

---

# Best-item Learning in Random Utility Models with Subset Choices

---

**Aadirupa Saha**  
aadirupa@iisc.ac.in  
Indian Institute of Science  
Bengaluru, India

**Aditya Gopalan**  
aditya@iisc.ac.in  
Indian Institute of Science  
Bengaluru, India

## Abstract

We consider the problem of PAC learning the most valuable item from a pool of  $n$  items using sequential, adaptively chosen plays of subsets of  $k$  items, when, upon playing a subset, the learner receives relative feedback sampled according to a general Random Utility Model (RUM) with independent noise perturbations to the latent item utilities. We identify a new property of such a RUM, termed the minimum advantage, that helps in characterizing the complexity of separating pairs of items based on their relative win/loss empirical counts, and can be bounded as a function of the noise distribution alone. We give a learning algorithm for general RUMs, based on pairwise relative counts of items and hierarchical elimination, along with a new PAC sample complexity guarantee of  $O(\frac{n}{c^2\epsilon^2} \log \frac{k}{\delta})$  rounds to identify an  $\epsilon$ -optimal item with confidence  $1 - \delta$ , when the worst case pairwise advantage in the RUM has sensitivity at least  $c$  to the parameter gaps of items. Fundamental lower bounds on PAC sample complexity show that this is near-optimal in terms of its dependence on  $n, k$  and  $c$ .

## 1 Introduction

Random utility models (RUMs) are a popular and well-established framework for studying behavioral choices by individuals and groups (Thurstone, 1927). In a RUM with finite alternatives or items, a distribution on the preferred alternative(s) is assumed to arise from a random utility drawn from a distribution for each

item, followed by rank ordering the items according to their utilities.

Perhaps the most widely known RUM is the Plackett-Luce or multinomial logit model (Plackett, 1975; Luce, 2012) which results when each item’s utility is sampled from an additive model with a Gumbel-distributed perturbation. It is unique in the sense of enjoying the property of independence of irrelevant attributes (IIA), which is often key in permitting efficient inference of Plackett-Luce models from data (Khetan and Oh, 2016). Other well-known RUMs include the probit model (Bliss, 1934) featuring random Gaussian perturbations to the intrinsic utilities, mixed logit, nested logit, etc.

A long line of work in statistics and machine learning focuses on estimating RUM properties from observed data (Soufiani et al., 2014; Zhao et al., 2018; Soufiani et al., 2013). Online learning or adaptive testing, on the other hand, has shown efficient ways of identifying the most attractive (i.e., highest utility) items in RUMs by learning from relative feedback from item pairs or more generally subsets (Szörényi et al., 2015; Saha and Gopalan, 2019; Jang et al., 2017). However, almost all existing work in this vein exclusively employs the Plackett-Luce model, arguably due to its very useful IIA property, and our understanding of learning performance in other, more general RUMs has been lacking. We take a step in this direction by framing the problem of sequentially learning the best item/items in general RUMs by adaptive testing of item subsets and observing relative RUM feedback. In the process, we uncover new structural properties in RUMs, including models with exponential, uniform, Gaussian (probit) utility distributions, and give algorithmic principles to exploit this structure, that permit provably sample-efficient online learning and allow us to go beyond Plackett-Luce.

**Our contributions:** We introduce a new property of a RUM, called the (pairwise) *advantage ratio*, which essentially measures the worst-case relative probabilities between an item pair across all possible contexts

(subsets) where they occur. We show that this ratio can be controlled (bounded below) as an affine function of the relative strengths of item pairs for RUMs based on several common centered utility distributions, e.g., exponential, Gumbel, uniform, Gamma, Weibull, normal, etc., even when the resulting RUM does not possess analytically favorable properties such as IIA.

We give an algorithm for sequentially and adaptively PAC (probably approximately correct) learning the best item from among a finite pool when, in each decision round, a subset of fixed size can be tested and top- $m$  rank ordered feedback from the RUM can be observed. The algorithm is based on the idea of maintaining pairwise win/loss counts among items, hierarchically testing subsets and propagating the surviving winners – principles that have been shown to work optimally in the more structured Plackett-Luce RUM (Szörényi et al., 2015; Saha and Gopalan, 2019).

In terms of performance guarantees, we derive a PAC sample complexity bound for our algorithm: when working with a pool of  $n$  items in total with subsets of size- $k$  chosen in each decision round, the algorithm terminates in  $O(\frac{n}{c^2\epsilon^2} \log \frac{k}{\delta})$  rounds where  $c$  is a lower bound on the advantage ratio’s sensitivity to intrinsic item utilities. This can in turn be shown to be a property of only the RUM’s perturbation distribution, independent of the subset size  $k$ . A novel feature of the guarantee is that, unlike existing sample complexity results for sequential testing in the Plackett-Luce model, it does not rely on specific properties like IIA which are not present in general RUMs. We also extend the result to cover top- $m$  rank ordered feedback, of which winner feedback ( $m = 1$ ) is a special case. Finally, we show that the sample complexity of our algorithm is order-wise optimal across RUMs having a given advantage ratio sensitivity  $c$ , by arguing an information-theoretic lower bound on the sample complexity of any online learning algorithm.

Our results and techniques represent a conceptual advance in the problem of online learning in general RUMs, moving beyond the Plackett-Luce model for the first time to the best of our knowledge.

**Related Work:** For classical multiarmed bandits setting, there is a well studied literature on PAC-arm identification problem (Even-Dar et al., 2006; Audibert and Bubeck, 2010; Kalyanakrishnan et al., 2012; Karnin et al., 2013; Jamieson et al., 2014), where the learner gets to see a noisy draw of absolute reward feedback of an arm upon playing a single arm per round. On the contrary, learning to identify the best item(s) with only relative preference information (ordinal as opposed to cardinal feedback) has seen steady progress since the introduction of the dueling bandit framework (Zoghi

et al., 2013) with pairs of items (size-2 subsets) that can be played, and subsequent work on generalisation to broader models both in terms of distributional parameters (Yue and Joachims, 2009; Gajane et al., 2015; Ailon et al., 2014; Zoghi et al., 2015) as well as combinatorial subset-wise plays (Mohajer et al., 2017; González et al., 2017; Saha and Gopalan, 2018b; Sui et al., 2017). There have been several developments on the PAC objective for different pairwise preference models, such as those satisfying stochastic triangle inequalities and strong stochastic transitivity (Yue and Joachims, 2011), general utility-based preference models (Urvoy et al., 2013), the Plackett-Luce model (Szörényi et al., 2015) and the Mallows model (Busa-Fekete et al., 2014a)]. Recent work has studied PAC-learning objectives other than identifying the single (near) best arm, e.g. recovering a few of the top arms (Busa-Fekete et al., 2013; Mohajer et al., 2017), or the true ranking of the items (Busa-Fekete et al., 2014b; Falahatgar et al., 2017). Some of the recent works also extended the PAC-learning objective with relative subsetwise preferences (Saha and Gopalan, 2018a; Chen et al., 2017, 2018; Saha and Gopalan, 2019; Ren et al., 2018).

However, none of the existing work considers strategies to learn efficiently in general RUMs with subset-wise preferences and to the best of our knowledge we are the first to address this general problem setup. In a different direction, there has been work on batch (non-adaptive) estimation in general RUMs, e.g., (Zhao et al., 2018; Soufiani et al., 2013); however, this does not consider the price of active learning and the associated exploration effort required as we study here. A related body of literature lies in dynamic assortment selection, where the goal is to offer a subset of items to customers in order to maximise expected revenue, which has been studied under different choice models, e.g. Multinomial-Logit (Talluri and Van Ryzin, 2004), Mallows and mixture of Mallows (Désir et al., 2016a), Markov chain-based choice models (Désir et al., 2016b), single transition model (Nip et al., 2017) etc., but again each of this work addresses a given and a very specific kind of choice model, and their objective is more suited to regret minimization type framework where playing every item comes with a associated cost.

## 2 Preliminaries

**Notation.** We denote by  $[n]$  the set  $\{1, 2, \dots, n\}$ . For any subset  $S \subseteq [n]$ , let  $|S|$  denote the cardinality of  $S$ . When there is no confusion about the context, we often represent (an unordered) subset  $S$  as a vector, or ordered subset,  $S$  of size  $|S|$  (according to, say, a fixed global ordering of all the items  $[n]$ ). In this case,  $S(i)$  denotes the item (member) at the  $i$ th position in subset  $S$ .  $\Sigma_S = \{\sigma \mid \sigma \text{ is a permutation over items of}$

$S\}$ , where for any permutation  $\sigma \in \Sigma_S$ ,  $\sigma(i)$  denotes the element at the  $i$ -th position in  $\sigma$ ,  $i \in [|S|]$ .  $\mathbf{1}(\varphi)$  is generically used to denote an indicator variable that takes the value 1 if the predicate  $\varphi$  is true, and 0 otherwise.  $x \vee y$  denotes the maximum of  $x$  and  $y$ , and  $Pr(A)$  is used to denote the probability of event  $A$ , in a probability space that is clear from the context.

## 2.1 Random Utility-based Discrete Choice Models

A discrete choice model specifies the relative preferences of two or more discrete alternatives in a given set. Random Utility Models (RUMs) are a widely-studied class of discrete choice models; they assume a (non-random) ground-truth utility score  $\theta_i \in \mathbb{R}$  for each alternative  $i \in [n]$ , and assign a distribution  $\mathcal{D}_i(\cdot|\theta_i)$  for scoring item  $i$ , where  $\mathbf{E}[\mathcal{D}_i | \theta_i] = \theta_i$ . To model a winning alternative given any set  $S \subseteq [n]$ , one first draws a random utility score  $X_i \sim \mathcal{D}_i(\cdot|\theta_i)$  for each alternative in  $S$ , and selects an item with the highest random score. More formally, the probability that an item  $i \in S$  emerges as the *winner* in set  $S$  is given by:

$$Pr(i|S) = Pr(X_i > X_j \quad \forall j \in S \setminus \{i\}) \quad (1)$$

In this paper, we assume that for each item  $i \in [n]$ , its random *utility score*  $X_i$  is of the form  $X_i = \theta_i + \zeta_i$ , where all the  $\zeta_i \sim \mathcal{D}$  are ‘noise’ random variables drawn independently from a probability distribution  $\mathcal{D}$ .

A 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(0, 1) distributions with location parameters 0 and scale parameter 1 (Azari et al., 2012), which results in score distributions  $Pr(X_i \in [x, x+dx]) = e^{-(x-\theta_i)} e^{-e^{-(x-\theta_i)}} dx$ ,  $\forall i \in [n]$ . Moreover, it can be shown that the probability that an alternative  $i$  emerges as the winner in any set  $S \ni i$  is simply proportional to its score parameter:  $Pr(i|S) = \frac{e^{\theta_i}}{\sum_{j \in S} e^{\theta_j}}$ .

Other families of discrete choice models can be obtained by imposing different probability distributions over the iid noise  $\zeta_i \sim \mathcal{D}$ ; e.g.,

1. *Exponential* noise:  $\mathcal{D}$  is the Exponential( $\lambda$ ) distribution ( $\lambda > 0$ ).
2. Noise from *Extreme value distributions*:  $\mathcal{D}$  is the Extreme-value-distribution( $\mu, \sigma, \xi$ ) ( $\mu \in \mathbb{R}, \sigma > 0, \xi \in \mathbb{R}$ ). Many well-known distributions fall in this class, e.g., *Frechet*, *Weibull*, *Gumbel*. For instance, when  $\chi = 0$ , this reduces to the *Gumbel*( $\mu, \sigma$ ) distribution.
3. *Uniform* noise:  $\mathcal{D}$  is the (continuous) Uniform( $a, b$ ) distribution ( $a, b \in \mathbb{R}, b > a$ ).

4. *Gaussian* or Frechet, Weibull, Gumbel noise:  $\mathcal{D}$  is the Gaussian( $\mu, \sigma$ ) distribution ( $\mu \in \mathbb{R}, \sigma > 0$ ).
5. *Gamma* noise:  $\mathcal{D}$  is the Gamma( $k, \xi$ ) distribution (where  $k, \xi > 0$ ).

Other distributions  $\mathcal{D}$  can alternatively be used for modelling the noise distribution, depending on desired tail properties, domain-specific information, etc.

Finally, we denote a RUM choice model, comprised of an instance  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)$  (with its implicit dependence on the noise distribution  $\mathcal{D}$ ) along with a playable subset size  $k \leq n$ , by RUM( $k, \boldsymbol{\theta}$ ).

## 3 Problem Setting

We consider the probably approximately correct (PAC) version of the sequential decision-making problem of finding the best item in a set of  $n$  items, by making only subset-wise comparisons.

Formally, the learner is given a finite set  $[n]$  of  $n > 2$  items or ‘arms’<sup>1</sup> along with a playable subset size  $k \leq n$ . At each decision round  $t = 1, 2, \dots$ , the learner selects a subset  $S_t \subseteq [n]$  of  $k$  distinct items, and receives (stochastic) feedback depending on (a) the chosen subset  $S_t$ , and (b) a RUM( $k, \boldsymbol{\theta}$ ) choice model with parameters  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)$  a priori unknown to the learner. The nature of the feedback can be of several types as described in Section 3.1. For the purposes of analysis, we assume, without loss of generality<sup>2</sup>, that  $\theta_1 > \theta_i \forall i \in [n] \setminus \{1\}$  for ease of exposition<sup>3</sup>. We define a *best item* to be one with the highest score parameter:  $i^* \in \operatorname{argmax}_{i \in [n]} \theta_i = \{1\}$ , under the assumptions above.

**Remark 1.** Under the assumptions above, it follows that item 1 is the Condorcet Winner (Zoghi et al., 2014) for the underlying pairwise preference model induced by RUM( $k, \boldsymbol{\theta}$ ).

### 3.1 Feedback models

We mean by ‘feedback model’ the information received (from the ‘environment’) once the learner plays a subset  $S \subseteq [n]$  of  $k$  items. Similar to different types of feedback models introduced earlier in the context of the specific Plackett-Luce RUM (Saha and Gopalan, 2019), we consider the following feedback mechanisms:

- **Winner of the selected subset (WI):** The environment returns a single item  $I \in S$ , drawn

<sup>1</sup>terminology borrowed from multi-armed bandits

<sup>2</sup>under the assumption that the learner’s decision rule does not contain any bias towards a specific item index

<sup>3</sup>The extension to the case where several items have the same highest parameter value is easily accomplished.

independently from the probability distribution  $Pr(I = i|S) = Pr(X_i > X_j, \forall j \in S \setminus \{i\}) \quad \forall i \in S, S \subseteq [n]$ .

- **Full ranking selected subset of items (FR):**

The environment returns a full ranking  $\sigma \in \Sigma_S$ , drawn from the probability distribution  $Pr(\sigma = \sigma|S) = \prod_{i=1}^{|S|} Pr(X_{\sigma(i)} > X_{\sigma(j)}, \forall j \in \{i+1, \dots, |S|\})$ ,  $\forall \sigma \in \Sigma_S$ . In fact, this is equivalent to picking  $\sigma(1)$  according to the winner feedback from  $S$ , then picking  $\sigma(2)$  from  $S \setminus \{\sigma(1)\}$  following the same feedback model, and so on, until all elements from  $S$  are exhausted, or, in other words, successively sampling  $|S|$  winners from  $S$  according to the  $RUM(k, \theta)$  model, without replacement.

