
Optimal, Efficient and Practical Algorithms for Assortment Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We address the problem of active online assortment optimization problem
2 with preference feedback, which is a framework for modeling user choices
3 and subsetwise utility maximization. The framework is useful in various
4 real-world applications including ad placement, online retail, recommender
5 systems, and fine-tuning language models, amongst many others. The prob-
6 lem, although has been studied in the past, lacks an intuitive and practical
7 solution approach with simultaneously efficient algorithm and optimal re-
8 gret guarantee. E.g., popularly used assortment selection algorithms often
9 require the presence of a ‘strong reference’ which is always included in the
10 choice sets, further they are also designed to offer the same assortments
11 repeatedly until the reference item gets selected—all such requirements
12 are quite unrealistic for practical applications. In this paper, we designed
13 efficient algorithms for the problem of regret minimization in assortment
14 selection with *Plackett Luce* (PL) based user choices. We designed a novel
15 concentration guarantee for estimating the score parameters of the PL model
16 using ‘*Pairwise Rank-Breaking*’, which builds the foundation of our proposed
17 algorithms. Moreover, our methods are practical, provably optimal, and
18 devoid of the aforementioned limitations of the existing methods. Empirical
19 evaluations corroborate our findings and outperform the existing baselines.

20 1 Introduction

21 Studies have shown that it is often easier, faster and less expensive to collect feedback on a
22 relative scale rather than asking ratings on an absolute scale. E.g., to understand the liking
23 for a given pair of items, say (A,B), it is easier for the users to answer preference-based
24 queries like: “Do you prefer Item A over B?”, rather than their absolute counterparts: “How
25 much do you score items A and B in a scale of [0-10]?”. Due to the widespread applicability
26 and ease of data collection with relative feedback, learning from preferences has gained much
27 popularity in the machine-learning community, especially the active learning literature which
28 has applications in Medical surveys, AI tutoring systems, Multi-player sports/games, or any
29 real-world systems that have ways to collect feedback in terms of preferences. The problem
30 is famously studied as the *Dueling-Bandit* (DB) problem in the active learning community
31 [41, 3, 45, 46, 44], which is an online learning framework for identifying a set of ‘good’ items
32 from a fixed decision-space (set of items) by querying preference feedback of actively chosen
33 item-pairs. Consequently, the generalization of Dueling-Bandits, with *subset-wise* preferences
34 has also been developed into an active field of research. For instance, applications like
35 Web search (e.g. Google, Bing, or even in some versions of ChatGPT), online shopping
36 (Amazon, App stores, Google Flights), recommender systems (e.g. Youtube, Netflix, Google
37 News/Maps, Spotify) typically involve users expressing preferences by choosing one result (or
38 a handful of results) from a subset of offered items and often the objective of the system is to

39 identify the ‘most-profitable’ subset to offer to their users. The problem, popularly termed
 40 as ‘Assortment Optimization’ is studied in many interdisciplinary literature, e.g. Online
 41 learning and bandits [10], Operations research [40, 2], Game theory [15], RLHF [20, 30], to
 42 name a few.

43 **Problem (Informal): Active Optimal Assortment (AOA)** Active Assortment Opti-
 44 mization (a.k.a. Utility Maximization with Subset Choices) [13, 2, 23, 22] is an active
 45 learning framework for finding the ‘optimal’ profit-maximizing subset. Formally, assume
 46 we have a decision set of $[K] := \{1, 2, \dots, K\}$ of K items, with each item being associated
 47 with the score (or utility) parameters $\theta := (\theta_1, \theta_2, \dots, \theta_K)$ (without loss of generality assume
 48 $\theta_1 \geq \theta_2 \geq \dots \geq \theta_K \geq 0$). At each round $t = 1, 2, \dots$, the learner or the algorithm gets to
 49 query an assortment (typically subsets containing up to m -items) $S_t \subseteq [K]$, upon which
 50 it gets to see some (noisy) relative preferences across the items in S_t , typically generated
 51 according to an underlying Plackett-Luce (PL) choice model with parameters θ (1). Further,
 52 to allow the event where no items are selected, we also model a No-Choice (NC) item, indexed
 53 by item-0, with PL parameter $\theta_0 \in \mathbb{R}_+$.

54 **(Objective 1.) Top- m :** identify the top- m item-set: $\{\theta_1, \dots, \theta_m\}$, for some $m \in [1, K]$.

55 **(Objective 2.) Wtd-Top- m :** A more general objective could also consider a weight (or
 56 price) $r_i \in \mathbb{R}_+$ associated with the item $i \in [K]$, and the goal could be to identify the
 57 assortment (subset) with maximum weighted utility ¹, as detailed in Sec. 2.

58 **Related Works and Limitations:** As stated above, the problem of AOA is fundamental
 59 in many practical scenarios, and thus widely studied in multiple research areas, including
 60 Online ML/learning theory and operations research.

61 • In the Online ML literature, the problem is well-studied as *Multi-Dueling Bandits* [39, 14],
 62 or *Battling Bandits* [35, 34, 11], which is an extension of the famous *Dueling Bandit* problem
 63 [46, 45]. The main limitation of this line of work is the lack of practical objectives, which either
 64 aim to identify the ‘best-item’ $1 (= \arg \max_{i \in [K]} \theta_i)$ within a PAC (probably approximately
 65 correct) framework [36, 16, 17, 31] or quantifying regret against the best items [35, 12]. Note
 66 the latter actually leads to the optimal subset choice of repeatedly selecting the optimal item,
 67 $\arg \max_i \theta_i$, m times, i.e. $(1, 1, \dots, 1)$, which is unrealistic from the viewpoint of real-world
 68 system design. Selecting an assortment of distinct top- m items (Top- m -AOA) or maximum
 69 expected utility (Wtd-Top- m -AOA) makes more sense.

70 • On the other hand, a similar line of the problem has been studied in operations research
 71 and dynamic assortment selection literature, where the goal is to offer a subset of items to
 72 the customers in order to maximize expected revenue. The problem has been studied under
 73 different user choice models, e.g. PL or Multinomial-Logit models [2], Mallows and mixture of
 74 Mallows [22], Markov chain-based choice models [23], single transition model [27] etc. While
 75 these works indeed consider a more practical objective of finding the best assortment (subset)
 76 with the highest expected utility for a regret minimization objective, (1) a major drawback
 77 in their approach lies in the algorithm design which *requires to keep on querying the same set*
 78 *multiple times*, e.g. [2, 29, 18, 1]. Such design techniques could be impractical to be deployed
 79 in real systems where users could easily get annoyed if the same items are shown again and
 80 again. For example, in ad-placement, music/movies/news/tweets/reels recommendations,
 81 offering the same assortment could increase user dissatisfaction and disengagement.

82 (2) The second major drawback of this line of work lies in the *structural assumption of*
 83 *their underlying choice models which requires the existence of a reference/default item, that*
 84 *needs to be part of every assortment S_t .* This leads to assuming a No-Choice item, typically
 85 denoted as item-0, which is a default choice of any assortment S_t . Further a stronger and
 86 more unrealistic assumption lies in the fact that they require to assume that the above pivot
 87 is stronger than the rest of the K items, i.e. $\theta_0 \geq \max_{i \in [K]} \theta_i$, i.e. the No-Choice (NC)
 88 action is the most likely outcome of any assortment S_t . This is often unrealistic, e.g., during
 89 user interactions with language models, or online shopping, or Route recommendation in
 90 GPS navigation, a NC action is highly improbable. Consequently, such assumption limits the
 91 use in real-systems. In the existing literature [2, 28, 1, 24], such assumptions are primarily

¹This is equivalent to finding the set with maximum expected revenue when r_i s represents the price of item i [2]

92 adapted solely for theoretical needs, precisely for maintaining concentration bounds of the
 93 PL parameters θ , and hence not well justified from a practical viewpoint. Some recent
 94 developments also generalized the AOA problem to linear MNL scores to incorporate large
 95 actions embedded in d -dimension [43, 42, 28], however, their approaches are either limited
 96 to the above restrictions or suffer sub-optimal regret guarantees without those assumptions
 97 (e.g. the regret bound of [28] is $O(d^{3/2}\sqrt{T})$ which is suboptimal by a d -factor). Considering
 98 the above limitations of the AOA literature, we set to answer two questions:

- 99 (1) Can we consider a general AOA model where the default item, like the NC item defined
 100 above, is not necessarily the strongest one, i.e. $\theta_0 \geq \max_{i \in [K]} \theta_i$?
 101 (2) Can we design a practical and regret optimal algorithm for the AOA framework, without
 102 needing to play the same repetitive actions and yet converge to the optimal assortment?

103 **Contributions** We answer these questions in the affirmative and present best of all
 104 scenarios. We design practical algorithms on practical AOA framework with practical
 105 objectives—Unlike the existing approaches of the AOA, literature [2, 18], we do not have to
 106 keep playing the same assortment multiple times, neither require a strongest default item
 107 (like NC satisfying $\theta_0 \geq \max_{i \in [K]} \theta_i$). Moreover, our objectives do not require us to converge
 108 to a multiset of replicated arms like $(1, 1, \dots, 1)$, but converge to the utility-maximizing set of
 109 distinct items. We list our contributions below:

110 **1. A General AOA Setup:** We work with a general problem of AOA for PL model,
 111 which requires no additional structural assumption of the θ parameters such as $\theta_0 \geq \max_i \theta_i$,
 112 unlike the existing works. We designed algorithms for two separate objectives Top- m and
 113 Wtd-Top- m as discussed above (Sec. 2).

114 **2. Practical, Efficient and Optimal Algorithm:** In Sec. 3, we give a practical,
 115 efficient and optimal algorithm for MNL Assortment (up to log factors and the magnitude of
 116 θ_{\max}). The regret bound of our algorithm AOA-RB_{PL} (Alg. 1) yields $\tilde{O}(\sqrt{KT})$ regret for
 117 both Top- m and Wtd-Top- m objective. Our algorithms use a novel parameter estimation
 118 technique for discrete choice models based on the concept of *Rank-Breaking* (RB) which is
 119 one of our key contributions towards designing the efficient and optimal algorithm. This
 120 enables our algorithm to perform optimally without requiring the No-Choice item to be
 121 the strongest. Appendix A details the key concept of our parameter estimation technique
 122 exploiting the concept of RB. Our resulting algorithm plays optimistically based on the UCB
 123 estimates of PL parameters and does not require repeating the same subset multiple times,
 124 justifying our title.

125 **3. Improvement with Adaptive Pivots:** In Sec. 4, we refine the performance of
 126 our algorithm by employing the novel idea of ‘adaptive pivots’ (a reference item) and
 127 proposed AOA-RB_{PL}-Adaptive. Performance-wise this removes the asymptotic dependence
 128 on $\theta_{\max} = \max_i \theta_i / \theta_0$ in the regret analysis. This enables the algorithm to work effectively
 129 in scenarios where the No-Choice item is less likely to be selected, i.e., $\theta_{\max} \gg 1$. This
 130 leads to a huge improvement in our experiments, especially in the range of low θ_0 , where
 131 AOA-RB_{PL}-Adaptive drastically outperforms over the existing baseline. Comparison of our
 132 regret bound with existing work is detailed in Table 1.

133 **4. Empirical Analysis.** Finally, we corroborate our theoretical results with empirical
 134 evaluations (Sec. 5), which certify our superior performance in the general AOA setups.

Work	Framework	Assume $\theta_0 = \theta_{\max} = 1$	Regret
Our (Alg. 1)	MNL model (Obj. 2)	No	$\sqrt{\min\{\theta_{\max}, K\}KT \log T}$
[2] (Thm 1)	MNL model (Obj. 2)	Yes	$\sqrt{KT \log T}$
[2] (Thm 4)	MNL model (Obj. 2)	No	$\sqrt{\theta_{\max}KT \log T}$
[1]	MNL model (Obj. 2)	Yes	$\sqrt{KT \log(mT)} + K \log^2(mT)$
[24]	MNL model with constraints (Obj. 2)	No	$\sqrt{\frac{KT}{\min_i r_i} \log T}$

Table 1: Our Contribution vs the Existing Results in the K -armed MNL-Assortment literature

135 It is also worth mentioning that our proposed algorithm and their respective regret analysis
 136 could be extended to any general random utility (RUM) based preference models [38, 37],
 137 as explained in Rem. 1. However, to keep the focus on the AOA problem and ease the
 138 presentation, we stick to the special case of MNL choice model based preferences.

139 2 Problem Setup

140 We write $[n] = \{1, 2, \dots, n\}$ and $\mathbf{1}\{\cdot\}$ denotes the indicator function. The symbol \lesssim , employed
 141 in the proof sketches, represents a coarse inequality.

142 We consider the sequential decision-making problem of Active Optimal Assortment (AOA),
 143 with preference/choice feedback. Formally, the learner is given $[K]$, a finite set of K items
 144 ($K > 2$). At each decision round $t = 1, 2, \dots$, the learner selects a subset $S_t \subseteq [K]$ of up to
 145 m items, and receives some (stochastic) feedback about the item preferences of S_t , drawn
 146 according to some unknown underlying Plackett-Luce (PL) choice model (1) with parameters
 147 $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_K) \in \mathbb{R}_+^K$. We assume $\theta_1 \geq \theta_2 \geq \dots \geq \theta_K$ without loss of generality. An
 148 interested reader may check App. A.1 for a detailed discussion on PL models. Given any
 149 assortment S_t we also consider the possibility of ‘no-selection’ of any items given an S_t .
 150 Following the literature of [2], we model this mathematically as a No-Choice (NC) item,
 151 indexed by item-0, and its corresponding PL utility parameter θ_0 . Unlike most existing
 152 literature on assortment selection, we are not assuming $\theta_0 \not\geq \max_{i \in [K]} \theta_i$. Further, since the
 153 PL model is scale independent, we set $\theta_0 = 1$ and scale the rest of the PL parameters.

