

Testing for the Bradley-Terry-Luce Model

Anuran Makur* and Japneet Singh†

*Department of Computer Science and †School of Electrical and Computer Engineering,
Purdue University, West Lafayette, IN 47907
Email: {amakur, sing1041}@purdue.edu

Abstract—The Bradley-Terry-Luce (BTL) model is one of the most widely used models for ranking a set of items given data about pairwise comparisons among them. While several studies in the literature have attempted to empirically test how accurately a BTL model can model some given pairwise comparison data, this work aims to develop a formal, computationally efficient hypothesis test to determine whether the BTL model accurately represents the data. Specifically, we first propose such a formal hypothesis test, establish an upper bound on the critical radius of the proposed test, and then provide a complementary lower bound on the critical radius. Our bounds prove the minimax optimality of the scaling of the critical radius with respect to the number of items (up to constant factors). Finally, we also take the first step towards characterizing the stability of rankings under the BTL model when there is a small model mismatch.

I. INTRODUCTION

Many applications, such as sports tournaments, consumer preference surveys, and political voting, generate data in the form of pairwise comparisons between a set of items or agents (e.g., choices, teams). These datasets are useful for performing various data analysis tasks, such as ranking the items, analyzing the skill level of a particular team over time, and examining market or sports competitiveness, cf. [1]–[15]. A popular modeling assumption to perform such learning and inference tasks is the *Bradley-Terry-Luce* (BTL) model [1]–[6]. The BTL model assigns a latent skill score $\alpha_i > 0$ to each item i , representing its relative merit compared to other items, and posits that the likelihood of i being preferred over an item j in a pairwise comparison is given by

$$\mathbb{P}(i \text{ is preferred over } j) = \frac{\alpha_i}{\alpha_i + \alpha_j}. \quad (1)$$

The BTL model is a natural consequence of the assumption of *independence of irrelevant alternatives* (IIA), which is widely used in economics and social choice theory [2]. Despite its widespread popularity, it is known that the IIA assumption does not hold well for various real-world datasets [16]–[18]. For example, the BTL model is oblivious to the “home-advantage effect” in sports, which refers to a home team’s possible advantage when playing against a visiting team (see, e.g., cricket [19], soccer [20]). Hence, several other models of pairwise comparisons have been proposed in the literature, e.g., modifications of the BTL model to incorporate home-advantage effect [21, Chapter 10], Thurstonian models [3], other generalizations of the BTL model [7], models of

rankings based on Borda scores [22], [23], and other non-parametric stochastically transitive models [24], [25].

Nevertheless, primarily because of its simplicity and interpretability, the BTL framework remains one of the most widely-used models. A large fraction of the associated results in the literature focuses on estimation of the skill score parameters of the BTL model. Some popular approaches include maximum likelihood estimation [6], [7], rank centrality (or Markov-chain-based) methods [10], [26], least-squares methods [13], and non-parametric methods [12], [24] (also see [8], [9] for Bayesian inference for BTL models). Once the parameters are estimated, they are then used for inference tasks, such as ranking items and learning skill distributions [14], [15]. An inherent assumption for such statistical analysis to be meaningful in real-world scenarios is that the BTL model accurately represents the given pairwise comparison data. Hence, it is important to understand precisely when the BTL assumption holds in a systematic manner.

A. Main Contributions

In contrast to the above works, we focus on the problem of testing whether the BTL assumption accurately models data generated from an underlying pairwise comparison model. First, we devise a notion of “distance” that allows us to quantify the deviation of a general pairwise comparison model from BTL models. We then use this distance to formally construct a hypothesis test to determine whether a BTL model accurately models the underlying data. We establish an upper bound on the minimax critical radius for this test. Furthermore, we also prove an information theoretic lower bound on the critical radius for this problem, thereby demonstrating the minimax optimality of the critical radius. Finally, we utilize the notion of distance mentioned above to analyze the stability of BTL assumptions in the context of rankings. More specifically, we investigate the deviation from the BTL condition that is sufficient for the ranking produced under the BTL assumption to differ from the classical Borda count ranking [27].

B. Related Literature

This work lies at the confluence of two fields of study: hypothesis testing and preference learning. The analysis of preference data, such as pairwise comparisons, has a rich history starting from the seminal works [1]–[6]. As mentioned earlier, the BTL model is one of the most well-studied models for pairwise comparison data [1]. It was first proposed by [6] as a method for estimating participants’ skill levels in chess

tournaments. Moreover, it is a special case of the Plackett-Luce model [2], [4], initially developed in mathematical psychology. We refer readers to [28], [29] for a comprehensive overview of different models of rankings. In the literature, many studies have focused on estimating parameters for the BTL model and characterizing the error bounds, cf. [7], [10], [11], [25], [26], [30]–[32]. For example, [11] presents non-asymptotic bounds for relative ℓ^∞ and ℓ^2 -norm estimation errors of normalized vectors of skill scores. We use some of these bounds in our arguments, although we need to make alterations since the proofs do not apply to general pairwise comparison models.