### 3.2 PAC Performance Objective: Correctness and Sample Complexity

For a  $RUM(k, \theta)$  instance with  $n \geq k$  arms, an arm  $i \in [n]$  is said to be  $\epsilon$ -optimal if  $\theta_i > \theta_1 - \epsilon$ . A sequential<sup>4</sup> learning algorithm that depends on feedback from an appropriate subset-wise feedback model is said to be  $(\epsilon, \delta)$ -PAC, for given constants  $0 < \epsilon \leq \frac{1}{2}, 0 < \delta \leq 1$ , if the following properties hold when it is run on any instance  $RUM(k, \theta)$ : (a) it stops and outputs an arm  $I \in [n]$  after a finite number of decision rounds (subset plays) with probability 1, and (b) the probability that its output  $I$  is an  $\epsilon$ -optimal arm in  $RUM(k, \theta)$  is at least  $1 - \delta$ , i.e.,  $Pr(I \text{ is } \epsilon\text{-optimal}) \geq 1 - \delta$ . Furthermore, by *sample complexity* of the algorithm, we mean the expected time (number of decision rounds) taken by the algorithm to stop when run on the instance  $RUM(k, \theta)$ .

## 4 Connecting Subsetwise preferences to Pairwise Scores

In this section, we introduce the key concept of Advantage ratio as a means to systematically relate subsetwise preference observations to pairwise scores in general RUMs.

Consider any set  $S \subseteq [n], |S| = k$ , and recall that the probability of item  $i$  winning in  $S$  is  $Pr(i|S) := Pr(X_i > X_j, \forall j \in [n] \setminus \{i\})$  for all  $i \in S, S \subseteq [n]$ . For any two items  $i, j \in [n]$ , let us denote  $\Delta_{ij} = (\theta_i - \theta_j)$ . Let us also denote by  $f(\cdot), F(\cdot)$  and  $\bar{F}(\cdot)$  the probability density function<sup>5</sup>, cumulative distribution func-

tion and complementary cumulative distribution function of the noise distribution  $\mathcal{D}$ , respectively; thus,  $F(x) = \int_{-\infty}^x f(x)dx$  for any  $x \in \text{Support}(\mathcal{D})$  and  $\bar{F}(x) = \int_x^{\infty} f(x)dx = 1 - F(x)$  for any  $x \in \text{Support}(\mathcal{D})$ .

We now introduce and analyse the *Advantage-Ratio* (Def. 1); we will see in Sec. 5.1 how this quantity helps us deriving an improved sample complexity guarantee for our  $(\epsilon, \delta)$ -PAC item identification problem.

**Definition 1** (Advantage ratio and Minimum advantage ratio). *Given any subsetwise preference model defined on  $n$  items, we define the advantage ratio of item  $i$  over item  $j$  within the subset  $S \subseteq [n], i, j \in S$  as  $Advantage\text{-Ratio}(i, j, S) = \frac{Pr(i|S)}{Pr(j|S)}$ .*

*Moreover, given a playable subset size  $k$ , we define the minimum advantage ratio,  $Min\text{-}AR$ , of item- $i$  over  $j$ , as the least advantage ratio of  $i$  over  $j$  across size- $k$  subsets of  $[n]$ , i.e.,*

$$Min\text{-}AR(i, j) = \min_{S \subseteq [n], |S|=k, S \ni i, j} \frac{Pr(i|S)}{Pr(j|S)}. \quad (2)$$

The key intuition here is that when  $Min\text{-}AR(i, j)$  does not equal 1, it serves as a distinctive measure for identifying item  $i$  and  $j$  separately irrespective of the context  $S$ . We specifically build on this intuition later in Sec. 5.1 to propose a new algorithm (Alg. 1) which finds the  $(\epsilon, \delta)$ -PAC best item relying on the unique distinctive property of the best-item  $\theta_1 > \theta_j \forall j \in [n] \setminus \{1\}$  (as described in Sec. 3).

The following result shows a variational lower bound, in terms of the noise distribution, for the minimum advantage ratio in a  $RUM(k, \theta)$  model with independent and identically distributed (iid) noise variables, that is often amenable to explicit calculation/bounding.

**Lemma 2** (Variational lower bound for the advantage ratio). *For any  $RUM(k, \theta)$  based subsetwise preference model and any item pair  $(i, j)$ ,*<sup>6</sup>

$$Min\text{-}AR(i, j) \geq \min_{z \in \mathbb{R}} \frac{Pr(X_i > \max(X_j, z))}{Pr(X_j > \max(X_i, z))}. \quad (3)$$

*Moreover for  $RUM(k, \theta)$  models one can show that for any triplet  $(i, j, S)$ ,  $Pr(X_i > \max(X_j, z)) = F(z - \theta_j)\bar{F}(z - \theta_i) + \int_{z-\theta_j}^{\infty} \bar{F}(x - \Delta_{ij})f(x)dx$ , which further lower bounds  $Min\text{-}AR(i, j)$  by:*

$$\min_{z \in \mathbb{R}} \frac{F(z - \theta_j)\bar{F}(z - \theta_i) + \int_{z-\theta_j}^{\infty} \bar{F}(x - \Delta_{ij})f(x)dx}{F(z - \theta_i)\bar{F}(z - \theta_j) + \int_{z-\theta_i}^{\infty} \bar{F}(x + \Delta_{ij})f(x)dx}.$$

The proof of the result appears in Appendix A.1. Fig. 1 shows a geometrical interpretation behind  $Min\text{-}AR(i, j)$ , under the joint realization of the pair of values  $(\zeta_i, \zeta_j)$ .

<sup>4</sup>We essentially mean a causal algorithm that makes present decisions using only past observed information at each time; the technical details for defining this precisely are omitted.

<sup>5</sup>We assume by default that all noise distributions have a density; the extension to more general noise distributions is left to future work.

<sup>6</sup>We assume  $\frac{0}{0}$  to be  $\infty$  in the right hand side of Eqn. 3.

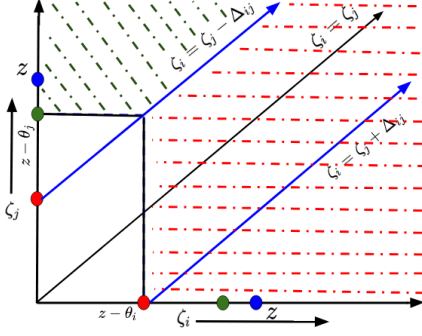


Figure 1: A two-dimensional geometrical interpretation for the quantity  $\text{Min-AR}(i, j)$ . Let  $z \in \mathbb{R}$  be a random variable denoting the max score observed for the rest of the items, i.e.  $\max_{a \in S \setminus \{i, j\}} X_a$ . Let the blue, green and red dot respectively denote the position of  $z$ ,  $z - \theta_j$  and  $z - \theta_i$ . With  $X_i = \theta_i + \zeta_i$ ,  $\forall i \in [n]$ , the green shaded region is where  $X_j > \max(X_i, z)$ , the red shaded region is where  $X_i > \max(X_j, z)$  i.e. item  $i$  is the winner, and the white rectangle is where  $\max(X_i, X_j) < z$  i.e. some other item wins. The shape of the green and red region varies as  $z$  moves on  $\mathbb{R}$  (in the hindsight this basically covers the realizations of all  $z$  over all possible subsets  $S$ )— $\text{Min-AR}(i, j)$  is attained at the particular  $z$  where the ratio of the mass of the red and green region is minimized (see Eqn. (3) for details).

**Remark 2.** Suppose  $\bar{S} := \arg \min_{|S|=k, i, j \in S} \frac{\Pr(i|S)}{\Pr(j|S)}$ . It is sufficient to consider the domain of  $z$  in the right hand side of (3) to be just the set  $\max_{r \in \bar{S} \setminus \{i, j\}} \theta_r + \text{support}(\mathcal{D})$ , as the proof of Lemma 2 brings out. However, for simplicity we use a smaller lower bound in Eqn. 3 and take  $z \in \mathbb{R}$ .

We next derive the  $\text{Min-AR}(i, j)$  values certain specific noise distributions:

**Lemma 3** (Analysing  $\text{Min-AR}$  for specific noise models). *Given a fixed item pair  $(i, j)$  such that  $\theta_i > \theta_j$ , the following bounds hold under the respective noise models in an iid RUM.*

1. *Exponential( $\lambda$ ):  $\text{Min-AR}(i, j) \geq e^{\Delta_{ij}} > 1 + \Delta_{ij}$  for Exponential noise with  $\lambda = 1$ .*
2. *Extreme value distribution( $\mu, \sigma, \chi$ ): For Gumbel( $\mu, \sigma$ ) ( $\chi = 0$ ) noise,  $\text{Min-AR}(i, j) = e^{\frac{\Delta_{ij}}{\sigma}} > 1 + \frac{\Delta_{ij}}{\sigma}$ .*
3. *Uniform( $a, b$ ):  $\text{Min-AR}(i, j) \geq 1 + \frac{2\Delta_{ij}}{b-a}$  for Uniform( $a, b$ ) noise ( $a, b \in \mathbb{R}, b > a$ , and  $\Delta_{ij} < \frac{a}{2}$ ).*
4. *Gamma( $k, \xi$ ):  $\text{Min-AR}(i, j) \geq 1 + \Delta_{ij}$  for Gamma( $2, 1$ ) noise.*
5. *Weibull( $\lambda, k$ ):  $\text{Min-AR}(i, j) \geq e^{\lambda \Delta_{ij}} > 1 + \lambda \Delta_{ij}$  for ( $k = 1$ ).*

6. *Normal  $\mathcal{N}(0, 1)$ :  $\exists c > 0$  such that, for  $\Delta_{ij}$  small enough (in a neighborhood of 0),  $\text{Min-AR}(i, j) \geq 1 + c\Delta_{ij}$ .*

The proof appears in Appendix A.2.

## 5 An optimal algorithm for the winner feedback model

In this section, we propose an algorithm (*Sequential-Pairwise-Battle*, Algorithm 1) for the  $(\epsilon, \delta)$ -PAC objective with winner feedback. We then analyse its correctness and sample complexity guarantee (Theorem 4) for any noise distribution  $\mathcal{D}$  (under a mild assumption of its being  $\text{Min-AR}$  bounded away from 1). Following this, we also prove a matching lower bound for the problem which shows that the sample complexity of Algorithm *Sequential-Pairwise-Battle* is unimprovable (up to a factor of  $\log k$ ).

### 5.1 The *Sequential-Pairwise-Battle* algorithm

Our algorithm is based on the simple idea of dividing the set of  $n$  items into sub-groups of size  $k$ , querying each subgroup ‘sufficiently enough’, retaining thereafter only the empirically ‘strongest item’ of each sub-group, and recursing on the remaining set of items until only one item remains.

More specifically, it starts by partitioning the initial item pool into  $G := \lceil \frac{n}{k} \rceil$  mutually exclusive and exhaustive sets  $\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_G$  such that  $\cup_{j=1}^G \mathcal{G}_j = S$  and  $\mathcal{G}_j \cap \mathcal{G}_{j'} = \emptyset$ ,  $\forall j, j' \in [G] | \mathcal{G}_j| = k$ ,  $\forall j \in [G - 1]$ . Each set  $\mathcal{G}_g$ ,  $g \in [G]$  is then queried for  $t = O\left(\frac{k}{\epsilon_\ell^2} \ln \frac{k}{\delta_\ell}\right)$  rounds, and only the ‘empirical winner’  $c_g$  of each group  $g$  is retained in a set  $S$ , rest are discarded. The algorithm next recurses the same procedure on the remaining set of surviving items, until a single item is left, which then is declared to be the  $(\epsilon, \delta)$  PAC-best item. Algorithm 1 presents the pseudocode in more detail.

**Key idea:** The primary novelty here is how the algorithm reasons about the ‘strongest item’ in each sub-group  $\mathcal{G}_g$ : It maintains the pairwise preferences of every item pair  $(i, j)$  in any sub-group  $\mathcal{G}_g$  and simply chooses the item that beats the rest of the items in the sub-group with a positive advantage of greater than  $\frac{1}{2}$  (alternatively, the item that wins maximum number of subset-wise plays). Our idea of maintaining pairwise preferences is motivated by a similar algorithm proposed in (Saha and Gopalan, 2019); however, their performance guarantee applies to only the very specific class of Plackett-Luce feedback models, whereas the novelty of our current analysis reveals the

**Algorithm 1** *Sequential-Pairwise-Battle*(Seq-PB)

---

```

1: Input:
2:   Set of items:  $[n]$ , Subset size:  $n \geq k > 1$ 
3:   Error bias:  $\epsilon > 0$ , Confidence parameter:  $\delta > 0$ 
4:   Noise model ( $\mathcal{D}$ ) dependent constant  $c > 0$ 
5: Initialize:
6:    $S \leftarrow [n]$ ,  $\epsilon_0 \leftarrow \frac{c\epsilon}{8}$ , and  $\delta_0 \leftarrow \frac{\delta}{2}$ 
7:   Divide  $S$  into  $G := \lceil \frac{n}{k} \rceil$  sets  $\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_G$  such
      that  $\cup_{j=1}^G \mathcal{G}_j = S$  and  $\mathcal{G}_j \cap \mathcal{G}_{j'} = \emptyset$ ,  $\forall j, j' \in [G]$ ,
      where  $|\mathcal{G}_j| = k$ ,  $\forall j \in [G-1]$ 
8:   If  $|\mathcal{G}_G| < k$ , then set  $\mathcal{R}_1 \leftarrow \mathcal{G}_G$  and  $G = G - 1$ 
9:   while  $\ell = 1, 2, \dots$  do
10:    Set  $S \leftarrow \emptyset$ ,  $\delta_\ell \leftarrow \frac{\delta_{\ell-1}}{2}$ ,  $\epsilon_\ell \leftarrow \frac{3}{4}\epsilon_{\ell-1}$ 
11:    for  $g = 1, 2, \dots, G$  do
12:      Play the set  $\mathcal{G}_g$  for  $t := \lceil \frac{k}{2\epsilon_\ell^2} \ln \frac{k}{\delta_\ell} \rceil$  rounds
13:       $w_i \leftarrow$  Number of times  $i$  won in  $t$  plays of  $\mathcal{G}_g$ ,
         $\forall i \in \mathcal{G}_g$ 
14:      Set  $c_g \leftarrow \arg \max_{i \in \mathcal{A}} w_i$  and  $S \leftarrow S \cup \{c_g\}$ 
15:    end for
16:     $S \leftarrow S \cup \mathcal{R}_\ell$ 
17:    if  $(|S| == 1)$  then
18:      Break (go out of the while loop)
19:    else if  $|S| \leq k$  then
20:       $S' \leftarrow$  Randomly sample  $k - |S|$  items from
         $[n] \setminus S$ , and  $S \leftarrow S \cup S'$ ,  $\epsilon_\ell \leftarrow \frac{c\epsilon}{2}$ ,  $\delta_\ell \leftarrow \delta$ 
21:    else
22:      Divide  $S$  into  $G := \lceil \frac{|S|}{k} \rceil$  sets  $\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_G$ ,
        such that  $\cup_{j=1}^G \mathcal{G}_j = S$ , and  $\mathcal{G}_j \cap \mathcal{G}_{j'} = \emptyset$ ,
         $\forall j, j' \in [G]$ , where  $|\mathcal{G}_j| = k$ ,  $\forall j \in [G-1]$ 
23:      If  $|\mathcal{G}_G| < k$ , then set  $\mathcal{R}_{\ell+1} \leftarrow \mathcal{G}_G$  and  $G = G - 1$ 
24:    end if
25:  end while
26: Output: The unique item left in  $S$ 

```

---

power of maintaining pairwise-estimates for more general  $\text{RUM}(k, \theta)$  subsetwise model (which includes the Plackett-Luce choice model as a special case). The pseudo code of *Sequential-Pairwise-Battle* is given in Alg. 1.