154 **Feedback model** The feedback model formulates the information received (from the
 155 ‘environment’) once the learner plays a subset $S_t \subseteq [K]$ of at most m items. Given S_t we
 156 consider the algorithm receives a winner feedback (or index of an item) $i_t \in S_t \cup \{0\}$, drawn
 157 according to the underlying PL choice model as:

$$\mathbb{P}(i_t = i | S_t) = \theta_i / (\theta_0 + \sum_{j \in S_t} \theta_j), \quad \forall i \in S_t. \quad (1)$$

158 We consider the following two objectives for the learner:

159 **1. Top- m -Objective.** One simple objective could be to identify the top- m item-set:
 160 $\{\theta_1, \dots, \theta_m\}$, for some $m \in [1, K]$. The performance of the learner can be captured by
 161 minimizing the following regret:

$$Reg_T^{\text{top}} := \sum_{t=1}^T \frac{\Theta_{S^*} - \Theta_{S_t}}{m}, \quad \text{where } S^* := \operatorname{argmax}_{S \subseteq [K]: |S|=m} \left\{ \Theta_S := \sum_{i \in S} \theta_i \right\}.$$

162 **2. Wtd-Top- m -Objective.** Here, each item- i is associated with a weight (for example
 163 price) $r_i \in \mathbb{R}_+$, and the goal is to identify the set of size at most m with maximum weighted
 164 utility. One could measure the regret of the learner as:

$$Reg_T^{\text{wtd}} := \sum_{t=1}^T (\mathcal{R}(S^*, \boldsymbol{\theta}) - \mathcal{R}(S_t, \boldsymbol{\theta})), \quad \text{where } \mathcal{R}(S, \boldsymbol{\theta}) := \sum_{i \in S} \frac{r_i \theta_i}{\theta_0 + \sum_{j \in S} \theta_j}, \quad \forall S \subseteq [K], \quad (2)$$

165 denotes $S^* := \operatorname{argmax}_{S \subseteq [K]: |S| \leq m} \mathcal{R}(S, \boldsymbol{\theta})$ is the optimal utility-maximizing subset. This
 166 objective corresponds to the standard objective in the MNL litterature [2].

167 3 A Practical and Efficient Algorithm for AOA with PL

168 In this section, we introduce our first algorithm, which works for both objectives.

169 3.1 Algorithm Design

170 At each time t , our algorithm (Alg. 1) maintains a pairwise preference matrix $\widehat{\mathbf{P}}_t \in [0, 1]^{n \times n}$,
 171 whose (i, j) -th entry $\widehat{p}_{ij,t}$ records the empirical probability of i having beaten j in a pairwise

172 duel, and a corresponding upper confidence bound $p_{ij,t}^{\text{ucb}}$. Let $[\tilde{K}] := [K] \cup \{0\}$. We define for
 173 each pair $(i, j) \in [\tilde{K}] \times [\tilde{K}]$,

$$p_{ij,t}^{\text{ucb}} := \hat{p}_{ij,t} + \sqrt{\frac{2\hat{p}_{ij,t}(1-\hat{p}_{ij,t})x}{n_{ij,t}}} + \frac{3x}{n_{ij,t}}, \quad \text{where } \hat{p}_{ij,t} := \frac{w_{ij,t}}{n_{ij,t}}, \quad (3)$$

174 where $w_{ij,t} = \sum_{s=1}^{t-1} \mathbb{1}\{i_s = i, j \in S_s\}$ denotes the number of pairwise wins of item- i over j
 175 and $n_{ij,t} = w_{ij,t} + w_{ji,t}$ being the number of times (i, j) has been compared. The above UCB
 176 estimates $p_{ij,t}^{\text{ucb}}$ are further used to design UCB estimates of the PL parameters θ_i as follows

$$\theta_{i,t}^{\text{ucb}} = p_{i0,t}^{\text{ucb}} / (1 - p_{i0,t}^{\text{ucb}})_+.$$

177 The estimates $\theta_{i,t}^{\text{ucb}}$ s are then used to select the set S_t , that maximizes the underlying objective.
 178 This optimization problem transforms into a static assortment optimization problem with
 179 upper confidence bounds $\theta_{i,t}^{\text{ucb}}$ as the parameters, and efficient solution methods for this case
 180 are available (see e.g., [7, 21, 32]).

Algorithm 1 AOA for PL model with RB (AOA-RB_{PL})

1: **input:** $x > 0$
 2: **init:** $\tilde{K} \leftarrow K + 1$, $[\tilde{K}] = [K] \cup \{0\}$, $\mathbf{W}_1 \leftarrow [0]_{\tilde{K} \times \tilde{K}}$
 3: **for** $t = 1, 2, 3, \dots, T$ **do**
 4: Set $\mathbf{N}_t = \mathbf{W}_t + \mathbf{W}_t^\top$, and $\hat{\mathbf{P}}_t = \frac{\mathbf{W}_t}{\mathbf{N}_t}$. Denote $\mathbf{N}_t = [n_{ij,t}]_{\tilde{K} \times \tilde{K}}$ and $\hat{\mathbf{P}}_t = [\hat{p}_{ij,t}]_{\tilde{K} \times \tilde{K}}$.
 5: Define for all i , $p_{ii,t}^{\text{ucb}} = \frac{1}{2}$ and for all $i, j \in [\tilde{K}], i \neq j$

$$p_{ij,t}^{\text{ucb}} = \hat{p}_{ij,t} + \left(\frac{2\hat{p}_{ij,t}(1-\hat{p}_{ij,t})x}{n_{ij,t}} \right)^{1/2} + \frac{3x}{n_{ij,t}}$$

 6: $\theta_{i,t}^{\text{ucb}} := p_{i0,t}^{\text{ucb}} / (1 - p_{i0,t}^{\text{ucb}})_+$
 7: $S_t \leftarrow \begin{cases} \text{Top-}m \text{ items from } \text{argsort}(\{\theta_{1,t}^{\text{ucb}}, \dots, \theta_{K,t}^{\text{ucb}}\}), \\ \quad \text{for Top-}m \text{ objective} \\ \text{argmax}_{S \subseteq [K] \mid |S| \leq m} \mathcal{R}(S, \theta_t^{\text{ucb}}), \\ \quad \text{for Wtd-Top-}m \text{ objective} \end{cases}$
 8: Play S_t
 9: Receive the winner $i_t \in [\tilde{K}]$ (drawn as per (1))
 10: Update: $\mathbf{W}_{t+1} = [w_{ij,t+1}]_{\tilde{K} \times \tilde{K}}$ s.t. $w_{i_t j, t+1} \leftarrow w_{i_t j, t} + 1 \quad \forall j \in S_t \cup \{0\}$
 11: **end for**

181 **3.2 Analysis: Concentration Lemmas**

182 We start the analysis by providing two technical lemmas, whose proofs are deferred to the
 183 appendix and that provide confidence bounds for the θ_i .

184 **Lemma 1.** *Let $T \geq 1$ and $x > 0$. Then, with probability at least $1 - 3KT e^{-x}$, for all $t \in [T]$
 185 and $i \in [K]$: $\theta_i \leq \theta_{i,t}^{\text{ucb}}$ atleast one of the following two inequalities is satisfied*

$$n_{i0,t} < 69x(\theta_0 + \theta_i) \quad \text{or} \quad \theta_{i,t}^{\text{ucb}} \leq \theta_i + 4(\theta_0 + \theta_i) \sqrt{\frac{2\theta_0\theta_i x}{n_{i0,t}}} + \frac{22x(\theta_0 + \theta_i)^2}{n_{i0,t}}.$$

186 The above lemma depends on $n_{i0,t}$ the number of times items i have been compared with
 187 item 0 up to round t . The latter is controlled using the following lemma:

188 **Lemma 2.** *Let $T \geq 1$ and $x > 0$. Then, with probability at least $1 - KTe^{-x}$: simultaneously
 189 for all $t \in [T]$ and $i \in [K]$*

$$\tau_{i,t} < 2x(\theta_0 + \Theta_{S^*})^2 \quad \text{or} \quad n_{i0,t} \geq \frac{(\theta_0 + \theta_i)\tau_{i,t}}{2(\theta_0 + \Theta_{S^*})}, \quad (4)$$

190 where $\tau_{i,t} = \sum_{s=1}^{t-1} \mathbb{1}\{i \in S_s\}$ denotes the number of rounds item i got selected before round t .

191 **3.3 Analysis: Top- m Objective:**

192 We are now ready to provide the regret upper bound for Algorithm 1 with Top- m objective.

193 **Theorem 3** (Top- m Objective). *Let $\theta_{\max} \geq 1$. Consider any instance of PL model on K*
 194 *items with parameters $\theta \in [0, \theta_{\max}]^K$, $\theta_0 = 1$. The regret of Alg. 1 with parameter $x = 2 \log T$*
 195 *is bounded as*

$$Reg_T^{\text{top}} = O(\theta_{\max}^{3/2} \sqrt{KT \log T}) \quad \text{when } T \rightarrow \infty.$$

196 The above rate of $\tilde{O}(KT)$ is optimal (up to log-factors), as a lower bound can be derived from
 197 standard multi-armed bandits [5, 6]. We only state here a sketch of the proof of Theorem 3.
 198 The detailed proof is deferred to the App. B.

199 *Proof Sketch of Theorem 3.* Let us define for any $S \subseteq [K]$,

$$\Theta_S = \sum_{i \in S} \theta_i, \quad \text{and} \quad \Theta_S^{\text{ucb}} := \sum_{i \in S} \theta_i^{\text{ucb}}.$$

200 Let \mathcal{E} be the high-probability event such that both Lemma 1 and 2 holds true. Then, $\mathbb{P}(\mathcal{E}) \geq$
 201 $1 - 4TK e^{-x}$. Let us first assume that \mathcal{E} holds true. Then, by Lemma 1, $\Theta_{S^*} \leq \Theta_{S^*}^{\text{ucb}} \leq \Theta_{S_t}^{\text{ucb}}$,
 202 which yields

$$Reg_T^{\text{top}} = \frac{1}{m} \sum_{t=1}^T \Theta_{S^*} - \Theta_{S_t} \leq \frac{1}{m} \sum_{t=1}^T \Theta_{S_t}^{\text{ucb}} - \Theta_{S_t} \lesssim \tau_0 + \frac{1}{m} \sum_{t=1}^T \sum_{i \in S_t} (\theta_i^{\text{ucb}} - \theta_i) \mathbb{1}\{\tau_{i,t} \geq \tau_0\},$$

203 where $\tau_0 = 138x(m+1)^2 \theta_{\max}^2$ corresponds to an exploration phase needed for the confidence
 204 upper bounds of Lem 1 and 2 to be satisfied. Then, noting that if \mathcal{E} holds true, we can show
 205 by Lemma 2, that $\mathbb{1}\{\tau_{i,t} \geq \tau_0\} \leq \mathbb{1}\{n_{i0,t} \geq 69x(\theta_0 + \theta_i)\}$. Therefore, we can apply Lemma 1
 206 that entails,

$$\begin{aligned} \frac{1}{m} \sum_{t=1}^T \sum_{i \in S_t} (\theta_i^{\text{ucb}} - \theta_i) \mathbb{1}\{\tau_{i,t} \geq \bar{n}_{i0}\} &\lesssim \frac{1}{m} \sum_{t=1}^T \sum_{i \in S_t} \left((\theta_0 + \theta_i) \sqrt{\frac{\theta_0 \theta_i x}{n_{i0,t}}} \mathbb{1}\{\tau_{i,t} \geq \tau_0\} \right) \\ &\stackrel{\text{Lem. 2}}{\lesssim} \frac{1}{m} \sum_{t=1}^T \sum_{i \in S_t} \theta_{\max}^{3/2} \sqrt{\frac{mx}{\tau_{i,t}}} \lesssim \frac{1}{m} \sum_{i=1}^K \theta_{\max}^{3/2} \sqrt{mx \tau_{i,t}} \lesssim \theta_{\max}^{3/2} \sqrt{xKT}. \end{aligned}$$

207 where we used $\sum_{i=1}^n 1/\sqrt{i} \leq 2\sqrt{n}$ and $\sum_i \tau_{i,t} = mT$ together with Jensen's inequality in the
 208 last inequality. We thus have under the event \mathcal{E} that $Reg_T^{\text{top}} \leq O(\theta_{\max}^{3/2} \sqrt{xKT})$ and the proof
 209 is concluded by taking the expectation with $x = 2 \log T$ to control $\mathbb{P}(\mathcal{E}^c)$. \square

210 **3.4 Analysis: Wtd-Top- m Objective**

211 We turn now to the analysis of the Wtd-Top- m objective (2). We start by stating a lemma
 212 from [2] that shows that the expected utility $\mathcal{R}(S^*, \theta)$ that corresponds to the optimal
 213 assortment $S^* = \arg\max_{S \subseteq [K], |S| \leq m} \mathcal{R}(S, \theta)$ is non-decreasing in the parameters θ .

