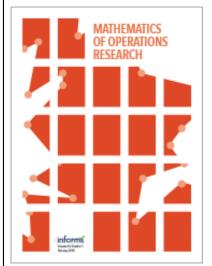
This article was downloaded by: [128.122.186.54] On: 20 October 2022, At: 08:24 Publisher: Institute for Operations Research and the Management Sciences (INFORMS) INFORMS is located in Maryland, USA



## **Mathematics of Operations Research**

Publication details, including instructions for authors and subscription information: <a href="http://pubsonline.informs.org">http://pubsonline.informs.org</a>

# Asymptotically Optimal Sequential Design for Rank Aggregation

Xi Chen, Yunxiao Chen, Xiaoou Li

#### To cite this article:

Xi Chen, Yunxiao Chen, Xiaoou Li (2022) Asymptotically Optimal Sequential Design for Rank Aggregation. Mathematics of Operations Research 47(3):2310-2332. <a href="https://doi.org/10.1287/moor.2021.1209">https://doi.org/10.1287/moor.2021.1209</a>

Full terms and conditions of use: <a href="https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions">https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions</a>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact permissions@informs.org.

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2022, INFORMS

Please scroll down for article-it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit <a href="http://www.informs.org">http://www.informs.org</a>



Vol. 47, No. 3, August 2022, pp. 2310-2332 ISSN 0364-765X (print), ISSN 1526-5471 (online)

## **Asymptotically Optimal Sequential Design for Rank Aggregation**

Xi Chen,<sup>a</sup> Yunxiao Chen,<sup>b</sup> Xiaoou Li<sup>c</sup>

<sup>a</sup> Stern School of Business, New York University, New York, New York 10013; <sup>b</sup> Department of Statistics, London School of Economics and Political Science, London WC2A 2AE, United Kingdom; School of Statistics, University of Minnesota, Minneapolis, Minnesota 55455 Contact: xc13@stern.nyu.edu, (b) https://orcid.org/0000-0002-9049-9452 (XC); y.chen186@lse.ac.uk (YC); lixx1766@umn.edu (XL)

Received: July 16, 2018

Revised: May 8, 2020; February 15, 2021

Accepted: June 11, 2021

Published Online in Articles in Advance:

January 5, 2022

MSC2020 Subject Classification: Primary:

62L10; secondary: 62L15

https://doi.org/10.1287/moor.2021.1209

Copyright: © 2022 INFORMS

Abstract. A sequential design problem for rank aggregation is commonly encountered in psychology, politics, marketing, sports, etc. In this problem, a decision maker is responsible for ranking K items by sequentially collecting noisy pairwise comparisons from judges. The decision maker needs to choose a pair of items for comparison in each step, decide when to stop data collection, and make a final decision after stopping based on a sequential flow of information. Because of the complex ranking structure, existing sequential analysis methods are not suitable. In this paper, we formulate the problem under a Bayesian decision framework and propose sequential procedures that are asymptotically optimal. These procedures achieve asymptotic optimality by seeking a balance between exploration (i.e., finding the most indistinguishable pair of items) and exploitation (i.e., comparing the most indistinguishable pair based on the current information). New analytical tools are developed for proving the asymptotic results, combining advanced change of measure techniques for handling the level crossing of likelihood ratios and classic large deviation results for martingales, which are of separate theoretical interest in solving complex sequential design problems. A mirror-descent algorithm is developed for the computation of the proposed sequential procedures.

Funding: X. Chen acknowledges support from the National Science Foundation (NSF) [Grant IIS-1845444]. Y. Chen acknowledges support from the National Academy of Education/Spencer Postdoctoral Fellowship. X. Li acknowledges support from the NSF [Grant DMS-1712657].

Supplemental Material: The online supplement is available at https://doi.org/10.1287/moor. 2021.1209.

Keywords: active sequential tests • asymptotically optimal policy • sequential analysis • rank aggregation

## 1. Introduction

This paper considers a sequential design problem for rank aggregation. In this problem, a decision maker is responsible for ranking K items by adaptively collecting the noisy outcome of pairwise comparison from judges. Sequential rank aggregation has a wide range of applications, including social choice (Saaty and Vargas [48]), sports (Elo [23]), search rankings (Page et al. [47]), etc. Pairwise comparison is the most popular approach for rank aggregation as sufficient evidence from cognitive psychology suggests that people make more accurate judgment when making pairwise comparisons (i.e., given a pair of items and asked to indicate which item is preferred to the other) as compared with multiwise comparison (Blumenthal [10]) and some applications, such as chess gaming, have a natural form of pairwise comparison.

In a rank aggregation problem, more comparisons usually lead to a more accurate global ranking. However, each comparison comes with some cost, for example, in crowdsourcing applications, a requester has to pay crowd workers a fixed amount of monetary reward for each labeled pair. Therefore, to design a cost-efficient ranking procedure, a decision maker faces the following three key challenges:

- How to adaptively decide the next pair of objects for comparison based on the collected information. The adaptive selection of pairs is important for saving the cost. For example, if we are confident that object 1 is ranked higher than 2 and object 2 is preferred over 3, there is no need to compare objects 1 and 3.
  - When to stop asking for more comparisons.
- 3. When stopping the comparison process, how to aggregate the pairwise comparisons to infer the global ranking.

Because of the wide applications of rank aggregation, there are several recent machine learning works devoted to the development of ranking algorithms with rigorous theoretical guarantees. For example, Negahban et al. [45], Hajek et al. [25], and Shah et al. [50] propose algorithms and establish the estimation error rates under the Bradley-Terry-Luce (BTL) model (Bradley and Terry [11], Luce [40]), the Thurstone [54] model, and a more general strong stochastic transitivity model (Ballinger and Wilcox [4], Morrison [43]). However, these works mainly focus on a static setting with either given or randomly drawn pairs. In contrast, under a sequential setting, we are interested in designing an adaptive pair-selection rule. Moreover, for recent active ranking works (e.g., Heckel et al. [26]), optimal stopping is usually not considered. For example, the commonly studied probably approximately correct sample complexity bound from the machine learning literature usually involves some large universal constants and cannot be directly used for an accurate stopping rule. Determining a right stopping time is critical for balancing accuracy and cost in many applications (e.g., ranking via crowdsourcing). Therefore, to address the challenge of optimal stopping, we adopt the sequential analysis framework from statistics that directly optimizes over the random stopping time. On the other hand, because of the complex structure of ranking aggregation, this problem cannot be formulated and solved by existing sequential adaptive design methods (Chernoff [20], Naghshvar and Javidi [44]).

Under a wide class of parametric comparison models (e.g., the BTL model; Bradley and Terry [11], Luce [40]), we develop new sequential analysis methods to conduct sequential experiments for pairwise comparisons and to balance the ranking accuracy and cost. We first formulate the problem under a general *Bayesian decision framework*. In particular, each item k is represented by a parameter  $\theta_k$ , which determines its underlying true rank among K items. For example, the parameter  $\theta_k$  can be viewed as the quality score for item k, and item i has a higher rank than item j if and only if  $\theta_i > \theta_j$ . The pairwise comparison of items i and j follows a probabilistic comparison model (e.g., Bradley and Terry [11], Luce [40], Thurstone [54]) parameterized by  $\theta_i$  and  $\theta_j$ . Under the Bayesian framework, the parameter vector for all product  $\theta$  is drawn from some prior distribution. A sequential procedure chooses a pair (i,j) for the next comparison in each stage and decides the stopping time T. Upon stopping, the final decision is to choose the global rank  $R := (R_1, \ldots, R_K)$  from the set of all permutations of  $\{1,2,\ldots,K\}$ . To measure the accuracy of a rank R, we adopt the widely used Kendall tau distance (Kendall and Gibbons [32]), which measures the number of inconsistent pairs between the decision R and the underlying true rank induced by the scores  $(\theta_1,\ldots,\theta_K)$ . Then, the loss function of this sequential design problem is defined by combining the cost of data collection and the Kendall tau distance:

$$\sum_{i < j} \{ I(\theta_i > \theta_j) I(R_i > R_j) + I(\theta_i < \theta_j) I(R_i < R_j) \} + cT, \tag{1}$$

where the constant c > 0 indicates the relative cost of each comparison and  $I(\cdot)$  denotes an indicator function. The goal is to optimize the expected loss in (1) over the pair-selection rule, stopping rule T, and final decision R (see Section 2 for more details). To justify the performance of the proposed policies, we adopt the notion of "asymptotic optimality" from Chernoff [20] (see Equation (7)) that is widely used in sequential analysis (Lai [35], Schwarz [49], Siegmund [51], Tartakovsky et al. [53]). Although finding an exact optimal policy is computationally intractable, we prove that the proposed policies are asymptotically optimal.

It is also worthwhile to note that, although, according to the final decision, our problem seems to be a multihypothesis sequential testing problem with adaptive experiment selection as considered in Naghshvar and Javidi [44], there exist fundamental differences. First, Naghshvar and Javidi [44] only consider simple hypotheses, and the ranking problem, when viewed as a multihypothesis testing problem, consists of composite hypotheses. Second, typically 0-1 loss is considered for measuring the decision accuracy in multihypothesis testing, and our problem has a more complex loss function based on the Kendall tau distance that is tailored to rank aggregation. Our problem is also a substantial generalization of a classic sequential test of two composite hypotheses (Kiefer and Sacks [33], Lai [34], Schwarz [49]). In particular, when the number of items is two (K = 2), our problem degenerates to testing two composite hypotheses without adaptive experiment selection.

## 1.1. Main Contribution

We summarize the main methodological and theoretical contributions of the paper as follows:

- Under a Bayesian decision framework and a large class of parametric pairwise comparison models, we derive an asymptotic lower bound (Theorem 1) for the Bayes risk of all possible sequential ranking policies. Note that the Bayes risk of the sequential rank aggregation problem, which combines the expected Kendall tau distance and the expected sample size, is more complex than that of the traditional sequential hypothesis testing problems (e.g., Chernoff [20], Kiefer and Sacks [33], Naghshvar and Javidi [44]).
- We propose two sequential ranking policies. In particular, we provide two choices of stopping rule and a class of randomized pair-selection rules. We quantify the expected Kendall tau and the sample size of the proposed methods (Theorems 2 and 3) and show that the Bayes risks match the asymptotic lower bound, which further implies that the proposed methods are asymptotically optimal (Corollary 1). Our randomized pair selection rule utilizes an epsilon-greedy strategy to balance the exploration (i.e., randomly selecting pairs to gain information about

the underlying parameters  $\{\theta_k\}_{k=1}^K$ ) and exploitation (i.e., choosing the best pair for comparison based on the current information). The exploration is critical for learning the rank, and the exploitation is critical for saving the sample size for comparison.

- For the exploration, we quantify the impact of the exploration rate on the estimation of model parameters and provide an exponential probability bound as an auxiliary result (Lemma 1).
- For the exploitation, we consider a randomized adaptive selection rule (see Section 3). Specifically, in each step, the probability of selecting each pair is obtained by solving a saddle point optimization problem. We further develop a mirror descent algorithm for solving the optimization (see Section 3.4).
- Technically, we develop new analytical tools for quantifying the level-crossing probability of a random function (e.g., likelihood function, martingale, or submartingale) double-indexed by model parameters and the sample size. As such a probability tends to zero, the problem falls into the rare-event analysis domain, in which an exact exponential decay rate is challenging to obtain. Traditional methods, such as the ones adopted in Naghshvar and Javidi [44] and Chernoff [20], are based on exponential change of measure of the log-likelihood ratio statistics and are not directly applicable to the ranking problem considered here. The method we use in the proof combines a mixture type of change of measure method recently proposed in Adler et al. [1], Li and Liu [37], and Li et al. [38] and large deviation results for martingales.

## 1.2. Related Works

Sequential hypothesis testing, initiated by the seminal works of Wald [57] and Wald and Wolfowitz [58], is an important area of statistics for processing data taken in a sequential experiment, in which the total number of observations is not fixed in advance. A sequential test is characterized by two components: (1) a stopping rule that decides when to stop the data-collection process and (2) a decision rule on choosing the hypothesis upon stopping. A large body of literature on sequential tests with two hypotheses has been developed, a partial list of which includes Schwarz [49], Hoeffding [27], and Lai [34]. Sequential testing with more than two hypotheses and sequential multiple testing have been extensively studied in recent decades (see, e.g., Dragalin et al. [21], Draglia et al. [22], Mei [42], Song and Fellouris [52], Xie and Siegmund [61]). For a comprehensive review on sequential analysis, we refer the readers to the surveys and books Siegmund [51], Lai [35], Hsiung et al. [29], Tartakovsky et al. [53], and references therein. In addition to optimizing over the stopping rule and final decision, Chernoff [20] first introduces the adaptive design into the sequential testing framework, followed by a large body of literature; see, for example, Albert [2], Kiefer and Sacks [33], Tsitovich [56], Naghshvar and Javidi [44], and Nitinawarat and Veeravalli [46]. Sequential analysis finds many applications in different disciplines, including clinical trials, educational testing, and industrial quality control (see, e.g., Bartroff and Lai [5], Bartroff et al. [6,7], Lai and Shih [36], Wang et al. [59], Ye et al. [62]).

Rank aggregation has been an active research problem in recent years (see, e.g., Chen et al. [16, 18], Chen and Suh [19], Garg and Johari [24], Hajek et al. [25], Kallus and Udell [31], Negahban et al. [45], Shah et al. [50], and references therein), and it finds many applications to social choice, tournament play, search rankings, advertisement placement, etc. With the advent of crowdsourcing services, one can easily ask crowd workers to conduct comparisons among a few objects in an online fashion at a low cost (Chen et al. [15, 17]). Therefore, active noisy sorting and ranking problems have received a lot of attention in recent years. For example, Braverman and Mossel [12], Braverman et al. [13], and Mao et al. [41] study the active sorting problem in which each query of (i, j) reveals the true ranking between i and j with a fixed probability  $1/2 + \gamma$  for some  $\gamma > 0$  regardless of the distance between i and j. In contrast, our model associates each item i with a preference score (aka utility)  $\theta_i$ . The comparison result between i and j would be based on the values of  $\theta_i$  and  $\theta_j$  according to some probabilistic model (e.g., see Equation (2)). Jamieson and Nowak [30] study the ranking problem with feature information for each item. Heckel et al. [26] investigates the active top-K ranking under a general class of nonparametric models and also establish a lower bound on the number of comparisons for parametric models. However, as we mention, although rank aggregation is extensively studied in the machine learning community, it has not been investigated under the sequential analysis framework, which incorporates the random stopping rule as a decision variable. The techniques developed in this work enable a sequential rank procedure with optimal stopping and adaptive design.

## 1.3. Paper Organization

The rest of the paper is organized as follows. In Section 2, we introduce the setup of the problem. Section 3 presents the proposed policies and the theoretical results and provides further discussions on the proof sketch and model misspecification. The concluding remarks are provided in Section 5. Technical proofs for the theorems are provided in Section 6. Proofs for all the lemmas are provided in the online supplement.