The following is our chief result; it proves correctness and a sample complexity bound for Algorithm 1.

**Theorem 4** (*Sequential-Pairwise-Battle: Correctness and Sample Complexity*). *Consider any iid subsetwise preference model  $\text{RUM}(k, \theta)$  based on a noise distribution  $\mathcal{D}$ , and suppose that for any item pair  $i, j$ , we have  $\text{Min-AR}(i, j) \geq 1 + \frac{4c\Delta_{ij}}{1-2c}$  for some  $\mathcal{D}$ -dependent constant  $c > 0$ . Then, Algorithm 1, with input constant  $c > 0$ , is an  $(\epsilon, \delta)$ -PAC algorithm with sample complexity  $O(\frac{n}{c^2\epsilon^2} \log \frac{k}{\delta})$ .*

The proof of the result appears in Appendix B.1.

**Remark 3.** *The linear dependence on the total num-*

*ber of items,  $n$ , is, in effect, indicates the price to pay for learning the  $n$  unknown model parameters  $\theta = (\theta_1, \dots, \theta_n)$  which decide the subsetwise winning probabilities of the  $n$  items. Remarkably, however, the theorem shows that the PAC sample complexity of the  $(\epsilon, \delta)$ -best item identification problem, with only winner feedback information from  $k$ -size subsets, is independent of  $k$  (except some mild logarithmic dependencies). One may expect to see improved sample complexity as the number of items being simultaneously tested in each round is large ( $k \geq 2$ ), but note that on the other side, the sample complexity could also worsen, since it is also harder for a good item to win and show itself in a few draws against a large population of  $k - 1$  other competitors – these effects roughly balance each other out, and the final sample complexity only depends on the total number of items  $n$  and the accuracy parameters  $(\epsilon, \delta)$ .*

Note that Lemma 3 gives specific values of the noise-model  $\mathcal{D}$  dependent constant  $c > 0$ , using which we can derive specific sample complexity bounds for certain noise models:

**Corollary 5** (Model specific correctness and sample complexity guarantees). *For the following representative noise distributions: Exponential(1), Gumbel( $\mu, \sigma$ ), Gamma(2, 1), Uniform( $a, b$ ), Weibull( $\lambda, 1$ ), Standard normal or Normal(0, 1), Seq-PB (Alg.1) finds an  $(\epsilon, \delta)$ -PAC item within sample complexity  $O(\frac{n}{\epsilon^2} \ln \frac{k}{\delta})$ .*

*Proof sketch.* The proof follows from the general performance guarantee of Seq-PB (Thm. 4) and Lem. 3. More specifically from Lem. 3 it follows that the value of  $c$  for these specific distributions are constant, which concludes the claim. For completeness the distribution-specific values of  $c$  are given in Appendix B.2.  $\square$

## 5.2 Sample Complexity Lower Bound

In this section we derive a sample complexity lower bound for any  $(\epsilon, \delta)$ -PAC algorithm for any  $\text{RUM}(k, \theta)$  model with  $\text{Min-AR}(i, j)$  strictly bounded away from 1 in terms of  $\Delta_{ij}$ . Our formal claim goes as follows:

**Theorem 6** (Sample Complexity Lower Bound for  $\text{RUM}(k, \theta)$  model). *Given  $\epsilon \in (0, \frac{1}{4}]$ ,  $\delta \in (0, 1]$ ,  $c > 0$  and an  $(\epsilon, \delta)$ -PAC algorithm  $A$  with winner item feedback, there exists a  $\text{RUM}(k, \theta)$  instance  $\nu$  with  $\text{Min-AR}(i, j) \geq 1 + 4c\Delta_{ij}$  for all  $i, j \in [n]$ , where the expected sample complexity of  $A$  on  $\nu$  is at least  $\Omega(\frac{n}{c^2\epsilon^2} \ln \frac{1}{2.4\delta})$ .*

The proof is given in Appendix B.3. It essentially involves a change of measure argument demonstrating a family of Plackett-Luce models (iid Gumbel noise), with the appropriate  $c$  value, that cannot easily be teased apart by any learning algorithm.

Comparing this result with the performance guarantee

of our proposed algorithm (Theorem 6) shows that the sample complexity of the algorithm is order-wise optimal (up to a  $\log k$  factor). Moreover, this result also shows that the IIA (independence of irrelevant attributes) property of the Plackett-Luce choice model is not essential for exploiting pairwise preferences via rank breaking, as was claimed in (Saha and Gopalan, 2019). Indeed, except for the case of *Gumbel* noise, none of the  $\text{RUM}(k, \theta)$  based models in Corollary 5 satisfies IIA, but they all respect the  $O\left(\frac{n}{\epsilon^2} \ln \frac{1}{\delta}\right)$   $(\epsilon, \delta)$ -PAC sample complexity guarantee.

**Remark 4.** For constant  $c = O(1)$ , the fundamental sample complexity bound of Theorem 6 resembles that of PAC best arm identification in the standard multi-armed bandit (MAB) problem (Even-Dar et al., 2006). Recall that our problem objective is exactly same as MAB, however our feedback model is very different since in MAB, the learner gets to see the noisy rewards/scores (i.e. the exact values of  $X_i$ , which can be seen as a noisy feedback of the true reward/score  $\theta_i$  of item- $i$ ), whereas here the learner only sees a  $k$ -wise relative preference feedback based on the underlying observed values of  $X_i$ , which is a more indirect way of giving feedback on the item scores, and thus intuitively our problem objective is at least as hard as that of MAB setup.

## 6 Results for Top- $m$ Ranking (TR) feedback model

We now address our  $(\epsilon, \delta)$ -PAC item identification problem for the case of more general, top- $m$  rank ordered feedback for the  $\text{RUM}(k, \theta)$  model, that generalises both the winner-item (WI) and full ranking (FR) feedback models.

**Top- $m$  ranking of items (TR- $m$ ):** In this feedback setting, the environment is assumed to return a ranking of only  $m$  items from among  $S$ , i.e., the environment first draws a full ranking  $\sigma$  over  $S$  according to  $\text{RUM}(k, \theta)$  as in **FR** above, and returns the first  $m$  rank elements of  $\sigma$ , i.e.,  $(\sigma(1), \dots, \sigma(m))$ . It can be seen that for each permutation  $\sigma$  on a subset  $S_m \subset S$ ,  $|S_m| = m$ , we must have  $\Pr(\sigma = \sigma|S) = \prod_{i=1}^m \Pr(X_{\sigma(i)} > X_{\sigma(j)}, \forall j \in \{i+1, \dots, m\})$ ,  $\forall \sigma \in \Sigma_S^m$ , where by  $\Sigma_S^m$  we denote the set of all possible  $m$ -length ranking of items in set  $S$ , it is easy to note that  $|S| = \binom{k}{m} m!$ . Thus, generating such a  $\sigma$  is also equivalent to successively sampling  $m$  winners from  $S$  according to the PL model, without replacement. It follows that **TR** reduces to **FR** when  $m = k = |S|$  and to **WI** when  $m = 1$ . Note that the idea for top- $m$  ranking feedback was introduced by (Saha and Gopalan, 2018a) but only for the specific Plackett Luce choice model.

### 6.1 Algorithm for top- $m$ ranking feedback

In this section, we extend the algorithm proposed earlier (Alg. 1) to handle feedback from the general top- $m$  ranking feedback model. We also show that we can achieve an  $\frac{1}{m}$ -factor improved sample complexity rate with top- $m$  ranking feedback (Thm. 7). We finally give a fundamental sample complexity bound (Thm. 8), which shows the optimality of our proposed algorithm mSeq-PB up to logarithmic factors.

**Main idea:** Same as Seq-PB, the algorithm proposed in this section (Alg. 2) in principle follows the same sequential elimination based strategy to find the nearest item of the  $\text{RUM}(k, \theta)$  model based on pairwise preferences. However, we use the idea of *rank breaking* (Soufiani et al., 2014; Saha and Gopalan, 2018a) to extract the pairwise preferences: formally, given any set  $S$  of size  $k$ , if  $\sigma \in \Sigma_S^m$ ,  $(S_m \subseteq S, |S_m| = m)$  denotes a possible top- $m$  ranking of  $S$ , then the *Rank-Breaking* subroutine considers each item in  $S$  to be beaten by its preceding items in  $\sigma$  in a pairwise sense. For instance, given a full ranking of a set of 4 elements  $S = \{a, b, c, d\}$ , say  $b \succ a \succ c \succ d$ , Rank-Breaking generates the set of 6 pairwise comparisons:  $\{(b \succ a), (b \succ c), (b \succ d), (a \succ c), (a \succ d), (c \succ d)\}$  etc.

As a whole, our new algorithm now again divides the set of  $n$  items into small groups of size  $k$ , say  $\mathcal{G}_1, \dots, \mathcal{G}_G$ ,  $G = \lceil \frac{n}{k} \rceil$ , and play each sub-group some  $t = O\left(\frac{k}{m\epsilon^2} \ln \frac{1}{\delta}\right)$  many rounds. Inside any fixed sub-group  $\mathcal{G}_g$ , after each round of play, it uses *Rank-Breaking* on the top- $m$  ranking feedback  $\sigma \in \Sigma_{\mathcal{G}_g}^m$ , to extract out  $\binom{m}{2} + (k-m)m$  many pairwise feedback, which is further used to estimate the empirical pairwise preferences  $\hat{p}_{ij}$  for each pair of items  $i, j \in \mathcal{G}_g$ . Based on these pairwise estimates it then only retains the strongest item of  $\mathcal{G}_g$  and recurse the same procedure on the set of surviving items, until just one item is left in the set. The complete algorithm is given in Alg. 2 (Appendix C.1).

Theorem 7 analyses the correctness and sample complexity bounds of mSeq-PB. Note that the sample complexity bound of mSeq-PB with top- $m$  ranking (TR) feedback model is  $\frac{1}{m}$ -times that of the WI model (Thm. 4). This is justified since intuitively revealing a ranking on  $m$  items in a  $k$ -set provides about  $m$  many WI feedback per round, which essentially leads to the  $m$ -factor improvement in the sample complexity.

**Theorem 7** (mSeq-PB(Alg. 2): Correctness and Sample Complexity). *Consider any  $\text{RUM}(k, \theta)$  subsetwise preference model based on noise distribution  $\mathcal{D}$  and suppose for any item pair  $i, j$ , we have  $\text{Min-AR}(i, j) \geq 1 + \frac{4c\Delta_{ij}}{1-2c}$  for some  $\mathcal{D}$ -dependent constant  $c > 0$ . Then mSeq-PB (Alg.2) with input constant*



$c > 0$  on top- $m$  ranking feedback model is an  $(\epsilon, \delta)$ -PAC algorithm with sample complexity  $O(\frac{n}{mc^2\epsilon^2} \log \frac{k}{\delta})$ .

(Proof is given in Appendix C.2.) Similar to Cor. 5, for the top- $m$  model again, we can derive specific sample complexity bounds for different noise distributions, e.g., *Exponential*, *Gumbel*, *Gaussian*, *Uniform*, *Gamma* etc., in this case as well.

## 6.2 Lower Bound: Top- $m$ ranking feedback

In this section, we analyze the fundamental limit of sample complexity lower bound for any  $(\epsilon, \delta)$ -PAC algorithm for RUM( $k, \theta$ ) model.

**Theorem 8** (Sample Complexity Lower Bound for RUM( $k, \theta$ ) model with TR- $m$  feedback). *Given  $\epsilon \in (0, \frac{1}{4}]$  and  $\delta \in (0, 1]$ , and an  $(\epsilon, \delta)$ -PAC algorithm  $A$  with winner item feedback, there exists a RUM( $k, \theta$ ) instance  $\nu$ , in which for any pair  $i, j \in [n]$   $\text{Min-AR}(i, j) \geq 1 + 4c\Delta_{ij}$ , where the expected sample complexity of  $A$  on  $\nu$  with top- $m$  ranking feedback has to be at least  $\Omega\left(\frac{n}{mc^2\epsilon^2} \ln \frac{1}{2.4\delta}\right)$  for  $A$  to be  $(\epsilon, \delta)$ -PAC.*

(The proof is given in Appendix C.3.) Similar to the case of winner feedback, comparing Theorem 7 with the above result shows that the sample complexity of mSeq-PB is orderwise optimal (up to logarithmic factors), for general case of top- $m$  ranking feedback as well.

## 7 Experiments

To complement our theoretical guarantees, we carry out some empirical simulations, as detailed below.

**RUM models.** We use the following 4 different noise models: **1.** Gumbel(0,1), **2.** Normal(0,1), **3.** Uniform(0,1), **4.** Exponential(1).

**Utility Scores.** Towards modelling different RUM based choice models, we combine the above noise models with the following 4 different ground utility scores ( $\theta$ ): **1.** b1:  $\theta_1 = 0.8, \theta_i = 0.6$ , otherwise. **2.** g1:  $\theta_1 = 0.8, \theta_i = 0.2$ , otherwise. **3.** geo:  $\theta_1 = 1, \frac{\theta_{i+1}}{\theta_i} = 0.9, \forall i \in [n]$ . **4.** arith:  $\theta_1 = 1, \theta_i - \theta_{i+1} = 0.01, \forall i \in [n]$ , with respectively  $n = 8, 16, 50$  and 100 items.

All reported performances are averaged across 50 runs. To the best of our knowledge no known algorithm address our problem setup for general RUM models, unfortunately we could not compare our method (Alg. 1 and 2) with any baseline. Given the above setup, we run two types the experiments to investigate:

**Success probability ( $1 - \delta$ ) (i.e. rate of correctness) vs sample complexity.** We set  $k = \frac{n}{2}, m = \frac{k}{2}$ ,

$\epsilon = \min_{i,j} |\theta_i - \theta_j|$  (i.e. the minimum pairwise gap among the utility scores) for each different environment. As expected from Thm 4, 6—with higher sample complexity, the success probability  $1 - \delta$  goes to 1 for each noise model showing that the algorithm is almost always correct with sufficient queries; similarly for lower number of observations the algorithms errors too often.

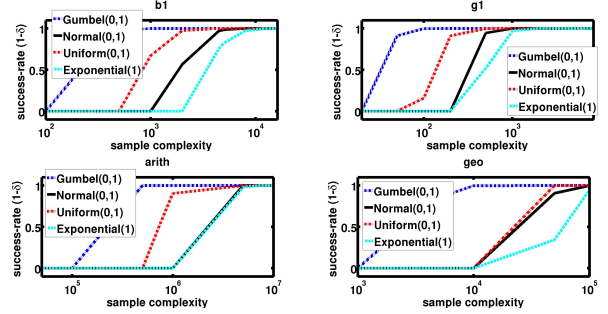


Figure 2: Success probability ( $1 - \delta$ ) vs sample complexity of Alg. 1 on different utility score-noise model combination

**Sample complexity vs length of rank-ordered feedback ( $m$ ).** We run these experiments on the *geo* dataset. Fig. 3 shows that the sample complexity seem to scale as  $O(\frac{1}{m})$  while  $\epsilon, \delta$  is kept fixed to 0.1 (validating the claim from Thm. 7).