214 **Lemma 4** (Lemma A.3 of [2]). *Assume $\theta_i^{\text{ucb}} \geq \theta_i$ for all $i \in [K]$, then $\mathcal{R}(S^*, \theta) \leq \mathcal{R}(S^*, \theta^{\text{ucb}})$.*

215 **Theorem 5** (Wtd-Top- m Objective). *Let $\theta_{\max} \geq 1$. Then, for any $\theta \in [0, \theta_{\max}]^K$ and*
 216 *weights $\mathbf{r} \in [0, 1]^K$, the weighted regret of AOA-RB $_{PL}$ (Alg. 1) with $x = 2 \log T$*

$$Reg_T^{\text{wtd}} = O(\sqrt{\theta_{\max} KT \log T}) \quad \text{when } T \rightarrow \infty.$$

217 The complete proof is postponed to App. B. The rate $\Omega(\sqrt{KT})$ is optimal as proved by the
 218 lower bound in [19] for MNL bandit problems for $\theta_{\max} = 1$. Our result recovers (up to a factor
 219 $\sqrt{\log T}$) the one of [2] when $\theta_{\max} = 1$. However, their algorithm relies on more sophisticated
 220 estimators that necessitate epochs repeating the same assortment until the No-Choice item
 221 is selected. Note for our problem setting, where it is possible to have $\theta_{\max} \gg \theta_0 = 1$, the
 222 length of these epochs could be of $O(K\theta_{\max})$, which could be potentially very large when
 223 $\theta_{\max} \gg 1$. This reduces the number of effective epochs, leading to poor estimation of the PL
 224 parameters. We see this tradeoff in our experiments (Sec. 5) where the MNL-UCB algorithm
 225 of [2] yields linear $O(T)$ regret for such choice of the problem parameters.

226 **Remark 1** (Beyond MNL Models). *Although, in this paper, we primarily focused on MNL*
 227 *based choice models, it is worth mentioning that our proposed algorithms can be generalized*
 228 *to more general random utility based models (RUMs) [9, 33] pursuing the ideas from [36]*
 229 *that extends the RB based parameter estimation technique to any RUM(θ) choice models.*
 230 *Our algorithms and analyses thus apply to any general RUM(θ) based choice models; we stick*
 231 *to the special case of MNL models in this paper for brevity and keep the main focus on the*
 232 *AOA problem and the related algorithmic novelties.*

233 *Proof sketch of Thm. 5.* Let \mathcal{E} be the high-probability event such that both Lemma 1 and 2
 234 are satisfied. Then,

$$\begin{aligned} \text{Reg}_T^{\text{wtd}} &= \sum_{t=1}^T \mathbb{E}[\mathcal{R}(S^*, \theta) - \mathcal{R}(S_t, \theta)] \lesssim \sum_{t=1}^T \mathbb{E}[(\mathcal{R}(S^*, \theta) - \mathcal{R}(S_t, \theta)) \mathbf{1}\{\mathcal{E}\}] + T\mathbb{P}(\mathcal{E}^c) \\ &\lesssim \sum_{t=1}^T \mathbb{E}[(\mathcal{R}(S_t, \theta_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \mathbf{1}\{\mathcal{E}\}] + T\mathbb{P}(\mathcal{E}^c) \end{aligned} \quad (5)$$

235 because $\mathcal{R}(S_t, \theta_t^{\text{ucb}}) \geq \mathcal{R}(S^*, \theta_t^{\text{ucb}}) \geq \mathcal{R}(S^*, \theta)$ under the event \mathcal{E} by Lemma 4. We now
 236 upper-bound the first term of the right-hand-side

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[(\mathcal{R}(S_t, \theta_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \mathbf{1}\{\mathcal{E}\}] &= \sum_{t=1}^T \mathbb{E}\left[\left(\sum_{i \in S_t} \frac{r_i \theta_{i,t}^{\text{ucb}}}{\theta_0 + \Theta_{S_t,t}} - \frac{r_i \theta_i}{\theta_0 + \Theta_{S_t}}\right) \mathbf{1}\{\mathcal{E}\}\right] \\ &\leq \sum_{t=1}^T \mathbb{E}\left[\left(\sum_{i \in S_t} \frac{r_i (\theta_{i,t}^{\text{ucb}} - \theta_i)}{\theta_0 + \Theta_{S_t}}\right) \mathbf{1}\{\mathcal{E}\}\right] \end{aligned}$$

237 Because $\Theta_{S_t,t}^{\text{ucb}} \geq \Theta_{S_t}$ under the event \mathcal{E} by Lemma 1. Then, using $r_i \leq 1$, we further upper-
 238 bound using an exploration parameter $\tau_0 = O(\log(T))$ so that the upper-confidence-bounds
 239 in Lemmas 1 and 2 are satisfied

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[(\mathcal{R}(S_t, \theta_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \mathbf{1}\{\mathcal{E}\}] &\leq \sum_{i=1}^K \mathbb{E}\left[\sum_{t=1}^T \left(\frac{|\theta_{i,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}}\right) \mathbf{1}\{i \in S_t, \mathcal{E}\}\right] \\ &\lesssim O(\tau_0) + \sum_{i=1}^K \mathbb{E}\left[\sum_{t=1}^T \frac{|\theta_{i,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}} \mathbf{1}\{i \in S_t, \tau_{i,t} \geq \tau_0, \mathcal{E}\}\right] \\ &\lesssim O(\tau_0) + \sum_{i=1}^K \sqrt{\sum_{t=1}^T \mathbb{E}\left[\frac{\theta_i \mathbf{1}\{i \in S_t\}}{\theta_0 + \Theta_{S_t}}\right]} \times \underbrace{\sqrt{\sum_{t=1}^T \mathbb{E}\left[\left(\frac{\theta_{i,t}^{\text{ucb}} - \theta_i}{\theta_0 + \Theta_{S_t}}\right)^2 \frac{\theta_0 + \Theta_{S_t}}{\theta_i} \mathbf{1}\{i \in S_t, \tau_{i,t} \geq \tau_0, \mathcal{E}\}\right]}}_{=: A_T(i)} \end{aligned} \quad (6)$$

240 where the last inequality is by Cauchy-Schwarz inequality. Now, the term $A_T(i)$ above may
 241 be upper-bounded using Lemmas 1 and 2,

$$\begin{aligned} A_T(i) &= \mathbb{E}\left[\frac{(\theta_{i,t}^{\text{ucb}} - \theta_i)^2}{\theta_i (\theta_0 + \Theta_{S_t})} \mathbf{1}\{i \in S_t, \tau_{i,t} \geq \tau_0, \mathcal{E}\}\right] \lesssim \sum_{t=1}^T \mathbb{E}\left[\frac{(\theta_0 + \theta_i)^2 x}{n_{i0,t} (\theta_0 + \Theta_{S_t})} \mathbf{1}\{i \in S_t\}\right] \\ &\lesssim \theta_{\max} x \sum_{t=1}^T \mathbb{E}\left[\frac{(\theta_0 + \theta_i) \mathbf{1}\{i \in S_t\}}{(\theta_0 + \Theta_{S_t}) n_{i0,t}}\right] = \theta_{\max} x \mathbb{E}\left[\sum_{t=1}^T \frac{\mathbf{1}\{i_t \in \{i, 0\}, i \in S_t\}}{n_{i0,t}}\right] \lesssim \theta_{\max} x \log T \end{aligned}$$

242 where in the last inequality we used that $\sum_{n=1}^T n^{-1} \leq 1 + \log T$. Substituting into (6),
 243 Jensen's inequality entails,

$$\sum_{t=1}^T \mathbb{E}[(\mathcal{R}(S_t, \theta_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \mathbf{1}\{\mathcal{E}\}] \lesssim O(\tau_0) + \mathbb{E}\left[\sqrt{\theta_{\max} x \log T} \sum_{i=1}^K \sqrt{\sum_{t=1}^T \frac{\theta_i \mathbf{1}\{i \in S_t\}}{\theta_0 + \Theta_{S_t}}}\right]. \quad (7)$$

244 The proof is finally concluded by applying Cauchy-Schwarz inequality which yields:

$$\sum_{i=1}^K \sqrt{\sum_{t=1}^T \frac{\theta_i \mathbb{1}\{i \in S_t\}}{\theta_0 + \Theta_{S_t}}} \leq \sqrt{K \sum_{t=1}^T \frac{\sum_{i=1}^K \theta_i \mathbb{1}\{i \in S_t\}}{\theta_0 + \Theta_{S_t}}} \leq \sqrt{KT}.$$

245 Finally, combining the above result with (5) and (7) concludes the proof

$$\text{Reg}_T^{\text{wt d}} \lesssim TP(\mathcal{E}^c) + O(\tau_0) + \sqrt{\theta_{\max} xKT \log T}.$$

246 Choosing $x = 2 \log T$ ensures $TP(\mathcal{E}^c) \leq O(1)$ and $\tau_0 \leq O(\log T)$. \square

247 4 Improved dependance on θ_{\max} with Adaptive Pivot Selection

248 A problem with Algorithm 1 stems from estimating all θ_i based on pairwise comparisons with
 249 item 0. When $\theta_{\max} \gg \theta_0 = 1$, item 0 may not be sampled enough as the winner, leading to
 250 poor estimators. This deficiency contributes to the suboptimal dependence on θ_{\max} observed
 251 in Theorems 3 and 5 and in prior work, such as [2]. We propose the following fix to optimize
 252 the pivot. For all $i, j \in [K] \cup \{0\}$ we define $\gamma_{ij} = \frac{\theta_i}{\theta_j}$, and the estimators:

$$\gamma_{ij,t}^{\text{ucb}} = p_{ij,t}^{\text{ucb}} / (1 - p_{ij,t}^{\text{ucb}})_+ \quad \text{and} \quad \gamma_{ii,t}^{\text{ucb}} = 1,$$

253 where $p_{ij,t}^{\text{ucb}}$ are defined in (3). For all rounds t , the algorithm AOA-RB_{PL}-Adaptive selects

$$S_t = \operatorname{argmax}_{|S| \leq m} \mathcal{R}(S, \hat{\theta}_t^{\text{ucb}}) \quad \text{where} \quad \hat{\theta}_{i,t}^{\text{ucb}} := \min_{j \in [K] \cup \{0\}} \gamma_{ij,t}^{\text{ucb}} \gamma_{j0,t}^{\text{ucb}}.$$

254 We offer below a regret bound that underscores the value of optimizing the pivot when
 255 $\theta_{\max} \gg K$. Note that while the algorithm and analysis are presented for the weighted
 256 objective with winner feedback only, it can be adapted to other objectives by replacing
 257 $\mathcal{R}(S, \theta)$ with the new objective in the analysis, as long as Lemma 4 remains valid.

258 **Theorem 6.** *Let $\theta_{\max} \geq 1$. For any $\theta \in [0, \theta_{\max}]^K$ and weights $\mathbf{r} \in [0, 1]^K$, the weighted
 259 regret of AOA-RB_{PL}-Adaptive is upper-bounded as*

$$\text{Reg}_T^{\text{wt d}} = O(\sqrt{\min\{\theta_{\max}, K\}KT \log T})$$

260 as $T \rightarrow \infty$ for the choice $x = 2 \log T$ (when defining $p_{ij,t}^{\text{ucb}}$).

261 Asymptotically, when θ_{\max} is constant, the regret is $O(K\sqrt{T} \log T)$, eliminating any depen-
 262 dence on θ_{\max} . This allows for handling scenarios where the No-Choice item is highly unlikely,
 263 which is not achievable in previous works such as [2, 1]. [2] did attempt in their Thm. 4 to
 264 relax the assumption of $\theta_{\max} = \theta_0$ and shows a bound of order $O(\max\{\theta_{\max}/\theta_0, 1\}^{1/2} \sqrt{KT})$,
 265 which unfortunately blows to ∞ as $\theta_0 \rightarrow 0$ or equivalently $\theta_{\max} \rightarrow \infty$, leading to a vac-
 266 uous bound. Here, lies the stark improvement and one of the key contributions, as also
 267 corroborated in our experimental evaluation Sec. 5 (Fig. 2).

268 The proof is deferred to the App. B, with a key step relying on selecting the pivot
 269 $j_t = \operatorname{argmax}_{j \in S_t \cup \{0\}} \theta_j$. The use of $|\hat{\theta}_{i,t}^{\text{ucb}} - \theta_i| \leq |\gamma_{ij_t,t}^{\text{ucb}} - \theta_i|$ provides confidence upper-
 270 bounds with an improved dependence on θ_{\max} , leveraging the fact that $\theta_{j_t} \geq \theta_i$. Due
 271 to the varying pivot over time, a telescoping argument introduces an additive factor \sqrt{K} .

272 5 Experiments

273 We provide here a synthetic experiments. All results are averaged across 100 runs. We
 274 evaluate the performance of our main algorithm AOA-RB_{PL}-Adaptive (Sec. 4), referred
 275 as ‘‘Our Alg-1 (Adaptive Pivot)’’, with the following two algorithms: AOA-RB_{PL} (Sec. 3)
 276 referred as ‘‘Our Alg-2 (No-Choice Pivot)’’, and MNL-UCB, the state-of-the-art algorithm
 277 for AOA ([2], Alg. 1).

278 **Different PL (θ) Environments.** We report our experiment results on two datasets with
 279 $K = 50$ items: (1) Arith50 with PL parameters $\theta_i = 1 - (i - 1)0.2$, $\forall i \in [50]$. (2) Bad50

280 with PL parameters $\theta_i = 0.6, \forall i \in [50] \setminus \{25\}$ and $\theta_{25} = 0.8$. For simplicity of computing
 281 the assortment choices S_t , we assume $r_i = 1, \forall i \in [K]$.

282 **(1). Averaged Regret with weak NC ($\theta_{\max}/\theta_0 \gg 1$) (Fig. 1):** In our first experiment,
 283 we set $m = 5$ and $\theta_0/\theta_{\max} = 0.01$ and report the average regret of the above three
 algorithms for our two objectives.