## 2. Problem Setup

We first introduce the comparison model and formulate the sequential ranking problem. Consider the task of inferring a global ranking over K items. Let  $A = \{(i,j): i,j \in \{1,\ldots,K\}, i < j\}$  be the set of pairs for comparison. At each time n ( $n = 1, 2, \ldots$ ), a pair  $a_n := (a_{n,1}, a_{n,2}) \in A$  is selected for comparison. For example,  $a_2 = (1,2)$  means that items 1 and 2 are compared at time 2. The comparison outcome is denoted by a random variable  $X_n \in \{0,1\}$ , where  $X_n = 1$  means item  $a_{n,1}$  is preferred to item  $a_{n,2}$ , and  $X_n = 0$  otherwise. The comparison outcome  $X_n$  is assumed to follow a ranking model, such as the widely used BTL (Bradley and Terry [11], Luce [40]) and Thurstone [54] models. Such a ranking model assumes that each item is associated with an unknown latent score  $\theta_i \in \mathbb{R}$  for  $i = 1, \ldots, K$ , where the global rank of the K items is given by the rank of  $\theta_1, \ldots, \theta_K$ . The distribution of  $X_n$  is determined by  $\theta_i$  and  $\theta_j$  when comparing pair (i,j). For example, given pair  $a_n := (a_{n,1}, a_{n,2})$ , the BTL model assumes that

$$\mathbb{P}(X_n = 1) = \frac{\exp(\theta_{a_{n,1}})}{\exp(\theta_{a_{n,1}}) + \exp(\theta_{a_{n,2}})};$$

$$\mathbb{P}(X_n = 0) = \frac{\exp(\theta_{a_{n,2}})}{\exp(\theta_{a_{n,1}}) + \exp(\theta_{a_{n,2}})}.$$
(2)

Under this model,  $\theta_{a_{n,1}} > \theta_{a_{n,2}}$  means that item  $a_{n,1}$  is preferred to item  $a_{n,2}$ , reflected by  $\mathbb{P}(X_n = 1) > 0.5$ . A common feature for many comparison models is that the distribution of the comparison between items i and j only depends on the pairwise differences  $\theta_i - \theta_j$ . Consequently, such models are not identifiable up to a location shift. To overcome this issue, we fix  $\theta_1 = 0$  and treat  $\theta = (\theta_2, \dots, \theta_K)$  as the unknown model parameters. The result of this paper applies to a wide class of comparison models, and thus, we denote the probability mass function of the comparison outcome x given pair a as  $f_{\theta}^a(x)$ . We point out that, although we focus on the case in which the distribution of the pairwise comparison only depends on  $\theta_{a_n,1} - \theta_{a_n,2}$ , our methods and results can be extended to more general cases without this requirement.

We now describe components in a sequential design for rank aggregation: an adaptive selection rule A, a stopping time T, and a decision rule R on the global rank. For the adaptive selection rule A, we consider the class of randomized adaptive selection rules, which contains deterministic selection rules as special cases. In particular, let  $A = \{\lambda_n : n = 1, 2, \ldots\}$ , where  $\lambda_n = (\lambda_n^{i,j})_{(i,j) \in \mathcal{A}} \in \Delta$  denotes the probability of selecting the pair (i, j). Here,  $\Delta = \{(\lambda^{i,j}) : \sum_{(i,j) \in \mathcal{A}} \lambda^{i,j} = 1, \ \lambda^{i,j} \ge 0\}$  is a probability simplex over K(K-1)/2 pairs. At each time n, a pair  $a_n$  is selected according to the categorical distribution with the parameter  $\lambda_n$ , where  $\lambda_n$  adapts to the filtration sigma algebra generated by the selected pairs and the observed outcomes, that is,  $\mathcal{F}_n = \sigma(X_1, \ldots, X_{n-1}, a_1, \ldots, a_{n-1})$ . The adaptive comparison process stops at time T, a stopping time with respect to the filtration  $\{\mathcal{F}_n\}_{n \ge 0}$ . It is worthwhile to note that the random stopping time T is also the number of samples being collected. Upon stopping, one needs to make a decision  $R := (R_1, \ldots, R_K)$ , the global ranking of the K items. For example, when K = 3, R = (3, 1, 2) means that one decides  $\theta_2 > \theta_3 > \theta_1$ . We further denote  $P_K$  as the set of permutations over  $\{1, \ldots, K\}$ , and thus,  $R \in P_K$ . The adaptive selection rule  $A = \{\lambda_n : n = 1, 2, \ldots\}$ , the stopping time T, and the decision R together form a sequential ranking policy, denoted by  $\pi = (A, T, R)$ .

The performance of a sequential ranking policy is measured via its ranking accuracy and the expected stopping time. Specifically, we measure the ranking accuracy by the Kendall tau distance (Kendall and Gibbons [32]), which is one of the most widely used measures for ranking consistency. More precisely, for each  $R = (R_1, ..., R_K) \in P_K$ , we convert it to the binary decisions over pairs  $\{R_{i,j} \in \{0,1\} : i,j \in \{1,...,K\}, i < j\}$ , where  $R_{i,j} = I(R_i < R_j)$ , and  $R_{i,j} = 1$  means that item i is preferred to item j. For example, if R = (3,1,2), we have  $R_{1,2} = 0$  and  $R_{2,3} = 1$ . The Kendall tau distance between R and the true ranking induced by  $(\theta_1, ..., \theta_K)$  is defined by

$$L_K(\{R_{i,j}\}) = \sum_{i < j} \{ I(\theta_i > \theta_j)(1 - R_{i,j}) + I(\theta_i < \theta_j)R_{i,j} \}.$$
(3)

On the other hand, the loss function associated with the random sample size *T* is defined as

$$L_c(T) = c \times T,\tag{4}$$

where the constant c > 0 indicates the *relative* cost of conducting one more pairwise comparison. The choice of c depends on the nature of the ranking problem. Generally, if obtaining each sample is expensive compared with the cost because of the inaccuracy of the ranking, then a large c is chosen and vice versa. Note that c is not a tuning parameter over which to optimize.

We define the risk associated with a sequential ranking policy under the Bayesian decision framework, in which the model parameter  $\theta$  is assumed to be random and follows a prior distribution. To avoid confusion, we write  $\Theta$  when  $\theta$  is viewed as random and denote by  $\rho(\theta)$  the prior density function of  $\Theta = (\Theta_2, ..., \Theta_K)$ . Recall that we have fixed  $\Theta_1 = 0$  to ensure identifiability. The Bayes risk combines the risks associated with the Kendall tau distance of the decision and the sampling cost

$$V_{c}(\rho, \pi) = \mathbb{E}^{\pi} \Big( L_{K}(\{R_{i,j}\}) + L_{c}(T) \Big)$$

$$= \mathbb{E}^{\pi} \left\{ \sum_{i < j} I(\Theta_{i} > \Theta_{j}) (1 - R_{i,j}) + I(\Theta_{i} < \Theta_{j}) R_{i,j} \right\} + c \mathbb{E}^{\pi} T,$$
(5)

where the expectation  $\mathbb{E}^{\pi}$  is taken under the policy  $\pi$  with respect to the randomness of the selected pairs, the observed comparison results, and the stopping time as well as the prior distribution  $\rho$ . Of particular interest is the minimum risk under the optimal sequential ranking policy given the prior distribution of  $\Theta$  and sampling cost c:

$$V_c^*(\rho) = \inf_{\pi} V_c(\rho, \pi). \tag{6}$$

For any given cost c, obtaining an analytical form of an optimal policy that achieves  $V^*(\rho, c)$  is typically infeasible. Following the literature of sequential analysis, a policy is usually evaluated by the notion of *asymptotic optimality* (Chernoff [20]). In particular, a policy  $\pi$  is said to be asymptotically optimal if

$$\lim_{c \to 0} \frac{V_c(\rho, \pi)}{V_c^*(\rho)} = 1,\tag{7}$$

that is, the Bayes risk of the policy matches the minimal Bayes risk asymptotically when the relative sampling cost converges to zero. It is worthwhile to note that, in the construction of our policy, we certainly allow the cost c to be nonzero. The notion of asymptotic optimality in (7) has been widely adopted in the sequential analysis literature as an optimality criterion (see, e.g., Chernoff [20], Kiefer and Sacks [33], Naghshvar and Javidi [44], Schwarz [49]). The limiting process  $c \to 0$  should be interpreted as the sample size n goes to infinity, which is a very common limiting process in statistical asymptotic theory. In asymptotic theory, letting n grow to infinity is only for the theoretical study of the properties of an estimator although, in practice, no data set has an infinite number of observations.

## 3. Sequential Policies and Asymptotic Optimality

In Section 3.1, we propose two sequential ranking policies:  $\pi_1$  and  $\pi_2$ . The asymptotic optimality of the two policies is presented in Section 3.2. Then, we provide the proof sketch in Section 3.3, the optimization algorithm for efficient computation in Section 3.4, and the discussions on model misspecification in Section 3.5.

## 3.1. Two Sequential Policies

We first introduce some notations. Let W be the support of the prior probability density function  $\rho$ , that is,  $W = \{\theta : \rho(\theta) > 0\}$ , where  $\overline{E}$  denotes the closure of a set E. We further define the set  $W_{i,j} = \{\theta : \theta_i \ge \theta_j\} \cap W$  for all  $i, j \in \{1, ..., K\}$ . It is worthwhile to note that  $W_{i,j}$  and  $W_{j,i}$  are different sets, and their union is the set W. Given a sequence of selected pairs  $a_1, ..., a_n$  and observed comparisons  $X_1, ..., X_n$ , the log-likelihood function is defined as

$$l_n(\boldsymbol{\theta}) = \sum_{i=1}^n \log f_{\boldsymbol{\theta}}^{a_i}(X_i),$$

and the corresponding maximum likelihood estimator  $\widehat{\boldsymbol{\theta}}^{(n)} = (\widehat{\boldsymbol{\theta}}_2^{(n)}, \dots, \widehat{\boldsymbol{\theta}}_K^{(n)})$  is

$$\widehat{\boldsymbol{\theta}}^{(n)} = \arg \sup_{\boldsymbol{\theta} \in W} l_n(\boldsymbol{\theta}). \tag{8}$$

In what follows, we present our proposed sequential policies in terms of the proposed stopping time T, selection rule A, and ranking decision R.

**3.1.1. Stopping Times.** We then introduce two stopping times based on the generalized likelihood ratio statistic:

$$T_1 = \inf \left\{ n > 1 : \sum_{(i,j) \in \mathcal{A}} \exp \left\{ - \left| \sup_{\theta \in W_{i,j}} l_n(\theta) - \sup_{\theta \in W_{j,i}} l_n(\theta) \right| \right\} \le e^{-h(c)} \right\}, \tag{9}$$

and

$$T_2 = \inf \left\{ n > 1 : \min_{(i,j) \in \mathcal{A}} \left| \sup_{\theta \in W_{i,j}} l_n(\theta) - \sup_{\theta \in W_{j,i}} l_n(\theta) \right| \ge h(c) \right\}, \tag{10}$$

where  $h(c) = |\log c|(1 + |\log c|^{-\alpha})$  for some constant  $\alpha \in (0,1)$  and c is the relative cost introduced in (4). We note that  $T_2$  is obtained by replacing the summation in  $T_1$  by maximization and taking log and minus on both sides. Intuitively, the stopping rule  $T_2$  stops when the likelihood can tell whether  $\theta_i \ge \theta_i$  or vice versa for each pair (i, j).

**3.1.2. Ranking Decision.** Upon stopping, the decision about the global rank is made according to the rank of maximum likelihood estimation (MLE) at the stopping time T ( $T = T_1$  or  $T_2$ ). That is,

$$R = r(\widehat{\boldsymbol{\theta}}^{(T)}),\tag{11}$$

where the function  $r(\theta)$ :  $\mathbb{R}^{K-1} \to P_K$  gives the rank of  $(0, \theta_2, \dots, \theta_K)$ . More precisely,  $r(\theta) = (r_1, \dots, r_K) \in P_K$ , satisfying  $\theta_{r_1} \ge \theta_{r_2} \ge \dots \ge \theta_{r_K}$ , where  $\theta_1 = 0$ .

**3.1.3. Randomized Selection Rule.** We proceed to the randomized selection rule A, which is obtained by solving an optimization program. For a given  $\theta$ , we define function  $D(\theta)$ ,

$$D(\boldsymbol{\theta}) = \max_{\lambda \in \Delta} \min_{\widetilde{\boldsymbol{\theta}} \in W: r(\widetilde{\boldsymbol{\theta}}) \neq r(\boldsymbol{\theta})} \sum_{(i,j)} \lambda^{i,j} D^{i,j} \Big( \boldsymbol{\theta} || \widetilde{\boldsymbol{\theta}} \Big), \tag{12}$$

where  $D^{i,j}(\boldsymbol{\theta} || \widetilde{\boldsymbol{\theta}})$  is the Kullback–Leibler (KL) divergence from  $f_{\widetilde{\boldsymbol{\theta}}}^{i,j}(\cdot)$  to  $f_{\boldsymbol{\theta}}^{i,j}(\cdot)$ , that is,

$$D^{i,j}(\boldsymbol{\theta}||\widetilde{\boldsymbol{\theta}}) := \sum_{x \in \{0,1\}} f_{\boldsymbol{\theta}}^{i,j}(x) \log \frac{f_{\boldsymbol{\theta}}^{i,j}(x)}{f_{\widetilde{\boldsymbol{\theta}}}^{i,j}(x)},$$

and  $f_{\theta}^{i,j}(x)$  denotes the probability mass function when the pair (i,j) is selected. We further define

$$\lambda^{*}(\boldsymbol{\theta}) = \arg\max_{\boldsymbol{\lambda} \in \Delta} \min_{\widetilde{\boldsymbol{\theta}} \in W: r(\widetilde{\boldsymbol{\theta}}) \neq r(\boldsymbol{\theta})} \sum_{(i,j)} \lambda^{i,j} D^{i,j}(\boldsymbol{\theta} || \widetilde{\boldsymbol{\theta}}), \tag{13}$$

and

$$\widehat{\lambda}_n = \left(\widehat{\lambda}_n^{i,j}\right) = \lambda^*(\widehat{\boldsymbol{\theta}}^{(n-1)}). \tag{14}$$

That is,  $\lambda^*(\theta)$  is the solution to the optimization problem (12), and  $\widehat{\lambda}_n$  is the solution to the optimization problem given the MLE based on the previous n-1 observations. The objective function in (12) is a weighted KL divergence for all pairs with the weights  $\lambda^{i,j}$ . The inner minimization problem is taken over all the parameter vector  $\widetilde{\theta} \in W$ , for which the induced rank  $r(\widetilde{\theta})$  is different from that of  $\theta$ . At each time n, given the MLE  $\widehat{\theta}^{(n-1)}$ , we compute  $\widehat{\lambda}_n$ , which is the maximizer of  $\lambda \in \Delta$  in  $D(\widehat{\theta}^{(n-1)})$ . We elaborate on the intuition behind the optimization in (12). First, for each  $\theta$ ,  $\sum_{(i,j)} \lambda^{i,j} D^{i,j}(\theta||\widetilde{\theta})$  gives the drift of the log-likelihood ratio statistics between  $f_{\theta}$  and  $f_{\overline{\theta}}$  under the model  $f_{\theta}$  and a randomized sampling scheme specified by  $\lambda$ , which is also the mutual information between  $f_{\theta}$  and  $f_{\overline{\theta}}$  when the pair is selected according to  $\lambda$ . Minimizing the inner part with respect to  $\widetilde{\theta}$  over the set  $\{\widetilde{\theta} \in W : r(\widetilde{\theta}) \neq r(\theta)\}$  provides a measure on the distinguishability of the rank of  $\theta$  under the sampling scheme  $\lambda$ . Second, if the true model parameter is  $\theta$ , we choose a sampling scheme  $\lambda$  such that it provides the highest distinguishability obtained by the first step. Thus, we perform maximization in the outer part of (12). Finally, as the true model parameter  $\theta$  is unknown, we replace  $\theta$  by the MLE based on the current information. In Section 3.4, we provide a mirror descent algorithm for solving (12).