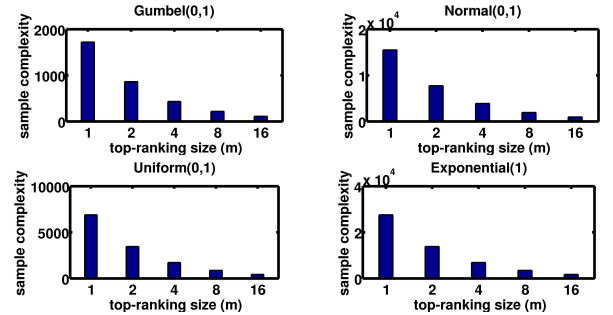


Figure 3: Sample complexity vs length of rank-ordered feedback ( $m$ ) of Alg. 2 on *geo* utility score for different RUM models

## 8 Conclusion and Future Directions

We have identified a new principle to learn with general subset-size preference feedback in general iid RUMs – rank breaking followed by pairwise comparisons. This is by extending the concept of pairwise advantage from the Plackett-Luce (PL) choice model to more general RUMs, and by showing that the IIA property that PL models enjoy is not essential for optimal sample complexity. Several interesting directions exist for future investigation, e.g., considering correlated noise models (more general RUMs), explicitly modeling item features or attributes, other metrics like regret for online utility optimization, and relative preference learning in time-correlated Markov Decision Processes.



**Acknowledgements.** This work was supported by the Qualcomm Innovation Fellowship IND-417067, 2019, and by a grant from the Robert Bosch Centre for Cyber-Physical Systems, Indian Institute of Science.

## References

- Ailon, N., Karnin, Z. S., and Joachims, T. (2014). Reducing dueling bandits to cardinal bandits. In *ICML*, volume 32, pages 856–864.
- Audibert, J.-Y. and Bubeck, S. (2010). Best arm identification in multi-armed bandits. In *COLT-23th Conference on Learning Theory-2010*, pages 13–p.
- Azari, H., Parkes, D., and Xia, L. (2012). Random utility theory for social choice. In *Advances in Neural Information Processing Systems*, pages 126–134.
- Bliss, C. I. (1934). The method of probits. *Science*.
- Busa-Fekete, R., Hüllermeier, E., and Szörényi, B. (2014a). Preference-based rank elicitation using statistical models: The case of mallows. In *Proceedings of The 31st International Conference on Machine Learning*, volume 32.
- Busa-Fekete, R., Szorenyi, B., Cheng, W., Weng, P., and Hüllermeier, E. (2013). Top-k selection based on adaptive sampling of noisy preferences. In *International Conference on Machine Learning*, pages 1094–1102.
- Busa-Fekete, R., Szörényi, B., and Hüllermeier, E. (2014b). Pac rank elicitation through adaptive sampling of stochastic pairwise preferences. In *AAAI*, pages 1701–1707.
- Chen, X., Gopi, S., Mao, J., and Schneider, J. (2017). Competitive analysis of the top-k ranking problem. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1245–1264. SIAM.
- Chen, X., Li, Y., and Mao, J. (2018). A nearly instance optimal algorithm for top-k ranking under the multinomial logit model. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2504–2522. SIAM.
- Désir, A., Goyal, V., Jagabathula, S., and Segev, D. (2016a). Assortment optimization under the mallows model. In *Advances in Neural Information Processing Systems*, pages 4700–4708.
- Désir, A., Goyal, V., Segev, D., and Ye, C. (2016b). Capacity constrained assortment optimization under the markov chain based choice model. *Operations Research*.
- Even-Dar, E., Mannor, S., and Mansour, Y. (2006). Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. *Journal of machine learning research*, 7(Jun):1079–1105.
- Falahatgar, M., Hao, Y., Orlitsky, A., Pichapati, V., and Ravindrakumar, V. (2017). Maxing and ranking with few assumptions. In *Advances in Neural Information Processing Systems*, pages 7063–7073.
- Gajane, P., Urvoy, T., and Clérôt, F. (2015). A relative exponential weighing algorithm for adversarial utility-based dueling bandits. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 218–227.
- González, J., Dai, Z., Damianou, A., and Lawrence, N. D. (2017). Preferential Bayesian optimization. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1282–1291. JMLR. org.
- Jamieson, K., Malloy, M., Nowak, R., and Bubeck, S. (2014). lil’ ucb : An optimal exploration algorithm for multi-armed bandits. In Balcan, M. F., Feldman, V., and Szepesvari, C., editors, *Proceedings of The 27th Conference on Learning Theory*, volume 35 of *Proceedings of Machine Learning Research*, pages 423–439. PMLR.
- Jang, M., Kim, S., Suh, C., and Oh, S. (2017). Optimal sample complexity of m-wise data for top-k ranking. In *Advances in Neural Information Processing Systems*, pages 1685–1695.
- Kalyanakrishnan, S., Tewari, A., Auer, P., and Stone, P. (2012). Pac subset selection in stochastic multi-armed bandits. In *ICML*, volume 12, pages 655–662.
- Karnin, Z., Koren, T., and Somekh, O. (2013). Almost optimal exploration in multi-armed bandits. In *International Conference on Machine Learning*, pages 1238–1246.
- Kaufmann, E., Cappé, O., and Garivier, A. (2016). On the complexity of best-arm identification in multi-armed bandit models. *The Journal of Machine Learning Research*, 17(1):1–42.
- Khetan, A. and Oh, S. (2016). Data-driven rank breaking for efficient rank aggregation. *Journal of Machine Learning Research*, 17(193):1–54.
- Luce, R. D. (2012). *Individual choice behavior: A theoretical analysis*. Courier Corporation.
- Mohajer, S., Suh, C., and Elmahdy, A. (2017). Active learning for top-k rank aggregation from noisy comparisons. In *International Conference on Machine Learning*, pages 2488–2497.
- Nip, K., Wang, Z., and Wang, Z. (2017). Assortment optimization under a single transition model.
- Plackett, R. L. (1975). The analysis of permutations. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 24(2):193–202.

- Popescu, P. G., Dragomir, S., Slusanschi, E. I., and Stanasila, O. N. (2016). Bounds for Kullback-Leibler divergence. *Electronic Journal of Differential Equations*, 2016.
- Ren, W., Liu, J., and Shroff, N. B. (2018). Pac ranking from pairwise and listwise queries: Lower bounds and upper bounds. *arXiv preprint arXiv:1806.02970*.
- Saha, A. and Gopalan, A. (2018a). Active ranking with subset-wise preferences. *International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- Saha, A. and Gopalan, A. (2018b). Battle of bandits. In *Uncertainty in Artificial Intelligence*.
- Saha, A. and Gopalan, A. (2019). PAC Battling Bandits in the Plackett-Luce Model. In *Algorithmic Learning Theory*, pages 700–737.
- Soufiani, H. A., Diao, H., Lai, Z., and Parkes, D. C. (2013). Generalized random utility models with multiple types. In *Advances in Neural Information Processing Systems*, pages 73–81.
- Soufiani, H. A., Parkes, D. C., and Xia, L. (2014). Computing parametric ranking models via rank-breaking. In *ICML*, pages 360–368.
- Sui, Y., Zhuang, V., Burdick, J. W., and Yue, Y. (2017). Multi-dueling bandits with dependent arms. *arXiv preprint arXiv:1705.00253*.
- Szörényi, B., Busa-Fekete, R., Paul, A., and Hüllermeier, E. (2015). Online rank elicitation for plackett-luce: A dueling bandits approach. In *Advances in Neural Information Processing Systems*, pages 604–612.
- Talluri, K. and Van Ryzin, G. (2004). Revenue management under a general discrete choice model of consumer behavior. *Management Science*, 50(1):15–33.
- Thurstone, L. L. (1927). A law of comparative judgment. *Psychological review*, 34(4):273.
- Urvoy, T., Clerot, F., Féraud, R., and Naamane, S. (2013). Generic exploration and k-armed voting bandits. In *International Conference on Machine Learning*, pages 91–99.
- Yue, Y. and Joachims, T. (2009). Interactively optimizing information retrieval systems as a dueling bandits problem. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 1201–1208. ACM.
- Yue, Y. and Joachims, T. (2011). Beat the mean bandit. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 241–248.
- Zhao, Z., Villamil, T., and Xia, L. (2018). Learning mixtures of random utility models. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Zoghi, M., Whiteson, S., and de Rijke, M. (2015). Mergerucb: A method for large-scale online ranker evaluation. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, pages 17–26. ACM.
- Zoghi, M., Whiteson, S., Munos, R., and de Rijke, M. (2013). Relative upper confidence bound for the k-armed dueling bandit problem. *arXiv preprint arXiv:1312.3393*.
- Zoghi, M., Whiteson, S., Munos, R., Rijke, M. d., et al. (2014). Relative upper confidence bound for the k-armed dueling bandit problem. In *JMLR Workshop and Conference Proceedings*, number 32, pages 10–18. JMLR.

# Supplementary for Best-item Learning in Random Utility Models with Subset Choices

## A Appendix for Section 4

### A.1 Proof of Lemma 2

**Lemma 2** (Variational lower bound for the advantage ratio). *For any RUM( $k, \theta$ ) based subsetwise preference model and any item pair  $(i, j)$ ,<sup>7</sup>*

$$\text{Min-AR}(i, j) \geq \min_{z \in \mathbb{R}} \frac{\Pr(X_i > \max(X_j, z))}{\Pr(X_j > \max(X_i, z))}. \quad (3)$$

Moreover for RUM( $k, \theta$ ) models one can show that for any triplet  $(i, j, S)$ ,  $\Pr(X_i > \max(X_j, z)) = F(z - \theta_j)\bar{F}(z - \theta_i) + \int_{z-\theta_j}^{\infty} \bar{F}(x - \Delta_{ij})f(x)dx$ , which further lower bounds Min-AR( $i, j$ ) by:

$$\min_{z \in \mathbb{R}} \frac{F(z - \theta_j)\bar{F}(z - \theta_i) + \int_{z-\theta_j}^{\infty} \bar{F}(x - \Delta_{ij})f(x)dx}{F(z - \theta_i)\bar{F}(z - \theta_j) + \int_{z-\theta_i}^{\infty} \bar{F}(x + \Delta_{ij})f(x)dx}.$$

*Proof.* Let us fix any subset  $S$  and two consider the items  $i, j \in S$  such that  $\theta_i > \theta_j$ . Recall that we also denote by  $\Delta_{ij} = (\theta_i - \theta_j)$ . Let us define a random variable  $X_r^S = \max_{r \in S \setminus \{i, j\}} X_r$  that denotes the maximum score value taken by the rest of the items in set  $S$ . Note that the support of  $X_r^S$ , say denoted by  $\text{supp}(X_r^S) = \max_{r \in S \setminus \{i, j\}} \theta_r + \text{supp}(\mathcal{D})$ .

Let us also denote  $\bar{S} := \arg \min_{S \subseteq [n] \mid |S|=k} \frac{\Pr(i|S)}{\Pr(j|S)}$ . We have:

$$\begin{aligned} \text{Min-AR}(i, j) &= \frac{\Pr(i|\bar{S})}{\Pr(j|\bar{S})} = \frac{\Pr(\{X_i > X_j\} \cap \{X_i > X_r \forall r \in \bar{S} \setminus \{i, j\}\})}{\Pr(\{X_j > X_i\} \cap \{X_j > X_r \forall r \in \bar{S} \setminus \{i, j\}\})} \\ &= \frac{\Pr(\{X_i > X_j\} \cap \{X_i > X_r^{\bar{S}}\})}{\Pr(\{X_j > X_i\} \cap \{X_j > X_r^{\bar{S}}\})} \\ &= \frac{\int_{\text{supp} X_r^{\bar{S}}} \Pr(\{X_i > x\} \cap \{X_i > X_j\}) f_{X_r^{\bar{S}}}(x) dx}{\int_{\text{supp} X_r^{\bar{S}}} \Pr(\{X_i > x\} \cap \{X_j > X_i\}) f_{X_r^{\bar{S}}}(x) dx} \\ &= \frac{\int_{\text{supp} X_r^{\bar{S}}} \Pr(\{X_i > x\} \cap \{X_j > X_i\}) \frac{\Pr(\{X_i > x\} \cap \{X_i > X_j\})}{\Pr(\{X_i > x\} \cap \{X_j > X_i\})} f_{X_r^{\bar{S}}}(x) dx}{\int_{\text{supp} X_r^{\bar{S}}} \Pr(\{X_i > x\} \cap \{X_j > X_i\}) f_{X_r^{\bar{S}}}(x) dx} \\ &> \min_{z \in \text{supp}(X_r^{\bar{S}})} \left[ \frac{\Pr(\{X_i > z\} \cap \{X_i > X_j\})}{\Pr(\{X_i > z\} \cap \{X_j > X_i\})} \right] \frac{\int_{\text{supp} X_r^{\bar{S}}} \Pr(\{X_i > x\} \cap \{X_j > X_i\}) f_{X_r^{\bar{S}}}(x) dx}{\int_{\text{supp} X_r^{\bar{S}}} \Pr(\{X_i > x\} \cap \{X_j > X_i\}) f_{X_r^{\bar{S}}}(x) dx} \\ &= \min_{z \in \text{supp}(X_r^{\bar{S}})} \frac{\Pr(\{X_i > \max(X_j, z)\})}{\Pr(\{X_j > \max(X_i, z)\})} \\ &> \min_{z \in \mathbb{R}} \frac{\Pr(\{X_i > \max(X_j, z)\})}{\Pr(\{X_j > \max(X_i, z)\})} \end{aligned}$$

Let us now introduce a random variable  $Y = \max(X_j, z)$ . Now owing to the ‘independent and identically distributed noise’ assumption of the RUM( $k, \theta$ ) model, we can further show that:

$$\Pr(X_i > \max(X_j, z)) = \Pr(X_i > Y) = \Pr(\{X_i > Y\} \cap \{Y = z\}) + \Pr(\{X_i > Y\} \cap \{Y > z\})$$

<sup>7</sup>We assume  $\frac{0}{0}$  to be  $\infty$  in the right hand side of Eqn. 3.

$$\begin{aligned}
 &= \Pr(\{X_i > z\} \mid \{Y = z\})\Pr(X_j < z) + \Pr(\{X_i > Y\} \cap \{Y > z\}) \\
 &= \Pr(\{\zeta_i + \theta_i > z\})\Pr(\zeta_j + \theta_j < z) + \Pr(\{X_i > X_j\} \cap \{X_j > z\}) \\
 &= \Pr(\{\zeta_i > z - \theta_i\})\Pr(\zeta_j < z - \theta_j) + \Pr(\{\zeta_i > \zeta_j - (\theta_i - \theta_j)\} \cap \{\zeta_j > z - \theta_j\}) \\
 &= F(z - \theta_j)\bar{F}(z - \theta_i) + \int_{z - \theta_j}^{\infty} \bar{F}(x - \Delta_{ij})f(x)dx,
 \end{aligned}$$

which proves the claim.  $\square$

## A.2 Proof of Lemma 3

**Lemma 3** (Analysing *Min-AR* for specific noise models). *Given a fixed item pair  $(i, j)$  such that  $\theta_i > \theta_j$ , the following bounds hold under the respective noise models in an iid RUM.*

1. *Exponential* $(\lambda)$ :  $\text{Min-AR}(i, j) \geq e^{\Delta_{ij}} > 1 + \Delta_{ij}$  for *Exponential* noise with  $\lambda = 1$ .
2. *Extreme value distribution* $(\mu, \sigma, \chi)$ : For *Gumbel* $(\mu, \sigma)$  ( $\chi = 0$ ) noise,  $\text{Min-AR}(i, j) = e^{\frac{\Delta_{ij}}{\sigma}} > 1 + \frac{\Delta_{ij}}{\sigma}$ .
3. *Uniform* $(a, b)$ :  $\text{Min-AR}(i, j) \geq 1 + \frac{2\Delta_{ij}}{b-a}$  for *Uniform* $(a, b)$  noise ( $a, b \in \mathbb{R}, b > a$ , and  $\Delta_{ij} < \frac{a}{2}$ ).
4. *Gamma* $(k, \xi)$ :  $\text{Min-AR}(i, j) \geq 1 + \Delta_{ij}$  for *Gamma* $(2, 1)$  noise.
5. *Weibull* $(\lambda, k)$ :  $\text{Min-AR}(i, j) \geq e^{\lambda\Delta_{ij}} > 1 + \lambda\Delta_{ij}$  for  $(k = 1)$ .
6. *Normal*  $\mathcal{N}(0, 1)$ :  $\exists c > 0$  such that, for  $\Delta_{ij}$  small enough (in a neighborhood of 0),  $\text{Min-AR}(i, j) \geq 1 + c\Delta_{ij}$ .