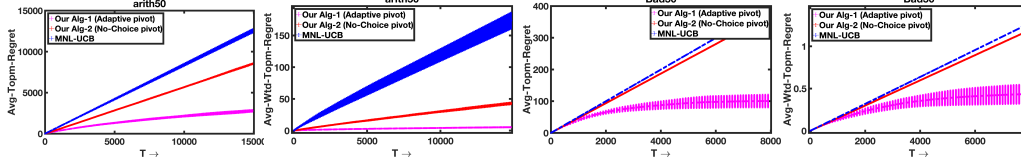
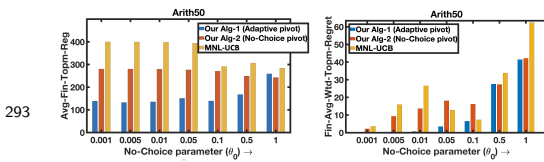


Figure 1: Averaged Regret for $m = 5, \theta_0 = 0.01$

284

285 Fig. 1 shows that our algorithm AOA-RB_{PL}-Adaptive (with adaptive pivot) significantly
 286 outperforms the other two algorithms, while our algorithm AOA-RB_{PL} with no-choice (NC)
 287 pivot still outperforms MNL-UCB.

288 **(2). Averaged Regret vs No-Choice PL Parameter (θ_{\max}/θ_0) (Fig. 2):** In this
 289 experiment, we evaluate the regret performance of our algorithm AOA-RB_{PL}-Adaptive. We
 290 report the experiment on Artith50 PL dataset and set the subsetsize $m = 5, \theta_{\max}/\theta_0 =$
 291 $\{1, 0.5, 0.1, 0.05, 0.01, 0.005, 0.001\}$. Fig. 2 shows the increase in the performance gap between
 292 our algorithm AOA-RB_{PL}-Adaptive (with adaptive pivot) with decreasing θ_0/θ_{\max} .



293

Figure 2: Comparative performance
 for varying $\theta_0/\theta_{\max}, m = 5$

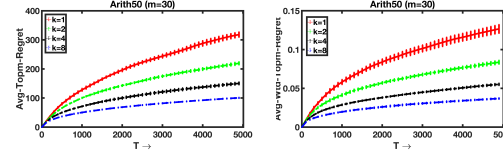


Figure 3: Tradoff: Averaged Regret vs
 length of the k rank-ordered feedback

294 **(3). Averaged Regret vs Length of the rank-ordered feedback (k) (Fig. 3):** We
 295 also run a thought experiment to understand the tradeoff between learning rate with k -length
 296 rank-ordered feedback, where given any assortment $S_t \subseteq [K]$ of size m , the learner gets to
 297 see the top- k draws ($k \leq m$) from the PL model without replacement. This is a stronger
 298 feedback than the winner (i.e. top-1 for $k = 1$) feedback and, as expected, we see in Fig. 3
 299 an improved regret (for both notions) when increasing k . The experiment are run on the
 300 Artith50 dataset with $m = 30$ and $k \in \{1, 2, 4, 8\}$.

301 6 Conclusion

302 We address the Active Optimal Assortment Selection problem with PL choice models, in-
 303 troducing a versatile framework (AOA) that eliminates the need for a strong default item,
 304 typically assumed as the No-Choice (NC) item in the existing literature. Our proposed
 305 algorithms employ a novel 'Rank-Breaking' technique to establish tight concentration guar-
 306 antees for estimating the score parameters of the PL model. Our approach stands out for
 307 its practicality and avoids the suboptimal practice of repeatedly selecting the same set of
 308 items until the default item prevails. This is beneficial when the default item's quality
 309 (θ_0) is significantly lower than the quality of the best item (θ_{\max}). Our algorithms are
 310 computationally efficient, optimal (up to log factors), and free from restrictive assumptions
 311 on the default item.

312 **Future Works.** Among many interesting questions to address in the future, it will be
 313 interesting to understand the role of the No-Choice (NC) item in the algorithm design,
 314 precisely, can we design efficient algorithms without the existence of NC items with a regret
 315 rate still linear in θ_{\max} ? Further, it will be interesting to extend our results to more general
 316 choice models beyond the PL model [18, 22, 23]. What is the tradeoff between the subsetsize
 317 m and the regret for such general choice models? Extending our results to large (potentially
 318 infinite) decision spaces and contextual settings would also be a very useful and practical
 319 contribution to the literature of assortment optimization.

References

- 320
- 321 [1] Shipra Agrawal, Vashist Avadhanula, Vineet Goyal, and Assaf Zeevi. Thompson sampling for
322 the mnl-bandit. In *Conference on learning theory*, pages 76–78. PMLR, 2017.
- 323 [2] Shipra Agrawal, Vashist Avadhanula, Vineet Goyal, and Assaf Zeevi. Mnl-bandit: A dynamic
324 learning approach to assortment selection. *Operations Research*, 67(5):1453–1485, 2019.
- 325 [3] Nir Ailon, Zohar Karnin, and Thorsten Joachims. Reducing dueling bandits to cardinal bandits.
326 In *International Conference on Machine Learning*, pages 856–864. PMLR, 2014.
- 327 [4] Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Exploration–exploitation tradeoff
328 using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19):1876–
329 1902, 2009.
- 330 [5] Peter Auer. Using upper confidence bounds for online learning. In *Foundations of Computer
331 Science, 2000. Proceedings. 41st Annual Symposium on*, pages 270–279. IEEE, 2000.
- 332 [6] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed
333 bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- 334 [7] Vashist Avadhanula, Jalaj Bhandari, Vineet Goyal, and Assaf Zeevi. On the tightness of
335 an lp relaxation for rational optimization and its applications. *Operations Research Letters*,
336 44(5):612–617, 2016.
- 337 [8] Hossein Azari, David Parkes, and Lirong Xia. Random utility theory for social choice. In
338 *Advances in Neural Information Processing Systems*, pages 126–134, 2012.
- 339 [9] Hossein Azari, David Parks, and Lirong Xia. Random utility theory for social choice. *Advances
340 in Neural Information Processing Systems*, 25, 2012.
- 341 [10] Viktor Bengs, Róbert Busa-Fekete, Adil El Mesaoudi-Paul, and Eyke Hüllermeier. Preference-
342 based online learning with dueling bandits: A survey. *Journal of Machine Learning Research*,
343 2021.
- 344 [11] Viktor Bengs, Róbert Busa-Fekete, Adil El Mesaoudi-Paul, and Eyke Hüllermeier. Preference-
345 based online learning with dueling bandits: A survey. *J. Mach. Learn. Res.*, 22:7–1, 2021.
- 346 [12] Viktor Bengs, Aadirupa Saha, and Eyke Hüllermeier. Stochastic contextual dueling bandits
347 under linear stochastic transitivity models. In *International Conference on Machine Learning*,
348 pages 1764–1786. PMLR, 2022.
- 349 [13] Gerardo Berbeglia and Gwenaél Joret. Assortment optimisation under a general discrete choice
350 model: A tight analysis of revenue-ordered assortments. *arXiv preprint arXiv:1606.01371*, 2016.
- 351 [14] Brian Brost, Yevgeny Seldin, Ingemar J. Cox, and Christina Lioma. Multi-dueling bandits and
352 their application to online ranker evaluation. *CoRR*, abs/1608.06253, 2016.
- 353 [15] Niladri S Chatterji, Aldo Pacchiano, Peter L Bartlett, and Michael I Jordan. On the theory of
354 reinforcement learning with once-per-episode feedback. *arXiv preprint arXiv:2105.14363*, 2021.
- 355 [16] Xi Chen, Sivakanth Gopi, Jieming Mao, and Jon Schneider. Competitive analysis of the top-k
356 ranking problem. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on
357 Discrete Algorithms*, pages 1245–1264. SIAM, 2017.
- 358 [17] Xi Chen, Yuanzhi Li, and Jieming Mao. A nearly instance optimal algorithm for top-k ranking
359 under the multinomial logit model. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM
360 Symposium on Discrete Algorithms*, pages 2504–2522. SIAM, 2018.
- 361 [18] Xi Chen, Chao Shi, Yining Wang, and Yuan Zhou. Dynamic assortment planning under nested
362 logit models. *Production and Operations Management*, 30(1):85–102, 2021.
- 363 [19] Xi Chen and Yining Wang. A note on a tight lower bound for mnl-bandit assortment selection
364 models. *arXiv preprint arXiv:1709.06109*, 2017.
- 365 [20] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
366 reinforcement learning from human preferences. *Advances in neural information processing
367 systems*, 30, 2017.

- 368 [21] James Davis, Guillermo Gallego, and Huseyin Topaloglu. Assortment planning under the
369 multinomial logit model with totally unimodular constraint structures. *Work in Progress*, 2013.
- 370 [22] Antoine Désir, Vineet Goyal, Srikanth Jagabathula, and Danny Segev. Assortment optimization
371 under the mallows model. In *Advances in Neural Information Processing Systems*, pages
372 4700–4708, 2016.
- 373 [23] Antoine Désir, Vineet Goyal, Danny Segev, and Chun Ye. Capacity constrained assortment
374 optimization under the markov chain based choice model. *Operations Research*, 2016.
- 375 [24] James A Grant and David S Leslie. Learning to rank under multinomial logit choice. *Journal*
376 *of Machine Learning Research*, 24(260):1–49, 2023.
- 377 [25] Minje Jang, Sunghyun Kim, Changho Suh, and Sewoong Oh. Optimal sample complexity of
378 m-wise data for top-k ranking. In *Advances in Neural Information Processing Systems*, pages
379 1685–1695, 2017.
- 380 [26] Ashish Khetan and Sewoong Oh. Data-driven rank breaking for efficient rank aggregation.
381 *Journal of Machine Learning Research*, 17(193):1–54, 2016.
- 382 [27] Kameng Nip, Zhenbo Wang, and Zizhuo Wang. Assortment optimization under a single
383 transition model. 2017.
- 384 [28] Min-hwan Oh and Garud Iyengar. Thompson sampling for multinomial logit contextual bandits.
385 *Advances in Neural Information Processing Systems*, 32, 2019.
- 386 [29] Mingdong Ou, Nan Li, Shenghuo Zhu, and Rong Jin. Multinomial logit bandit with linear
387 utility functions. *arXiv preprint arXiv:1805.02971*, 2018.
- 388 [30] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,
389 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to
390 follow instructions with human feedback. *Advances in Neural Information Processing Systems*,
391 35:27730–27744, 2022.
- 392 [31] Wenbo Ren, Jia Liu, and Ness B Shroff. PAC ranking from pairwise and listwise queries: Lower
393 bounds and upper bounds. *arXiv preprint arXiv:1806.02970*, 2018.
- 394 [32] Paat Rusmevichientong, Zuo-Jun Max Shen, and David B Shmoys. Dynamic assortment
395 optimization with a multinomial logit choice model and capacity constraint. *Operations*
396 *research*, 58(6):1666–1680, 2010.
- 397 [33] Aadirupa Saha and Suprovat Ghoshal. Exploiting correlation to achieve faster learning rates in
398 low-rank preference bandits. In *International Conference on Artificial Intelligence and Statistics*,
399 pages 456–482. PMLR, 2022.
- 400 [34] Aadirupa Saha and Aditya Gopalan. Active ranking with subset-wise preferences. *International*
401 *Conference on Artificial Intelligence and Statistics (AISTATS)*, 2018.
- 402 [35] Aadirupa Saha and Aditya Gopalan. Combinatorial bandits with relative feedback. In *Advances*
403 *in Neural Information Processing Systems*, 2019.
- 404 [36] Aadirupa Saha and Aditya Gopalan. PAC Battling Bandits in the Plackett-Luce Model. In
405 *Algorithmic Learning Theory*, pages 700–737, 2019.
- 406 [37] Aadirupa Saha and Aditya Gopalan. Best-item learning in random utility models with subset
407 choices. In *International Conference on Artificial Intelligence and Statistics*, pages 4281–4291.
408 PMLR, 2020.
- 409 [38] Hossein Azari Soufiani, David C Parkes, and Lirong Xia. Computing parametric ranking models
410 via rank-breaking. In *ICML*, pages 360–368, 2014.
- 411 [39] Yanan Sui, Vincent Zhuang, Joel Burdick, and Yisong Yue. Multi-dueling bandits with
412 dependent arms. In *Conference on Uncertainty in Artificial Intelligence, UAI’17*, 2017.
- 413 [40] Kalyan Talluri and Garrett Van Ryzin. Revenue management under a general discrete choice
414 model of consumer behavior. *Management Science*, 50(1):15–33, 2004.
- 415 [41] Yisong Yue, Josef Broder, Robert Kleinberg, and Thorsten Joachims. The k -armed dueling
416 bandits problem. *Journal of Computer and System Sciences*, 78(5):1538–1556, 2012.

- 417 [42] Yu-Jie Zhang and Masashi Sugiyama. Online (multinomial) logistic bandit: Improved regret
418 and constant computation cost. *Advances in Neural Information Processing Systems*, 36, 2024.
- 419 [43] Zihan Zhang and Xiangyang Ji. Regret minimization for reinforcement learning by evaluating
420 the optimal bias function. In *Advances in Neural Information Processing Systems*, pages
421 2827–2836, 2019.
- 422 [44] Masrour Zoghi, Zohar S Karnin, Shimon Whiteson, and Maarten De Rijke. Copeland dueling
423 bandits. In *Advances in Neural Information Processing Systems*, pages 307–315, 2015.
- 424 [45] Masrour Zoghi, Shimon Whiteson, Remi Munos, Maarten de Rijke, et al. Relative upper
425 confidence bound for the k -armed dueling bandit problem. In *JMLR Workshop and Conference*
426 *Proceedings*, number 32, pages 10–18. JMLR, 2014.
- 427 [46] Masrour Zoghi, Shimon A Whiteson, Maarten De Rijke, and Remi Munos. Relative confidence
428 sampling for efficient on-line ranker evaluation. In *Proceedings of the 7th ACM international*
429 *conference on Web search and data mining*, pages 73–82. ACM, 2014.

