem of ranking  $N \geq K$  items using partial rankings over  $K$  candidates, but under a context-free setting. In contrast, we study a stochastic contextual setting, where contexts are drawn from an unknown (and fixed) distribution.

In the active learning setting, the algorithm jointly selects the context  $x_t$ —which is no longer given but actively chosen—and the assortment  $S_t$  by maximizing the average uncertainty objective. For computational efficiency, we employ the arbitrary reference action strategy described in Equation (H.1). (Note that one may alternatively use the reference action selection method from Equation (4), which selects  $\bar{a}_t$  to maximize uncertainty.)

$$(x_t, S_t) = \operatorname{argmax}_{x \in \mathcal{X}} \operatorname{argmax}_{\substack{S \in \mathcal{S} \\ \bar{a}_t \in S}} \frac{1}{|S|} \sum_{a \in S} \|\phi(x, a) - \phi(x, \bar{a}_t)\|_{H_t^{-1}}, \quad \text{for any } \bar{a}_t \in \mathcal{A}. \quad (\text{H.3})$$

The rest of the algorithm proceeds in the same manner as Algorithm 1. With the above context-assortment selection strategy, M-AUP0 achieves the following bound on the worst-case suboptimality gap, matching the order established in Theorem 1 (and in Theorem E.1):

**Theorem H.2.** *Under the same setting as Theorem 1 and E.1, with probability at least  $1 - \delta$ , M-AUP0 achieves the following worst-case suboptimality gap:*

$$\mathbf{WorstSubOpt}(T) = \begin{cases} \tilde{O}\left(\frac{d}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} + \frac{d^2 K^2}{\kappa T}\right), & (\text{PL loss}) \\ \tilde{O}\left(\frac{d}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} + \frac{d^2}{\kappa T}\right). & (\text{RB loss}) \end{cases}$$

*Proof of Theorem H.2.* We present only the proof using the PL loss (2), as extending it to the RB loss case (E.1) follows similarly to the extension from Theorem 1 to Theorem E.1.

By the definition of the worst-case suboptimality gap, we have

$$\begin{aligned} \mathbf{WorstSubOpt}(T) &= \max_{x \in \mathcal{X}} \left[ (\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)))^\top \theta^* \right] \\ &\leq \max_{x \in \mathcal{X}} \left[ (\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)))^\top (\theta^* - \hat{\theta}_{T+1}) \right] \\ &\quad (\hat{\pi}_T(x) = \operatorname{argmax}_{a \in \mathcal{A}} \phi(x, a)^\top \hat{\theta}_{T+1}) \\ &= \frac{1}{T} \sum_{t=1}^T \max_{x \in \mathcal{X}} \left[ (\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)))^\top (\theta^* - \hat{\theta}_{T+1}) \right]. \end{aligned}$$

We adopt the same definitions for  $\mathcal{T}_0$  (Equation (D.3)),  $\mathcal{T}^w$  (Equation (D.1)), and  $\Lambda_t$  (Equation (D.2)) as in Theorem 1. Then, we have

$$\begin{aligned} &\frac{1}{T} \sum_{t=1}^T \max_{x \in \mathcal{X}} \left[ (\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)))^\top (\theta^* - \hat{\theta}_{T+1}) \right] \\ &= \frac{1}{T} \sum_{t \in \mathcal{T}_0} \max_{x \in \mathcal{X}} \left[ (\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)))^\top (\theta^* - \hat{\theta}_{T+1}) \right] \\ &\quad + \frac{1}{T} \sum_{t \in \mathcal{T}^w \cap (\mathcal{T}_0)^c} \max_{x \in \mathcal{X}} \left[ (\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)))^\top (\theta^* - \hat{\theta}_{T+1}) \right] \\ &\quad + \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \max_{x \in \mathcal{X}} \left[ (\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)))^\top (\theta^* - \hat{\theta}_{T+1}) \right] \\ &\leq \frac{8B}{\log(2)T} d \log \left( 1 + \frac{2K}{\log(2)\lambda} \right) + \frac{48\sqrt{2}BK^2}{\kappa T} \beta_{T+1}(\delta)^2 d \log \left( 1 + \frac{2KT}{d\lambda} \right) \\ &\quad (\text{Lemma D.4 and D.5}) \\ &\quad + \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \max_{x \in \mathcal{X}} \left[ (\phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)))^\top (\theta^* - \hat{\theta}_{T+1}) \right]. \end{aligned} \quad (\text{H.4})$$