*Proof.* We can derive the  $\text{Min-AR}(i, j)$  values for the following distributions by simply applying the lower bound formula stated in Thm. 2  $\left( \min_{z \in \mathbb{R}} \frac{F(z - \theta_j)\bar{F}(z - \theta_i) + \int_{z - \theta_j}^{\infty} \bar{F}(x - \Delta_{ij})f(x)dx}{\bar{F}(z - \theta_i)\bar{F}(z - \theta_j) + \int_{z - \theta_i}^{\infty} \bar{F}(x + \Delta_{ij})f(x)dx} \right)$  along with their specific density functions as stated below for each specific distributions:

### 1. Exponential noise:

When the noise distribution  $\mathcal{D}$  is *Exponential* $(1)$ , i.e.  $\zeta_i, \zeta_j \stackrel{iid}{\sim} \text{Exponential}(1)$  note that:  $f(x) = e^{-x}$ ,  $F(x) = 1 - e^{-x}$ , and  $\text{support}(\mathcal{D}) = [0, \infty)$ .

### 2. Gumbel noise:

When the noise distribution  $\mathcal{D}$  is *Gumbel* $(\mu, \sigma)$ , i.e.  $\zeta_i, \zeta_j \stackrel{iid}{\sim} \text{Gumbel}(\mu, \sigma)$  note that:  $f(x) = e^{-\frac{(x-\mu)}{\sigma}} e^{-e^{-\frac{(x-\mu)}{\sigma}}}$ ,  $F(x) = e^{-e^{-\frac{(x-\mu)}{\sigma}}}$ , and  $\text{support}(\mathcal{D}) = (-\infty, \infty)$ .

### 3. Uniform noise case:

When the noise distribution  $\mathcal{D}$  is *Uniform* $(a, b)$ , i.e.  $\zeta_i, \zeta_j \stackrel{iid}{\sim} \text{Uniform}(a, b)$  note that:  $f(x) = \frac{1}{b-a}$ ,  $F(x) = \frac{x-a}{b-a}$ , and  $\text{support}(\mathcal{D}) = [a, b]$ .

### 4. Gamma noise:

When the noise distribution  $\mathcal{D}$  is *Gamma* $(k, \xi)$ , with  $k = 2$  and  $\xi = 1$ , i.e.  $\zeta_i, \zeta_j \stackrel{iid}{\sim} \text{Gamma}(2, 1)$  note that:  $f(x) = xe^{-x}$ ,  $F(x) = 1 - e^{-x} - xe^{-x}$ , and  $\text{support}(\mathcal{D}) = [0, \infty)$ .

### 5. Weibull noise:

When the noise distribution  $\mathcal{D}$  is *Weibull* $(\lambda, k)$ , with  $k = 1$ , i.e.  $\zeta_i, \zeta_j \stackrel{iid}{\sim} \text{Weibull}(\lambda, 1)$  note that:  $f(x) = \frac{1}{\lambda} e^{-\frac{x}{\lambda}}$ ,  $F(x) = 1 - e^{-\frac{x}{\lambda}}$ , and  $\text{support}(\mathcal{D}) = [0, \infty)$ .

**6. Gaussian noise.** (sketch) Note that Gaussian distributions do not have closed form CDFs and are difficult to compute in general, so we propose a different line of analysis specifically for the Gaussian noise case: Take the noise distribution to be standard normal, i.e.,  $\zeta_i, \zeta_j \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ , with density  $f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$ . When  $X_i = \theta_i + \zeta_i$  and  $X_j = \theta_j + \zeta_j$  with  $\Delta_{ij} = \theta_i - \theta_j > 0$ , we describe a method to find a lower bound for the quantity

$$\inf_{z \in \mathbb{R}} \frac{\Pr(X_i > \max(X_j, z))}{\Pr(X_j > \max(X_i, z))}.$$

First, note that by translation, we can take  $\theta_j = 0$  and  $\theta_i = \Delta$  without loss of generality. Doing so allows us to write

$$Pr(X_i > \max(X_j, z)) = F(z)(1 - F(z - \Delta)) + \int_z^\infty (1 - F(y - \Delta))f(y)dy = g(\Delta, z),$$

where

$$g(a, b) = Pr_{(U,V) \sim \mathcal{N}(a,1) \times \mathcal{N}(0,1)} [U > \max(V, b)].$$

It follows that

$$Pr(X_j > \max(X_i, z)) = F(z - \Delta)(1 - F(z)) + \int_{z-\Delta}^\infty (1 - F(y + \Delta))f(y)dy = g(-\Delta, z - \Delta).$$

With this notation, we wish to minimize the ratio  $\frac{g(\Delta, z)}{g(-\Delta, z - \Delta)}$  over  $z \in \mathbb{R}$ .

Notice that  $g(0, z) = \frac{1 - F^2(z)}{2}$ , and  $\frac{\partial g(\Delta, z)}{\partial \Delta} = F(z)f(z - \Delta) + \int_z^\infty f(y - \Delta)f(y)dy$ . Hence, up to first order, for  $\Delta$  small enough, we have<sup>8</sup>

$$\begin{aligned} \frac{g(\Delta, z)}{g(-\Delta, z - \Delta)} &\approx \frac{g(0, z) + \Delta \frac{\partial g(\Delta, z)}{\partial \Delta} \big|_{\Delta=0}}{g(0, z - \Delta) - \Delta \frac{\partial g(\Delta, z - \Delta)}{\partial \Delta} \big|_{\Delta=0}} \\ &= \frac{\frac{1}{2} - \frac{F^2(z)}{2} + \Delta F(z)f(z) + \Delta \int_z^\infty f(y)^2 dy}{\frac{1}{2} - \frac{F^2(z - \Delta)}{2} - \Delta F(z - \Delta)f(z - \Delta) - \Delta \int_{z-\Delta}^\infty f(y)^2 dy} \\ &\equiv \frac{h_1(z)}{h_2(z)}, \quad \text{say.} \end{aligned}$$

Differentiating the above ratio w.r.t.  $z$  and equating it to 0 to find its minimum, we obtain the condition

$$\begin{aligned} h_1'(z_*)h_2(z_*) &= h_1(z_*)h_2'(z_*) \\ \Leftrightarrow (1 - \Delta)F(z_*)f'(z_*)h_2(z_*) &= (1 + \Delta)F(z_* - \Delta)f'(z_* - \Delta)h_1(z_*). \end{aligned} \quad (4)$$

The solution  $z_*$  to (4) is 0 for  $\Delta = 0$ . Since the Gaussian density is infinitely smooth, it follows that there exists a universal constant  $c_1 > 0$  such that the solution  $z_*$  to (4), for a general small  $\Delta$  is,  $z_* = c_1\Delta$  up to first order. This implies that

$$\begin{aligned} \frac{h_1(z_*)}{h_2(z_*)} &= \frac{\frac{1}{2} - \frac{F^2(c_1\Delta)}{2} + \Delta F(c_1\Delta)f(c_1\Delta) + \Delta \int_{c_1\Delta}^\infty f(y)^2 dy}{\frac{1}{2} - \frac{F^2(c_1\Delta - \Delta)}{2} - \Delta F(-c_1\Delta)f(-c_1\Delta) - \Delta \int_{-c_1\Delta}^\infty f(y)^2 dy} \\ &\approx 1 + c\Delta, \end{aligned}$$

for a universal constant  $c > 0$ . This concludes the argument.  $\square$

## B Appendix for Section 5.1

### B.1 Proof of Theorem 4

**Theorem 4** (*Sequential-Pairwise-Battle: Correctness and Sample Complexity*). *Consider any iid subsetwise preference model  $RUM(k, \theta)$  based on a noise distribution  $\mathcal{D}$ , and suppose that for any item pair  $i, j$ , we have  $Min-AR(i, j) \geq 1 + \frac{4c\Delta_{ij}}{1-2c}$  for some  $\mathcal{D}$ -dependent constant  $c > 0$ . Then, Algorithm 1, with input constant  $c > 0$ , is an  $(\epsilon, \delta)$ -PAC algorithm with sample complexity  $O(\frac{n}{c^2\epsilon^2} \log \frac{k}{\delta})$ .*

*Proof.* We start by analyzing the required sample complexity of *Sequential-Pairwise-Battle*. Note that at any iteration  $\ell$ , any set  $\mathcal{G}_\ell$  is played for exactly  $t = \frac{k}{2\epsilon_\ell^2} \ln \frac{k}{\delta_\ell}$  many number of rounds. Also, since the algorithm

<sup>8</sup>The argument can be made rigorous using the Taylor expansion up to 2nd order.

discards exactly  $k - 1$  items from each set  $\mathcal{G}_g$ , the maximum number of iterations possible is  $\lceil \ln_k n \rceil$ . Now at any iteration  $\ell$ , since  $G = \left\lfloor \frac{|S_\ell|}{k} \right\rfloor < \frac{|S_\ell|}{k}$ , the total sample complexity the for iteration is at most  $\frac{|S_\ell|}{k} t \leq \frac{n}{2k^{\ell-1}\epsilon_\ell^2} \ln \frac{k}{\delta_\ell}$ , as  $|S_\ell| \leq \frac{n}{k^\ell}$  for all  $\ell \in [\lceil \ln_k n \rceil]$ . Also note that for all but last iteration  $\ell \in [\lceil \ln_k n \rceil]$ , we have  $\epsilon_\ell = \frac{c\epsilon}{8} \left(\frac{3}{4}\right)^{\ell-1}$ , and  $\delta_\ell = \frac{\delta}{2^{\ell+1}}$ . Moreover, for the last iteration  $\ell = \lceil \ln_k n \rceil$ , the sample complexity is clearly  $t = \frac{2k}{c^2\epsilon^2} \ln \frac{2k}{\delta}$ , as in this case  $\epsilon_\ell = \frac{c\epsilon}{2}$ , and  $\delta_\ell = \frac{\delta}{2}$ , and  $|S| = k$ . Thus, the total sample complexity of Algorithm 1 is given by

$$\begin{aligned} \sum_{\ell=1}^{\lceil \ln_k n \rceil} \frac{|S_\ell|}{2\epsilon_\ell^2} \ln \frac{k}{\delta_\ell} &\leq \sum_{\ell=1}^{\infty} \frac{n}{2k^\ell \left(\frac{c\epsilon}{8} \left(\frac{3}{4}\right)^{\ell-1}\right)^2} k \ln \frac{k2^{\ell+1}}{\delta} + \frac{2k}{c^2\epsilon^2} \ln \frac{2k}{\delta} \\ &\leq \frac{64n}{2c^2\epsilon^2} \sum_{\ell=1}^{\infty} \frac{16^{\ell-1}}{(9k)^{\ell-1}} \left( \ln \frac{k}{\delta} + (\ell+1) \right) + \frac{2k}{c^2\epsilon^2} \ln \frac{2k}{\delta} \\ &\leq \frac{32n}{c^2\epsilon^2} \ln \frac{k}{\delta} \sum_{\ell=1}^{\infty} \frac{4^{\ell-1}}{(9k)^{\ell-1}} (3\ell) + \frac{2k}{c^2\epsilon^2} \ln \frac{2k}{\delta} = O\left(\frac{n}{c^2\epsilon^2} \ln \frac{k}{\delta}\right) \text{ [for any } k > 1], \end{aligned}$$

and this proves the sample complexity bound of Theorem 4. We next prove the  $(\epsilon, \delta)$ -PAC property of *Sequential-Pairwise-Battle*.

Consider any fixed subgroup  $\mathcal{G}$  of size  $k$ , such that two items  $a, b \in \mathcal{G}$ . Now suppose we denote by  $Pr(\{ab\}|\mathcal{G}) = Pr(a|\mathcal{G}) + Pr(b|\mathcal{G})$  the probability that either  $a$  or  $b$  wins in the subset  $\mathcal{G}$ . Then the probability that  $a$  wins in  $\mathcal{G}$  given either  $a$  or  $b$  won in  $\mathcal{G}$  is given by  $p_{ab|\mathcal{G}} := \frac{Pr(a|\mathcal{G})}{Pr(\{ab\}|\mathcal{G})} = \frac{Pr(a|\mathcal{G})}{Pr(a|\mathcal{G}) + Pr(b|\mathcal{G})}$  — this quantity in a way models the pairwise preference of  $a$  over  $b$  in the set  $\mathcal{G}$ . Note that as long as  $\theta_a > \theta_b$ ,  $p_{ab|\mathcal{G}} > \frac{1}{2}$ , for any  $\mathcal{G}$  (since  $Pr(a|\mathcal{G}) > Pr(b|\mathcal{G})$ ). We in fact now introduce the notation  $p_{ab} := \min_{\mathcal{G} \subseteq [n], |\mathcal{G}|=k} p_{ab|\mathcal{G}}$ .