Unfortunately, directly using  $\lambda_n$  in the selection rule A as the choice probability does not guarantee asymptotic optimality. This is because  $\widehat{\lambda}_n$  does not guarantee sufficient exploration of all item pairs, which may lead to the imbalance between the exploration and exploitation for the sequential procedure. To fix this issue, we combine  $\widehat{\lambda}_n$  with an  $\epsilon$ -greedy approach, which is widely used in balancing exploration and exploitation in multiarmed bandit and decision-making problems (see, e.g., Watkins [60]). Specifically, an exploration probability  $p \in (0,1)$  is chosen, which is typically small and may be chosen depending on the value of the relative sampling cost c. At each time n, with probability p, we select the next pair uniformly from A. With probability 1-p, the next pair is selected according to the categorical distribution specified by  $\widehat{\lambda}_n$ . In other words, for each pair (i,j), the choice

probability of the selection rule at time n is given by

$$\lambda_n^{i,j} = p \frac{2}{K(K-1)} + (1-p)\widehat{\lambda}_n^{i,j}.$$

**Remark 1.** We clarify that the proposed " $\epsilon$ -greedy" algorithm is one of the asymptotically optimal exploration methods, and there may be other exploration methods with similar theoretical properties. For example, the  $\epsilon$ -greedy algorithm with the exploration probability decaying at a rate  $n^{-\beta}$  when the sample size is n may be asymptotically optimal for a range of  $\beta > 0$ . The theoretical properties of these additional exploration methods is an interesting problem and worth further investigation.

We call the preceding selection rule  $A_p$ , where the subscript emphasizes its dependence on the exploration rate p. The two proposed sequential ranking policies are defined by  $\pi_1 := (A_p, T_1, R)$  and  $\pi_2 := (A_p, T_2, R)$ . The proposed sequential ranking policies are summarized in Algorithm 1, where the prior information of  $\Theta$  is only utilized through its support W in steps 1 and 2. Algorithm 1 is an iterative algorithm, which runs in  $T_1$  (or  $T_2$ ) iterations, where  $T_1$  (or  $T_2$ ) is a data-dependent stopping time. The major computational complexity for each iteration arises from solving two optimization problems in steps 1 and 2. Step 1 is a standard maximum likelihood estimation, which depends on the structure of the loss function I and the constraint I0. The computation for solving (13) is discussed in Section 3.4. The proofs of the theoretical results are provided in Section 6.

## Algorithm 1 (Sequential Ranking Policy)

**Input:** The probability mass (density) function  $f_{\theta}^{a}(x)$  for any pair  $a \in \mathcal{A}$ , the probability  $p \in (0,1)$  in  $\epsilon$ -greedy, and the support W of  $\rho(\theta)$ .

**Initialization:** Uniformly sample a pair  $a_1$  at random and observe the comparison outcome  $X_1$ .

**Iterate:** For  $n = 2, 3, \ldots$  until the stopping time T in (9) (or (10)) is reached.

1. Compute the MLE based on the previous n-1 comparisons:

$$\widehat{\boldsymbol{\theta}}^{(n-1)} = \arg \sup_{\boldsymbol{\theta} \in W} l_{n-1}(\boldsymbol{\theta}).$$

2. Compute

$$\widehat{\lambda}_{n} = \arg\max_{\lambda \in \Delta} \min_{\widetilde{\boldsymbol{\theta}} \in W: r(\widetilde{\boldsymbol{\theta}}^{(n-1)})} \sum_{(i,j) \in \mathcal{A}} \lambda^{i,j} D^{i,j} (\widehat{\boldsymbol{\theta}}^{(n-1)} || \widetilde{\boldsymbol{\theta}}).$$

$$(15)$$

- 3. Flip a coin with heads probability p.
  - If the outcome is heads, select the pair  $a_n$  uniformly at random over all pairs from A.
  - Otherwise, select the pair  $a_n$  according to the categorical distribution specified by  $\lambda_n$ .
- 4. Observe the comparison result  $X_n$  and update the likelihood function  $l_n(\theta)$ .

**Output:** The rank  $R = r(\widehat{\boldsymbol{\theta}}^{(T)})$ , that is, the global rank induced by  $\widehat{\boldsymbol{\theta}}^{(T)}$ .

## 3.2. Asymptotic Optimality

This section contains the main results of the paper, including (1) a lower bound on the risk of a general sequential ranking procedure and (2) theoretical analysis of the proposed procedures, which leads to their asymptotic optimality. The asymptotic optimality of the proposed method is established through the following theorems, which are introduced later in this section. Theorem 1 provides an asymptotic lower bound for the Bayes risk of an arbitrary sequential ranking policy. Theorems 2 and 3 provide asymptotic upper bounds for the proposed procedures in terms of their expected Kendall tau and expected stopping time, respectively. These upper bounds together lead to an asymptotic upper bound for the Bayes risk of the proposed procedures that matches the lower bound in Theorem 1. As the asymptotic lower and upper bounds match, we conclude that the proposed method is asymptotically optimal in Corollary 1. As a by-product, an exponential deviation bound for the MLE over a time window is also obtained in Lemma 1. The assumptions for our results are described and discussed.

**3.2.1. Notations.** Throughout the rest of the paper, we write  $a_c = O(b_c)$  for two sequences  $a_c$  and  $b_c$  if  $|a_c|/|b_c|$  is bounded uniformly in  $\theta$  as  $c \to 0$ . Similarly, we write  $a_c = \Omega(b_c)$  if  $a_c > 0$ ,  $b_c > 0$ , and  $b_c = O(a_c)$ . We also write  $a_c = o(b_c)$  if  $a_c/b_c \to 0$  uniformly in  $\theta$ . The norm  $\|\cdot\|$  indicates the  $\ell_2$  vector norm. Throughout the paper, we use the uppercase Greek letter  $\Theta$  to indicate the *random* score parameter and the lowercase Greek letter  $\theta$  to denote a *deterministic* vector.

**3.2.2. Main Results.** We first describe the assumptions. For technical needs, we make some regularity conditions on the prior distribution  $\rho(\theta)$ . Recall that we have fixed  $\theta_1 = 0$  and let  $\theta = (\theta_2, \dots, \theta_K) \in \mathbb{R}^{K-1}$  be the unknown model parameter.

**Assumption 1.** The support  $W := \overline{\{\theta \in \mathbb{R}^{K-1} : \rho(\theta) > 0\}}$  is a compact set in  $\mathbb{R}^{K-1}$ , where  $\overline{E}$  denotes the closure of a set E. In addition, for any permutation  $\sigma \in P_K$ ,  $(\{\theta \in \mathbb{R}^{K-1} : r(\theta) = \sigma\} \cap W)^{\circ} \neq \emptyset$ , where  $E^{\circ}$  denotes the interior of a set E.

**Assumption 2.** There exists a constant  $\delta_b > 0$  such that, for all s > 0 and  $\theta \in W$ ,  $m(B(\theta, s) \cap W) \ge \min{\{\delta_b s^{K-1}, 1\}}$ , where  $B(\theta, s)$  denotes the open ball centered at  $\theta$  with radius s and  $m(\cdot)$  denotes the Lebesgue measure.

**Assumption 3.** The function  $\log f_{\theta}^{a}(x)$  is continuously differentiable in  $\theta$  for all x uniformly. That is,

$$\sup_{\theta \in W, a \in \mathcal{A}, x} \|\nabla_{\theta} \log f_{\theta}^{a}(x)\| < \infty$$

In addition,  $\inf_{\theta \in W, a \in A, x} f_{\theta}^{a}(x) > 0$ .

**Assumption 4.** The Kullback-Leibler divergence satisfies  $\inf_{\theta,\widetilde{\theta} \in W: r(\widetilde{\theta}) \neq r(\theta)} \max_{(i,j)} D^{i,j}(\theta || \widetilde{\theta}) > 0$ .

**Assumption 5.** The prior density satisfies  $\inf_{\theta \in W^{\circ}} \rho(\theta) > 0$  and  $\sup_{\theta \in W} \rho(\theta) < \infty$ .

We provide some remarks on the regularity assumptions. Assumption 1 requires the prior distribution for  $\Theta$  to have a bounded support, which has a nonempty interior for each rank. Assumption 2 avoids the support W being singular. Assumption 3 requires the smoothness of the likelihood function. It also requires that the comparison probability is bounded away from zero and one. Assumption 4 requires that there is no tie in the support of the prior distribution. This is a standard assumption in sequential analysis, which corresponds to the classic "indifference zone" assumption in sequential hypothesis testing (Kiefer and Sacks [33], Lorden [39], Schwarz [49]). In particular, the indifference zone condition assumes that the null and alternative hypotheses are separated in the sense that the Kullback–Leibler divergence between the two hypotheses is positive, and if the true model parameter is in between the two hypotheses, then it is considered to be indifferent for selecting the null and alternative hypothesis. For example, for any  $\delta > 0$ ,  $\kappa > 0$ , the set

$$W = \{\theta : \|\theta\| \le \kappa \text{ and } \forall i \ne j \text{ such that } |\theta_i - \theta_i| \ge \delta\}$$
 (16)

satisfies Assumptions 1, 2, and 4. Assumption 5 requires the prior distribution to have a positive density function (bounded from zero) over the support. For instance, for the set W described in (16), the uniform prior over W satisfies Assumption 5. In addition, with such a uniform prior over W, the BTL model defined in (2) satisfies Assumptions 3 and 4. It is worthwhile to note that these technical assumptions are mainly for the theoretical development, and the proposed adaptive ranking policies are applicable in practice regardless of the conditions on W.

Recall the definition of  $D(\theta)$  in (12). We further define

$$t_c(\boldsymbol{\theta}) = \frac{|\log c|}{D(\boldsymbol{\theta})}.$$
 (17)

Note that, under Assumption 4,  $t_c(\theta)$  is always finite. Intuitively, for small c,  $|\log c|/\{\min_{\widetilde{\theta} \in W: r(\widetilde{\theta}) \neq r(\theta)} \sum_{(i,j)} \lambda^{i,j} D^{i,j}(\theta||\widetilde{\theta})\}$  is approximately the smallest expected sample for the simple-against-simple hypothesis testing problem  $H_0: X_n \sim f_{\widetilde{\theta}}^{a_n}$  against  $H_1: X_n \sim f_{\widetilde{\theta}}^{a_n}$  for some  $r(\widetilde{\theta}) \neq r(\theta)$ , where  $a_n$  is sampled from  $\lambda$ . Note that  $t_c(\theta) = |\log c|/D(\theta) = \inf_{\lambda \in \triangle} [|\log c|/\{\min_{\widetilde{\theta} \in W: r(\widetilde{\theta}) \neq r(\theta)} \sum_{(i,j)} \lambda^{i,j} D^{i,j}(\theta||\widetilde{\theta})\}]$ . Thus,  $t_c(\theta)$  is approximately the smallest expected sample size for distinguishing the global rank of  $\theta$  from other ranks with an adaptive selection step. We formalize these heuristic arguments in the following Theorems 1–3.

We first present a lower bound on the minimal Bayes risk  $V_c^*(\rho)$  defined in (6).

**Theorem 1.** *Under Assumptions* 1–5, we have

$$\liminf_{c\to 0} \frac{V_c^*(\rho)}{c\mathbb{E}t_c(\Theta)} \ge 1,$$

where  $\mathbb{E}t_c(\Theta) = \int_W t_c(\theta) \rho(\theta) d\theta$ .

Recall the definition in (7) that a policy  $\pi$  is said to be asymptotically optimal if  $V_c(\pi, \rho) = (1 + o(1))V_c^*(\rho)$  as  $c \to 0$ . Thus, to show a policy  $\pi$  is indeed asymptotically optimal, we only need to show that  $V_c(\pi, \rho) = (1 + o(1))c\mathbb{E}t_c(\Theta)$  as  $c \to 0$ , according to Theorem 1. We proceed to show that the proposed sequential ranking

method is asymptotically optimal. In Section 3.1, we propose two policies  $\pi_1 = (A_p, T_1, R)$ ,  $\pi_2 = (A_p, T_2, R)$ . Their risks consist of two parts, the expected Kendall tau and the expected sample size.

**Assumption 6.** For each  $\theta, \theta' \in W$  and  $\theta \neq \theta'$ , there exists  $a \in A$  that can distinguish  $\theta$  and  $\theta'$ . That is,  $\sum_{a \in A} D^a(\theta || \theta') > 0$  for  $\theta, \theta' \in W$ . In addition, there is a constant  $\delta > 0$  such that  $\sum_{a \in A} D^a(\theta || \theta') \geq \delta || \theta - \theta' ||^2$ .

Assumption 6 requires the identifiability of the model, which is critical for the consistency of the MLE. For the BTL model described in (2), Assumption 6 is satisfied after fixing  $\theta_1 = 0$ . In what follows, Theorems 2 and 3 provide asymptotic upper bounds for the expected Kendall tau and expected stopping time of the proposed method, respectively.

**Theorem 2.** Under Assumptions 1–6, we consider a policy  $\pi_l = (A, T_l, R)$  (l = 1, 2), where we choose  $p \propto |\log c|^{-\frac{1}{2} + \delta_0}$  for some  $\delta_0$  satisfying  $0 < \delta_0 < \frac{1}{2}$  in Algorithm 1 and  $R = \{R_{i,j}\}$ . Then,

$$\mathbb{E}L_K(\{R_{i,j}\}) = O(c)$$
 for  $l = 1, 2$ .

**Theorem 3.** Under Assumptions 1–6, we consider a policy  $\pi_l = (A, T_l, R)$  (l = 1, 2), where we choose  $p \propto |\log c|^{-\frac{1}{2} + \delta_0}$  for some  $\delta_0$  satisfying  $0 < \delta_0 < \frac{1}{2}$  in Algorithm 1 and  $R = \{R_{i,j}\}$ . Then,

$$\limsup_{c\to 0} \frac{\mathbb{E}T_l}{\mathbb{E}t_c(\Theta)} \le 1 \text{ for } l = 1, 2.$$

Combining this with the asymptotic lower bound on the minimal Bayes risk in Theorem 1 and noticing that  $\lim_{c\to 0} \mathbb{E}t_c(\Theta) = \infty$ , we arrive at the asymptotic optimality of the proposed policies.

**Corollary 1.** Under Assumptions 1–6, if we choose  $p \propto |\log c|^{-\frac{1}{2}+\delta_0}$  for some  $\delta_0$  satisfying  $0 < \delta_0 < \frac{1}{2}$ , then  $\pi_l = (A_p, T_l, R)$ , l = 1, 2, are asymptotically optimal policies.

**3.2.3. Consistency of MLE.** An auxiliary result obtained in deriving the upper bound for the expected sample size is the following exponential bound for the MLE over a time window.