To further bound the last term of Equation (H.4), we get

$$\begin{aligned}
& \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \max_{x \in \mathcal{X}} \left[ \left( \phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)) \right)^\top \left( \theta^* - \hat{\theta}_{T+1} \right) \right] \\
& \leq \frac{1}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \max_{x \in \mathcal{X}} \left[ \left\| \phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)) \right\|_{H_{T+1}^{-1}} \left\| \theta^* - \hat{\theta}_{T+1} \right\|_{H_{T+1}} \right] \\
& \quad \text{(Hölder's ineq.)} \\
& \leq \frac{\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \max_{x \in \mathcal{X}} \left[ \left\| \phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)) \right\|_{H_t^{-1}} \right] \\
& \quad (H_{T+1} \geq H_t \text{ and Corollary D.1, with prob. } 1 - \delta)
\end{aligned}$$

Then, for any arbitrary  $\bar{a}_t \in \mathcal{A}$ , we have

$$\begin{aligned}
& \frac{\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \max_{x \in \mathcal{X}} \left[ \left\| \phi(x, \pi^*(x)) - \phi(x, \hat{\pi}_T(x)) \right\|_{H_t^{-1}} \right] \\
& \leq \frac{\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \max_{x \in \mathcal{X}} \left[ \left\| \phi(x, \pi^*(x)) - \phi(x, \bar{a}_t) \right\|_{H_t^{-1}} + \left\| \phi(x, \hat{\pi}_T(x)) - \phi(x, \bar{a}_t) \right\|_{H_t^{-1}} \right] \\
& \leq \frac{4\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{1}{|S_t|} \sum_{a \in S_t} \left\| \phi(x_t, a) - \phi(x_t, \bar{a}_t) \right\|_{H_t^{-1}}. \quad ((x_t, S_t) \text{ selection rule, Eqn. (H.3)})
\end{aligned}$$

Hence, we further obtain

$$\begin{aligned}
& \frac{4\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{1}{|S_t|} \sum_{a \in S_t} \left\| \phi(x_t, a) - \phi(x_t, \bar{a}_t) \right\|_{H_t^{-1}} \\
& \leq \frac{4\beta_{T+1}(\delta)}{T} \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \frac{1}{|S_t|} \sum_{a \in S_t} \left\| \phi(x_t, a) - \phi(x_t, \bar{a}_t) \right\|_{H_t^{-1}} \quad \text{(Eqn. (H.3))} \\
& \leq \frac{4\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \left( \frac{1}{|S_t|} \right)^2 |S_t|} \sqrt{\sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \sum_{a \in S_t} \left\| \phi(x_t, a) - \phi(x_t, \bar{a}_t) \right\|_{H_t^{-1}}^2} \\
& \quad \text{(Cauchy-Schwartz ineq.)} \\
& = \frac{4\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \left( \frac{1}{|S_t|} \right)^2 |S_t|} \sqrt{50 \sum_{t \notin \mathcal{T}_0 \cup \mathcal{T}^w} \sum_{a \in S_t} \left\| \phi(x_t, a) - \phi(x_t, \bar{a}_t) \right\|_{\Lambda_t^{-1}}^2} \\
& \quad \text{(Lemma D.2, with prob. } 1 - \delta) \\
& \leq \frac{30\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} \sqrt{\sum_{t=1}^T \min \left\{ 1, \sum_{a \in S_t} \left\| \phi(x_t, a) - \phi(x_t, \bar{a}_t) \right\|_{\Lambda_t^{-1}}^2 \right\}} \\
& \leq \frac{30\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} \sqrt{2d \log \left( 1 + \frac{2KT}{d\lambda} \right)} \quad \text{(Lemma D.3)} \\
& = \mathcal{O} \left( \frac{\beta_{T+1}(\delta)}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} \cdot \sqrt{d \log(KT)} \right). \quad \text{(H.5)}
\end{aligned}$$