**Lemma 9.** *For any item pair  $i, j \in [n]$  and any set  $S \subseteq [n]$ , if their advantage ratio  $\frac{Pr(i|S)}{Pr(j|S)} \geq 1 + \alpha$ , for some  $\alpha > 0$ , then pairwise preference of item  $i$  over  $j$  in set  $S$   $p_{ij|S} > \frac{1}{2} + \frac{\alpha}{4}$ .*

*Proof.* Note that

$$\begin{aligned} \frac{Pr(i|S)}{Pr(j|S)} \geq 1 + \alpha &\implies \frac{Pr(i|S) - Pr(j|S)}{Pr(j|S)} \geq \alpha \\ \implies p_{ij|S} - 0.5 &= \frac{Pr(i|S) - Pr(j|S)}{2(Pr(i|S) + Pr(j|S))} \geq \frac{\alpha Pr(j|S)}{2(Pr(j|S) + Pr(j|S))} = \frac{\alpha}{4}, \end{aligned}$$

which concludes the proof.  $\square$

**Corollary 10.** *For any item pair  $i, j \in [n]$ , if  $Min-AR(i, j) \geq 1 + \alpha$  for some  $\alpha > 0$ , then  $p_{ij} > \frac{1}{2} + \frac{\alpha}{4}$ .*

*Proof.* The proof directly follows from Lem . 9 by using subset  $S = \min_{S \subseteq [n], |S|=k} Min-AR(i, j)$ .  $\square$

Let us denote the set of surviving items  $S$  at the beginning of phase  $\ell$  as  $S_\ell$ . We now claim the following crucial lemma which shows at any phase  $\ell$ , the best (the one with highest  $\theta$  parameter) item retained in  $S_{\ell+1}$  can not be too bad in comparison to the best item of  $S_\ell$ . The formal claim goes as follows:

**Lemma 11.** *At any iteration  $\ell$ , for any  $\mathcal{G}_g$ , if  $i_g := \arg \max_{i \in \mathcal{G}_g} \theta_i$ , then with probability at least  $(1 - \delta_\ell)$ ,  $\theta_{c_g} > \theta_{i_g} - \frac{\epsilon_\ell}{c}$ .*

*Proof.* Let us define  $\hat{p}_{ij} = \frac{w_i}{w_i + w_j}$ ,  $\forall i, j \in \mathcal{G}_g, i \neq j$ . Then clearly  $\hat{p}_{c_g i_g} \geq \frac{1}{2}$ , as  $c_g$  is the empirical winner in  $t$  rounds, i.e.  $c_g \leftarrow \arg \max_{i \in \mathcal{G}_g} w_i$ . Moreover  $c_g$  being the empirical winner of  $\mathcal{G}_g$  we also have  $w_{c_g} \geq \frac{t}{k}$ , and thus  $w_{c_g} + w_{r_g} \geq \frac{t}{k}$  as well. Let  $n_{ij} := w_i + w_j$  denotes the number of pairwise comparisons of item  $i$  and  $j$  in  $t$  rounds,

$i, j \in \mathcal{G}_g$ . Clearly  $0 \leq n_{ij} \leq t$ . Then let us analyze the probability of a ‘bad event’ where  $c_g$  is indeed such that  $\theta_{c_g} < \theta_{i_g} - \frac{\epsilon_\ell}{c}$ .

This implies that the advantage ratio of  $i_g$  and  $c_g$  in  $\mathcal{G}$  is  $\frac{Pr(i_g|\mathcal{G})}{Pr(c_g|\mathcal{G})} \geq 1 + 4\epsilon_\ell$ .

But now by Lem. 9 this further implies  $p_{i_g c_g|\mathcal{G}} \geq \frac{1}{2} + \epsilon_\ell$ . But since  $c_g$  beats  $i_g$  empirically in the subgroup  $\mathcal{G}$ , this implies  $\hat{p}_{c_g i_g} > \frac{1}{2}$ . The following argument shows that this is even unlikely to happen, more formally with probability  $(1 - \delta_\ell/k)$ :

$$\begin{aligned} & Pr\left(\{\hat{p}_{c_g i_g} \geq \frac{1}{2}\}\right) \\ &= Pr\left(\{\hat{p}_{c_g i_g} \geq \frac{1}{2}\} \cap \{n_{c_g i_g} \geq \frac{t}{k}\}\right) + Pr\left(\{n_{c_g i_g} < \frac{t}{k}\}\right) Pr\left(\{\hat{p}_{c_g i_g} \geq \frac{1}{2}\} \mid \{n_{c_g i_g} < \frac{t}{k}\}\right) \\ &= Pr\left(\{\hat{p}_{c_g i_g} - \epsilon_\ell \geq \frac{1}{2} - \epsilon_\ell\} \cap \{n_{c_g i_g} \geq \frac{t}{k}\}\right) \\ &\leq Pr\left(\{\hat{p}_{c_g i_g} - p_{c_g i_g|\mathcal{G}} \geq \epsilon_\ell\} \cap \{n_{c_g i_g} \geq \frac{t}{k}\}\right) \\ &\leq \exp\left(-2\frac{t}{k}(\epsilon_\ell)^2\right) = \frac{\delta_\ell}{k}. \end{aligned}$$

where the first inequality holds as  $p_{c_g i_g|\mathcal{G}} < \frac{1}{2} - \epsilon_\ell$ , and the second inequality follows from Hoeffdings lemma. Now taking the union bound over all  $\epsilon_\ell$ -suboptimal elements  $i'$  of  $\mathcal{G}_g$  (i.e.  $\theta_{i'} < \theta_{i_g} - \epsilon_\ell$ ), we get:

$$Pr\left(\left\{\exists i' \in \mathcal{G}_g \mid p_{i' i_g} < \frac{1}{2} - \epsilon_\ell, \text{ and } c_g = i'\right\}\right) \leq \frac{\delta_\ell}{k} \left|\left\{\exists i' \in \mathcal{G}_g \mid p_{i' i_g} < \frac{1}{2} - \epsilon_\ell, \text{ and } c_g = i'\right\}\right| \leq \delta_\ell,$$

as  $|\mathcal{G}_g| = k$ , and the claim follows henceforth.  $\square$

Let us denote the single element remaining in  $S$  at termination by  $r \in [n]$ . Also note that for the last iteration  $\ell = \lceil \ln_k n \rceil$ , since  $\epsilon_\ell = \frac{\epsilon}{2}$ , and  $\delta_\ell = \frac{\delta}{2}$ , applying Lemma 11 on  $S$ , we get that  $Pr\left(\theta_r < \theta_{i_g} - \frac{\epsilon}{2}\right) \leq \frac{\delta}{2}$ .

Without loss of generality we assume the best item of the RUM( $k, \theta$ ) model is  $\theta_1$ , i.e.  $\theta_1 > \theta_i \forall i \in [n] \setminus \{1\}$ . Now for any iteration  $\ell$ , let us define  $g_\ell \in [G]$  to be the index of the set that contains *best item* of the entire set  $S_\ell$ , i.e.  $\arg \max_{i \in S_\ell} \theta_i \in \mathcal{G}_{g_\ell}$ . Then applying Lemma 11, with probability at least  $(1 - \delta_\ell)$ ,  $\theta_{c_{g_\ell}} > \theta_{i_{g_\ell}} - \epsilon_\ell/c$ . Note that initially, at phase  $\ell = 1$ ,  $i_{g_\ell} = 1$ . Then, for each iteration  $\ell$ , applying Lemma 11 recursively to  $\mathcal{G}_{g_\ell}$ , we finally get  $\theta_r > \theta_1 - \left(\frac{\epsilon}{8} + \frac{\epsilon}{8}\left(\frac{3}{4}\right) + \dots + \frac{\epsilon}{8}\left(\frac{3}{4}\right)^{\lceil \ln_k n \rceil}\right) - \frac{\epsilon}{2} \geq \theta_1 - \frac{\epsilon}{8}\left(\sum_{i=0}^{\infty} \left(\frac{3}{4}\right)^i\right) - \frac{\epsilon}{2} \geq \theta_1 - \epsilon$ . Thus assuming the algorithm does not fail in any of the iteration  $\ell$ , we finally have that  $p_{r*1} > \frac{1}{2} - \epsilon$ —this shows that the final item output by Seq-PB is  $\epsilon$  optimal.

Finally since at any phase  $\ell$ , the algorithm fails with probability at most  $\delta_\ell$ , the total failure probability of the algorithm is at most  $\left(\frac{\delta}{4} + \frac{\delta}{8} + \dots + \frac{\delta}{2^{\lceil \ln_k n \rceil}}\right) + \frac{\delta}{2} \leq \delta$ . This concludes the correctness of the algorithm showing that it indeed satisfies the  $(\epsilon, \delta)$ -PAC objective.  $\square$

## B.2 Proof of Corollary 5

*Proof.* The proof essentially follows from the general performance guarantee of Seq-PB (Thm. 4) and Lem. 3. More specifically from Lem. 3 it follows that the value of  $c$  for these specific distributions are constant, which concludes the claim. For completeness the distribution-specific values of  $c$  are given below:

1.  $c = 0.25$  for Exponential noise with  $\lambda = 1$
2.  $c = \frac{0.25}{\sigma}$  for  $Gumbel(\mu, \sigma)$
3.  $c = \frac{0.5}{(b-a)}$  for  $Uniform(a, b)$
4.  $c = \frac{1}{4}$  for  $Gamma(2, 1)$



5.  $c = \frac{\lambda}{4}$  for *Weibull*( $\lambda, 1$ )
6.  $c = \frac{1}{3}$  *Normal*  $\mathcal{N}(0, 1)$ , etc.

□

### B.3 Proof of Theorem 6

Before proving the lower bound result we state a key lemma from (Kaufmann et al., 2016) which is a general result for proving information theoretic lower bound for bandit problems:

Consider a multi-armed bandit (MAB) problem with  $n$  arms or actions  $\mathcal{A} = [n]$ . At round  $t$ , let  $A_t$  and  $Z_t$  denote the arm played and the observation (reward) received, respectively. Let  $\mathcal{F}_t = \sigma(A_1, Z_1, \dots, A_t, Z_t)$  be the sigma algebra generated by the trajectory of a sequential bandit algorithm up to round  $t$ .

**Lemma 12** (Lemma 1, (Kaufmann et al., 2016)). *Let  $\nu$  and  $\nu'$  be two bandit models (assignments of reward distributions to arms), such that  $\nu_i$  (resp.  $\nu'_i$ ) is the reward distribution of any arm  $i \in \mathcal{A}$  under bandit model  $\nu$  (resp.  $\nu'$ ), and such that for all such arms  $i$ ,  $\nu_i$  and  $\nu'_i$  are mutually absolutely continuous. Then for any almost-surely finite stopping time  $\tau$  with respect to  $(\mathcal{F}_t)_t$ ,*

$$\sum_{i=1}^n \mathbf{E}_{\nu}[N_i(\tau)] KL(\nu_i, \nu'_i) \geq \sup_{\mathcal{E} \in \mathcal{F}_{\tau}} kl(Pr_{\nu}(\mathcal{E}), Pr_{\nu'}(\mathcal{E})),$$

where  $kl(x, y) := x \log(\frac{x}{y}) + (1-x) \log(\frac{1-x}{1-y})$  is the binary relative entropy,  $N_i(\tau)$  denotes the number of times arm  $i$  is played in  $\tau$  rounds, and  $Pr_{\nu}(\mathcal{E})$  and  $Pr_{\nu'}(\mathcal{E})$  denote the probability of any event  $\mathcal{E} \in \mathcal{F}_{\tau}$  under bandit models  $\nu$  and  $\nu'$ , respectively.

We now proceed to proof our lower bound result of Thm. 6.

**Theorem 6** (Sample Complexity Lower Bound for RUM( $k, \theta$ ) model). *Given  $\epsilon \in (0, \frac{1}{4}]$ ,  $\delta \in (0, 1]$ ,  $c > 0$  and an  $(\epsilon, \delta)$ -PAC algorithm  $A$  with winner item feedback, there exists a RUM( $k, \theta$ ) instance  $\nu$  with  $Min-AR(i, j) \geq 1 + 4c\Delta_{ij}$  for all  $i, j \in [n]$ , where the expected sample complexity of  $A$  on  $\nu$  is at least  $\Omega(\frac{n}{c^2\epsilon^2} \ln \frac{1}{2.4\delta})$ .*

*Proof.* In order to apply the change of measure based lemma Lem. 12, we constructed the following specific instances of the RUM( $k, \theta$ ) model for our purpose and assume  $\mathcal{D}$  to be the *Gumbel*(0, 1) noise:

$$\text{True Instance } (\nu^1) : \theta_j^1 = 1 - \epsilon, \forall j \in [n] \setminus \{1\}, \text{ and } \theta_1^1 = 1,$$

Note the only  $\epsilon$ -optimal arm in the true instance is arm 1. Now for every suboptimal item  $a \in [n] \setminus \{1\}$ , consider the modified instances  $\nu^a$  such that:

$$\text{Instance-a } (\nu^a) : \theta_j^a = 1 - 2\epsilon, \forall j \in [n] \setminus \{a, 1\}, \theta_1^a = 1 - \epsilon, \text{ and } \theta_a^a = 1.$$

For any problem instance  $\nu^a$ ,  $a \in [n] \setminus \{1\}$ , the probability distribution associated with arm  $S \in \mathcal{A}$  is given by

$$\nu_S^a \sim \text{Categorical}(p_1, p_2, \dots, p_k), \text{ where } p_i = Pr(i|S), \forall i \in [k], \forall S \in \mathcal{A},$$

where  $Pr(i|S)$  is as defined in Section 3.1. Note that the only  $\epsilon$ -optimal arm for **Instance-a** is arm  $a$ . Now applying Lemma 12, for any event  $\mathcal{E} \in \mathcal{F}_{\tau}$  we get,

$$\sum_{\{S \in \mathcal{A} : a \in S\}} \mathbf{E}_{\nu^1}[N_S(\tau_A)] KL(\nu_S^1, \nu_S^a) \geq kl(Pr_{\nu}(\mathcal{E}), Pr_{\nu'}(\mathcal{E})). \quad (5)$$

The above result holds from the straightforward observation that for any arm  $S \in \mathcal{A}$  with  $a \notin S$ ,  $\nu_S^1$  is same as  $\nu_S^a$ , hence  $KL(\nu_S^1, \nu_S^a) = 0$ ,  $\forall S \in \mathcal{A}$ ,  $a \notin S$ . For notational convenience, we will henceforth denote  $S^a = \{S \in \mathcal{A} : a \in S\}$ .

Now let us analyse the right hand side of (5), for any set  $S \in S^a$ .

**Case-1:** First let us consider  $S \in S^a$  such that  $1 \notin S$ . Note that in this case:

$$\nu_S^1(i) = \frac{1}{k}, \text{ for all } i \in S$$

On the other hand, for problem **Instance-a**, we have that:

$$\nu_S^a(i) = \begin{cases} \frac{e^1}{(k-1)e^{1-2\epsilon} + e^1} & \text{when } S(i) = a, \\ \frac{e^{1-2\epsilon}}{(k-1)e^{1-2\epsilon} + e^1}, & \text{otherwise} \end{cases}$$

Now using the following upper bound on  $KL(\mathbf{p}_1, \mathbf{p}_2) \leq \sum_{x \in \mathcal{X}} \frac{p_1^2(x)}{p_2(x)} - 1$ ,  $\mathbf{p}_1$  and  $\mathbf{p}_2$  be two probability mass functions on the discrete random variable  $\mathcal{X}$  (Popescu et al., 2016) we get:

$$\begin{aligned} KL(\nu_S^1, \nu_S^a) &\leq (k-1) \frac{(k-1)e^{1-2\epsilon} + e^1}{k^2(e^{1-2\epsilon})} + \frac{(k-1)e^{1-2\epsilon} + e^1}{k^1 e^1} - 1 \\ &= \frac{(k-1)}{k^2} \left( e^\epsilon - e^{-\epsilon} \right)^2 = \frac{(k-1)}{k^2} e^{-2\epsilon} (e^\epsilon - 1)^2 \leq \frac{\epsilon^2}{k} \text{ for any } \epsilon \in \left[ 0, \frac{1}{2} \right] \end{aligned}$$

**Case-2:** Now let us consider the remaining set in  $S^a$  such that  $S \ni 1, a$ . Similar to the earlier case in this case we get that:

$$\nu_S^a(i) = \begin{cases} \frac{e^1}{(k-1)e^{1-\epsilon} + e^1} & \text{when } S(i) = 1, \\ \frac{e^{1-\epsilon}}{(k-1)e^{1-\epsilon} + e^1}, & \text{otherwise} \end{cases}$$

On the other hand, for problem **Instance-a**, we have that:

$$\nu_S^a(i) = \begin{cases} \frac{e^{1-\epsilon}}{(k-2)e^{1-2\epsilon} + e^{1-\epsilon} + e^1} & \text{when } S(i) = 1, \\ \frac{e^1}{(k-2)e^{1-2\epsilon} + e^{1-\epsilon} + e^1} & \text{when } S(i) = a, \\ \frac{e^{1-2\epsilon}}{(k-2)e^{1-2\epsilon} + e^{1-\epsilon} + e^1}, & \text{otherwise} \end{cases}$$

Now using the previously mentioned upper bound on the KL divergence, followed by some elementary calculations one can show that for any  $[0, \frac{1}{4}]$ :

$$KL(\nu_S^1, \nu_S^a) \leq \frac{8\epsilon^2}{k}$$

Thus combining the above two cases we can conclude that for any  $S \in S^a$ ,  $KL(\nu_S^1, \nu_S^a) \leq \frac{8\epsilon^2}{k}$ , and as argued above for any  $S \notin S^a$ ,  $KL(\nu_S^1, \nu_S^a) = 0$ .