**Lemma 1.** Let  $m \ge n$  and let  $\varepsilon_{\lambda,m,n}$  be a sequence of real numbers such that  $\min_{n \le t \le m,(i,j)} \lambda_t^{i,j} \ge \varepsilon_{\lambda,m,n}$ . In addition, let  $\delta_{m,n}$  be a sequence of positive numbers such that  $n\varepsilon_{\lambda,m,n}\delta_{m,n}^2 \to \infty$  as  $n \to \infty$ . Then,

$$\mathbb{P}_{\boldsymbol{\theta}}\left(\sup_{n\leq t\leq m}\|\widehat{\boldsymbol{\theta}}^{(t)}-\boldsymbol{\theta}\|\geq \delta_{m,n}\right)\leq e^{-\Omega(n\epsilon_{\lambda,m,n}^2\delta_{m,n}^4)}\times O(m^K),$$

where we denote  $\mathbb{P}_{\theta}(\cdot)$  the conditional probability  $\mathbb{P}(\cdot|\Theta=\theta)$  and  $\widehat{\theta}^{(t)}$  is the MLE defined in (8). Moreover, this upper bound is uniform for  $\theta \in W$ .

The proof is provided in the online supplement. From this lemma, we can derive exponential upper bounds concerning the uniform consistency of  $\widehat{\theta}^{(t)}$ . In particular, if we let  $\delta_{m,n}$  be a fixed positive constant and  $\varepsilon_{\lambda,m,n}^2 \gg m^{-1}\log m$  as  $m \to \infty$ , then we can show  $\sup_{t>n} \|\widehat{\theta}^{(t)} - \theta\| \to 0$  in probability as  $n \to \infty$  with additional steps.

## 3.3. Proof Strategy

We briefly explain the proof strategy for each of the main theorems. Theorem 1 provides a lower bound on  $V(\rho,\pi)$  for an arbitrary policy  $\pi=(A,T,R)$  by discussing two cases:  $\mathbb{E}L_K(R) \geq c |\log c|^2$  and  $\mathbb{E}L_K(R) < c |\log c|^2$ . For the first case, Theorem 1 is easily justified. The main technicalities are in the second case, in which the main step is to develop an upper bound for the probability  $\mathbb{P}(T \leq (1-\delta)\mathbb{E}t_c(\Theta))$  for any constant  $\delta > 0$ . Heuristically, we argue that, whenever  $\mathbb{E}L_K(R)$  is small, it implies that the likelihood ratios between the conditional probability measures of data given that  $\Theta$  has different ranking patterns are relatively large, which cannot be achieved with a relatively small sample size T. The rigorous proof for this heuristic statement is done through a change-of-measure argument and a large deviation bound for the likelihood ratio.

The proof of Theorem 2 is based on the analysis of the expected Kendall tau under the stopping times  $T_1$  and  $T_2$ . The analysis under  $T_2$  is based on the equation

$$\mathbb{E}L_K(R) = \sum_{i,j} \int_{\boldsymbol{\theta} \in W_{j,i}} \mathbb{P}_{\boldsymbol{\theta}} \left( \sup_{\widetilde{\boldsymbol{\theta}} \in W_{j,i}} l_{T_2}(\widetilde{\boldsymbol{\theta}}) - \sup_{\boldsymbol{\theta}' \in W_{i,j}} l_{T_2}(\boldsymbol{\theta}') > h(c) \right) \rho(\boldsymbol{\theta}) d\boldsymbol{\theta},$$

followed by developing an upper bound for the probability  $\mathbb{P}_{\theta}(\sup_{\widetilde{\theta} \in W_{i,i}} l_{T_2}(\widetilde{\theta}) - \sup_{\theta' \in W_{i,i}} l_{T_2}(\theta') > h(c))$ , where  $h(c) = |\log c|(1 + |\log c|^{-\alpha})$  is slightly larger than  $|\log c|$ . Intuitively, thanks to the  $\varepsilon$ -greedy algorithm and the stopping time, a sufficient amount of information has been collected upon stopping so that the error probability is

well controlled. The analysis under  $T_1$  is similar, and we omit the details here. To prove Theorem 3, we first note that  $p \propto |\log c|^{-\frac{1}{2} + \delta_0}$  for some positive  $\delta_0$  in the  $\varepsilon$ -greedy algorithm. Thus, we can apply Lemma 1 and show that the MLE  $\widehat{\theta}^{(t)}$  is consistent with an exponential error bound. Roughly, this justifies that  $\widehat{\lambda}_n$  defined in (14) is close to  $\lambda^*(\theta)$  given  $\Theta = \theta$ . Thus, the expected sample size  $\mathbb{E}(T_i|\Theta = \theta)$  approximates the one given by the selection rule  $\lambda^*(\theta)$  that can be further approximated by  $h(c)/D(\theta) = (1 + o(1))t_c(\theta)$ , where we recall that  $t_c(\theta) = |\log c|/D(\theta)$  is defined in (17). We can then justify Theorem 3 by taking the expectation with respect to the prior distribution of  $\Theta$  on both sides.

## 3.4. Optimization in Algorithm 1

In this section, we show that the key optimization problem in (13) can be solved efficiently using the mirror descent algorithm (see, e.g., Beck and Teboulle [8]).

Algorithm 2 (Mirror Descent Algorithm for Solving Equation (13))

**Input:** The MLE estimator  $\theta$  and total number of iterations m.

**Initialization:** A starting point  $\lambda^0 \in \Delta$  and a constant  $c_0 > 0$ .

**Iterate:** For t = 1, 2, ..., m:

Compute the maximizer:

$$\widetilde{\boldsymbol{\theta}}^{0}(\lambda^{t-1}) \in \underset{\widetilde{\boldsymbol{\theta}} \in W: r(\widetilde{\boldsymbol{\theta}}) \neq r(\boldsymbol{\theta})}{\arg \max} - \sum_{(i,j)} \lambda^{i,j,t-1} D^{i,j}(\boldsymbol{\theta} || \widetilde{\boldsymbol{\theta}})$$

- 2. Compute the subgradient  $\mathbf{g}(\lambda^{t-1})$ , where  $\mathbf{g}(\lambda^{t-1})_{i,j} = -D^{i,j}(\boldsymbol{\theta}||\widetilde{\boldsymbol{\theta}}^0(\lambda^{t-1}))$
- 3. Update for  $\lambda^t$ :

$$\lambda^{t} = \arg\min_{\lambda \in \Lambda} \left\{ \eta_{t} \langle \mathbf{g}(\lambda^{t-1}), \lambda \rangle + D(\lambda || \lambda^{t-1}) \right\}, \tag{18}$$

 $\lambda^t = \arg\min_{\lambda \in \Delta} \big\{ \eta_t \langle \mathbf{g}(\lambda^{t-1}), \lambda \rangle + D(\lambda || \lambda^{t-1}) \big\},$  where  $\eta_t = \frac{c_0}{\sqrt{t}}$  and  $D(\lambda || \lambda^{t-1})$  is the KL divergence between  $\lambda$  and  $\lambda^{t-1}$ , that is,  $D(\lambda || \lambda^{t-1}) = \sum_{i,j} \lambda_{i,j} \log \frac{\lambda^{i,j}}{\lambda^{i,j,l-1}}$ 

**Output:** The solution  $\widehat{\lambda} = \frac{1}{m} \sum_{t=1}^{m} \lambda^{t}$ .

Let us first consider the inner optimization problem

$$\widetilde{\boldsymbol{\theta}}^{0}(\lambda) \in \underset{\widetilde{\boldsymbol{\theta}} \in W: r(\widetilde{\boldsymbol{\theta}}) \neq r(\boldsymbol{\theta})}{\arg \max} - \sum_{(i,j)} \lambda^{i,j} D^{i,j}(\boldsymbol{\theta} || \widetilde{\boldsymbol{\theta}}), \tag{19}$$

in step 1 of Algorithm 2. We clarify that, in this optimization,  $\theta$  is fixed,  $\widetilde{\theta}$  is the decision variable with which we want to optimize, and the resulting  $\widetilde{\theta}^0(\lambda)$  depends on  $\theta$  and  $\lambda$ . For almost all the popular comparison models, the objective function  $-\sum_{(i,j)} \lambda^{i,j} D^{i,j}(\theta \| \widetilde{\theta})$  is smooth in  $\widetilde{\theta}$ . Moreover, the objective function is also concave in  $\widetilde{\theta}$  for comparison models in an exponential family form (e.g., the BTL model in (2)). When the support  $\{\widetilde{\theta} \in W : r(\widetilde{\theta}) \neq 0\}$  $r(\theta)$  can be written as the union of a finite number of convex sets (see Equation (20)), (19) can be obtained by solving finite maximization problems, each with a smooth concave objective function constrained in a convex set. Therefore, from now on, we assume that the inner optimization problem can be solved.

We then discuss the outer optimization problem

$$\min_{\lambda \in \triangle} h(\lambda), \ h(\lambda) = \max_{\widetilde{\boldsymbol{\theta}} \in W: r(\widetilde{\boldsymbol{\theta}}) \neq r(\boldsymbol{\theta})} \phi(\lambda, \widetilde{\boldsymbol{\theta}}), \ \phi(\lambda, \widetilde{\boldsymbol{\theta}}) = -\sum_{(i,j)} \lambda^{i,j} D^{i,j}(\boldsymbol{\theta} || \widetilde{\boldsymbol{\theta}}).$$

When  $\phi(\lambda, \hat{\theta})$  is a continuous and bounded function and the set W is compact, further noting that  $\phi(\lambda, \hat{\theta})$  is convex in  $\lambda$  for every  $\widetilde{\theta}$ ,  $h(\lambda)$  is a convex function in  $\lambda$ , by Danskin's Theorem (see Bertsekas [9, proposition B.25]). Moreover, for a given  $\lambda$ , let  $\widetilde{\boldsymbol{\theta}}^0(\lambda) \in \arg\max_{\widetilde{\boldsymbol{\theta}} \in W: r(\widetilde{\boldsymbol{\theta}}) \neq r(\boldsymbol{\theta})} \phi(\lambda, \widetilde{\boldsymbol{\theta}})$  be one of the maximizers. Then, by Danskin's theorem,  $\mathbf{g}(\lambda)$  with  $\mathbf{g}(\lambda)_{i,j} = -D^{i,j}(\boldsymbol{\theta}||\widetilde{\boldsymbol{\theta}}^{0}(\lambda))$  is a subgradient of  $h(\lambda)$  as used in step 2 of Algorithm 2.

Finally, (18) in step 3 of the algorithm has a closed-form solution obtained by writing down the Karush–Kuhn–Tucker condition. That is,

$$\lambda^{i,j,t} = \frac{1}{C} \lambda^{i,j,t-1} \exp\left(-\eta_t \mathbf{g}(\lambda^{t-1})_{i,j}\right),\,$$

where  $\lambda^{i,j,t}$  is the (i,j)th component of  $\lambda^t$  and the normalization constant  $C = \sum_{i,j} \lambda^{i,j,t-1} \exp\left(-\eta_t \mathbf{g}(\lambda^{t-1})_{i,j}\right)$ .

From Beck and Teboulle [8] or Bubeck [14, theorem 4.2], we have the following convergence rate for Algorithm 2.

**Proposition 1** (Beck and Teboulle [8]). Assuming the inner optimization in (19) can be solved exactly, the mirror descent algorithm in Algorithm 2 is guaranteed to converge to the optimal solution at the rate of  $O(\sqrt{1/t})$ . That is, when  $t = O(1/\epsilon^2)$ , we have  $h(\widehat{\lambda}) - \min_{\lambda \in \Delta} h(\lambda) \le \epsilon$ .

We clarify that, for W defined in the example (16), it is a union over exponentially many convex sets. Thus, the proposed method requires exponential computational time for such a W. On the other hand, it is possible to have a fully polynomial computational-time algorithm if a misspecified  $\widetilde{W}$  is adopted (see (20) in the next section).

## 3.5. Model Misspecification

In practice, the support W of the prior distribution  $\rho(\cdot)$  may be unknown. In this case, we may choose

$$\widetilde{W} = \bigcup_{(i,j)} \widetilde{W}_{i,j} \text{ and } \widetilde{W}_{i,j} = \{ \boldsymbol{\theta} : \theta_i \ge \theta_j \} \cap \{ \boldsymbol{\theta} : |\theta_i| \le M, 2 \le i \le K \}$$
 (20)

in the sequential ranking policy for some reasonable positive constant M. With this misspecified support of  $\rho(\cdot)$ , the resulting policy may not achieve the asymptotic lower bound of the Bayes risk presented in Theorem 1 because of the incomplete information. On the other hand, the Bayes risk of the resulting ranking procedure can still achieve the same order of the minimal Bayes risk as  $c \to 0$ . That is,  $\limsup_{c \to 0} V_c(\rho, \pi)/V_c^*(\rho)$  is finite but greater than one. The following assumption is made to guarantee that the function  $f_{\theta}^a(x)$  has similar regularity on  $\widetilde{W}$  as on W. This assumption is mild. For example, it is satisfied for W,  $\widetilde{W}$ , and  $f_{\theta}^a(x)$  described in (16), (20), and (2), respectively.

**Assumption 7.**  $\sup_{\theta \in \widetilde{W}, a \in \mathcal{A}, x} \|\nabla_{\theta} \log f_{\theta}^{a}(x)\| < \infty$ ,  $\inf_{\theta \in \widetilde{W}, a \in \mathcal{A}, x} f_{\theta}^{a}(x) > 0$ , and  $\inf_{\theta \in W, \widetilde{\theta} \in \widetilde{W}: r(\widetilde{\theta}) \neq r(\theta)} \max_{(i,j)} D^{i,j}(\theta \| \widetilde{\theta}) > 0$ . In addition, there is a constant  $\delta > 0$  such that  $\sum_{a \in \mathcal{A}} D^{a}(\theta \| \theta') \ge \delta \|\theta - \theta'\|^{2}$  for all  $\theta \in W$  and  $\theta' \in \widetilde{W}$ .

**Theorem 4.** If we replace W by  $\widetilde{W}$  and replace  $W_{i,j}$  by  $\widetilde{W}_{i,j}$  (defined in (20)) in (9), (10), and (13) as well as in Algorithm 1 and adopt the resulting policy  $\pi_l = (A, T_l, R)$  (l = 1, 2) with  $p \propto |\log c|^{-\frac{1}{2} + \delta_0}$  for some  $\delta_0$  satisfying  $0 < \delta_0 < \frac{1}{2}$ , then under Assumptions 1, 2, 5, and 7,

$$\limsup_{c\to 0} \frac{V_c(\rho,\pi)}{V_c^*(\rho)} \leq \frac{\mathbb{E}\{1/\widetilde{D}(\Theta)\}}{\mathbb{E}\{1/D(\Theta)\}},$$

where  $D(\theta)$  is defined in (12) and  $\widetilde{D}(\theta)$  is defined as

$$\widetilde{D}(\boldsymbol{\theta}) = \max_{\boldsymbol{\lambda} \in \boldsymbol{\Lambda}} \inf_{\widetilde{\boldsymbol{\theta}} \in \widetilde{\boldsymbol{W}}: r(\widetilde{\boldsymbol{\theta}}) \neq r(\boldsymbol{\theta})} \sum_{(i,j)} \boldsymbol{\lambda}^{i,j} D^{i,j}(\boldsymbol{\theta} || \widetilde{\boldsymbol{\theta}}).$$

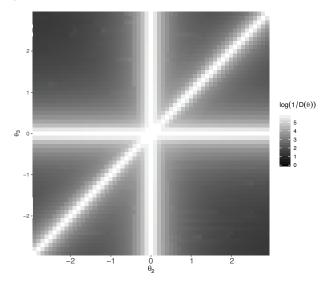
To obtain Theorem 4, we perform a similar analysis as those for Theorems 2 and 3. Although  $\widetilde{W}$  violates the separation property required by Assumption 4, a similar proof strategy still applies under Assumption 7. Roughly, this is because the expected sample size  $\mathbb{E}(T_l|\Theta=\theta)$  is now approximated by  $|\log c|/\widetilde{D}(\theta)$  and  $\widetilde{D}(\theta)>0$  for  $\theta\in W$ . Note that, to have  $\widetilde{D}(\theta)>0$ , we only need the support W to have the separation property and  $\widetilde{W}$  can contain ties among the parameters.