430 **NeurIPS Paper Checklist**

431 **1. Claims**

432 Question: Do the main claims made in the abstract and introduction accurately
433 reflect the paper’s contributions and scope?

434 Answer: [\[Yes\]](#)

435 Justification: In the abstract, we list the main claims of this paper in a general
436 fashion. Then, in the introduction we state them in more detail. They accurately
437 reflect the paper’s contribution and scope.

438 Guidelines:

- 439 • The answer NA means that the abstract and introduction do not include the
440 claims made in the paper.
- 441 • The abstract and/or introduction should clearly state the claims made, including
442 the contributions made in the paper and important assumptions and limitations.
443 A No or NA answer to this question will not be perceived well by the reviewers.
- 444 • The claims made should match theoretical and experimental results, and reflect
445 how much the results can be expected to generalize to other settings.
- 446 • It is fine to include aspirational goals as motivation as long as it is clear that
447 these goals are not attained by the paper.

448 **2. Limitations**

449 Question: Does the paper discuss the limitations of the work performed by the
450 authors?

451 Answer: [\[Yes\]](#)

452 Justification: We discuss the limitations and assumptions of our work throughout
453 the paper. Additional limitations are highlighted in the discussion.

454 Guidelines:

- 455 • The answer NA means that the paper has no limitation while the answer No
456 means that the paper has limitations, but those are not discussed in the paper.
- 457 • The authors are encouraged to create a separate "Limitations" section in their
458 paper.
- 459 • The paper should point out any strong assumptions and how robust the results
460 are to violations of these assumptions (e.g., independence assumptions, noiseless
461 settings, model well-specification, asymptotic approximations only holding
462 locally). The authors should reflect on how these assumptions might be violated
463 in practice and what the implications would be.
- 464 • The authors should reflect on the scope of the claims made, e.g., if the approach
465 was only tested on a few datasets or with a few runs. In general, empirical
466 results often depend on implicit assumptions, which should be articulated.
- 467 • The authors should reflect on the factors that influence the performance of the
468 approach. For example, a facial recognition algorithm may perform poorly when
469 image resolution is low or images are taken in low lighting. Or a speech-to-text
470 system might not be used reliably to provide closed captions for online lectures
471 because it fails to handle technical jargon.
- 472 • The authors should discuss the computational efficiency of the proposed algo-
473 rithms and how they scale with dataset size.
- 474 • If applicable, the authors should discuss possible limitations of their approach
475 to address problems of privacy and fairness.
- 476 • While the authors might fear that complete honesty about limitations might
477 be used by reviewers as grounds for rejection, a worse outcome might be that
478 reviewers discover limitations that aren’t acknowledged in the paper. The
479 authors should use their best judgment and recognize that individual actions in
480 favor of transparency play an important role in developing norms that preserve
481 the integrity of the community. Reviewers will be specifically instructed to not
482 penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: The assumptions can be found in the Problem Setup section and in the paragraphs before the theorems and remarks. We provide proof sketches in the main text and complete proofs in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: experimental details including algorithms and setups are clearly provided. In addition, the main contribution of the paper is theoretical and synthetic experiments are mostly provided as an illustration.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.

- 535 (c) If the contribution is a new model (e.g., a large language model), then there
536 should either be a way to access this model for reproducing the results or a
537 way to reproduce the model (e.g., with an open-source dataset or instructions
538 for how to construct the dataset).
- 539 (d) We recognize that reproducibility may be tricky in some cases, in which
540 case authors are welcome to describe the particular way they provide for
541 reproducibility. In the case of closed-source models, it may be that access to
542 the model is limited in some way (e.g., to registered users), but it should be
543 possible for other researchers to have some path to reproducing or verifying
544 the results.

545 5. Open access to data and code

546 Question: Does the paper provide open access to the data and code, with sufficient
547 instructions to faithfully reproduce the main experimental results, as described in
548 supplemental material?

549 Answer: [No]

550 Justification: experimental setups are synthetic and can easily be reproduced.

551 Guidelines:

- 552 • The answer NA means that paper does not include experiments requiring code.
- 553 • Please see the NeurIPS code and data submission guidelines ([https://nips.
554 cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 555 • While we encourage the release of code and data, we understand that this might
556 not be possible, so “No” is an acceptable answer. Papers cannot be rejected
557 simply for not including code, unless this is central to the contribution (e.g., for
558 a new open-source benchmark).
- 559 • The instructions should contain the exact command and environment needed
560 to run to reproduce the results. See the NeurIPS code and data submis-
561 sion guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>)
562 for more details.
- 563 • The authors should provide instructions on data access and preparation, in-
564 cluding how to access the raw data, preprocessed data, intermediate data, and
565 generated data, etc.
- 566 • The authors should provide scripts to reproduce all experimental results for
567 the new proposed method and baselines. If only a subset of experiments are
568 reproducible, they should state which ones are omitted from the script and why.
- 569 • At submission time, to preserve anonymity, the authors should release
570 anonymized versions (if applicable).
- 571 • Providing as much information as possible in supplemental material (appended
572 to the paper) is recommended, but including URLs to data and code is permitted.

573 6. Experimental Setting/Details

574 Question: Does the paper specify all the training and test details (e.g., data splits,
575 hyperparameters, how they were chosen, type of optimizer, etc.) necessary to
576 understand the results?

577 Answer: [Yes]

578 Justification: All details are provided to reproduce the experiments.

579 Guidelines:

- 580 • The answer NA means that the paper does not include experiments.
- 581 • The experimental setting should be presented in the core of the paper to a level
582 of detail that is necessary to appreciate the results and make sense of them.
- 583 • The full details can be provided either with the code, in appendix, or as
584 supplemental material.

585 7. Experiment Statistical Significance

586 Question: Does the paper report error bars suitably and correctly defined or other
587 appropriate information about the statistical significance of the experiments?

588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639

Answer: [Yes]

Justification: [Yes]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We do not see any potential negative social impact of this work and it follows the NeurIPS code of ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

640 **10. Broader Impacts**

641 Question: Does the paper discuss both potential positive societal impacts and
642 negative societal impacts of the work performed?

643 Answer: [NA]

644 Justification: This work addresses the problem of designing efficient and optimal
645 algorithms for different assortment selection problems with MNL models. Our
646 work is purely theoretical and studies a fundamental mathematical optimization
647 framework that is unrelated to societal considerations

648 Guidelines:

- 649 • The answer NA means that there is no societal impact of the work performed.
- 650 • If the authors answer NA or No, they should explain why their work has no
651 societal impact or why the paper does not address societal impact.
- 652 • Examples of negative societal impacts include potential malicious or unintended
653 uses (e.g., disinformation, generating fake profiles, surveillance), fairness consid-
654 erations (e.g., deployment of technologies that could make decisions that unfairly
655 impact specific groups), privacy considerations, and security considerations.
- 656 • The conference expects that many papers will be foundational research and
657 not tied to particular applications, let alone deployments. However, if there
658 is a direct path to any negative applications, the authors should point it out.
659 For example, it is legitimate to point out that an improvement in the quality
660 of generative models could be used to generate deepfakes for disinformation.
661 On the other hand, it is not needed to point out that a generic algorithm for
662 optimizing neural networks could enable people to train models that generate
663 Deepfakes faster.
- 664 • The authors should consider possible harms that could arise when the technology
665 is being used as intended and functioning correctly, harms that could arise when
666 the technology is being used as intended but gives incorrect results, and harms
667 following from (intentional or unintentional) misuse of the technology.
- 668 • If there are negative societal impacts, the authors could also discuss possible
669 mitigation strategies (e.g., gated release of models, providing defenses in addition
670 to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a
671 system learns from feedback over time, improving the efficiency and accessibility
672 of ML).

673 **11. Safeguards**

674 Question: Does the paper describe safeguards that have been put in place for
675 responsible release of data or models that have a high risk for misuse (e.g., pretrained
676 language models, image generators, or scraped datasets)?

677 Answer: [NA]

678 Justification: [NA]

679 Guidelines:

- 680 • The answer NA means that the paper poses no such risks.
- 681 • Released models that have a high risk for misuse or dual-use should be released
682 with necessary safeguards to allow for controlled use of the model, for example
683 by requiring that users adhere to usage guidelines or restrictions to access the
684 model or implementing safety filters.
- 685 • Datasets that have been scraped from the Internet could pose safety risks. The
686 authors should describe how they avoided releasing unsafe images.
- 687 • We recognize that providing effective safeguards is challenging, and many papers
688 do not require this, but we encourage authors to take this into account and
689 make a best faith effort.

690 **12. Licenses for existing assets**

691 Question: Are the creators or original owners of assets (e.g., code, data, models),
692 used in the paper, properly credited and are the license and terms of use explicitly
693 mentioned and properly respected?

694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

Supplementary: Optimal, Efficient and Practical Algorithms for Assortment Optimization

761
762

A Preliminaries: Some Useful Concepts for PL choice models

A.1 Plackett-Luce (PL): A Discrete Choice Model

765 A discrete choice model specifies the relative preferences of two or more discrete alternatives
766 in a given set. A widely studied class of discrete choice models is the class of *Random*
767 *Utility Models* (RUMs), which assume a ground-truth utility score $\theta_i \in \mathbb{R}$ for each alternative
768 $i \in [n]$, and assign a conditional distribution $\mathcal{D}_i(\cdot|\theta_i)$ for scoring item i . To model a winning
769 alternative given any set $S \subseteq [n]$, one first draws a random utility score $X_i \sim \mathcal{D}_i(\cdot|\theta_i)$ for
770 each alternative in S , and selects an item with the highest random score.

One widely used RUM is the *Multinomial-Logit (MNL)* or *Plackett-Luce model (PL)*, where the \mathcal{D}_i s are taken to be independent Gumbel distributions with parameters θ'_i [8], i.e., with probability densities

$$\mathcal{D}_i(x_i|\theta'_i) = e^{-(x_i-\theta'_i)} e^{-e^{-(x_i-\theta'_i)}}, \quad \theta'_i \in \mathbb{R}, \forall i \in [n].$$

771 Moreover assuming $\theta'_i = \ln \theta_i$, $\theta_i > 0 \forall i \in [n]$, it can be shown in this case the probability
772 that an alternative i emerges as the winner in the set $S \ni i$ becomes: $\mathbb{P}(i|S) = \frac{\theta_i}{\sum_{j \in S} \theta_j}$.

773 Other families of discrete choice models can be obtained by imposing different probability
774 distributions over the utility scores X_i , e.g. if $(X_1, \dots, X_n) \sim \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{\Lambda})$ are jointly normal
775 with mean $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)$ and covariance $\boldsymbol{\Lambda} \in \mathbb{R}^{n \times n}$, then the corresponding RUM-based
776 choice model reduces to the *Multinomial Probit (MNP)*.

A.2 Rank Breaking

778 *Rank breaking* (RB) is a well-understood idea involving the extraction of pairwise comparisons
779 from (partial) ranking data, and then building pairwise estimators on the obtained pairs by
780 treating each comparison independently [26, 25], e.g., a winner a sampled from among a, b, c is
781 rank-broken into the pairwise preferences $a \succ b$, $a \succ c$. We use this idea to devise estimators
782 for the pairwise win probabilities $p_{ij} = \mathbb{P}(i|\{i, j\}) = \theta_i/(\theta_i + \theta_j)$ for our problem setting.
783 We used the idea of RB in both our algorithms (AOA-RB_{PL} and AOA-RB_{PL}-Adaptive) to
784 update the pairwise win-count estimates $w_{i,j,t}$ for all the item pairs $(i, j) \in [K] \times [K]$, which
785 is further used for deriving the empirical pairwise preference estimates $\hat{p}_{i,j,t}$, at any time t .