Note that the only  $\epsilon$ -optimal arm for any **Instance-a** is arm  $a$ , for all  $a \in [n]$ . Now, consider  $\mathcal{E}_0 \in \mathcal{F}_\tau$  be an event such that the algorithm  $A$  returns the element  $i = 1$ , and let us analyse the left hand side of (5) for  $\mathcal{E} = \mathcal{E}_0$ . Clearly,  $A$  being an  $(\epsilon, \delta)$ -PAC algorithm, we have  $Pr_{\nu^1}(\mathcal{E}_0) > 1 - \delta$ , and  $Pr_{\nu^a}(\mathcal{E}_0) < \delta$ , for any suboptimal arm  $a \in [n] \setminus \{1\}$ . Then we have

$$kl(Pr_{\nu^1}(\mathcal{E}_0), Pr_{\nu^a}(\mathcal{E}_0)) \geq kl(1 - \delta, \delta) \geq \ln \frac{1}{2.4\delta} \quad (6)$$

where the last inequality follows from (Kaufmann et al., 2016) (Eqn. 3).

Now applying (5) for each modified bandit **Instance- $\nu^a$** , and summing over all suboptimal items  $a \in [n] \setminus \{1\}$  we get,

$$\sum_{a=2}^n \sum_{\{S \in \mathcal{A} | a \in S\}} \mathbf{E}_{\nu^1}[N_S(\tau_A)] KL(\nu_S^1, \nu_S^a) \geq (n-1) \ln \frac{1}{2.4\delta}. \quad (7)$$

Using the upper bounds on  $KL(\nu_S^1, \nu_S^a)$  as shown above, the right hand side of (7) can be further upper bounded as:

$$\begin{aligned} \sum_{a=2}^n \sum_{\{S \in \mathcal{A} | a \in S\}} \mathbf{E}_{\nu^1}[N_S(\tau_A)] KL(\nu_S^1, \nu_S^a) &\leq \sum_{S \in \mathcal{A}} \mathbf{E}_{\nu^1}[N_S(\tau_A)] \sum_{\{a \in S | a \neq 1\}} \frac{8\epsilon^2}{k} \\ &= \sum_{S \in \mathcal{A}} \mathbf{E}_{\nu^1}[N_S(\tau_A)] k - (\mathbf{1}(1 \in S)) \frac{8\epsilon^2}{k} \leq \sum_{S \in \mathcal{A}} \mathbf{E}_{\nu^1}[N_S(\tau_A)] 8\epsilon^2. \end{aligned} \quad (8)$$

Finally noting that  $\tau_A = \sum_{S \in \mathcal{A}} [N_S(\tau_A)]$ , combining (7) and (8), we get

$$(8\epsilon^2) \mathbf{E}_{\nu^1}[\tau_A] = \sum_{S \in \mathcal{A}} \mathbf{E}_{\nu^1}[N_S(\tau_A)] (8\epsilon^2) \geq (n-1) \ln \frac{1}{2.4\delta}. \quad (9)$$

Now note that as derived in Lem. 3, for *Gumbel*(0,1) noise, we have shown that for any pair  $i, j \in [n]$ ,  $\text{Min-AR}(i, j) = e^{\Delta_{ij}} > 1 + \Delta_{ij} = 1 + 4\frac{1}{4}\Delta_{ij} \implies$  the value of the noise dependent constant  $c$  can be taken to be  $c = \frac{1}{4}$ . Thus rewriting Eqn. 9 we get  $\mathbf{E}_{\nu^1}[\tau_A] \geq \frac{(n-1)}{8\epsilon^2} \ln \frac{1}{2.4\delta} = \frac{(n-1)}{128c^2\epsilon^2} \ln \frac{1}{2.4\delta}$ . The above construction shows the existence of a problem instance of  $\text{RUM}(k, \theta)$  model where any  $(\epsilon, \delta)$ -PAC algorithm requires at least  $\Omega(\frac{n}{c^2\epsilon^2} \ln \frac{1}{2.4\delta})$  samples to ensure correctness of its performance, concluding our proof.  $\square$

**Remark 5.** *It is worth noting that our lower bound analysis is essentially in spirit the same as the one proposed by (Saha and Gopalan, 2019) for the Plackett-Luce model. However note that, their PAC objective is quite different than the one considered in our case—precisely their model is positive scale invariant, unlike ours which is shift invariant w.r.t the model parameters  $\theta$ . Moreover our setting aims to find a  $\epsilon$ -best item in additive sense (i.e. to find an item  $i$  whose score difference w.r.t to the best item 1 is at most  $\epsilon > 0$  or  $\theta_1 - \theta_i < \epsilon$ ), as opposed to the  $(\epsilon, \delta)$ -PAC objective considered in (Saha and Gopalan, 2019) which seeks to find a multiplicative- $\epsilon$ -best item (i.e. to find an item  $i$  which matches the score of the best item up to  $\epsilon$ -factor or  $\theta_i > \epsilon\theta_1$ ). Therefore the problem instance construction for proving a suitable lower bound these two setups are very different where lies the novelty of our current lower bound analysis.*

## C Appendix for Section 6

### C.1 Pseudo code of Sequential-Pairwise-Battle for top- $m$ ranking feedback (mSeq-PB)

The description is given in Algorithm 2.

### C.2 Proof of Theorem 7

**Theorem 7** (mSeq-PB(Alg. 2): Correctness and Sample Complexity). *Consider any  $\text{RUM}(k, \theta)$  subsetwise preference model based on noise distribution  $\mathcal{D}$  and suppose for any item pair  $i, j$ , we have  $\text{Min-AR}(i, j) \geq 1 + \frac{4c\Delta_{ij}}{1-2c}$  for some  $\mathcal{D}$ -dependent constant  $c > 0$ . Then mSeq-PB (Alg.2) with input constant  $c > 0$  on top- $m$  ranking feedback model is an  $(\epsilon, \delta)$ -PAC algorithm with sample complexity  $O(\frac{n}{mc^2\epsilon^2} \log \frac{k}{\delta})$ .*

---

**Algorithm 2** *Sequential-Pairwise-Battle* (TR- $m$  feedback)
 

---

```

1: Input:
2:   Set of items:  $[n]$ , and subset size:  $k > 2$  ( $n \geq k \geq m$ )
3:   Error bias:  $\epsilon > 0$ , and confidence parameter:  $\delta > 0$ 
4:   Noise model ( $\mathcal{D}$ ) dependent constant  $c > 0$ 
5: Initialize:
6:    $S \leftarrow [n]$ ,  $\epsilon_0 \leftarrow \frac{c\epsilon}{8}$ , and  $\delta_0 \leftarrow \frac{\delta}{2}$ 
7:   Divide  $S$  into  $G := \lceil \frac{n}{k} \rceil$  sets  $\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_G$  such that  $\cup_{j=1}^G \mathcal{G}_j = S$  and  $\mathcal{G}_j \cap \mathcal{G}_{j'} = \emptyset$ ,  $\forall j, j' \in [G]$ ,  $|\mathcal{G}_j| = k$ ,  $\forall j \in [G-1]$ . If  $|\mathcal{G}_G| < k$ , then set  $\mathcal{R}_1 \leftarrow \mathcal{G}_G$  and  $G = G - 1$ .
8: while  $\ell = 1, 2, \dots$  do
9:   Set  $S \leftarrow \emptyset$ ,  $\delta_\ell \leftarrow \frac{\delta_{\ell-1}}{2}$ ,  $\epsilon_\ell \leftarrow \frac{3}{4}\epsilon_{\ell-1}$ 
10:  for  $g = 1, 2, \dots, G$  do
11:    Initialize pairwise (empirical) win-count  $w_{ij} \leftarrow 0$ , for each item pair  $i, j \in \mathcal{G}_g$ 
12:    for  $\tau = 1, 2, \dots, t$  ( $t := \lceil \frac{4k}{m\epsilon_\ell^2} \ln \frac{2k}{\delta_\ell} \rceil$ ) do
13:      Play the set  $\mathcal{G}_g$  (one round of battle)
14:      Receive: The top- $m$  ranking  $\sigma_\tau \in \Sigma_{\mathcal{G}}^m$ 
15:      Update win-count  $w_{ij}$  of each item pair  $i, j \in \mathcal{G}_g$  applying Rank-Breaking on  $\sigma_\tau$ 
16:    end for
17:    Define  $\hat{p}_{i,j} = \frac{w_{ij}}{w_{ij} + w_{ji}}$ ,  $\forall i, j \in \mathcal{G}_g$ 
18:    If  $\exists$  any  $i \in \mathcal{G}_g$  such that  $\hat{p}_{ij} + \frac{\epsilon_\ell}{2} \geq \frac{1}{2}$ ,  $\forall j \in \mathcal{G}_g$ , then set  $c_g \leftarrow i$ , else select  $c_g \leftarrow$  uniformly at random from  $\mathcal{G}_g$ , and set  $S \leftarrow S \cup \{c_g\}$ 
19:  end for
20:   $S \leftarrow S \cup \mathcal{R}_\ell$ 
21:  if ( $|S| == 1$ ) then
22:    Break (go out of the while loop)
23:  else if  $|S| \leq k$  then
24:     $S' \leftarrow$  Randomly sample  $k - |S|$  items from  $[n] \setminus S$ , and  $S \leftarrow S \cup S'$ ,  $\epsilon_\ell \leftarrow \frac{c\epsilon}{2}$ ,  $\delta_\ell \leftarrow \delta$ 
25:  else
26:    Divide  $S$  into  $G := \lceil \frac{|S|}{k} \rceil$  sets  $\mathcal{G}_1, \dots, \mathcal{G}_G$  such that  $\cup_{j=1}^G \mathcal{G}_j = S$ ,  $\mathcal{G}_j \cap \mathcal{G}_{j'} = \emptyset$ ,  $\forall j, j' \in [G]$ ,  $|\mathcal{G}_j| = k$ ,  $\forall j \in [G-1]$ . If  $|\mathcal{G}_G| < k$ , then set  $\mathcal{R}_{\ell+1} \leftarrow \mathcal{G}_G$  and  $G = G - 1$ .
27:  end if
28: end while
29: Output: The unique item left in  $S$ 
    
```

---

*Proof.* Same as the proof of Thm. 4, we start by analyzing the required sample complexity of the algorithm. Note that at any iteration  $\ell$ , any set  $\mathcal{G}_g$  is played for exactly  $t = \frac{4k}{m\epsilon_\ell^2} \ln \frac{2k}{\delta_\ell}$  many number of times. Also since the algorithm discards away exactly  $k - 1$  items from each set  $\mathcal{G}_g$ , hence the maximum number of iterations possible is  $\lceil \ln_k n \rceil$ . Now at any iteration  $\ell$ , since  $G = \lceil \frac{|S_\ell|}{k} \rceil < \frac{|S_\ell|}{k}$ , the total sample complexity for iteration  $\ell$  is at most  $\frac{|S_\ell|}{k} t \leq \frac{4n}{mk^{\ell-1}\epsilon_\ell^2} \ln \frac{2k}{\delta_\ell}$ , as  $|S_\ell| \leq \frac{n}{k^\ell}$  for all  $\ell \in [\lceil \ln_k n \rceil]$ . Also note that for all but last iteration  $\ell \in [\lceil \ln_k n \rceil]$ ,  $\epsilon_\ell = \frac{\epsilon}{8} \left(\frac{3}{4}\right)^{\ell-1}$ , and  $\delta_\ell = \frac{\delta}{2^{\ell-1}}$ . Moreover for the last iteration  $\ell = \lceil \ln_k n \rceil$ , the sample complexity is clearly  $t = \frac{4k}{mc^2(\epsilon/2)^2} \ln \frac{4k}{\delta}$ , as in this case  $\epsilon_\ell = \frac{c\epsilon}{2}$ , and  $\delta_\ell = \frac{\delta}{2}$ , and  $|S| = k$ . Thus the total sample complexity of Algorithm 2 is given by

$$\begin{aligned}
 \sum_{\ell=1}^{\lceil \ln_k n \rceil} \frac{|S_\ell|}{m(\epsilon_\ell/2)^2} \ln \frac{2k}{\delta_\ell} &\leq \sum_{\ell=1}^{\infty} \frac{4n}{mc^2 k^\ell \left(\frac{\epsilon}{8} \left(\frac{3}{4}\right)^{\ell-1}\right)^2} k \ln \frac{k 2^{\ell+1}}{\delta} + \frac{16k}{mc^2 \epsilon^2} \ln \frac{4k}{\delta} \\
 &\leq \frac{256n}{mc^2 \epsilon^2} \sum_{\ell=1}^{\infty} \frac{16^{\ell-1}}{(9k)^{\ell-1}} \left( \ln \frac{k}{\delta} + (\ell+1) \right) + \frac{16k}{mc^2 \epsilon^2} \ln \frac{4k}{\delta}
 \end{aligned}$$

$$\leq \frac{256n}{mc^2\epsilon^2} \ln \frac{k}{\delta} \sum_{\ell=1}^{\infty} \frac{4^{\ell-1}}{(9k)^{\ell-1}} (3\ell) + \frac{16k}{mc^2\epsilon^2} \ln \frac{4k}{\delta} = O\left(\frac{n}{mc^2\epsilon^2} \ln \frac{k}{\delta}\right) \text{ [for any } k > 1].$$

We are now only left with proving the  $(\epsilon, \delta)$ -PAC correctness of the algorithm. We used the same notations as introduced in the proof of Thm. 4.