## 4. Numerical Examples

## 4.1. Behavior of $D(\Theta)$

Our main results suggest that the oracle risk  $V_c^*(\rho) \approx c|\log c|\mathbb{E}\{1/D(\Theta)\}$  when cost c is close to zero under the assumptions required by Theorems 2 and 3. The quantity  $1/D(\theta)$  can be naturally viewed as a measure of difficulty for the rank aggregation task when the true parameter vector is  $\theta$ . In what follows, we numerically investigate the behavior of  $1/D(\theta)$ .

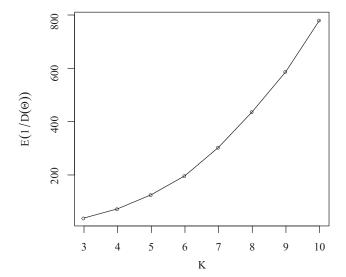
**Figure 1.** A level plot for the value of  $\log(1/D(\theta))$  as a function of  $\theta_2$  (*x*-axis) and  $\theta_3$  (*y*-axis), where K = 3 and  $W = \{\theta : ||\theta|| \le 3$ ,  $\theta_1 = 0$ , and  $\forall i \ne j$  such that  $|\theta_i - \theta_j| \ge 0.1$ }.



We first show the value of  $1/D(\theta)$  as a function of  $\theta$ , when the number of items K = 3. The support W of the prior distribution is chosen according to (16) that satisfies  $W = \{\theta : ||\theta|| \le 3, \theta_1 = 0, \text{ and } \forall i \ne j \text{ such that } |\theta_i - \theta_j| \ge 0.1\}$ . Figure 1 provides a level plot for the value of  $\log(1/D(\theta))$  as a function of  $\theta_2$  and  $\theta_3$ . As we can see, the value of  $1/D(\theta)$  becomes larger when the values of  $\theta_1$ ,  $\theta_2$ , and  $\theta_3$  are closer to each other and becomes smaller when they are more distinct.

We further show how the value of  $\mathbb{E}(1/D(\Theta))$  depends on the number of items K. For each choice of K, the support W is chosen as (16) with  $\theta_1=0$ ,  $\kappa=3$ , and  $\delta=0.1$ . Figure 2 shows that the value of  $\mathbb{E}(1/D(\Theta))$  is an increasing function of K, where  $\mathbb{E}(1/D(\Theta))$  is approximated by 2,000 Monte Carlo simulations. As we can see from Figure 2,  $\mathbb{E}(1/D(\Theta))$  increases with K, suggesting that the rank aggregation task becomes more difficult, on average, when the number of items becomes larger.

**Figure 2.** The value of  $\mathbb{E}(1/D(\Theta))$  as a function of K, where  $K = 3, 4, \dots, 10$ . For each choice of K, the support  $W = \{\theta : ||\theta|| \le 3, \theta_1 = 0, \text{ and } \forall i \ne j \text{ such that } |\theta_i - \theta_j| \ge 0.1\}$ , and  $\Theta$  follows a uniform distribution on W. Each  $\mathbb{E}(1/D(\Theta))$  is computed by 2,000 Monte Carlo simulations.



**Table 1.** Comparison between adaptive selection and random selection rules under a fixed-length stopping criterion. Each cell gives the averaged Kendall tau distance/0–1 loss for global ranking based on 1,000 independent simulations.

		Kendall tau		0–1 loss				
Sample size	20	40	60	20	40	60		
Adaptive selection	0.217	0.115	0.075	0.195	0.113	0.074		
Random selection	0.226	0.137	0.114	0.210	0.137	0.111		

## 4.2. Effectiveness of Adaptive Selection

We now show the power of the proposed adaptive selection rule by comparing it with a random selection rule that randomly picks a pair of items in each iteration. For each selection rule, we stop data collection once a fixed number of observations are collected, and sample sizes 20, 40, and 60 are considered. In the adaptive selection method, we set p=0.2 for the  $\epsilon$ -greedy strategy. The adaptive selection is implemented using Algorithm 2 with the number of iterations m=200,  $\lambda^{i,j,0}=2/(K(K-1))$ , and  $c_0=1$ . Note that the random selection method is essentially an off-line approach. The comparison is conducted under a model with K=3,  $W=\{\theta:\|\theta\|\leq 3,\ \theta_1=0,$  and  $\forall i\neq j$  such that  $|\theta_i-\theta_j|\geq 0.1\}$ , and the prior distribution  $\rho$  is a uniform distribution on W. For each selection rule and each sample size, 1,000 independent simulations are conducted. Two performance metrics are considered, including the Kendall tau distance (3) and the 0–1 loss for the recovery of global ranking that indicates whether the global ranking of  $\theta$  is completely recovered.

The results are given in Table 1 on the averaged Kendall tau distance and the averaged 0–1 loss for global ranking. As we can see, for each sample size, both the average Kendall tau distance and the average 0–1 loss for global ranking are smaller when applying the adaptive selection rule. The advantage of adaptive over random selection becomes more substantial as the observation size increases.

Under the current simulation setting, collecting one additional sample takes about six seconds, which is mainly because of solving optimization Problem (15) in Algorithm 1. Note that the complexity of solving (15) depends on the number of disconnected regions that the support W has, which grows exponentially with W. Therefore, for large values of W, it is suggested to simplify the computation by using the misspecified support W in (20), which can be written as the union of W half-planes.

## 4.3. Effectiveness of Adaptive Stopping

We further assess the effectiveness of the two stopping rules. The same model as before is used, that is, K = 3,  $W = \{\theta : ||\theta|| \le 3, \theta_1 = 0, \text{ and } \forall i \ne j \text{ such that } |\theta_i - \theta_j| \ge 0.1\}$ , and the prior distribution  $\rho$  is a uniform distribution on W. For the proposed adaptive stopping rules, we set  $h(c) = |\log c|(1 + |\log c|^{-0.5})$ , where  $\log c = -0.25$ , -0.5, -0.75, -1, -1.25, and -1.5 are considered. The proposed adaptive selection rule is used with  $p = 0.2 \times |\log c|^{-\frac{1}{4}}$ . For each stopping rule and each value of c, 1,000 independent simulations are conducted for which the averaged sample size, the Kendall tau distance, and the Bayes risk (5) are recorded as shown in Tables 2 and 3.

We then compare these adaptive stopping rules with the fixed-length stopping rule. More precisely, for each value of *c* and each adaptive stopping rule, we consider a policy with the same adaptive selection rule and the sample size fixed to be the corresponding averaged sample size. The averaged Kendall tau distance is also obtained based on 1,000 independent simulations and is reported in Tables 2 and 3.

Comparing each adaptive stopping rule with the corresponding fixed-length stopping rule, we see that the adaptive stopping rule gives substantially smaller averaged Kendall tau distances for all choices of *c*. It suggests

**Table 2.** Comparison between the proposed stopping rule  $T_1$  and a fixed-length stopping rule with the same adaptive selection rule. For both methods, the averaged Kendall tau distances are given, each of which is computed based on 1,000 independent simulations. For stopping rule  $T_1$ , the Bayes risks are also given as a linear combination of the Kendall tau distance and sampling cost.

		Kendall tau					Bayes risk						
Sample size	31	46	60	76	91	110	$\log(c)$	-0.25	-0.5	-0.75	-1	-1.25	-1.5
$T_1$ Fixed length	0.107 0.207	0.057 0.133	0.039 0.100	0.019 0.069	0.017 0.059	0.014 0.052	$T_1$	23.9	27.8	28.5	28.3	26.2	24.5

**Table 3.** Comparison between the proposed stopping rule  $T_2$  and a fixed-length stopping rule with the same adaptive selection rule. For both methods, the averaged Kendall tau distances are given, each of which is computed based on 1,000 independent simulations. For stopping rule  $T_2$ , the Bayes risks are also given as a linear combination of the Kendall tau distance and sampling cost.

		Kendall tau						Bayes risk						
Sample size	19	35	50	64	88	105	$\log(c)$	-0.25	-0.5	-0.75	-1	-1.25	-1.5	
$T_2$ Fixed length	0.190 0.207	0.112 0.133	0.057 0.100	0.029 0.069	0.020 0.059	0.014 0.052	$T_2$	15.3	21.3	23.8	23.7	25.3	23.5	

that the adaptive stopping rules lead to more accurate ranking aggregation results than the nonadaptive stopping rule.

Comparing the results in Tables 2 and 3, it seems that stopping rule  $T_1$  has slightly better performance than  $T_2$  in terms of Kendall's tau distance when the value of c is large. For example, the averaged Kendall tau distance for  $T_1$  is 0.107 when the averaged sample size is 31, and that for  $T_2$  is 0.112 when the averaged sample size is 35. Similarly,  $T_1$  achieves an averaged Kendall tau distance 0.057 when the averaged sample size is 46, and  $T_2$  achieves the same value with an averaged sample size of 50. However, as c decays (e.g., when  $\log(c) = -1.25$ , -1.5), the two procedures have similar performance in terms of the averaged sample size and Kendall tau distance. Regarding Bayes risks, we see that, for each value of c, the Bayes risks of  $T_2$  tend to be smaller than those of  $T_1$ . This is because sampling cost is the dominant term in the Bayes risk. As  $T_2$  tends to stop slightly earlier than  $T_1$ , its Bayes risks tend to be smaller. The difference in the corresponding Bayes risks becomes smaller when c decays. When  $\log(c) = -1.25$ , -1.5, the Bayes risks of the two methods are quite close to each other. It is worth pointing out that the difference in the finite sample performance when c is relatively large may depend on the choice of h(c), and the two stopping times are asymptotically equivalent when c goes to zero.

## 5. Concluding Remarks

In this paper, we consider the sequential design of rank aggregation with adaptive pairwise comparison. This problem is not only of practical importance because of its wide applications in fields such as psychology, politics, marketing, and sports, but it is also of theoretical significance in sequential analysis. Because of the more complex structure of the ranking problem than the hypothesis testing problems, no existing sequential analysis framework is suitable. We formulate the problem under a Bayesian decision framework and develop asymptotically optimal policies. Compared with the existing Bayesian sequential hypothesis testing problems, the problem solved in this paper is technically more challenging because of the more structured risk function. Novel technical tools are developed to solve this problem, and they are of separate theoretical interest in solving complex sequential design problems.

The current work may be extended in several directions. First, an even larger class of comparison models may be considered. The models considered in the current paper all assume the judges to be homogeneous; that is, the comparison outcome does not depend on who the judge is. It is of interest to consider the heterogeneity of the judges by incorporating judge-specific random effects into the comparison models and develop corresponding sequential designs. Second, different risk structures are incorporated into the sequential ranking designs to account for practical needs in different applications. For example, we consider other metrics for assessing the ranking accuracy (e.g., based on the accuracy of identifying the set of top *K* items) and nonuniform costs for different judges.

The results for the pairwise comparison problem can be extended to the case for multiple choices by extending the BTL model in (2) to the multinomial logit model (Train [55]). More specifically, given an L-tuple at time n,  $a_n = (a_{n,1}, \ldots, a_{n,L})$ , the annotator chooses  $X_n \in \{1, \ldots, L\}$  following the distribution  $\mathbb{P}(X_n = k) = \frac{\exp(\theta_{a_{n,k}})}{\sum_{k'=1}^{L} \exp(\theta_{a_{n,k'}})}$ .

However, additional challenges arise from solving the corresponding optimization problem in (13) that incurs higher complexity as a result of exploring more combinations of choices. For example, if there are K items and L choices presented to the annotator each time, we need to solve an optimization problem involving  $\binom{K}{L}$  combinations. It is worth further investigation on how to reduce the computational burden while keeping a certain optimality.

## 6. Proof of Theorems

In this section, we present the proofs of Theorems 1–3. The proofs for the lemmas are delayed to the online supplement. Throughout the proofs, we use the constants  $\delta_{\rho} = \inf_{\theta \in W} \rho(\theta) > 0$  and  $\sup_{\theta \in W, x, a \in \mathcal{A}} |\nabla f_{\theta}^{a}(x)| \le \kappa_{0}$ . According to Assumptions 5 and 3, these two constants are finite.