Plugging Equation (H.5) into Equation (H.4), and setting  $\beta_{T+1}(\delta) = \mathcal{O}(B\sqrt{d \log(KT)} + B\sqrt{\lambda})$ , with probability at least  $1 - 3\delta$ , we have

$$\mathbf{SubOpt}(T) = \tilde{\mathcal{O}} \left( \frac{d}{T} \sqrt{\sum_{t=1}^T \frac{1}{|S_t|}} + \frac{d^2 K^2}{\kappa T} \right).$$

Substituting  $\delta \leftarrow \frac{\delta}{3}$ , we conclude the proof of Theorem H.2.  $\square$



## I Numerical Experiments

### I.1 Synthetic Data

**Setup.** In the synthetic data experiment, we sample the true but unknown parameter  $\theta^* \in \mathbb{R}^d$  from a  $d$ -dimensional standard normal distribution, i.e.,  $\theta^* \sim \mathcal{N}(0, I_d)$ , and then normalize it to ensure  $\|\theta^*\|_2 \leq 1$ . We consider four different types of context sets  $\mathcal{X}$ :

1. **Instance 1 (Stochastic contexts):** For each  $x \in \mathcal{X}$ , the feature vectors  $\phi(x, \cdot)$  are sampled from a standard normal distribution and then normalized to satisfy  $\|\phi(x, \cdot)\|_2 \leq 1$ . Here,  $|\mathcal{X}| = 100$ .
2. **Instance 2 (Non-contextual):** A single shared context is used for all rounds, i.e.,  $\mathcal{X} = \{x_1\}$  and  $|\mathcal{X}| = 1$ . The corresponding feature vectors  $\phi(x_1, \cdot)$  are sampled from a standard normal distribution and then normalized to satisfy  $\|\phi(x_1, \cdot)\|_2 \leq 1$ .
3. **Instance 3 (Hard-to-learn contexts):** For each  $x \in \mathcal{X}$ , the feature vectors  $\phi(x, \cdot)$  are constructed such that most of them are approximately orthogonal to the true parameter  $\theta^*$ . Here,  $|\mathcal{X}| = 100$ .
4. **Instance 4 (Skewed stochastic contexts):** For each  $x \in \mathcal{X}$ , the feature vectors  $\phi(x, \cdot)$  are sampled in a skewed or biased manner and then normalized to satisfy  $\|\phi(x, \cdot)\|_2 \leq 1$ . Here,  $|\mathcal{X}| = 100$ .

Additionally, we set the feature dimension to  $d = 5$  and the number of available actions to  $|\mathcal{A}| = N = 100$ . The suboptimality gap is measured every 25 rounds. All results are averaged over 20 independent runs with different random seeds, and standard errors are reported to indicate variability. The experiments are run on a Xeon(R) Gold 6226R CPU @ 2.90GHz (16 cores).

**Baselines.** For the baselines, in addition to DopeWolfe [77] introduced in the main paper, we also compare the performance of our algorithm, M-AUP0, against a uniform random assortment selection strategy, Uniform, as defined in Equation (H.2).

Thekumparampil et al. [77] propose a D-optimal design approach for the Plackett-Luce objective to efficiently select informative subsets of items for comparison. Recognizing the computational complexity inherent in this method, they introduce a randomized Frank-Wolfe algorithm, named DopeWolfe, which approximates the optimal design by solving linear maximization sub-problems on randomly chosen variables. This approach reduces computational overhead while maintaining effective learning performance. However, their approach is specifically tailored to the single-context setting (e.g., **Instance 2**) and may not generalize well to the multiple-context scenarios (e.g., **Instances 1, 3, and 4**). While their original implementation updates the model parameters using a maximum likelihood estimation (MLE) procedure, we instead adopt an online update strategy to ensure a fair comparison across all methods. For sampling size parameter  $R$ , we set  $R = \min\{\binom{N}{K}, 100, 000\}$ .

The uniform random assortment selection strategy, Uniform, selects  $K$  actions uniformly at random from the available action set  $\mathcal{A}$  at each round, without utilizing any uncertainty or reward-based information. This approach can be effective when the feature representations are sufficiently diverse (e.g., **Instances 1, 2, and 4**), but may perform poorly when the diversity parameter  $\lambda_0$  in Assumption H.1 is very small (e.g., **Instance 3**).