A.3 Parameter Estimation with PL based preference data

787 **Lemma 7** (Pairwise win-probability estimates for the PL model [34]). *Consider a Plackett-*
788 *Luce choice model with parameters $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)$, and fix two items $i, j \in [n]$. Let*
789 *S_1, \dots, S_T be a sequence of (possibly random) subsets of $[n]$ of size at least 2, where T is*
790 *a positive integer, and i_1, \dots, i_T a sequence of random items with each $i_t \in S_t$, $1 \leq t \leq T$,*
791 *such that for each $1 \leq t \leq T$, (a) S_t depends only on S_1, \dots, S_{t-1} , and (b) i_t is distributed*
792 *as the Plackett-Luce winner of the subset S_t , given $S_1, i_1, \dots, S_{t-1}, i_{t-1}$ and S_t , and (c)*
793 *$\forall t : \{i, j\} \subseteq S_t$ with probability 1. Let $n_i(T) = \sum_{t=1}^T \mathbb{P}(i_t = i)$ and $n_{ij}(T) = \sum_{t=1}^T \mathbb{P}(\{i_t \in$
794 *$\{i, j\}\})$. Then, for any positive integer v , and $\eta \in (0, 1)$,**

$$\mathbb{P}\left(\frac{n_i(T)}{n_{ij}(T)} - \frac{\theta_i}{\theta_i + \theta_j} \geq \eta, n_{ij}(T) \geq v\right) \leq e^{-2v\eta^2},$$

$$\mathbb{P}\left(\frac{n_i(T)}{n_{ij}(T)} - \frac{\theta_i}{\theta_i + \theta_j} \leq -\eta, n_{ij}(T) \geq v\right) \leq e^{-2v\eta^2}.$$

795 **B Omitted Proofs from Sec. 3 and Sec. 4**

796 **B.1 A concentration bounds for the $p_{ij,t}$**

797 We first prove below a concentration inequality based on Bernstein's inequality for the
798 estimators $p_{ij,t}$.

799 **Lemma 8.** *Let $(i, j) \in [K] \times [K]$. Let $T \geq 1$ and $x > 0$. Then, with probability at least
800 $1 - 3Te^{-x}$,*

$$p_{ij} \leq p_{ij,t}^{\text{ucb}} \leq p_{ij} + 2\sqrt{\frac{2p_{ij}(1-p_{ij})x}{n_{ij,t}}} + \frac{11x}{n_{ij,t}}, \quad (8)$$

801 *simultaneously for all $t \in [T]$.*

802 *Proof of Lemma 8.* Let $T \geq 1$, $x > 0$ and $i, j \in [K]$. Applying Thm. 1 of [4], with probability
803 at least $1 - \beta(x, T)$, we get simultaneously for all $t \in [T]$,

$$|\widehat{p}_{ij,t} - p_{ij}| \leq \sqrt{\frac{2\widehat{p}_{ij,t}(1-\widehat{p}_{ij,t})x}{n_{ij,t}}} + \frac{3x}{n_{ij,t}}, \quad (9)$$

804 where $\beta(x, T) = 3 \inf_{1 < \alpha \leq 3} \min \left\{ \frac{\log T}{\log \alpha}, T \right\} e^{-x/\alpha} \leq 3Te^{-x}$. Note that the inequality holds
805 true although $n_{ij,t}$ is a random variable. This, shows the first inequality

$$p_{ij} \leq p_{ij,t}^{\text{ucb}}.$$

806 For the second inequality, (9) implies

$$\begin{aligned} p_{ij,t}^{\text{ucb}} &= \widehat{p}_{ij,t} + \sqrt{\frac{2\widehat{p}_{ij,t}(1-\widehat{p}_{ij,t})x}{n_{ij,t}}} + \frac{3x}{n_{ij,t}} \\ &\leq p_{ij} + 2\sqrt{\frac{2\widehat{p}_{ij,t}(1-\widehat{p}_{ij,t})x}{n_{ij,t}}} + \frac{6x}{n_{ij,t}}. \end{aligned} \quad (10)$$

807 Furthermore, because $x \mapsto x(1-x)$ is 1-Lipschitz on $[0, 1]$, we have

$$\begin{aligned} |\widehat{p}_{ij,t}(1-\widehat{p}_{ij,t}) - p_{ij}(1-p_{ij})| &\leq |\widehat{p}_{ij,t} - p_{ij}| \\ &\stackrel{(9)}{\leq} \sqrt{\frac{2\widehat{p}_{ij,t}(1-\widehat{p}_{ij,t})x}{n_{ij,t}}} + \frac{3x}{n_{ij,t}}. \end{aligned}$$

808 Therefore,

$$\begin{aligned} \widehat{p}_{ij,t}(1-\widehat{p}_{ij,t}) &\leq p_{ij}(1-p_{ij}) + \sqrt{\frac{2\widehat{p}_{ij,t}(1-\widehat{p}_{ij,t})x}{n_{ij,t}}} + \frac{3x}{n_{ij,t}} \\ &\leq \left(\sqrt{p_{ij}(1-p_{ij})} + \sqrt{\frac{3x}{n_{ij,t}}} \right)^2, \end{aligned}$$

809 which yields

$$\sqrt{\widehat{p}_{ij,t}(1-\widehat{p}_{ij,t})} \leq \sqrt{p_{ij}(1-p_{ij})} + \sqrt{\frac{3x}{n_{ij,t}}}. \quad (11)$$

810 Plugging back into (10), we get

$$p_{ij,t}^{\text{ucb}} \leq 2\sqrt{\frac{2p_{ij}(1-p_{ij})x}{n_{ij,t}}} + \frac{11x}{n_{ij,t}}.$$

811 □

812 **B.2 Proof of Lemma 1**

813 *Proof.* Let $i \in [K]$ and $x > 0$. Then, by a union bound on Lemma 8 and 2, with probability
 814 at least $1 - 4Te^{-x}$, (8) and (4) hold true for all $t \in [T]$. We consider this high-probability
 815 event in the rest of the proof. Define the function $f : x \mapsto x/(1-x)_+$ on $[0, 1]$ (with the
 816 convention $f(1) = +\infty$), so that $\theta_{i,t}^{\text{ucb}} = f(p_{i0,t}^{\text{ucb}})$ and $\theta_i = f(p_{i0})$. Because f is non-decreasing,
 817 and $p_{i0,t}^{\text{ucb}} \geq p_{i0}$ by (8), we have

$$\theta_{i,t}^{\text{ucb}} \geq \theta_i. \quad (12)$$

818 Furthermore, denote

$$\Delta_{i,t} := 2\sqrt{\frac{2p_{ij}(1-p_{ij})x}{n_{i0,t}}} + \frac{11x}{n_{i0,t}} = 2\sqrt{\frac{2\theta_0\theta_i x}{(\theta_0 + \theta_i)^2 n_{i0,t}}} + \frac{11x}{n_{i0,t}}. \quad (13)$$

819 In the rest of the proof we assume, $n_{i0,t} \geq 69x(\theta_0 + \theta_i)$. Then, using that $\theta_0\theta_i \leq \theta_0 + \theta_i$
 820 since $\theta_0 = 1$, it implies

$$(\theta_0 + \theta_i)\Delta_{i,t} \leq 2\sqrt{\frac{2\theta_0\theta_i x}{n_{i0,t}}} + \frac{11x(\theta_0 + \theta_i)}{n_{i0,t}} \leq \frac{1}{2},$$

821 and

$$p_{i0} + \Delta_{i,t} = \frac{\theta_i}{\theta_0 + \theta_i} + \Delta_{i,t} \leq \frac{\theta_i + 1/2}{\theta_i + 1} < 1.$$

822 Thus, because f is non-decreasing

$$\begin{aligned} \theta_{i,t}^{\text{ucb}} - \theta_i &= f(p_{i0,t}^{\text{ucb}}) - f(p_{i0}) \\ &\stackrel{(8)}{\leq} f(p_{i0} + \Delta_{i,t}) - f(p_{i0}) \\ &= \frac{p_{i0} + \Delta_{i,t}}{1 - p_{i0} - \Delta_{i,t}} - \frac{p_{i0}}{1 - p_{i0}} \\ &= \frac{\Delta_{i,t}}{(1 - p_{i0})(1 - p_{i0} - \Delta_{i,t})} \\ &= \frac{(\theta_0 + \theta_i)^2 \Delta_{i,t}}{1 - (\theta_0 + \theta_i)\Delta_{i,t}} \\ &\leq 2(\theta_0 + \theta_i)^2 \Delta_{i,t} \\ &\stackrel{(13)}{\leq} 4(\theta_0 + \theta_i) \sqrt{\frac{2\theta_0\theta_i x}{n_{i0,t}}} + \frac{22x(\theta_0 + \theta_i)^2}{n_{i0,t}}, \end{aligned}$$

823 which concludes the proof. \square

824 **B.3 Proof of Lemma 2**

825 *Proof.* Let $T \geq 1$ and $i \in [K]$. Recall that $\tau_{i,t} = \sum_{s=1}^{t-1} \mathbb{1}\{i \in S_s\}$ is the number of times i
 826 was played at the start of round t and $n_{i0,t} = \sum_{s=1}^{t-1} \mathbb{1}\{i_t \in \{i, 0\}, i \in S_t\}$ is the number of
 827 times i or 0 won up to round t when played together. When i is played the probability of 0
 828 or i to win is

$$\mathbb{P}(i_t \in \{i, 0\} | S_t) = \frac{\theta_0 + \theta_i}{\theta_0 + \Theta_{S_t}} \geq \frac{\theta_0 + \theta_i}{\theta_0 + \Theta_{S^*}}.$$

829 Therefore, applying Chernoff-Hoeffding inequality together with a union bound (to deal with
 830 the fact that $\tau_{i,t}$ is random), we have with probability at least $1 - Te^{-x}$

$$n_{i0,t} \geq \frac{\theta_0 + \theta_i}{\theta_0 + \Theta_{S^*}} \tau_{i,t} - \sqrt{\frac{\tau_{i,t} x}{2}}$$

831 simultaneously for all $t \in [T]$. Noting that

$$\frac{\theta_0 + \theta_i}{\theta_0 + \Theta_{S^*}} \tau_{i,t} - \sqrt{\frac{\tau_{i,t} x}{2}} \geq \frac{\theta_0 + \theta_i}{2(\theta_0 + \Theta_{S^*})} \tau_{i,t}$$

832 if $\tau_{i,t} \geq 2x(\theta_0 + \Theta_{S^*})^2 \geq \frac{2x(\theta_0 + \Theta_{S^*})^2}{(\theta_0 + \theta_i)^2}$ concludes the proof. \square

833 **B.4 Proof of Theorem 3**

834 *Proof.* Let us define for any $S \subseteq [K]$,

$$\Theta_S = \sum_{i \in S} \theta_i, \quad \text{and} \quad \Theta_S^{\text{ucb}} := \sum_{i \in S} \theta_i^{\text{ucb}}.$$

835 Let \mathcal{E} be the high-probability event such that both Lemma 1 and 2 holds true. Then,
 836 $\mathbb{P}(\mathcal{E}) \geq 1 - 4TKe^{-x}$. Let us first assume that \mathcal{E} holds true. Then, by Lemma 1,

$$\begin{aligned} \text{Reg}_T^{\text{top}} &= \frac{1}{m} \sum_{t=1}^T \Theta_{S^*} - \Theta_{S_t} \\ &\leq \frac{1}{m} \sum_{t=1}^T \min \left\{ \Theta_{S^*}, \Theta_{S_t}^{\text{ucb}} - \Theta_{S_t} \right\} \leftarrow \text{because } \Theta_{S^*} \leq \Theta_{S^*}^{\text{ucb}} \leq \Theta_{S_t}^{\text{ucb}} \text{ under the event } \mathcal{E} \\ &= \frac{1}{m} \sum_{t=1}^T \min \left\{ \Theta_{S^*}, \sum_{i \in S_t} \theta_{i,t}^{\text{ucb}} - \theta_i \right\} \\ &\leq \frac{1}{m} \Theta_{S^*} \sum_{i=1}^K \bar{\tau}_{i0} + \frac{1}{m} \sum_{t=1}^T \sum_{i \in S_t} (\theta_{i,t}^{\text{ucb}} - \theta_i) \mathbb{1}\{\tau_{i,t} \geq \bar{\tau}_{i0}\} \end{aligned}$$

837 where $\bar{\tau}_{i0} = 2x(\theta_0 + \Theta_{S^*}) \max\{\theta_0 + \Theta_{S^*}, 69\} \leq 138x(m+1)^2 \theta_{\max}^2$, where $\theta_{\max} := \max_i \theta_i$.
 838 Then, noting that if \mathcal{E} holds true, by Lemma 2, we also have $n_{i0,t} \geq \frac{1}{2(\theta_0 + \Theta_{S^*})} (\theta_0 + \theta_i) \tau_{i,t}$,
 839 which yields

$$\mathbb{1}\{\tau_{i,t} \geq \bar{\tau}_{i0}\} \leq \mathbb{1}\{n_{i0,t} \geq 69x(\theta_0 + \theta_i)\}.$$

840 Therefore, we can apply Lemma 1 that entails,

$$\begin{aligned} &\frac{1}{m} \sum_{t=1}^T \sum_{i \in S_t} (\theta_{i,t}^{\text{ucb}} - \theta_i) \mathbb{1}\{\tau_{i,t} \geq \bar{\tau}_{i0}\} \\ &\stackrel{\text{Lem. 1}}{\leq} \frac{1}{m} \sum_{t=1}^T \sum_{i \in S_t} \left(4(\theta_0 + \theta_i) \sqrt{\frac{2\theta_0 \theta_i x}{n_{i0,t}}} + \frac{22x(\theta_0 + \theta_i)^2}{n_{i0,t}} \right) \mathbb{1}\{n_{i0,t} \geq 69x(\theta_0 + \theta_i)\} \\ &\stackrel{\text{Lem 2}}{\leq} \frac{1}{m} \sum_{t=1}^T \sum_{i \in S_t} \left(8 \sqrt{\frac{(\theta_0 + \Theta_{S^*})(\theta_0 + \theta_i) \theta_0 \theta_i x}{\tau_{i,t}}} + \frac{44x(\theta_0 + \Theta_{S^*})(\theta_0 + \theta_i)}{\tau_{i,t}} \right) \\ &\leq \frac{1}{m} \sum_{i=1}^K 16 \sqrt{(\theta_0 + \Theta_{S^*})(\theta_0 + \theta_i) \theta_0 \theta_i x \tau_{i,T}} + 44x(\theta_0 + \Theta_{S^*}) \sum_{i=1}^K (\theta_0 + \theta_i) (1 + \log(\tau_{i,T})), \end{aligned}$$