## 6.1. Proof for Theorem 1

Let  $\varepsilon = c|\log c|^2$ . For an arbitrary policy  $\pi = (A, T, R)$  and a prior probability density function  $\rho$ , there are two possibilities: either  $\mathbb{E}L_K(R) \ge \varepsilon$  or  $\mathbb{E}L_K(R) < \varepsilon$ . For the first case, we can see  $V(\rho, \pi) \ge \varepsilon \ge (1 + o(1))c\mathbb{E}t_c(\Theta)$ . For the second case, we have

$$V(\pi, \rho) = \mathbb{E}L_K(R) + c\mathbb{E}T \ge c\mathbb{E}T.$$

Therefore, to prove the theorem, it is sufficient to show that

$$\liminf_{c\to 0} \frac{c\mathbb{E}T}{c\mathbb{E}t_c(\Theta)} \ge 1,$$

or equivalently, for each  $\delta > 0$ , there exists a positive constant  $c_0 > 0$  such that, for  $c < c_0$ ,

$$\mathbb{E}T \ge (1 - \delta)\mathbb{E}t_c(\Theta).$$

Let  $t_{c,\delta}(\theta) = (1 - 2\delta/3)t_c(\theta)$  for each  $\delta > 0$ . Then, we arrive at a lower bound

$$\begin{split} \mathbb{E}T &\geq \mathbb{E}[TI(T > t_{c,\delta}(\Theta))] \\ &\geq \int \rho(\theta) t_{c,\delta}(\theta) \mathbb{P}_{\theta}(T > t_{c,\delta}(\theta)) d\theta \\ &= \mathbb{E}t_{c,\delta}(\Theta) - \int \rho(\theta) t_{c,\delta}(\theta) \mathbb{P}_{\theta}(T \leq t_{c,\delta}(\theta)) d\theta \\ &\geq \mathbb{E}t_{c,\delta}(\Theta) - t_{\max,\delta} \mathbb{P}(T \leq t_{c,\delta}(\Theta)), \end{split}$$

where we define  $t_{\max,\delta} = \max_{\theta \in W} t_{c,\delta}(\theta)$  and recall that  $\mathbb{P}_{\theta}$  represents for the conditional probability  $\mathbb{P}(\cdot|\Theta = \theta)$ . According to Assumption 4, we have  $t_{\max,\delta} = O(|\log c|) = O(\mathbb{E}t_c(\Theta))$ . Therefore, it is sufficient to show

$$\mathbb{P}(T \le t_{c,\delta}(\Theta)) = o(1).$$

We proceed to an upper bound for  $\mathbb{P}(T \le t_{c,\delta}(\Theta))$ . We abuse the notation a little and write  $U_r = \{\theta : r(\theta) = r\}$ , the set of parameters that gives the rank r. Then, we have

$$\mathbb{P}(T \le t_{c,\delta}(\Theta)) = \sum_{r \in P_K} \mathbb{P}(T \le t_{c,\delta}(\Theta), \Theta \in U_r)$$

$$= O(1) \times \max_{r \in P_K} \mathbb{P}(T \le t_{c,\delta}(\Theta), \Theta \in U_r). \tag{21}$$

We proceed to an upper bound for  $\mathbb{P}(T \leq t_{c,\delta}(\Theta), \Theta \in U_r)$  for each  $r \in P_K$ . Define an event

$$B_r = \left\{ \frac{\mathbb{P}(\Theta \in U_r | \mathcal{F}_T)}{\max_{(i,j):W_{i,j} \cap U_r = \emptyset} \mathbb{P}(\Theta \in W_{i,j} | \mathcal{F}_T)} > \frac{c^{\frac{\delta}{10}}}{\varepsilon} \right\},\tag{22}$$

where  $\mathcal{F}_n = \sigma(X_1, ..., X_n, a_1, ..., a_n)$  denotes the  $\sigma$ -algebra generated by  $X_1, ..., X_n$  and  $a_1, ..., a_n$ . We split the probability

$$\mathbb{P}(T \le t_{c,\delta}(\Theta), \Theta \in U_r)$$

$$= \mathbb{P}(T \le t_{c,\delta}(\Theta), \Theta \in U_r, B_r) + \mathbb{P}(T \le t_{c,\delta}(\Theta), \Theta \in U_r, B_r^c),$$

which can be bounded from above by

$$\mathbb{P}(T \le t_{c,\delta}(\Theta), \Theta \in U_r) \le \mathbb{P}(T \le t_{c,\delta}(\Theta), \Theta \in U_r, B_r) + \mathbb{P}(\Theta \in U_r, B_r^c). \tag{23}$$

We establish upper bounds for the two terms on the right-hand side of the preceding inequality separately. The next lemma, whose proof is presented in the online supplement, provides an upper bound for the second term.

**Lemma 2.** For all  $r \in P_K$ , if  $\mathbb{E}L_K(R) \leq \varepsilon$ , then

$$\mathbb{P}(\Theta \in U_r, B_r^c) \leq \left(1 + \frac{c^{\frac{\delta}{10}}}{\varepsilon}\right) \varepsilon.$$

We proceed to the first term  $\mathbb{P}(T \leq t_{c,\delta}(\Theta), \Theta \in U_r, B_r)$  on the right-hand side of (23). Then,

$$\mathbb{P}(T \le t_{c,\delta}(\Theta), \Theta \in U_r, B_r) = \int_{U_r} \mathbb{P}_{\theta}(T \le t_{c,\delta}(\theta), B_r) \rho(\theta) d\theta. \tag{24}$$

Recall the definition of the event  $B_r$  in (22), and we have

$$B_r \cap \left\{ T \leq t_{c,\delta}(\boldsymbol{\theta}) \right\} \subset \left\{ \max_{1 \leq t \leq t_{c,\delta}(\boldsymbol{\theta})} \frac{\mathbb{P}(\boldsymbol{\Theta} \in U_r | \mathcal{F}_t)}{\max_{W_{i,j} \cap U_r = \emptyset} \mathbb{P}(\boldsymbol{\Theta} \in W_{i,j} | \mathcal{F}_t)} > \frac{c^{\frac{\delta}{10}}}{\varepsilon} \right\}.$$

Consequently,

$$\mathbb{P}_{\theta}(T \le t_{c,\delta}(\theta), B_r) \le \mathbb{P}_{\theta}\left(\max_{1 \le t \le t_{c,\delta}(\theta)} \frac{\mathbb{P}(\Theta \in U_r | \mathcal{F}_t)}{\max_{W_{i,j} \cap U_r = \emptyset} \mathbb{P}(\Theta \in W_{i,j} | \mathcal{F}_t)} > \frac{c^{\frac{\delta}{10}}}{\varepsilon}\right). \tag{25}$$

We proceed to an upper bound for the preceding display. For each  $\theta$ , we define a random sequence  $\{\theta_t^* : 1 \le t \le t_{c,\delta}(\theta)\}$  as follows:

$$\boldsymbol{\theta}_{t}^{*} = \underset{\widetilde{\boldsymbol{\theta}} \in W: r(\widetilde{\boldsymbol{\theta}}) \neq r(\boldsymbol{\theta})}{\arg \min} \sum_{n=1}^{t} \sum_{i,j} \lambda_{n}^{i,j} D^{i,j}(\boldsymbol{\theta} || \widetilde{\boldsymbol{\theta}}).$$

Intuitively,  $\theta_t^*$  is the score parameter that is most difficult to distinguish from  $\theta$  at time t among those that have different rank with  $\theta$  given that item selection rules  $\lambda_1, \ldots, \lambda_n$  have been adopted. We further choose the index process  $(i_t^*, j_t^*)$  to be such that  $\theta_t^* \in W_{i_t^*, j_t^*}$  but  $\theta \notin W_{i_t^*, j_t^*}$ . If there are multiple (i, j)'s satisfying this, then we choose  $(i_t^*, j_t^*)$  arbitrarily from them. From the definition, we know  $\theta_t^*$  and  $(i_t^*, j_t^*)$  are adapted to  $\sigma(\lambda_1, \ldots, \lambda_t)$  and, thus, are adapted to  $\mathcal{F}_{t-1}$ . We use the next lemma to transform the probability in (25) to a probability based on a martingale parameterized by  $\theta$ .

**Lemma 3.** For each  $\theta' \in U_r$ , define a martingale with respect to the filtration  $\{\mathcal{F}_n : n \geq 1\}$  and probability measure  $\mathbb{P}_{\theta}$  as follows:

$$M_{t}(\boldsymbol{\theta}') = l_{t}^{\vec{a}}(\boldsymbol{\theta}') - l_{t}^{\vec{a}}(\boldsymbol{\theta}_{t}^{*}) - \sum_{n=1}^{t} \sum_{(i,j)} \lambda_{n}^{i,j} D^{i,j}(\boldsymbol{\theta} || \boldsymbol{\theta}_{t}^{*}) + \sum_{n=1}^{t} \sum_{(i,j)} \lambda_{n}^{i,j} D^{i,j}(\boldsymbol{\theta} || \boldsymbol{\theta}'),$$

where  $l_t^{\vec{a}}(\theta) = \log \prod_{i=1}^t f_{\theta}^{a_i}(X_i)$ . Then, there exists a positive constant  $c_0 > 0$  such that, for  $0 < c < c_0$ ,

$$\mathbb{P}_{\theta} \left( \max_{1 \leq t \leq t_{c,\delta}(\theta)} \frac{\mathbb{P}(\Theta \in U_r | \mathcal{F}_t)}{\max_{W_{i,j} \cap U_r = \emptyset} \mathbb{P}(\Theta \in W_{i,j} | \mathcal{F}_t)} > \frac{c^{\frac{\delta}{10}}}{\varepsilon} \right) \\
\leq \mathbb{P}_{\theta} \left( \max_{1 \leq t \leq t_{c,\delta}(\theta), \ \theta' \in U_r} M_t(\theta') \geq \frac{\delta}{2} |\log c| \right).$$
(26)

According to this lemma, to find an upper bound for (25), it is sufficient to find an upper bound for the right-hand side of (26), which is the probability that a stochastic process indexed by  $\theta'$  and t goes above a certain level. In this paper, we use the following two lemmas repeatedly to handle this type of level-crossing probability. The first one is the Azuma–Hoeffding inequality proved by Azuma [3] and Hoeffding [28].

**Lemma 4** (Azuma-Hoeffding Inequality). Let  $M_n$  be a martingale with respect to the filtration  $\{\mathcal{F}_n : n = 1, 2, ..\}$ . Let  $X_n = M_n - M_{n-1}$ . Assume that  $X_n$  is bounded and  $X_n \in [a_n, b_n]$ , where  $a_n$  and  $b_n$  are deterministic constants. Then, for each t > 0, we have

$$\mathbb{P}\left(\max_{1 \le m \le n} M_m \ge t\right) \le \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right).$$

The next lemma is the key lemma that allows us to derive the level-crossing probability by aggregating marginal tail bounds of a random field. Its proof is given in the online supplement.

**Lemma 5.** Let  $\{\zeta(\theta): \theta \in W\}$  be a random field over a compact set  $U \subset \mathbb{R}^K$  that satisfies Assumption 2. Let  $\beta(\theta,b)$  be defined as follows:

$$\beta(\boldsymbol{\theta}, b) = \mathbb{P}(\zeta(\boldsymbol{\theta}) \ge b),$$

where  $\mathbb{P}$  is a probability measure and we assume that  $\zeta(\cdot)$  has a continuous sample path almost surely under  $\mathbb{P}$ . Assume that  $\zeta(\cdot)$  has a Lipschitz-continuous sample path in the sense that there exists a constant  $\kappa_L$  such that, for all  $\theta, \theta' \in W$ ,

$$|\zeta(\theta) - \zeta(\theta')| \le \kappa_L ||\theta - \theta'||$$
 almost surely under  $\mathbb{P}$ .

Then, we have that, for all positive  $\gamma$ ,

$$\mathbb{P}\Big(\max_{\theta \in W} \zeta(\theta) \geq b\Big) \leq \int_{W} \beta(\theta, b - \gamma) d\theta \times \frac{\kappa_{L}^{K-1}}{\gamma^{K-1} \delta_{b}},$$

where  $\delta_b$  is the constant defined in Assumption 2.

Set  $n := t_{c,\delta}(\theta)$ ,  $t := \frac{\delta}{2} |\log c| - 1$ ,  $M_n := M_n(\theta')$ , and  $a_n = -b_n := 2 \max_{x,a \in \mathcal{A}, \theta \in W} |\log f_{\theta,x}^a(x)|$  in Lemma 4, and we have, for each  $\theta'$ ,

$$\mathbb{P}_{\theta}\left(\max_{1\leq n\leq t_{c,\delta}(\theta)} M_n(\theta') \geq \frac{\delta}{2}|\log c| - 1\right) \leq \exp\left(-\frac{2\left(\frac{\delta}{2}|\log c| - 1\right)^2}{t_{c,\delta}(\theta)a_1^2}\right).$$

According to Assumptions 1 and 3, we have  $a_1 < \infty$ , and consequently,

$$\mathbb{P}_{\theta}\left(\max_{1\leq n\leq t_{c,\delta}(\theta)} M_n(\theta') \geq \frac{\delta}{2}|\log c| - 1\right) \leq \exp\left(-\Omega(\delta^2|\log c|)\right). \tag{27}$$

Note that, for  $\theta'$ ,  $\widetilde{\theta} \in U_r$ ,

$$\begin{split} \max_{1 \leq n \leq t_{c,\delta}(\boldsymbol{\theta})} M_n(\boldsymbol{\theta}') - \max_{1 \leq n \leq t_{c,\delta}(\boldsymbol{\theta})} M_n(\widetilde{\boldsymbol{\theta}}) \\ \leq \max_{1 \leq n \leq t_{c,\delta}(\boldsymbol{\theta})} |M_n(\boldsymbol{\theta}') - M_n(\widetilde{\boldsymbol{\theta}})| \\ \leq t_{c,\delta}(\boldsymbol{\theta}) \kappa_0 ||\boldsymbol{\theta}' - \widetilde{\boldsymbol{\theta}}||, \end{split}$$

where  $\kappa_0 = 4\sup_{a \in \mathcal{A}, \theta' \in W, x} |\nabla \log f_{\theta}^a(x)| < \infty$  denotes the Lipschitz constant of  $M_1(\theta')$ . Therefore,  $M_n(\theta')$  is a Lipschitz-continuous random field in  $\theta'$ . The preceding display and (27), together with Lemma 5, give

$$\mathbb{P}_{\theta} \left( \max_{1 \leq n \leq t_{c,\delta}(\theta), \, \theta' \in U_r} M_n(\theta') \geq \frac{\delta}{2} |\log c| \right)$$

$$\leq \exp\left(-\Omega(\delta^2 |\log c|)\right) m(U_r) \frac{t_{c,\delta}(\theta)^{K-1} \kappa_0^{K-1}}{\delta_b}$$

$$= \exp\left(-\Omega(\delta^2 |\log c|)\right) \times O(|\log c|^{K-1}),$$

where we recall that  $m(\cdot)$  denotes the Lebesgue measure. The preceding inequality and (24)–(26) give

$$\mathbb{P}(T \le t_{c,\delta}(\Theta), \Theta \in U_r, B_r) \le \exp\left(-\Omega(\delta^2|\log c|)\right) \times O(|\log c|^{K-1}).$$

Combine this with Lemma 2 and (23), and we have

$$\mathbb{P}(T \leq t_{c,\delta}(\Theta), \Theta \in U_r) \leq \left(1 + \frac{c^{\frac{\delta}{10}}}{\varepsilon}\right) \varepsilon + \exp\left(-\Omega(\delta^2 |\log c|)\right) \times O(|\log c|^{K-1}).$$

Combine the preceding display with (21), and we have

$$\mathbb{P}(T \leq t_{c,\delta}(\Theta)) \leq O(1) \times \left\{ \left(1 + \frac{c^{\frac{\delta}{10}}}{\varepsilon}\right) \varepsilon + \exp\left(-\Omega(\delta^2 |\log c|)\right) \times O(|\log c|^{K-1}) \right\}.$$

Therefore,  $\mathbb{P}(T \leq t_{c,\delta}(\Theta)) = o(1)$  as  $c \to 0$ . This completes the proof.