We start by making a crucial observation that at any phase, for any subgroup  $\mathcal{G}_g$ , the strongest item of the  $\mathcal{G}_g$  gets picked in the top- $m$  ranking quite often. More formally:

**Lemma 13.** *Consider any particular set  $\mathcal{G}_g$  at any phase  $\ell$ , and let us denote by  $q_i$  as the number of times any item  $i \in \mathcal{G}_g$  appears in the top- $m$  rankings when items in the set  $\mathcal{G}_g$  are queried for  $t$  rounds. Then if  $i_g := \arg \max_{i \in \mathcal{G}_g} \theta_i$ , then with probability at least  $\left(1 - \frac{\delta_\ell}{2k}\right)$ , one can show that  $q_{i_g} > (1 - \eta) \frac{mt}{k}$ , for any  $\eta \in \left(\frac{3}{32\sqrt{2}}, 1\right]$ .*

*Proof.* Fix any iteration  $\ell$  and a set  $\mathcal{G}_g$ ,  $g \in 1, 2, \dots, G$ . Define  $i_g^\tau := \mathbf{1}(i \in \sigma_\tau)$  as the indicator variable if  $i^{\text{th}}$  element appeared in the top- $m$  ranking at iteration  $\tau \in [t]$ . Recall the definition of TR feedback model (Sec. 3.1). Using this we get  $\mathbf{E}[i_g^\tau] = \Pr(\{i_g \in \sigma\}) = \Pr(\exists j \in [m] \mid \sigma(j) = i_g) = \sum_{j=1}^m \Pr(\sigma(j) = i_g) = \sum_{j=0}^{m-1} \frac{1}{k-j} \geq \frac{m}{k}$ , as  $\Pr(\{i_g | S\}) \geq \frac{1}{|S|}$  for any  $S \subseteq [\mathcal{G}_g]$  ( $i_g := \arg \max_{i \in \mathcal{G}_g} \theta_i$  being the best item of set  $\mathcal{G}_g$ ). Hence  $\mathbf{E}[q_{i_g}] = \sum_{\tau=1}^t \mathbf{E}[i_g^\tau] \geq \frac{mt}{k}$ . Now applying Chernoff-Hoeffdings bound for  $w_{i_g}$ , we get that for any  $\eta \in \left(\frac{3}{32}, 1\right]$ ,

$$\begin{aligned} \Pr\left(q_{i_g} \leq (1 - \eta)\mathbf{E}[q_{i_g}]\right) &\leq \exp\left(-\frac{\mathbf{E}[q_{i_g}]\eta^2}{2}\right) \leq \exp\left(-\frac{mt\eta^2}{2k}\right) \\ &= \exp\left(-\frac{2\eta^2}{\epsilon_\ell^2} \ln\left(\frac{2k}{\delta_\ell}\right)\right) = \exp\left(-\frac{(\sqrt{2}\eta)^2}{\epsilon_\ell^2} \ln\left(\frac{2k}{\delta_\ell}\right)\right) \\ &\leq \exp\left(-\ln\left(\frac{2k}{\delta_\ell}\right)\right) \leq \frac{\delta_\ell}{2k}, \end{aligned}$$

where the second last inequality holds as  $\eta \geq \frac{3}{32\sqrt{2}}$  and  $\epsilon_\ell \leq \frac{3}{32}$ , for any iteration  $\ell \in [\ln n]$ ; in other words for any  $\eta \geq \frac{3}{32\sqrt{2}}$ , we have  $\frac{\sqrt{2}\eta}{\epsilon_\ell} \geq 1$  which leads to the second last inequality. Thus we finally derive that with probability at least  $\left(1 - \frac{\delta_\ell}{2k}\right)$ , one can show that  $q_{i_g} > (1 - \eta)\mathbf{E}[q_{i_g}] \geq (1 - \eta) \frac{mt}{k}$ , and the proof follows henceforth.  $\square$

In particular, fixing  $\eta = \frac{1}{2}$  in Lemma 13, we get that with probability at least  $\left(1 - \frac{\delta_\ell}{2}\right)$ ,  $q_{i_g} > (1 - \frac{1}{2})\mathbf{E}[q_{i_g}] > \frac{mt}{2k}$ . Note that, for any round  $\tau \in [t]$ , whenever an item  $i \in \mathcal{G}_g$  appears in the top- $m$  set  $\mathcal{G}_{gm}^\tau$ , then the rank breaking update ensures that every element in the top- $m$  set gets compared with rest of the  $k - 1$  elements of  $\mathcal{G}_g$ . Based on this observation, we now prove that for any set  $\mathcal{G}_g$ , a near-best ( $\epsilon_\ell$ -optimal of  $i_g$ ) is retained as the winner  $c_g$  with probability at least  $\left(1 - \frac{\delta_\ell}{2}\right)$ . More formally:

**Lemma 14.** *Consider any particular set  $\mathcal{G}_g$  at any iteration  $\ell$ . Let  $i_g \leftarrow \arg \max_{i \in \mathcal{G}_g} \theta_i$ , then with probability at least  $\left(1 - \delta_\ell\right)$ ,  $\theta_{c_g} > \theta_{i_g} - \frac{\epsilon_\ell}{c}$ .*

*Proof.* With top- $m$  ranking feedback, the crucial observation lies in that at any round  $\tau \in [t]$ , whenever an item  $i \in \mathcal{G}_g$  appears in the top- $m$  ranking  $\sigma_\tau$ , then the rank breaking update ensures that every element in the top- $m$  set gets compared to each of the rest  $k - 1$  elements of  $\mathcal{G}_g$  - it defeats to every element preceding item in  $\sigma \in \Sigma_{\mathcal{G}_{gm}}$ , and wins over the rest. If  $n_{ij} = w_{ij} + w_{ji}$  denotes the number of times item  $i$  and  $j$  are compared after rank-breaking, for  $i, j \in \mathcal{G}_g$ ,  $n_{ij} = n_{ji}$ , and from Lemma 13 with  $\eta = \frac{1}{2}$  we have that  $n_{i_g j} \geq \frac{mt}{2k}$  with probability at least  $\left(1 - \delta_\ell/2k\right)$ . Given the above arguments in place, for any item  $j \in \mathcal{G}_g \setminus \{i_g\}$ , by Hoeffdings inequality:

$$\Pr\left(\left\{\hat{p}_{ji_g} - p_{ji_g|\mathcal{G}_g} > \frac{\epsilon_\ell}{2}\right\} \cap \left\{n_{ji_g} \geq \frac{mt}{2k}\right\}\right) \leq \exp\left(-2\frac{mt}{2k}(\epsilon_\ell/2)^2\right) \leq \frac{\delta_\ell}{2k},$$

Now consider any item  $j$  such that  $\theta_{i_g} - \theta_j > \epsilon_\ell/c$ , then we have  $\frac{Pr(i_g|\mathcal{G}_g)}{Pr(j|\mathcal{G}_g)} > 1 + 4\epsilon_\ell$ , which by Lem. 9 implies  $p_{i_g j|\mathcal{G}_g} > \frac{1}{2} + \epsilon_\ell$ , or equivalently  $p_{j i_g|\mathcal{G}_g} < \frac{1}{2} - \epsilon_\ell$ .

But since we show that for any item  $j \in \mathcal{G}_g \setminus \{1\}$ , with high probability  $(1 - \delta_\ell/2k)$ , we have  $\hat{p}_{j i_g} - p_{j i_g|\mathcal{G}_g} < \frac{\epsilon_\ell}{2}$ . Taking union bound above holds true for any  $j \in \mathcal{G}_g \setminus \{1\}$  with probability at least  $(1 - \delta/2)$ . Combining with the above claim of  $p_{j i_g|\mathcal{G}_g} < \frac{1}{2} - \epsilon_\ell$ , this further implies  $\hat{p}_{j i_g} + \frac{\epsilon_\ell}{2} < p_{j i_g|\mathcal{G}_g} + \epsilon_\ell < \frac{1}{2}$ . Thus no such  $\epsilon_\ell$  suboptimal item can be picked as  $c_g$  for any subgroup  $\mathcal{G}_g$ , at any phase  $\ell$ .

On the other hand, following the same chain of arguments note that  $\hat{p}_{i_g j} - p_{i_g j|\mathcal{G}_g} > -\frac{\epsilon_\ell}{2} \implies \hat{p}_{i_g j} + \frac{\epsilon_\ell}{2} > p_{i_g j|\mathcal{G}_g} > \frac{1}{2}$  for all  $j \in \mathcal{G}_g$ ,  $i_g$  is a valid candidate for  $c_g$  always, or in other case some other  $\epsilon_\ell$ -suboptimal item  $j$  (such  $\theta_j > \theta_{i_g} - \epsilon_\ell$ ) can be chosen as  $c_g$ . This concludes the proof.  $\square$

The correctness-claim now follows using a similar argument as given for the proof of Thm. 4. We add the details below for the sake of completeness: Without loss of generality, we assume the best item of the RUM( $k, \theta$ ) model is  $\theta_1$ , i.e.  $\theta_1 > \theta_i \forall i \in [n] \setminus \{1\}$ . Now for any iteration  $\ell$ , let us define  $g_\ell \in [G]$  to be the index of the set that contains *best item* of the entire set  $S_\ell$ , i.e.  $\arg \max_{i \in S_\ell} \theta_i \in \mathcal{G}_{g_\ell}$ . Then applying Lemma 14, with probability at least  $(1 - \delta_\ell)$ ,  $\theta_{c_{g_\ell}} > \theta_{i_{g_\ell}} - \epsilon_\ell/c$ . Note that initially, at phase  $\ell = 1$ ,  $i_{g_\ell} = 1$ . Then, for each iteration  $\ell$ , applying Lemma 14 recursively to  $\mathcal{G}_{g_\ell}$ , we finally get  $\theta_r > \theta_1 - \left(\frac{\epsilon}{8} + \frac{\epsilon}{8} \left(\frac{3}{4}\right) + \dots + \frac{\epsilon}{8} \left(\frac{3}{4}\right)^{\lfloor \ln_k n \rfloor}\right) - \frac{\epsilon}{2} \geq \theta_1 - \frac{\epsilon}{8} \left(\sum_{i=0}^{\infty} \left(\frac{3}{4}\right)^i\right) - \frac{\epsilon}{2} \geq \theta_1 - \epsilon$ . Thus assuming the algorithm does not fail in any of the iteration  $\ell$ , we finally have that  $p_{r^*1} > \frac{1}{2} - \epsilon$ —this shows that the final item output by Seq-PB is  $\epsilon$  optimal.

Finally note that since at each iteration  $\ell$ , the algorithm fails with probability at most  $\delta_\ell(1/2 + \frac{1}{2k}) \leq \delta_\ell$ , the total failure probability of the algorithm is at most  $\left(\frac{\delta}{4} + \frac{\delta}{8} + \dots + \frac{\delta}{2^{\lceil \frac{n}{k} \rceil}}\right) + \frac{\delta}{2} \leq \delta$ . This shows the correctness of the algorithm, concluding the proof.  $\square$

### C.3 Proof of Theorem 8

The proof proceeds almost same as the proof of Thm. 6, the only difference lies in the analysis of the KL-divergence terms with top- $m$  ranking feedback.

Consider the exact same set of RUM( $k, \theta$ ) instances,  $\{\nu^a\}_{a=1}^n$  we constructed for Thm. 6. It is now interesting to note that how the top- $m$  ranking feedback affects the KL-divergence analysis, precisely the KL-divergence shoots up by a factor of  $m$  which in fact triggers an  $\frac{1}{m}$  reduction in regret learning rate. We show this below formally.

Note that for top- $m$  ranking feedback for any problem instance  $\nu^a$ ,  $a \in [n]$ , each  $k$ -set  $S \subseteq [n]$  is associated to  $\binom{k}{m}(m!)$  number of possible outcomes, each representing one possible ranking of set of  $m$  items of  $S$ , say  $S_m$ . Also the probability of any permutation  $\sigma \in \Sigma_S^m$  is given by  $p_S^a(\sigma) = Pr_{\nu^a}(\sigma|S)$ , where  $Pr_{\nu^a}(\sigma|S)$  is as defined for top- $m$  ranking feedback for RUM( $k, \theta$ ) problem instance  $\nu^a$  (see Sec. 6). More formally, for problem **Instance-a**, we have that:

$$\begin{aligned} p_S^a(\sigma) &= Pr_{\nu^a}(\sigma = \sigma|S) = \prod_{i=1}^m Pr(X_{\sigma(i)} > X_{\sigma(j)}, \forall j \in \{i+1, \dots, m\}), \forall \sigma \in \Sigma_S^m \\ &= Pr_{\nu^a}(\sigma = \sigma|S) = \prod_{i=1}^m Pr(\zeta_{\sigma(i)} > \zeta_{\sigma(j)} - (\theta_{\sigma(i)} - \theta_{\sigma(j)}^a), \forall j \in \{i+1, \dots, m\}), \forall \sigma \in \Sigma_S^m \end{aligned}$$

As also argued in the proof of Thm. 6, note that for any top- $m$  ranking of  $\sigma \in \Sigma_S^m$ ,  $KL(p_S^1(\sigma), p_S^a(\sigma)) = 0$  for any set  $S \not\ni a$ . Hence while comparing the KL-divergence of instances  $\nu^1$  vs  $\nu^a$ , we need to focus only on sets containing  $a$  (recall we denote this as  $S^a$ ). Applying the chain rule for KL-divergence, we now get:

$$\begin{aligned} KL(p_S^1, p_S^a) &= KL(p_S^1(\sigma_1), p_S^a(\sigma_1)) + KL(p_S^1(\sigma_2 | \sigma_1), p_S^a(\sigma_2 | \sigma_1)) + \dots \\ &\quad + KL(p_S^1(\sigma_m | \sigma(1 : m-1)), p_S^a(\sigma_m | \sigma(1 : m-1))), \end{aligned} \tag{10}$$

where we abbreviate  $\sigma(i)$  as  $\sigma_i$  and  $KL(P(Y | X), Q(Y | X)) := \sum_x Pr(X = x) [KL(P(Y | X = x), Q(Y | X = x))]$  denotes the conditional KL-divergence. Moreover it is easy to note that for any  $\sigma \in \Sigma_S^m$  such that  $\sigma(i) = a$ , we have  $KL(p_S^1(\sigma_{i+1} | \sigma(1:i)), p_S^a(\sigma_{i+1} | \sigma(1:i))) = 0$ , for all  $i \in [m]$ .

Now using the KL divergence upper bounds, as derived in the proof of Thm. 6, we have than

$$KL(p_S^1(\sigma_1), p_S^a(\sigma_1)) \leq \frac{\Delta_a'^2}{\frac{8\epsilon^2}{k}}.$$

One can potentially use the same line of argument to upper bound the remaining KL divergence terms of (10) as well. More formally note that for all  $i \in [m-1]$ , we can show that:

$$\begin{aligned} & KL(p_S^1(\sigma_{i+1} | \sigma(1:i)), p_S^a(\sigma_{i+1} | \sigma(1:i))) \\ &= \sum_{\sigma' \in \Sigma_S^i} Pr(\sigma') KL(p_S^1(\sigma_{i+1} | \sigma(1:i)) = \sigma', p_S^a(\sigma_{i+1} | \sigma(1:i)) = \sigma') \leq \frac{8\epsilon^2}{k} \end{aligned}$$

Thus applying above in (10) we get:

$$KL(p_S^1, p_S^a) = KL(p_S^1(\sigma_1) + \dots + KL(p_S^1(\sigma_m | \sigma(1:m-1)), p_S^a(\sigma_m | \sigma(1:m-1)))) \leq \frac{8m\epsilon^2}{k}. \quad (11)$$

Eqn. (11) precisely gives the main result to derive Thm. 8. Note that it shows an  $m$ -factor blow up in the KL-divergence terms owing to top- $m$  ranking feedback. The rest of the proof can be derived by following exactly the same argument used in 6, which yields to the desired sample complexity lower bound.