**Performance measure.** Since computing the exact suboptimality gap is challenging under a general distribution  $\rho$ , we instead evaluate the *average realized regret*, which serves as a slightly relaxed proxy for the suboptimality gap.

$$\begin{aligned} \text{SubOpt}(T) &\leq \frac{1}{T} \sum_{t=1}^T (\phi(x_t, \pi^*(x_t)) - \phi(x_t, \hat{\pi}_T(x_t)))^\top \theta^* + \underbrace{\tilde{\mathcal{O}}\left(\frac{1}{\sqrt{T}}\right)}_{\text{incurred by MDS terms}} \\ &\leq \underbrace{\frac{1}{T} \sum_{t=1}^T (\phi(x_t, \pi^*(x_t)) - \phi(x_t, \pi_t(x_t)))^\top \theta^*}_{=: \text{average realized regret}} + \tilde{\mathcal{O}}\left(\frac{1}{\sqrt{T}}\right), \end{aligned}$$

where we define  $\pi_t(x) := \arg\max_a \phi(x, a)^\top \hat{\theta}_t$ , and let  $\hat{\pi}_T$  denote the best policy among  $\{\pi_t\}_{t=1}^T$ , possibly selected using a validation set.

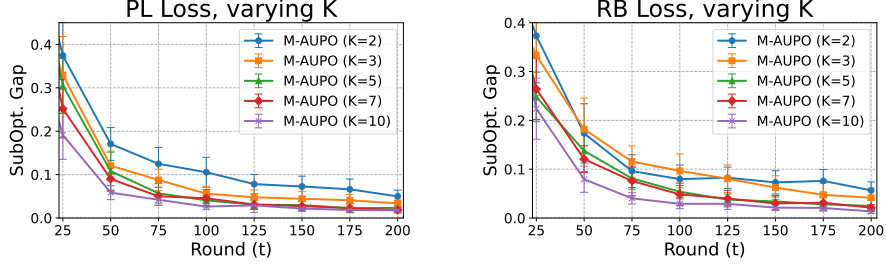


Figure I.1: Suboptimality gap of M-AUPO on Instance 1 under varying  $K$ , evaluated using PL loss (left) and RB loss (right).

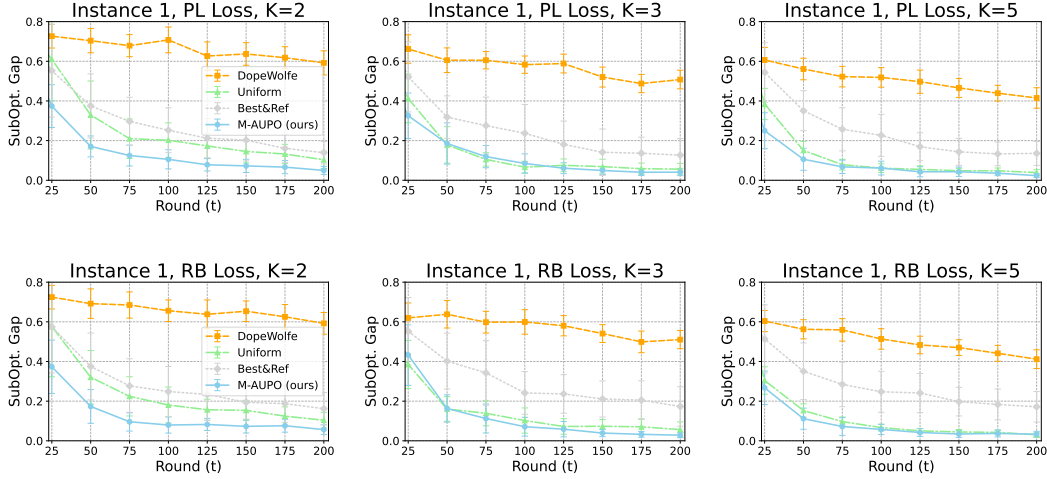


Figure I.2: Performance comparisons for Instance 1 (Stochastic contexts) with  $K = 2, 3$ , and  $5$ , evaluated under the PL loss (first row) and RB loss (second row).