#### 6.2. Proof of Theorem 2

We start with the stopping time  $T_2$ . With the decision rule D defined in (11), the expected Kendall tau at the stopping time  $T_2$  is

$$\mathbb{E}L_{K}(R) = \mathbb{E}\sum_{(i,j)} I(\Theta_{i} < \Theta_{j}) R_{i,j}$$

$$= \int_{W(i,j): \theta \notin W_{j,i}} \mathbb{P}_{\theta} \left( \sup_{\widetilde{\theta} \in W_{j,i}} l_{T_{2}}(\widetilde{\theta}) > \sup_{\theta' \in W_{i,j}} l_{T_{2}}(\theta') \right) \rho(\theta) d\theta$$

$$= \int_{W \theta \notin W_{j,i}} \mathbb{P}_{\theta} \left( \sup_{\widetilde{\theta} \in W_{j,i}} l_{T_{2}}(\widetilde{\theta}) - \sup_{\theta' \in W_{i,j}} l_{T_{2}}(\theta') > h(c) \right) \rho(\theta) d\theta, \tag{28}$$

where we write  $l_t(\theta) = \sum_{n=1}^t \log f_{\theta}^{a_n}(X_n)$  as the log-likelihood function. Equation (28) is bounded from above by

$$\mathbb{E}L_{K}(R) \leq \sup_{\boldsymbol{\theta} \in W} \rho(\boldsymbol{\theta}) \times m(W) \times \frac{K(K-1)}{2}$$

$$\times \sup_{\boldsymbol{\theta} \in W} \max_{(i,j): \boldsymbol{\theta} \notin W_{j,i}} \mathbb{P}_{\boldsymbol{\theta}} \left\{ \sup_{\widetilde{\boldsymbol{\theta}} \in W_{j,i}} l_{T_{2}}(\widetilde{\boldsymbol{\theta}}) - l_{T_{2}}(\boldsymbol{\theta}) > h(c) \right\}. \tag{29}$$

To obtain the preceding inequality, we use the fact that  $\sup_{\theta' \in W_{i,j}} l_{T_2}(\theta') \ge l_{T_2}(\theta)$  for (i, j) such that  $\theta \notin W_{j,i}$  and  $\sup_{\theta \in W} \rho(\theta) < \infty$  according to Assumption 5. We split the probability

$$\mathbb{P}_{\theta} \left\{ \sup_{\widetilde{\boldsymbol{\theta}} \in W_{ji}} l_{T_{2}}(\widetilde{\boldsymbol{\theta}}) - l_{T_{2}}(\boldsymbol{\theta}) > h(c) \right\}$$

$$\leq \mathbb{P}_{\theta} \left\{ \sup_{\widetilde{\boldsymbol{\theta}} \in W_{ji}} l_{T_{2}}(\widetilde{\boldsymbol{\theta}}) - l_{T_{2}}(\boldsymbol{\theta}) > h(c) \text{ and } T_{2} \leq \tau \right\} + \mathbb{P}_{\theta}(T_{2} \geq \tau). \tag{30}$$

We clarify that  $\theta'$ ,  $\theta$ , and  $\widetilde{\theta}$  are deterministic vectors here. The second term on the right-hand side of the preceding display is controlled by the next lemma.

**Lemma 6.** If  $\tau = \Omega(|\log c|^3)$ , then

$$\mathbb{P}_{\theta}(T_i \ge \tau) \le c^2 \quad (i = 1, 2).$$

We proceed to an upper bound of the first term on the right-hand side of (30). Define a stopping time  $T_2 \wedge \tau = \min(T_2, \tau)$ , and then we have

$$\mathbb{P}_{\theta} \left\{ \sup_{\widetilde{\boldsymbol{\theta}} \in W_{ji}} l_{T_{2}}(\widetilde{\boldsymbol{\theta}}) - l_{T_{2}}(\boldsymbol{\theta}) > h(c) \text{ and } T_{2} \leq \tau \right\}$$

$$\leq \mathbb{P}_{\theta} \left\{ \sup_{\widetilde{\boldsymbol{\theta}} \in W_{ji}} l_{T_{2} \wedge \tau}(\widetilde{\boldsymbol{\theta}}) - l_{T_{2} \wedge \tau}(\boldsymbol{\theta}) > h(c) \right\}.$$

Now, we consider the random field  $\eta(\widetilde{\boldsymbol{\theta}}) = l_{T_2 \wedge \tau}(\widetilde{\boldsymbol{\theta}}) - l_{T_2 \wedge \tau}(\boldsymbol{\theta})$  for  $\widetilde{\boldsymbol{\theta}} \in W_{ji}$ . We proceed to an upper bound for  $\mathbb{P}_{\boldsymbol{\theta}}(\sup_{\widetilde{\boldsymbol{\theta}} \in W_{ii}} \eta(\widetilde{\boldsymbol{\theta}}) > h(c))$  through Lemma 5. We first note that  $\eta(\widetilde{\boldsymbol{\theta}})$  is a Lipschitz-continuous function:

$$|\eta(\widetilde{\boldsymbol{\theta}}) - \eta(\widetilde{\boldsymbol{\theta}}')| \le |l_{T \wedge \tau}(\widetilde{\boldsymbol{\theta}}) - l_{T \wedge \tau}(\widetilde{\boldsymbol{\theta}}')| \le \tau \kappa_0 ||\widetilde{\boldsymbol{\theta}} - \widetilde{\boldsymbol{\theta}}'||. \tag{31}$$

We further obtain the marginal tail probability of  $\eta(\tilde{\theta})$  through the next lemma.

**Lemma 7.** For all  $\theta \neq \theta$  and all constant A > 0, we have

$$\mathbb{P}_{\boldsymbol{\theta}}\Big(l_{T\wedge\tau}(\widetilde{\boldsymbol{\theta}})-l_{T\wedge\tau}(\boldsymbol{\theta})\geq A\Big)\leq e^{-A}.$$

We take A = h(c) - 1 in the preceding lemma and obtain

$$\mathbb{P}_{\theta}(\eta(\widetilde{\boldsymbol{\theta}}) \ge h(c) - 1) \le e^{-h(c) + 1}$$
.

Combining the preceding display with (31) and Lemma 5, we arrive at

$$\mathbb{P}_{\theta} \left( \sup_{\widetilde{\theta} \in W_{ji}} \eta(\widetilde{\theta}) > h(c) \right) \le O(\tau^{K-1} e^{-h(c)}). \tag{32}$$

We combine (32), (29) and Lemma 6 and arrive at

$$\mathbb{P}_{\theta} \left\{ \sup_{\widetilde{\theta} \in W_{ji}} l_{T_2}(\widetilde{\theta}) - l_{T_2}(\theta) > h(c) \right\}$$

$$\leq O(c^2) + O(e^{-|\log c| - |\log c|^{1-\alpha} + (K-1)\log \tau})$$

$$= O(c^2) + O(ce^{-|\log c|^{1-\alpha} + 3(K-1)\log|\log c|})$$

$$= o(c).$$

This completes our analysis for  $T_2$ . We proceed to the analysis of the policy  $\pi_1$  and the stopping time  $T_1$ . According to the definition of  $T_1$  in (10), we can see that, upon stopping,

$$\begin{split} \max_{(i,j):1 \leq i < j \leq K} & \exp \left[ \min \left\{ \sup_{\widetilde{\boldsymbol{\theta}} \in W_{i,j}} l_{T_1}(\widetilde{\boldsymbol{\theta}}) - \sup_{\boldsymbol{\theta} \in W} l_{T_1}(\boldsymbol{\theta}), \sup_{\widetilde{\boldsymbol{\theta}} \in W_{j,i}} l_{T_1}(\widetilde{\boldsymbol{\theta}}) - \sup_{\boldsymbol{\theta} \in W} l_{T_1}(\boldsymbol{\theta}) \right\} \\ & \leq \sum_{(i,j):1 \leq i < j \leq K} \exp \left[ \min \left\{ \sup_{\widetilde{\boldsymbol{\theta}} \in W_{i,j}} l_{T_1}(\widetilde{\boldsymbol{\theta}}) - \sup_{\boldsymbol{\theta} \in W} l_{T_1}(\boldsymbol{\theta}), \sup_{\widetilde{\boldsymbol{\theta}} \in W_{j,i}} l_{T_1}(\widetilde{\boldsymbol{\theta}}) - \sup_{\boldsymbol{\theta} \in W} l_{T_1}(\boldsymbol{\theta}) \right\} \right] \leq e^{-h(c)}. \end{split}$$

Taking the logarithm and rearranging terms in the preceding display, we have

$$\min_{1 \le i < j \le K} \left[ \sup_{\boldsymbol{\theta} \in W} l_n(\boldsymbol{\theta}) - \min \left\{ \sup_{\widetilde{\boldsymbol{\theta}} \in W_{i,j}} l_n(\widetilde{\boldsymbol{\theta}}), \sup_{\widetilde{\boldsymbol{\theta}} \in W_{j,i}} l_n(\widetilde{\boldsymbol{\theta}}) \right\} \right] \ge h(c). \tag{33}$$

With (33), we can follow similar derivations as those for (29) and arrive at

$$\mathbb{E}L_{K}(\overline{D}_{T_{1}}) \leq \sup_{\theta' \in W} \rho(\theta)m(W) \times \frac{K(K-1)}{2} \sup_{\theta \in W} \max_{(i,j): \theta \notin W_{j,i}} \mathbb{P}_{\theta} \left[ \sup_{\widetilde{\theta} \in W_{ji}} l_{T_{1}}(\widetilde{\theta}) - l_{T_{1}}(\theta) > h(c) \right].$$

The rest of the proof is similar to that for the stopping time  $T_2$ . We omit the details.

## 6.3. Proof of Theorem 3

Let  $\delta$  be an arbitrary positive number; then, we can find an upper bound for the expectation of a stopping time T as follows:

$$\mathbb{E}T = \sum_{m=0}^{\infty} \mathbb{E}[TI(m(1+\delta)t_{c}(\Theta) \leq T < (m+1)(1+\delta)t_{c}(\Theta))]$$

$$\leq (1+\delta)\mathbb{E}t_{c}(\Theta) + \sum_{m=1}^{\infty} \mathbb{E}[TI(m(1+\delta)t_{c}(\Theta) \leq T < (m+1)(1+\delta)t_{c}(\Theta))]$$

$$\leq (1+\delta)\mathbb{E}t_{c}(\Theta)$$

$$+ (1+\delta)\max_{\theta \in W} t_{c}(\theta) \sum_{m=1}^{\infty} (m+1)\mathbb{P}(m(1+\delta)t_{c}(\Theta) \leq T < (m+1)(1+\delta)t_{c}(\Theta))$$

$$\leq (1+\delta)\mathbb{E}t_{c}(\Theta) + (1+\delta)\max_{\theta \in W} t_{c}(\theta) \sum_{m=1}^{\infty} (m+1)\max_{\theta \in W} \mathbb{P}_{\theta}(m(1+\delta)t_{c}(\theta) \leq T < (m+1)(1+\delta)t_{c}(\theta)). \tag{34}$$

We proceed to an upper bound for the probability in the preceding sum for  $T = T_i$  (i = 1, 2). We start with  $T = T_2$ . We split the probability for  $m \ge 1$ :

$$\mathbb{P}_{\boldsymbol{\theta}}(m(1+\delta)t_{c}(\boldsymbol{\theta}) \leq T_{2} < (m+1)(1+\delta)t_{c}(\boldsymbol{\theta}))$$

$$\leq \mathbb{P}_{\boldsymbol{\theta}}\left(m(1+\delta)t_{c}(\boldsymbol{\theta}) \leq T_{2} < (m+1)(1+\delta)t_{c}(\boldsymbol{\theta}), \max_{m(1+\delta)\delta_{2}t_{c}(\boldsymbol{\theta})\leq t\leq m(1+\delta)t_{c}(\boldsymbol{\theta})} \|\widehat{\boldsymbol{\theta}}^{(t)} - \boldsymbol{\theta}\| \leq |\log c|^{-\delta_{1}}\right)$$

$$+ \mathbb{P}_{\boldsymbol{\theta}}\left(\max_{m(1+\delta)\delta_{2}t_{c}(\boldsymbol{\theta})\leq t\leq m(1+\delta)t_{c}(\boldsymbol{\theta})} \|\widehat{\boldsymbol{\theta}}^{(t)} - \boldsymbol{\theta}\| \geq |\log c|^{-\delta_{1}}\right), \tag{35}$$

where we choose  $\delta_1 = \frac{\delta_0}{8}$  and  $\delta_2 = |\log c|^{-\delta_0/2}$ , and  $\delta_0$  is defined in the selection rule where we recall that  $p \propto |\log c|^{-\frac{1}{2}+\delta_0}$ . The second term on the preceding display is bounded from above according to Lemma 1, in which we set  $n := m(1+\delta)\delta_2 t_c(\theta)$ ,  $m := m(1+\delta)t_c(\theta)$ ,  $\varepsilon_{\lambda} = \Omega(|\log c|^{-\frac{1}{2}+\delta_0})$  and  $\delta_{m,n} = |\log c|^{-\delta_1}$  and arrive at

$$\mathbb{P}_{\theta} \left( \max_{m(1+\delta)\delta_{2}t_{c}(\theta) \leq t \leq m(1+\delta)t_{c}(\theta)} ||\widehat{\theta}^{(t)} - \theta|| \geq |\log c|^{-\delta_{1}} \right) \\
\leq e^{-\Omega(m(1+\delta)\delta_{2}t_{c}(\theta)|\log c|^{-4\delta_{1}}|\log c|^{-1+2\delta_{0}})} \times O(m^{K-1}|\log c|^{K-1}) \\
= e^{-\Omega(m|\log c|^{2\delta_{0}-4\delta_{1}}\delta_{2})} O(m^{K-1}|\log c|^{K-1}) \\
= e^{-\Omega(m|\log c|^{\delta_{0}})} O(m^{K-1}|\log c|^{K-1}).$$
(36)

We proceed to the first term on the right-hand side of (35). For  $m \ge 1$ , we can see that  $T_2 > m(1+\delta)t_c(\theta)$  implies that there exists (i,j) such that  $|\sup_{\widetilde{\theta} \in W_{i,j}} l_n(\widetilde{\theta}) - \sup_{\theta' \in W_{j,i}} l_n(\theta')| \le h(c)$  for  $n = (1+\delta)mt_c(\theta)$ . Without loss of generality, we assume that  $\theta \in W_{i,j}$ ; then,  $T_2 > m(1+\delta)t_c(\theta)$  further implies  $l_n(\theta) - \sup_{\theta' \in W_{j,i}} l_n(\theta') \le h(c)$ . Therefore, an upper bound for the first term on the right-hand side of (35) is

$$\mathbb{P}_{\theta}\left(m(1+\delta)t_{c}(\theta) \leq T_{2} \leq (m+1)(1+\delta)t_{c}(\theta), \max_{m(1+\delta)\delta_{2}t_{c}(\theta)\leq t\leq m(1+\delta)t_{c}(\theta)} \|\widehat{\boldsymbol{\theta}}^{(n)} - \boldsymbol{\theta}\| \leq |\log c|^{-\delta_{1}}\right) \\
\leq \mathbb{P}_{\theta}\left(l_{n}(\theta) - \sup_{\theta' \in W_{j,i}} l_{n}(\theta') \leq h(c), \max_{m(1+\delta)\delta_{2}t_{c}(\theta)\leq t\leq m(1+\delta)t_{c}(\theta)} \|\widehat{\boldsymbol{\theta}}^{(n)} - \boldsymbol{\theta}\| \leq |\log c|^{-\delta_{1}}\right), \tag{37}$$

We present an upper bound for the preceding display in the next lemma.

**Lemma 8.** If the strategy  $\lambda^*(\widehat{\boldsymbol{\theta}}^{(t)})$  is adopted with probability 1 - o(1) uniformly for  $mt_c(\boldsymbol{\theta})(1 + \delta)\delta_2 \le t \le m(1 + \delta)t_c(\boldsymbol{\theta})$ . Then,

$$\mathbb{P}_{\boldsymbol{\theta}} \left( l_n(\boldsymbol{\theta}) - \sup_{\boldsymbol{\theta}' \in W_{j,i}} l_n(\boldsymbol{\theta}') \le h(c), \max_{m(1+\delta)\delta_2 t_c(\boldsymbol{\theta}) \le t \le m(1+\delta)t_c(\boldsymbol{\theta})} \|\widehat{\boldsymbol{\theta}}^{(n)} - \boldsymbol{\theta}\| \le |\log c|^{-\delta_1} \right)$$

$$\le e^{-\Omega(m|\log c|)} \times O(|\log c|^{K-1} m^{K-1}),$$

where  $n = (1 + \delta)mt_c(\boldsymbol{\theta})$ .