841 where we used $\sum_{i=1}^n 1/\sqrt{i} \leq 2\sqrt{n}$ and $\sum_{i=1}^n i^{-1} \leq 1 + \log n$. We thus have

$$\begin{aligned} \text{Reg}_T^{\text{top}} &\leq 138x(m+1)^2 K \theta_{\max}^3 + \frac{1}{m} \sum_{i=1}^K 16\theta_{\max}^{3/2} \sqrt{(m+1)x\tau_{i,T}} \\ &\quad + 44x(m+1)(1 + \theta_{\max})^2 \sum_{i=1}^K (1 + \log(\tau_{i,T})) \\ &\leq 138x(m+1)^2 K \theta_{\max}^3 + 16\theta_{\max}^{3/2} \sqrt{2xKT} + 88x(m+1)K\theta_{\max}^2 \left(1 + \log \left(\frac{mT}{K} \right) \right). \end{aligned}$$

842 Therefore,

$$\begin{aligned} \mathbb{E}[\text{Reg}_T^{\text{top}}] &\leq 12\sqrt{2}xmK\theta_{\max}^3 + 16\theta_{\max}^{3/2}\sqrt{2xKT} + 88xmK\theta_{\max}^2 \left(1 + \log \left(\frac{mT}{K} \right) \right) \\ &\quad + 4mKT^2 e^{-x} \theta_{\max}. \end{aligned}$$

843 Choosing $x = 2 \log T$ concludes the proof. \square

844 **B.5 Proof of Theorem 5**

845 *Proof.* Let \mathcal{E} be the high-probability event such that Lemma 1 and 2 are satisfied, so that
 846 $\mathbb{P}(\mathcal{E}) \geq 1 - 4KT e^{-x}$. Then, denoting $x \wedge y := \min\{x, y\}$,

$$\begin{aligned} \text{Reg}_T^{\text{wd}} &= \sum_{t=1}^T \mathbb{E}[\mathcal{R}(S^*, \theta) - \mathcal{R}(S_t, \theta)] \\ &= \sum_{t=1}^T \mathbb{E}[(\mathcal{R}(S^*, \theta) - \mathcal{R}(S_t, \theta))\mathbb{1}\{\mathcal{E}\} + (\mathcal{R}(S^*, \theta) - \mathcal{R}(S_t, \theta))\mathbb{1}\{\mathcal{E}^c\}] \\ &\leq \sum_{t=1}^T \mathbb{E}\left[\left((\mathcal{R}(S_t, \theta_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge \mathcal{R}(S^*, \theta)\right)\mathbb{1}\{\mathcal{E}\} + \mathcal{R}(S^*, \theta)\mathbb{1}\{\mathcal{E}^c\}\right] \end{aligned} \quad (14)$$

847 because $\mathcal{R}(S_t, \theta_t^{\text{ucb}}) \geq \mathcal{R}(S^*, \theta_t^{\text{ucb}}) \geq \mathcal{R}(S^*, \theta)$ under the event \mathcal{E} by Lemma 4. Then, using
 848 $\mathcal{R}(S^*, \theta) \leq \max_i r_i \leq 1$, we get

$$\begin{aligned} \text{Reg}_T^{\text{wd}} &\leq \sum_{t=1}^T \mathbb{E}\left[\left((\mathcal{R}(S_t, \theta_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge 1\right)\mathbb{1}\{\mathcal{E}\} + \mathbb{1}\{\mathcal{E}^c\}\right] \\ &\leq 4T^2 K e^{-x} + \sum_{t=1}^T \mathbb{E}\left[\left((\mathcal{R}(S_t, \theta_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right]. \end{aligned}$$

849 Let us upper-bound the second term of the right-hand-side

$$\begin{aligned} &\sum_{t=1}^T \mathbb{E}\left[\left((\mathcal{R}(S_t, \theta_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right] \\ &= \sum_{t=1}^T \mathbb{E}\left[\left(\left(\sum_{i \in S_t} \frac{r_i \theta_{i,t}^{\text{ucb}}}{\theta_0 + \Theta_{S_t,t}} - \frac{r_i \theta_i}{\theta_0 + \Theta_{S_t}}\right) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right] \\ &\leq \sum_{t=1}^T \mathbb{E}\left[\left(\left(\sum_{i \in S_t} \frac{r_i (\theta_{i,t}^{\text{ucb}} - \theta_i)}{\theta_0 + \Theta_{S_t}}\right) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right] && \text{because } \Theta_{S_t,t}^{\text{ucb}} \geq \Theta_{S_t} \text{ under } \mathcal{E} \\ &\leq \sum_{t=1}^T \mathbb{E}\left[\left(\left(\sum_{i \in S_t} \frac{|\theta_{i,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}}\right) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right] && \text{because } r_i \leq 1 \\ &\leq \sum_{i=1}^K \mathbb{E}\left[\sum_{t=1}^T \left(\frac{|\theta_{i,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}} \wedge 1\right)\mathbb{1}\{i \in S_t\}\mathbb{1}\{\mathcal{E}\}\right] \\ &\leq 138xm^2 K \theta_{\max}^2 + \sum_{i=1}^K \mathbb{E}\left[\sum_{t=1}^T \frac{|\theta_{i,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}} \mathbb{1}\{i \in S_t, \tau_{i,t} \geq 138x(m+1)^2 \theta_{\max}^2\}\mathbb{1}\{\mathcal{E}\}\right] \\ &\leq 138xm^2 K \theta_{\max}^2 + \sum_{i=1}^K \sqrt{\sum_{t=1}^T \mathbb{E}\left[\frac{(\frac{\theta_0}{m} + \theta_i)\mathbb{1}\{i \in S_t\}}{\theta_0 + \Theta_{S_t}}\right]} \\ &\quad \times \underbrace{\sqrt{\sum_{t=1}^T \mathbb{E}\left[\left(\frac{|\theta_{i,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}}\right)^2 \frac{\theta_0 + \Theta_{S_t}}{\frac{\theta_0}{m} + \theta_i} \mathbb{1}\{i \in S_t, \tau_{i,t} \geq 138x(m+1)^2 \theta_{\max}^2\}\mathbb{1}\{\mathcal{E}\}\right]}}_{=: A_T(i)} \end{aligned} \quad (16)$$

850 where the last inequality is by Cauchy-Schwarz inequality. Now, the term $A_T(i)$ above may
 851 be upper-bounded as follows

$$A_T(i) := \sum_{t=1}^T \mathbb{E}\left[\left(\frac{|\theta_{i,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}}\right)^2 \frac{\theta_0 + \Theta_{S_t}}{\frac{\theta_0}{m} + \theta_i} \mathbb{1}\{i \in S_t, \tau_{i,t} \geq 138x(m+1)^2 \theta_{\max}^2\}\mathbb{1}\{\mathcal{E}\}\right]$$

$$= \mathbb{E} \left[\frac{(\theta_{i,t}^{\text{ucb}} - \theta_i)^2}{\left(\frac{\theta_0}{m} + \theta_i\right)\theta_0 + \Theta_{S_t}} \mathbf{1}\{i \in S_t, \tau_{i,t} \geq 138x(m+1)^2\theta_{\max}^2\} \mathbf{1}\{\mathcal{E}\} \right].$$

852 Now, since under the event \mathcal{E} by Lemma 2, $\tau_{i,t} \geq 138x(m+1)^2\theta_{\max}^2$ implies

$$n_{i0,t} \geq 69x(\theta_0 + \theta_i)(m+1)\theta_{\max} \geq 69x(\theta_0 + \theta_i).$$

853 Therefore, we can apply Lemma 1, which further upper-bounds

$$\begin{aligned} A_T(i) &\leq \sum_{t=1}^T \mathbb{E} \left[\left(\frac{2^6(\theta_0 + \theta_i)^2 x}{n_{i0,t}} + \frac{2(22x)^2(\theta_0 + \theta_i)^4}{n_{i0,t}^2 \left(\frac{\theta_0}{m} + \theta_i\right)} \right) \right. \\ &\quad \left. \times \frac{\mathbf{1}\{i \in S_t, \tau_{i,t} \geq 138x(m+1)^2\theta_{\max}^2\} \mathbf{1}\{\mathcal{E}\}}{\theta_0 + \Theta_{S_t}} \right] \\ &\leq \sum_{t=1}^T \mathbb{E} \left[\left(\frac{2^6(\theta_0 + \theta_i)^2 x}{n_{i0,t}} + \frac{15x(\theta_0 + \theta_i)^3}{n_{i0,t}\theta_{\max}(\theta_0 + m\theta_i)} \right) \times \frac{\mathbf{1}\{i \in S_t\} \mathbf{1}\{\mathcal{E}\}}{\theta_0 + \Theta_{S_t}} \right] \end{aligned}$$

854 where we used $n_{i0,t} \geq 69x(\theta_0 + \theta_i)m\theta_{\max}$ in the last inequality. Then, we get

$$\begin{aligned} A_T(i) &\leq \sum_{t=1}^T \mathbb{E} \left[\left(\frac{(\theta_0 + \theta_i)^2 x}{n_{i0,t}} + \frac{30x(\theta_0 + \theta_i)}{n_{i0,t}} \right) \times \frac{\mathbf{1}\{i \in S_t\} \mathbf{1}\{\mathcal{E}\}}{\theta_0 + \Theta_{S_t}} \right] \\ &\leq (94 + 64\theta_i)x \sum_{t=1}^T \mathbb{E} \left[\frac{(\theta_0 + \theta_i) \mathbf{1}\{i \in S_t\}}{(\theta_0 + \Theta_{S_t})n_{i0,t}} \right] \\ &= (94 + 64\theta_i)x \mathbb{E} \left[\sum_{t=1}^T \frac{\mathbf{1}\{i_t \in \{i, 0\}, i \in S_t\}}{n_{i0,t}} \right] \\ &= (94 + 64\theta_i)x \mathbb{E} [1 + \log(n_{i0}(T))] \\ &\leq 158\theta_{\max}x(1 + \log T). \end{aligned}$$

855 Substituting into (16), we then obtain using Cauchy-Schwarz inequality,

$$\begin{aligned} &\sum_{t=1}^T \mathbb{E} \left[\left((\mathcal{R}(S_t, \theta_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge 1 \right) \mathbf{1}\{\mathcal{E}\} \right] \\ &\leq 138xm^2K\theta_{\max}^2 + 13\sqrt{\theta_{\max}x(1 + \log T)} \sum_{i=1}^K \sqrt{\sum_{t=1}^T \mathbb{E} \left[\frac{(\frac{\theta_0}{m} + \theta_i) \mathbf{1}\{i \in S_t\}}{\theta_0 + \Theta_{S_t}} \right]} \\ &\leq 138xm^2K\theta_{\max}^2 + 13\sqrt{\theta_{\max}x(1 + \log T)} \sqrt{\mathbb{E} \left[K \sum_{t=1}^T \frac{\sum_{i=1}^K (\frac{\theta_0}{m} + \theta_i) \mathbf{1}\{i \in S_t\}}{\theta_0 + \Theta_{S_t}} \right]} \\ &= 138xm^2K\theta_{\max}^2 + 13\sqrt{\theta_{\max}x(1 + \log T)KT}. \end{aligned}$$

856 Finally, replacing into Inequality (15) yields

$$\text{Reg}_T^{\text{rtd}} \leq 4T^2K e^{-x} + 138xm^2K\theta_{\max}^2 + 13\sqrt{\theta_{\max}x(1 + \log T)KT}.$$

857 Choosing $x = 2 \log T$ concludes the proof. \square

858 B.6 Proof of Theorem 6

859 The proof follows the one of Theorem 5, except that the concentration lemmas should be
 860 generalized to any pairs (i, j) instead of only with respect to item 0, whose proofs are left
 861 to the reader and closely follows the one of Lemma 1 and 2. For simplicity, this proof is
 862 performed up to universal multiplicative constants, using the rough inequality \lesssim .

863 **Lemma 9.** Let $T \geq 1$ and $x > 0$. Then, with probability at least $1 - 3K(K+1)Te^{-x}$,
864 simultaneously for all $t \in [T]$ and $i \neq j$ in $[K]$: $\gamma_{ij} := \frac{\theta_i}{\theta_j} \leq \gamma_{ij,t}^{\text{ucb}}$ and one of the following
865 two inequalities is satisfied

$$n_{ij,t} < 69x(1 + \gamma_{ij}) \quad \text{or} \quad \gamma_{ij,t}^{\text{ucb}} \leq \gamma_{ij} + 4(\gamma_{ij} + 1) \sqrt{\frac{2\gamma_{ij}x}{n_{ij,t}}} + \frac{22x(\gamma_{ij} + 1)^2}{n_{ij,t}}.$$

866 **Lemma 10.** Let $T \geq 1$ and $x > 0$. Then, with probability at least $1 - 3K(K+1)Te^{-x}$,
867 simultaneously for all $t \in [T]$ and $i \in [K]$: $\hat{\theta}_{i,t}^{\text{ucb}} := \min_j \gamma_{ij,t}^{\text{ucb}} \gamma_{j0,t}^{\text{ucb}} \geq \theta_i$ and for all j one of
868 the following two inequalities is satisfied

$$n_{ij,t} \lesssim x(1 + \gamma_{ij}) \quad \text{or} \quad n_{j0,t} \lesssim x(1 + \theta_j)^2 \theta_j^{-1}$$

869 or

$$\gamma_{ij,t}^{\text{ucb}} \gamma_{j0,t}^{\text{ucb}} - \theta_i \lesssim \sqrt{(\gamma_{ij} + 1)\theta_i x} \left(\sqrt{\frac{(\theta_i + \theta_j)}{n_{ij,t}}} + \sqrt{\frac{(1 + \theta_j)}{n_{j0,t}}} \right) + (\gamma_{ij} + 1) \frac{(\theta_i + \theta_j)x}{n_{ij,t}} + \frac{\gamma_{ij}(1 + \theta_j)^2 x}{n_{j0,t}}.$$