**Results.** The Figure I.1 show the suboptimality gap of M-AUPO under both the PL loss (2) and RB loss (E.1) as the maximum assortment size  $K$  varies. The results clearly show that performance improves as  $K$  increases, supporting our theoretical findings.

We present performance comparisons in Figures I.2 through I.5, corresponding to Instances 1 through 4, respectively. Overall, our algorithm, M-AUPO, consistently outperforms the baseline methods. While DopeWolfe also considers the selection of  $K$  actions from  $N$  actions, it treats each context  $x$  independently and is specifically designed for the context-free setting (i.e., a singleton context). As a result, DopeWolfe cannot leverage information sharing across varying contexts and performs poorly in our setting. The only exception is in Instance 2 (Figure I.3), a special case of the non-contextual setting, where M-AUPO performs slightly worse than DopeWolfe. This is an expected outcome, as DopeWolfe leverages a D-optimal design strategy, which is known to be highly effective in the single-context setting. However, it is important to note that DopeWolfe completely fails in more general contextual scenarios (Figures I.2, I.4, and I.5), and its computational cost is significantly higher than that of our approach (see Table I.1).

The uniform random assortment selection strategy, Uniform, demonstrates competitive performance—though still worse than M-AUPO—in Instances 1, 2, and 4, as illustrated in Figures I.2, I.3, and I.5, respectively. However, in Instance 3 (Figure I.4), where the diversity parameter  $\lambda_0$  is very small due to most feature vectors lying within a hyperplane, Uniform performs significantly worse, as discussed in Appendix H.2.

Moreover, we observe that the suboptimality gap decreases with increasing  $K$  across all algorithms. For both M-AUPO and Uniform, this trend is consistent with our theoretical results (Theorems 1, E.1,

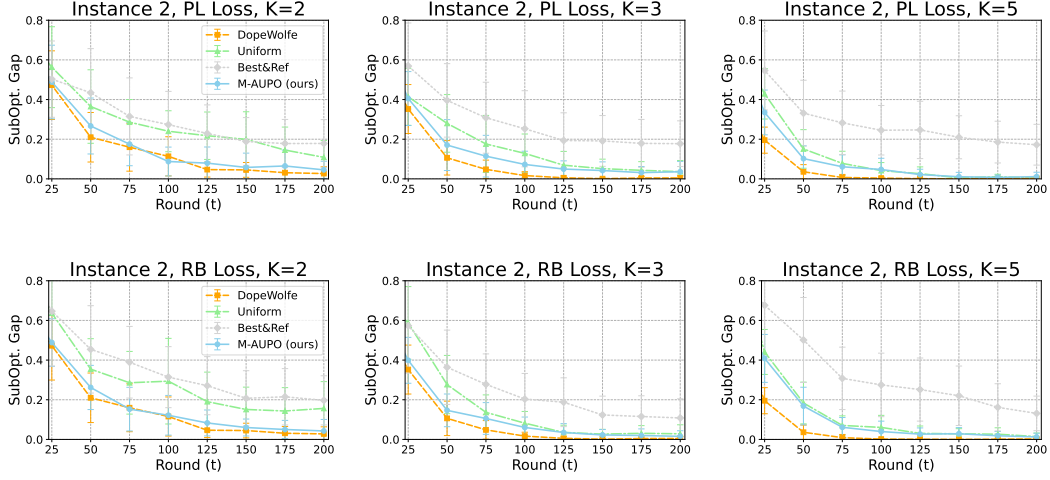


Figure I.3: Performance comparisons for Instance 2 (Non-contextual) with  $K = 2, 3$ , and  $5$ , evaluated under the PL loss (first row) and RB loss (second row).