We combine the preceding lemma with (36) and (35), we arrive at

$$\mathbb{P}_{\theta}(m(1+\delta)t_{c}(\theta) \leq T_{2} < (m+1)(1+\delta)t_{c}(\theta)) \leq (e^{-\Omega(m|\log c|)} + e^{-\Omega(m|\log c|^{\delta_{0}})}) \times O(m^{K-1}|\log c|^{K-1}).$$

This together with (34) gives

$$\begin{split} \mathbb{E}T_2 &\leq (1+\delta)\mathbb{E}t_c(\Theta) \\ &+ O(|\log c|) \times \sum_{m=1}^{\infty} (m+1) \Big\{ (e^{-\Omega(m|\log c|)} + e^{-\Omega(m|\log c|^{\delta_0})}) \times O(m^{K-1}|\log c|^{K-1}) \Big\} \\ &\leq (1+\delta)\mathbb{E}t_c(\Theta) + o(|\log c|). \end{split}$$

This completes our analysis for  $T_2$ . We proceed to the analysis of  $T_1$ . We can see that the event  $T_1 > n$  implies that

$$\sum_{(i,j)} \exp \left[ \min \left\{ \sup_{\widetilde{\boldsymbol{\theta}} \in W_{i,j}} l_n(\widetilde{\boldsymbol{\theta}}) - \sup_{\boldsymbol{\theta} \in W} l_n(\boldsymbol{\theta}), \sup_{\widetilde{\boldsymbol{\theta}} \in W_{j,i}} l_n(\widetilde{\boldsymbol{\theta}}) - \sup_{\boldsymbol{\theta} \in W} l_n(\boldsymbol{\theta}) \right\} \right] > e^{-h(c)},$$

which further implies that

$$K(K-1)\max_{(i,j)} \exp \left[ \min \left\{ \sup_{\widetilde{\boldsymbol{\theta}} \in W_{i,j}} l_n(\widetilde{\boldsymbol{\theta}}) - \sup_{\boldsymbol{\theta} \in W} l_n(\boldsymbol{\theta}), \sup_{\widetilde{\boldsymbol{\theta}} \in W_{j,i}} l_n(\widetilde{\boldsymbol{\theta}}) - \sup_{\boldsymbol{\theta} \in W} l_n(\boldsymbol{\theta}) \right\} \right] > e^{-h(c)}.$$

Simplifying the preceding display, we can see it is equivalent to there existing (i, j) such that

$$\left|\sup_{\widetilde{\boldsymbol{\theta}}\in W_{i,j}}l_n(\widetilde{\boldsymbol{\theta}}) - \sup_{\boldsymbol{\theta}\in W_{j,i}}l_n(\boldsymbol{\theta})\right| \le h(c) + \log K(K-1).$$

The analysis is similar for the stopping time  $T_1$  to that of  $T_2$  by replacing h(c) by  $h(c) + \log K(K-1)$  in the derivation following (37). We omit the details.

## 6.4. Proof of Theorem 4

First, to distinguish between the sequential method with and without model misspecification, we use the notation "–" over a method (e.g., the sequential ranking rule  $\overline{\pi}_l = (\overline{A}, \overline{T}_l, \overline{R})$  and the MLE  $\overline{\theta}^{(t)}$ ) to indicate that it is based on the algorithm with the misspecified support  $\widetilde{W}$  of the prior distribution  $\rho(\cdot)$ . The proof of Theorem 4 follows similar arguments as those of Corollary 1. That is, we show the following modified version of Theorems 2 and 3, whose proofs are provided in the online supplement.

**Proposition 2.** Following the sequential ranking rules  $\overline{\pi}_l = (\overline{A}, \overline{T}_l, \overline{R})$  (for l = 1, 2), we have

$$\mathbb{E}L_K(\{\overline{R}_{i,i}\}) = O(c).$$

## **Proposition 3.**

$$\limsup_{c\to 0} \frac{\mathbb{E}\overline{T}_l}{\mathbb{E}\widetilde{t}_c(\Theta)} \leq 1,$$

where we define  $\widetilde{t}_c(\theta) = \frac{|\log(c)|}{\widetilde{D}(\theta)}$  and  $\widetilde{D}(\theta) = \max_{\lambda \in \Delta} \min_{\widetilde{\theta} \in \widetilde{W}: r(\widetilde{\theta}) \neq r(\theta)} \sum_{(i,j)} \lambda^{i,j} D^{i,j}(\theta || \widetilde{\theta}).$ 

Combining Theorems 2 and 3 with Propositions 2 and 3, we arrive at

$$\limsup_{c \to 0} \frac{V_c(\rho, \pi)}{V_c^*(\rho)} \le \lim_{c \to 0} \frac{O(c) + c\mathbb{E}\widetilde{t}_c(\Theta)}{O(c) + c\mathbb{E}t_c(\Theta)} = \lim_{c \to 0} \frac{O(c) + c|\log c|\mathbb{E}\{1/\widetilde{D}(\Theta)\}}{O(c) + c|\log c|\mathbb{E}\{1/D(\Theta)\}} = \frac{\mathbb{E}\{1/\widetilde{D}(\Theta)\}}{\mathbb{E}\{1/D(\Theta)\}}$$

## Acknowledgments

The authors thank the anonymous associate editor and two anonymous referees for many useful suggestions and feedback that greatly improved the paper. The authors also thank Dr. Jingchen Liu and Dr. Zhiliang Ying for helpful discussions.

#### **Endnote**

<sup>1</sup> The computation time is evaluated based on our implementation of the proposed method in R version 3.6.1 on a standard desktop PC with Intel(R) Core(TM) i5-5300 @2.3 GHZ.

## References

- [1] Adler RJ, Blanchet JH, Liu J (2012) Efficient Monte Carlo for high excursions of Gaussian random fields. *Ann. Appl. Probab.* 22(3): 1167–1214.
- [2] Albert AE (1961) The sequential design of experiments for infinitely many states of nature. Ann. Math. Statist. 32(3):774–799.
- [3] Azuma K (1967) Weighted sums of certain dependent random variables. Tohoku Math. J. (2). 19(3):357-367.
- [4] Ballinger TP, Wilcox NT (1997) Decisions, error and heterogeneity. Econom. J. (London). 107(443):1090–1105.
- [5] Bartroff J, Lai TL (2008) Efficient adaptive designs with mid-course sample size adjustment in clinical trials. *Statist. Medicine* 27(10): 1593–1611.
- [6] Bartroff J, Finkelman M, Lai TL (2008) Modern sequential analysis and its applications to computerized adaptive testing. Psychometrika 73:473–486.
- [7] Bartroff J, Lai TL, Shih MC (2013) Sequential Experimentation in Clinical Trials (Springer, New York).
- [8] Beck A, Teboulle M (2003) Mirror descent and nonlinear projected subgradient methods for convex optimization. *Oper. Res. Lett.* 31(3): 167–175.
- [9] Bertsekas D (1999) Nonlinear Programming (Athena Scientific).
- [10] Blumenthal AL (1977) The Process of Cognition (Prentice Hall/Pearson Education).

- [11] Bradley R, Terry M (1952) Rank analysis of incomplete block designs: I. the method of paired comparisons. Biometrika. 39(3/4):324–345.
- [12] Braverman M, Mossel E (2009) Sorting from noisy information. Preprint, submitted October 7, https://arxiv.org/abs/0910.1191.
- [13] Braverman M, Mao J, Weinberg MS (2016) Parallel algorithms for select and partition with noisy comparisons. Proc. *Annual Sympos. Theory Comput.*
- [14] Bubeck S (2015) Convex optimization: Algorithms and complexity. Foundations Trends Machine Learn. 8(3-4):231-357.
- [15] Chen X, Jiao K, Lin Q (2016) Bayesian decision process for cost-efficient dynamic ranking via crowdsourcing. J. Machine Learn. Res. 17(217):1–40.
- [16] Chen X, Li Y, Mao J (2018) An instance optimal algorithm for top-K ranking under the multinomial logit model. ACM-SIAM Sympos. Discrete Algorithms.
- [17] Chen X, Bennett PN, Collins-Thompson K, Horvitz E (2013) Pairwise ranking aggregation in a crowdsourced setting. *Proc. ACM Internat. Conf. Web Search Data Mining*.
- [18] Chen X, Gopi S, Mao J, Schneider J (2018) Optimal instance adaptive algorithm for the top-k ranking problem. *IEEE Trans. Inform. Theory* 64(9):6139–6160.
- [19] Chen Y, Suh C (2015) Spectral MLE: Top-k rank aggregation from pairwise comparisons. Proc. Internat. Conf. Machine Learn.
- [20] Chernoff H (1959) Sequential design of experiments. Ann. Math. Statist. 30(3):755-770.
- [21] Dragalin VP, Tartakovsky AG, Veeravalli VV (2000) Multihypothesis sequential probability ratio tests. II. Accurate asymptotic expansions for the expected sample size. *IEEE Trans. Inform. Theory* 46(4):1366–1383.
- [22] Draglia V, Tartakovsky AG, Veeravalli VV (1999) Multihypothesis sequential probability ratio tests. I. Asymptotic optimality. IEEE Trans. Inform. Theory 45(7):2448–2461.
- [23] Elo AE (1978) The Rating of Chessplayers, Past, and Present (Arco Publishing).
- [24] Garg N, Johari R (2019) Designing optimal binary rating systems. Proc. 22nd Internat. Conf. Artificial Intelligence Statist.
- [25] Hajek B, Oh S, Xu J (2014) Minimax-optimal inference from partial rankings. Proc. Adv. Neural Inform. Processing Systems.
- [26] Heckel R, Shah NB, Ramchandran K, Wainwright MJ (2019) Active ranking from pairwise comparisons and when parametric assumptions do not help. *Ann. Statist.* 47(6):3099–3126.
- [27] Hoeffding W (1960) Lower bounds for the expected sample size and the average risk of a sequential procedure. *Ann. Math. Statist.* 31(2): 352–368.
- [28] Hoeffding W (1963) Probability inequalities for sums of bounded random variables. J. Amer. Statist. Assoc. 58(301):13–30.
- [29] Hsiung AC, Ying ZL, Zhang CH, eds. (2004) Random Walk, Sequential Analysis and Related Topics: A Festschrift in Honor of Yuan-Shih Chow (World Scientific).
- [30] Jamieson K, Nowak R (2011) Active ranking using pairwise comparisons. Proc. 24th Internat. Conf. Neural Inform. Processing Systems, 2240–2248
- [31] Kallus N, Udell M (2020) Dynamic assortment personalization in high dimensions. Oper. Res. 68(4):1020-1037.
- [32] Kendall M, Gibbons JD (1990) Rank Correlation Methods, 5th ed. (Charles Griffin).
- [33] Kiefer J, Sacks J (1963) Asymptotically optimum sequential inference and design. Ann. Math. Statist. 34(3):705–750.
- [34] Lai TL (1988) Nearly optimal sequential tests of composite hypotheses. Ann. Statist. 16(2):856-886.
- [35] Lai TL (2001) Sequential analysis: Some classical problems and new challenges. Statista Sinica 11(2):303–351.
- [36] Lai TL, Shih MC (2004) Power, sample size and adaptation considerations in the design of group sequential clinical trials. *Biometrika* 91(3):507–528.
- [37] Li X, Liu J (2015) Rare-event simulation and efficient discretization for the supremum of Gaussian random fields. *Adv. Appl. Probab.* 47(3): 787–816.
- [38] Li X, Liu J, Ying Z (2018) Chernoff index for Cox test of separate parametric families. Ann. Statist. 46(1):1-29.
- [39] Lorden G (1976) 2-SPRT's and the modified Kiefer-Weiss problem of minimizing an expected sample size. Ann. Statist. 4(2):281–291.
- [40] Luce RD (1959) Individual Choice Behavior: A Theoretical Analysis (Wiley, New York).
- [41] Mao C, Weed J, Rigollet P (2018) Minimax rates and efficient algorithms for noisy sorting. Proc. Algorithmic Learn. Theory.
- [42] Mei Y (2010) Efficient scalable schemes for monitoring a large number of data streams. Biometrika 97(2):419-433.
- [43] Morrison HW (1963) Testable conditions for triads of paired comparison choices. Psychometrika 28:369–390.
- [44] Naghshvar M, Javidi T (2013) Active sequential hypothesis testing. Ann. Statist. 41(6):2703–2738.
- [45] Negahban S, Oh S, Sha D (2017) Rank centrality: Ranking from pair-wise comparisons. Oper. Res. 65(1):266-287.
- [46] Nitinawarat S, Veeravalli VV (2015) Controlled sensing for sequential multihypothesis testing with controlled Markovian observations and non-uniform control cost. Sequential Anal. 34(1):1–24.
- [47] Page L, Brin S, Motwani R, Winograd T (1999) The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab.
- [48] Saaty TL, Vargas LG (2012) The possibility of group choice: Pairwise comparisons and merging functions. Soc. Choice Welfare 38(3): 481–496.
- [49] Schwarz G (1962) Asymptotic shapes of Bayes sequential testing regions. Ann. Math. Statist. 33(1):224-236.
- [50] Shah NB, Balakrishnan S, Guntuboyina A, Wainright MJ (2017) Stochastically transitive models for pairwise comparisons: Statistical and computational issues. *IEEE Trans. Inform. Theory* 63(2):934–959.
- [51] Siegmund D (1985) Sequential Analysis: Tests and Confidence Intervals (Springer, New York).
- [52] Song Y, Fellouris G (2017) Asymptotically optimal, sequential, multiple testing procedures with prior information on the number of signals. *Electronic J. Statist.* 11(1):338–363.
- [53] Tartakovsky A, Nikiforov I, Basseville M (2014) Sequential Analysis: Hypothesis Testing and Changepoint Detection (Chapman and Hall/CRC).
- [54] Thurstone LL (1927) A law of comparative judgement. Psych. Rev. 34(4):273–286.
- [55] Train K (2009) Discrete Choice Methods with Simulation (Cambridge University Press).
- [56] Tsitovich I (1985) Sequential design of experiments for hypothesis testing. Theory Probab. Appl. 29(4):814-817.
- [57] Wald A (1945) Sequential tests of statistical hypotheses. Ann. Math. Statist. 16(2):117–186.
- [58] Wald A, Wolfowitz J (1948) Optimum character of the sequential probability ratio test. Ann. Math. Statist. 19(3):326–339.

- [59] Wang S, Lin H, Chang HH, Douglas J (2016) Hybrid computerized adaptive testing: From group sequential design to fully sequential design. J. Ed. Measurement 53(1):45–62.
- [60] Watkins CJCH (1989) Learning from delayed rewards. Unpublished PhD thesis, Cambridge University, Cambridge, UK.
- [61] Xie Y, Siegmund DO (2013) Sequential multi-sensor change-point detection. Ann. Statist. 41(2):670–692.
- [62] Ye S, Fellouris G, Culpepper S, Douglas J (2016) Sequential detection of learning in cognitive diagnosis. *British J. Math. Statist. Psych.* 69(2):139–158.