870 *Proof of Lemma 10.* The proof follows from Lemma 9. If $n_{ij,t} > Cx(1 + \gamma_{ij})$ and $n_{j0,t} >$
871 $Cx(1 + \theta_j)$ for some large enough constant C , we have

$$\gamma_{ij,t}^{\text{ucb}} \leq \gamma_{ij} + 4(\gamma_{ij} + 1) \sqrt{\frac{2\gamma_{ij}x}{n_{ij,t}}} + \frac{22x(\gamma_{ij} + 1)^2}{n_{ij,t}}$$

872 and

$$\gamma_{j0,t}^{\text{ucb}} \leq \gamma_{j0} + 4(\gamma_{j0} + 1) \sqrt{\frac{2\gamma_{j0}x}{n_{j0,t}}} + \frac{22x(\gamma_{j0} + 1)^2}{n_{j0,t}} \leq 2\gamma_{j0}.$$

873 This implies,

$$\begin{aligned} \gamma_{ij,t}^{\text{ucb}} \gamma_{j0,t}^{\text{ucb}} - \theta_i &= \gamma_{ij,t}^{\text{ucb}} \gamma_{j0,t}^{\text{ucb}} - \gamma_{ij} \gamma_{j0} = (\gamma_{ij,t}^{\text{ucb}} - \gamma_{ij}) \gamma_{j0,t}^{\text{ucb}} + \gamma_{ij} (\gamma_{j0,t}^{\text{ucb}} - \gamma_{j0}) \\ &\leq 2(\gamma_{ij,t}^{\text{ucb}} - \gamma_{ij}) \gamma_{j0} + \gamma_{ij} (\gamma_{j0,t}^{\text{ucb}} - \gamma_{j0}) \\ &\leq 8\gamma_{j0}(\gamma_{ij} + 1) \sqrt{\frac{2\gamma_{ij}x}{n_{ij,t}}} + \frac{44x\gamma_{j0}(\gamma_{ij} + 1)^2}{n_{ij,t}} \\ &\quad + 4\gamma_{ij}(\gamma_{j0} + 1) \sqrt{\frac{2\gamma_{j0}x}{n_{j0,t}}} + \frac{22x\gamma_{ij}(\gamma_{j0} + 1)^2}{n_{j0,t}}. \end{aligned}$$

874 Replacing $\gamma_{ij} = \theta_i/\theta_j$ and $\gamma_{j0} = \theta_j$ concludes the proof. \square

875 **Lemma 11.** Let $T \geq 1$ and $x > 0$. Then, with probability at least $1 - K(K+1)Te^{-x}$

$$\tau_{ij,t} < 2x \frac{(\theta_0 + \Theta_{S^*})^2}{\theta_i + \theta_j} \quad \text{or} \quad n_{ij,t} \geq \frac{(\theta_i + \theta_j)\tau_{ij,t}}{2(\theta_0 + \Theta_{S^*})}, \quad (17)$$

876 where $\tau_{ij,t} := \sum_{s=1}^{t-1} \mathbf{1}\{\{i, j\} \subseteq S_s\}$ simultaneously for all $t \in [T]$ and $i \neq j \in [K]$.

877 *Proof of Theorem 6.* Let \mathcal{E} be the high-probability event of Lemmas 10 and 11 are satisfied,
878 so that $\mathbb{P}(\mathcal{E}) \geq 1 - 4K^2Te^{-x}$. First, note that since we have under the event \mathcal{E} , $\hat{\theta}_t^{\text{ucb}} \leq \theta_t^{\text{ucb}}$,
879 our procedure also satisfies the regret upper-bound

$$\text{Reg}_T^{\text{wt}} \leq O(\sqrt{\theta_{\max}} K T \log T)$$

880 of Theorem 5. Indeed, all upper-bounds of the proof of Theorem 5 remain valid upper-bounds
881 except the probability of the event \mathcal{E}^c which is $O(T^{-1})$ for $x = 2 \log T$.

882 Let us now prove that we also have $R_T \leq O(K\sqrt{T} \log T)$ with no asymptotic dependence on
883 θ_{\max} when $T \rightarrow \infty$.

884 Then,

$$\begin{aligned}
\text{Reg}_T^{\text{wtd}} &= \sum_{t=1}^T \mathbb{E}[\mathcal{R}(S^*, \theta) - \mathcal{R}(S_t, \theta)] \\
&= \sum_{t=1}^T \mathbb{E}[(\mathcal{R}(S^*, \theta) - \mathcal{R}(S_t, \theta))\mathbb{1}\{\mathcal{E}\} + (\mathcal{R}(S^*, \theta) - \mathcal{R}(S_t, \theta))\mathbb{1}\{\mathcal{E}^c\}] \\
&\leq \sum_{t=1}^T \mathbb{E}\left[\left((\mathcal{R}(S_t, \hat{\theta}_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge \mathcal{R}(S^*, \theta)\right)\mathbb{1}\{\mathcal{E}\} + \mathcal{R}(S^*, \theta)\mathbb{1}\{\mathcal{E}^c\}\right].
\end{aligned} \tag{18}$$

885 Then, using $\mathcal{R}(S^*, \theta) \leq \max_i r_i \leq 1$, we get

$$\begin{aligned}
\text{Reg}_T^{\text{wtd}} &\leq \sum_{t=1}^T \mathbb{E}\left[\left((\mathcal{R}(S_t, \hat{\theta}_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge 1\right)\mathbb{1}\{\mathcal{E}\} + \mathbb{1}\{\mathcal{E}^c\}\right] \\
&\leq 4T^2 K(K+1)^2 e^{-x} + \sum_{t=1}^T \mathbb{E}\left[\left((\mathcal{R}(S_t, \hat{\theta}_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right].
\end{aligned} \tag{19}$$

886 Follow the proof of Theorem 5, we upper-bound the second term of the right-hand-side
887 of (19):

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{E}\left[\left((\mathcal{R}(S_t, \hat{\theta}_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right] \\
&= \sum_{t=1}^T \mathbb{E}\left[\left(\left(\min_{j \in [K]} \sum_{i \in S_t} \frac{r_i \hat{\theta}_{i,t}^{\text{ucb}}}{1 + \sum_{j \in S_t} \hat{\theta}_{j,t}^{\text{ucb}}} - \frac{r_i \theta_i}{1 + \sum_{j \in S_t} \theta_j}\right) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right] \\
&\leq \sum_{t=1}^T \mathbb{E}\left[\left(\left(\sum_{i \in S_t} \frac{r_i (\hat{\theta}_{i,t}^{\text{ucb}} - \theta_i)}{\theta_0 + \Theta_{S_t}}\right) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right] \quad \text{because } \sum_{i \in S_t} \hat{\theta}_{i,t}^{\text{ucb}} \geq \Theta_{S_t} \text{ under } \mathcal{E} \\
&\leq \sum_{t=1}^T \mathbb{E}\left[\left(\left(\sum_{i \in S_t} \frac{|\hat{\theta}_{i,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}}\right) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right] \quad \text{because } r_i \leq 1 \\
&\leq \sum_{i=1}^K \mathbb{E}\left[\sum_{t=1}^T \left(\frac{|\hat{\theta}_{i,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}} \wedge 1\right)\mathbb{1}\{i \in S_t\}\mathbb{1}\{\mathcal{E}\}\right] \\
&\leq \sum_{i=1}^K \mathbb{E}\left[\sum_{t=1}^T \left(\frac{|\gamma_{ij_t,t}^{\text{ucb}} \gamma_{j_t,0,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}} \wedge 1\right)\mathbb{1}\{i \in S_t\}\mathbb{1}\{\mathcal{E}\}\right]
\end{aligned} \tag{20}$$

888 where $j_t = \operatorname{argmax}_{j \in S_t \cup \{0\}} \theta_j$, where the last inequality is by definition of $\hat{\theta}_{i,t}^{\text{ucb}}$. Now,
889 from Lemma 10, paying an additive exploration cost to ensure that $n_{ij,t} \gtrsim x(1 + \gamma_{ij})$ and
890 $n_{j0,t} \gtrsim x(1 + \theta_j)^2 \theta_j$ for all $j \in S_t$ such that $\theta_j \geq \theta_0$. From Lemma 11, this is satisfied if for
891 some constant $C > 0$

$$\tau_{ij,t} > Cm^2 \theta_{\max}^2 x.$$

892 Such a condition can be wrong for a couple $(i, j) \in S_t^2$ at most during $CK^2 m^2 \theta_{\max}^2 x =$
893 $O(\log T)$ rounds (since $\tau_{ij,t}$ increases then). Thus, for C large enough,

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{E}\left[\left((\mathcal{R}(S_t, \hat{\theta}_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge 1\right)\mathbb{1}\{\mathcal{E}\}\right] \\
&\leq O(\log T) + \sum_{i=1}^K \mathbb{E}\left[\sum_{t=1}^T \frac{|\gamma_{ij_t,t}^{\text{ucb}} \gamma_{j_t,0,t}^{\text{ucb}} - \theta_i|}{\theta_0 + \Theta_{S_t}} \mathbb{1}\{i \in S_t, \tau_{ij_t,t} \wedge \tau_{j_t,t} \geq Cxm^2 \theta_{\max}^2\} \mathbb{1}\{\mathcal{E}\}\right]
\end{aligned}$$

$$\begin{aligned}
&\lesssim O(\log T) + \sum_{i=1}^K \mathbb{E} \left[\sum_{t=1}^T \left(\sqrt{(\gamma_{ij_t} + 1)\theta_i} x \left(\sqrt{\frac{(\theta_i + \theta_{j_t})}{n_{ij_t,t}}} + \sqrt{\frac{(1 + \theta_j)}{n_{j_t,0,t}}} \right) \right. \right. \\
&\quad \left. \left. + (\gamma_{ij_t} + 1) \frac{(\theta_i + \theta_{j_t})x}{n_{ij_t,t}} + \frac{\gamma_{ij_t}(1 + \theta_{j_t})^2 x}{n_{j_t,0,t}} \right) \frac{\mathbb{1}\{i \in S_t\}}{\theta_0 + \Theta_{S_t}} \right] \\
&\leq O(\log T) + \sum_{i=1}^K \mathbb{E} \left[\sum_{t=1}^T \sqrt{(\gamma_{ij_t} + 1)\theta_i} x \left(\sqrt{\frac{(\theta_i + \theta_{j_t})}{n_{ij_t,t}}} + \sqrt{\frac{(1 + \theta_{j_t})}{n_{j_t,0,t}}} \right) \frac{\mathbb{1}\{i \in S_t\}}{\theta_0 + \Theta_{S_t}} \right]
\end{aligned}$$

894 where the last inequality is because using that $\{i, j_t, 0\} \subseteq S_t$, we have

$$\mathbb{E} \left[\sum_{t=1}^T \frac{1 + \theta_{j_t}}{(1 + \Theta_{S_t})n_{j_t,0,t}} \right] = \mathbb{E} \left[\sum_{t=1}^T \sum_{j=1}^K \frac{\mathbb{1}\{i_t \in \{j, 0\}\}}{n_{j,0,t}} \mathbb{1}\{j = j_t\} \right] \leq K(1 + \log T).$$

895 and

$$\mathbb{E} \left[\sum_{t=1}^T \frac{\theta_i + \theta_{j_t}}{(1 + \Theta_{S_t})n_{ij_t,t}} \right] = \mathbb{E} \left[\sum_{t=1}^T \sum_{j=1}^K \frac{\mathbb{1}\{i_t \in \{j, i\}\}}{n_{j,0,t}} \mathbb{1}\{j = j_t\} \right] \leq K(1 + \log T).$$

896 Then, by Cauchy-Schwarz inequality we further get

$$\begin{aligned}
&\sum_{t=1}^T \mathbb{E} \left[\left((\mathcal{R}(S_t, \hat{\theta}_t^{\text{ucb}}) - \mathcal{R}(S_t, \theta)) \wedge 1 \right) \mathbb{1}\{\mathcal{E}\} \right] \\
&\lesssim O(\log T) + \sum_{i=1}^K \sqrt{\mathbb{E} \left[\sum_{t=1}^T \frac{(\gamma_{ij_t} + 1)\theta_i \mathbb{1}\{i \in S_t\} x}{\theta_0 + \Theta_{S_t}} \right]} \tag{21}
\end{aligned}$$

$$\begin{aligned}
&\quad \times \sqrt{\mathbb{E} \left[\sum_{t=1}^T \left(\frac{(\theta_i + \theta_{j_t})}{n_{ij_t,t}} + \frac{(1 + \theta_{j_t})}{n_{j_t,0,t}} \right) \frac{\mathbb{1}\{i \in S_t\}}{\theta_0 + \Theta_{S_t}} \right]} \\
&\lesssim O(\log T) + \sum_{i=1}^K \sqrt{\mathbb{E} \left[\sum_{t=1}^T \frac{(\gamma_{ij_t} + 1)\theta_i \mathbb{1}\{i \in S_t\} x}{\theta_0 + \Theta_{S_t}} \right]} \sqrt{K \log T} \\
&\lesssim O(\log T) + \sum_{i=1}^K \sqrt{\mathbb{E} \left[\sum_{t=1}^T \frac{\theta_i \mathbb{1}\{i \in S_t\} x}{\theta_0 + \Theta_{S_t}} \right]} \sqrt{K \log T} \text{ (because } \gamma_{ij_t} \leq 1 \text{ by definition of } j_t) \\
&\leq O(K\sqrt{T}x \log T) = O(K\sqrt{T} \log T), \tag{22}
\end{aligned}$$

897 where the last inequality is by Jensen's inequality and the equality by setting $x = 2 \log T$ to
898 control the probability that \mathcal{E}^c occurs. This concludes the proof. \square