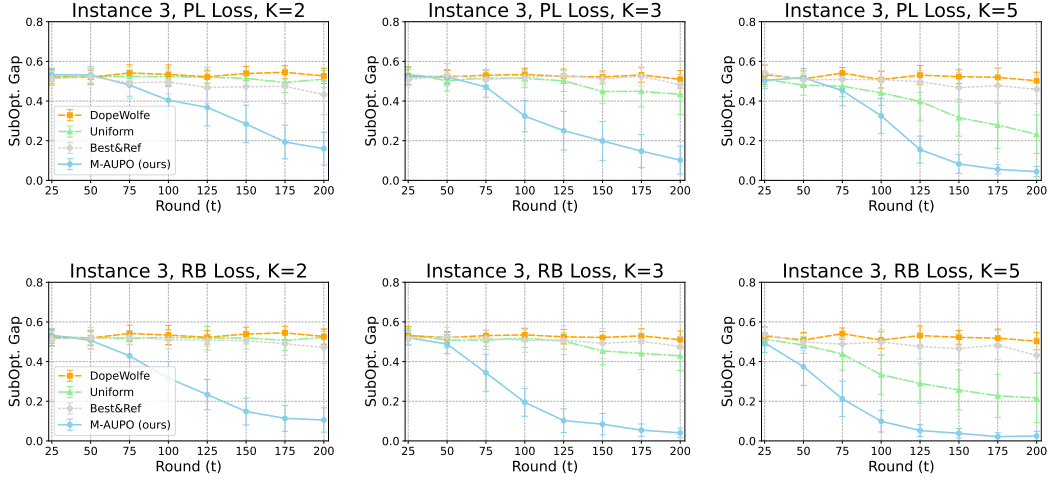


Figure I.4: Performance comparisons for Instance 3 (Hard-to-learn contexts) with  $K = 2, 3$ , and  $5$ , evaluated under the PL loss (first row) and RB loss (second row).

$K$	DopeWolfe	Uniform	M-AUPO (ours)
2	7.28 s	0.10 s	1.94 s
3	99.6 s	0.18 s	2.37 s
5	150.5 s	0.35 s	2.94 s
7	218.8 s	0.58 s	4.17 s
10	331.1 s	0.99 s	4.50 s

Table I.1: Runtime comparison over 200 rounds (seconds)

and H.1). In contrast, the improvement observed for DopeWolfe suggests that its current theoretical guarantees may be loose. This indicates that tighter bounds might be achievable by incorporating some of the techniques introduced in our work.

Table I.2 presents the average assortment size  $|S_t|$  of M-AUPO for various values of the maximum assortment size  $K$ . In most cases, the algorithm selects the full  $K$  actions, i.e.,  $|S_t| = K$ . An

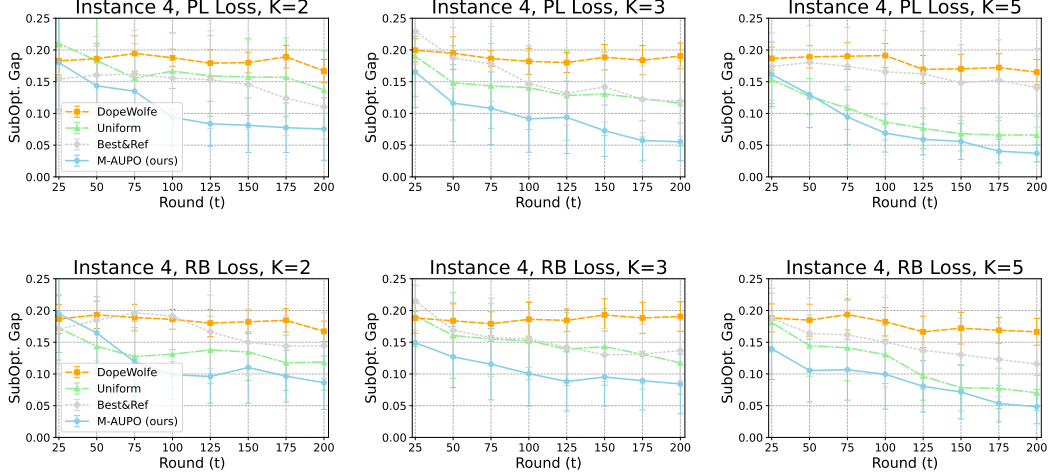


Figure I.5: Performance comparisons for Instance 4 (Skewed stochastic contexts) with  $K = 2, 3$ , and  $5$ , evaluated under the PL loss (first row) and RB loss (second row).

exception occurs when  $K$  is large (e.g., 30 or more), which may be impractical in real-world applications due to the increased annotation burden on human labelers.

$K$	2	3	5	7	10	30	50
PL loss, $ S_t $	2.00	3.00	5.00	7.00	10.00	18.31	18.69
RB loss, $ S_t $	2.00	3.00	5.00	7.00	10.00	18.39	18.40

Table I.2: Assortment size  $|S_t|$  of M-AUPO with varying maximum size  $K$  in the synthetic data experiment

## I.2 Real-World Dataset

**Setup.** In our real-world dataset experiments, we evaluate performance on two widely used benchmark datasets: TREC Deep Learning (TREC-DL)<sup>3</sup> and NECTAR<sup>4</sup>. The TREC-DL dataset provides 100 candidate answers for each query, offering a rich and diverse set of responses suitable for learning from listwise feedback. In contrast, the NECTAR dataset presents a more concise setup, with only 7 candidate answers per question. From each dataset, we randomly sample  $|\mathcal{X}| = 5000$  prompts (i.e., questions), each paired with its corresponding set of candidate actions—100 for TREC-DL and 7 for NECTAR.

We use the Gemma-2B<sup>5</sup> language model [76] to construct the feature representation  $\phi(x, a)$ . To obtain  $\phi(x, a)$ , we first concatenate the input prompt  $x$  and the candidate response  $a$  into a single sequence, which is then fed into Gemma-2B. The resulting feature vector is extracted from the last hidden layer of the model and has a dimensionality of  $d = 2048$ . We then apply  $\ell_1$  normalization to enhance numerical stability and ensure consistent scaling. For each round  $t$ , we sample the context index from an exponential distribution with rate  $\lambda = 0.1$ , which assigns higher probability to smaller indices and thus biases the selection toward earlier contexts. To generate ranking feedback and evaluate the suboptimality gap, we use the Mistral-7B<sup>6</sup> reward model [31] as the ground-truth reward function, denoted by  $r_{\theta^*}$ .

We measure the suboptimality gap every 2,500 rounds throughout the training process and report the average performance over 10 independent runs, each with a different random seed. Along with the

<sup>3</sup><https://microsoft.github.io/msmarco/TREC-Deep-Learning>

<sup>4</sup><https://huggingface.co/datasets/berkeley-nest/Nectar>

<sup>5</sup><https://huggingface.co/google/gemma-2b-it>

<sup>6</sup><https://huggingface.co/Ray2333/reward-model-Mistral-7B-instruct-Unified-Feedback>

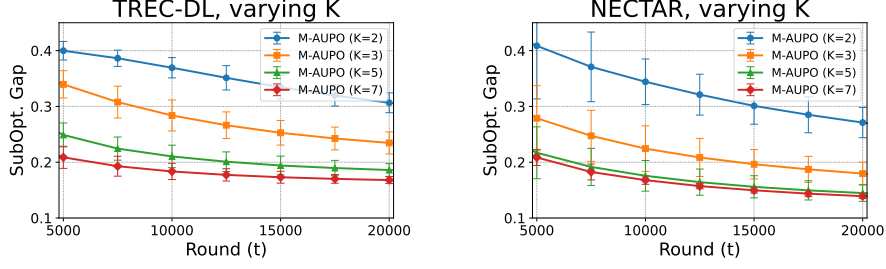


Figure I.6: Real-world dataset experiment: suboptimality gap of M-AUP0 under varying  $K$  on the TREC-DL dataset (left) and the NECTAR dataset (right). The results are rescaled to align the performances between the two datasets.

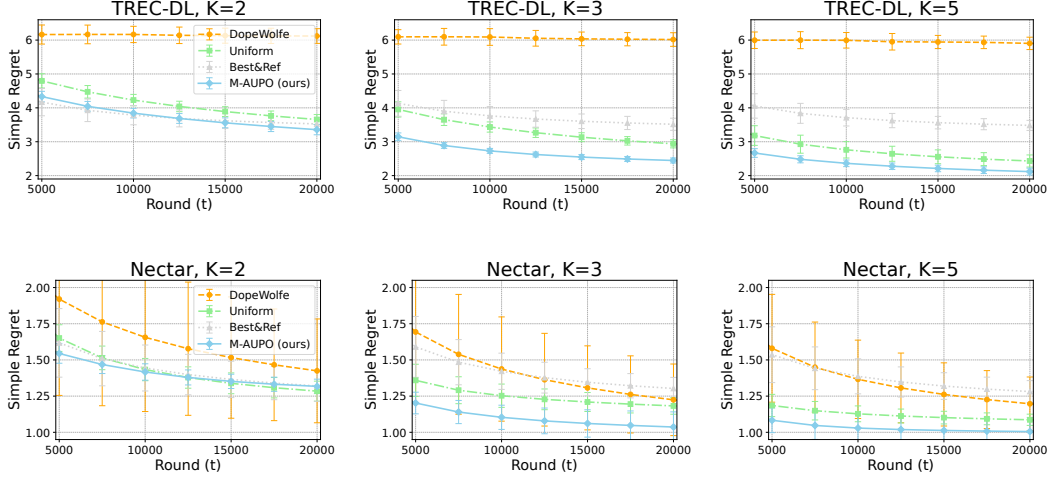


Figure I.7: Performance comparisons on the TREC-DL dataset (top row) and the NECTAR dataset (bottom row) for varying values of  $K = 2, 3$ , and  $5$ .

average, we also include the standard error to indicate variability across runs. In these experiments, we report results under the PL loss only, since the performance difference between PL and RB losses is minimal, as demonstrated in the synthetic data experiments. The experiments are conducted on a Xeon(R) Gold 6226R CPU @ 2.90GHz (16 cores) and a single GeForce RTX 3090 GPU.

**Baselines.** We use the same set of baselines as in the synthetic data experiments. For DopeWolfe [77], we set the sampling size parameter  $R$  as  $R = \min\{\binom{N}{K}, 1000\}$ . Although a small value of  $R \leq 1000$  may introduce significant approximation error—since the theoretically minimal-error choice is  $R = \mathcal{O}(\binom{N}{K})$ —we adopt this smaller value in our experiment to reduce computational overhead.

**Performance measure.** We measure the average realized regret as in the synthetic experiment (Appendix I.1).

**Results.** The Figure I.6 show the suboptimality gap of M-AUP0 under the PL loss on two real-world datasets as the maximum assortment size  $K$  varies. Consistent with our theoretical findings, the performance improves as  $K$  increases.

We present performance comparisons in Figure I.7. Our algorithm, M-AUP0, consistently outperforms all baselines by a significant margin. As in the synthetic data experiments, the suboptimality gap for all methods decreases as  $K$  increases. Notably, DopeWolfe performs particularly poorly on the TREC-DL dataset. This may be attributed to the use of a small sampling size  $R$ , which is insufficient compared to the full subset space of size  $\binom{N}{K} = \mathcal{O}(N^K) \gg 1000 \geq R$ . This result highlights an important practical limitation of DopeWolfe: despite its use of approximate optimization to

reduce runtime, the method still depends on combinatorial sampling to perform well, which becomes computationally infeasible in large-scale settings. In contrast, our algorithm, M-AUP0, maintains strong performance while requiring only  $\tilde{O}(NK)$  computational cost, making it significantly more scalable and practical for real-world applications.

Table I.3 reports the actual assortment size  $|S_t|$  selected by M-AUP0 on both datasets. In the TREC-DL experiment,  $|S_t|$  is nearly equal to  $K$  for all values of  $K$ , as the number of available actions is large ( $N = 100$ ). In contrast, in the NECTAR experiment, where the number of available actions is much smaller ( $N = 7$ ), the actual assortment size  $|S_t|$  is often smaller than  $K$ , especially when  $K = N$ . This reduction occurs because the limited action space constrains the potential informativeness of larger assortments—for example, it becomes difficult to achieve high average uncertainty when there are too few actions to choose from.

$K$	2	3	5	7
TREC-DL dataset, $ S_t $	2.00	3.00	4.99	6.95
NECTAR dataset, $ S_t $	2.00	2.99	4.31	4.74

Table I.3: Assortment size  $|S_t|$  of M-AUP0 with varying maximum size  $K$  in the real-world dataset experiment

## J Limitations

In this paper, we primarily focus on the online PbRL setting, where contexts are drawn stochastically from a fixed distribution. We also consider the active learning variant in Appendix H.3. However, we do not explore the offline setting [94], which may involve a different set of challenges. As a result, it remains an open question whether similar improvements—such as better performance with larger  $K$ —can be achieved in the offline setting. We view this as a promising direction for future research and leave it as an open problem.