Hypothesis testing also has a rich history in statistics, ranging from Pearson’s χ^2 -test [33] to non-parametric tests [34]. Yet, to the best of our knowledge, no study has developed rigorous hypothesis tests to determine the validity of the BTL assumption in the literature. The minimax perspective of hypothesis testing, which we specialize in our setting, was initially proposed by [35]. It is worth mentioning that recently, [36] analyzed two-sample testing on pairwise comparison data, and [37] derived lower bounds for testing the IIA assumption given general preference data. For the special case of pairwise comparisons, the lower bounds in [37] agree with ours in terms of the high-level scaling law of the critical radius. However, our hypothesis testing problem is formulated differently to [37] and [37] does not provide upper bounds.

Furthermore, investigating the stability of the BTL assumption is another interesting question in the literature. For example, [22] provided empirical evidence that the BTL assumption is not very robust to changes in the pairwise comparison matrix. So, in this work, we also take the first steps towards rigorously characterizing the stability of rankings under the BTL model.

II. FORMAL SETUP AND DECISION RULE

A. Notational Preliminaries

We briefly collect some notation here that is used throughout this work. Let $\mathbf{1}_n \in \mathbb{R}^n$ be the column vector with all entries equal to 1, where we drop the subscript when it is clear from context, and $[n] \triangleq \{1, \dots, n\}$. Furthermore, for any vector $x \in \mathbb{R}^n$, $\text{diag}(x) \in \mathbb{R}^{n \times n}$ is the diagonal matrix with x along its principal diagonal. For a vector $x \in \mathbb{R}^n$, $\|x\|_2$ denotes its ℓ^2 -norm and x^p denotes the entry-wise p th power of x , i.e., $x^p = [x_1^p, x_2^p, \dots, x_n^p]^T$. For any matrix $A \in \mathbb{R}^{n \times n}$, $\|A\|_2$ and $\|A\|_F$ denote the operator norm and Frobenius norm of A , respectively. For a strictly positive vector $\pi \in \mathbb{R}^n$, we define a Hilbert space on \mathbb{R}^n with inner product $\langle x, y \rangle_\pi = \sum_{i=1}^n \pi_i x_i y_i$ and the corresponding vector and matrix norms $\|x\|_\pi = \sqrt{\langle x, x \rangle_\pi}$ and $\|A\|_\pi = \sup_{\|x\|_\pi=1} \|Ax\|_\pi$ and $\|A\|_{\pi, F} = (\sum_i \sum_j \pi_j A_{ij}^2)^{1/2}$.

B. Formal Model and Goal

We begin by introducing general pairwise comparison models. Consider a set of n agents, indexed by $[n]$ with $n \in \mathbb{N} \setminus \{1\}$, that engage in a tournament consisting of several pairwise comparisons. This scenario is ubiquitous in many real-world

applications. For example, in a sports tournament, $[n]$ represents the teams or players that play pairwise games with each other, and in discrete choice models from economics, $[n]$ represents alternatives that an individual may choose from. Several probabilistic models exist in the literature to capture such pairwise comparison settings, e.g., BTL model [1], [2], [5], Thurstonian models [3], and non-parametric models [24], [25], and all of them turn out to be specializations of the following general pairwise comparison model.

Definition 1 (Pairwise Comparison Model): For any pair of agents $i, j \in [n]$, $i \neq j$, let $p_{ij} \in (0, 1)$ denote the probability that agent j beats agent i in a “ i vs. j ” pairwise comparison. We refer to the collection of $n(n-1)$ parameters $\{p_{ij} : i, j \in [n], i \neq j\}$ as a *pairwise comparison model*.

Specifically, we consider an asymmetric setting where an “ i vs. j ” comparison may have a different meaning to a “ j vs. i ” comparison. Such asymmetric settings are commonly observed in sports like cricket, football, etc. [19]. Hence, such a model can be aptly summarized by a matrix $P \in \mathbb{R}^{n \times n}$ with

$$P_{ij} = \begin{cases} p_{ij}, & i \neq j, \\ \frac{1}{2}, & i = j, \end{cases} \quad (2)$$

where we have set $P_{ii} = \frac{1}{2}$ for notational convenience.

In our analysis, we will find it convenient to assign a time-homogenous Markov chain (or row stochastic matrix) on the finite state space $[n]$ to any pairwise comparison model. This canonical assignment is defined next.

Definition 2 (Canonical Markov Chain): For any pairwise comparison model $\{p_{ij} \in (0, 1) : i, j \in [n], i \neq j\}$ with matrix $P \in \mathbb{R}^{n \times n}$, its *canonical Markov chain* is given by the row stochastic matrix $S \in \mathbb{R}^{n \times n}$, where

$$S_{ij} = \begin{cases} \frac{p_{ij}}{n}, & i \neq j, \\ 1 - \frac{1}{n} \sum_{k \in [n] \setminus \{i\}} p_{ik}, & i = j. \end{cases}$$

As noted earlier, the most well-known specialization of the pairwise comparison model in Definition 1 is the BTL model defined below [1], [2], [5].

Definition 3 (BTL Model [1], [2], [5]): A pairwise comparison model $\{p_{ij} \in (0, 1) : i, j \in [n], i \neq j\}$ is known as a BTL (or multinomial logit) model if there exist skill score parameters $\alpha_i > 0$ for every agent $i \in [n]$ such that:

$$\forall i, j \in [n], i \neq j, \quad p_{ij} = \frac{\alpha_j}{\alpha_i + \alpha_j}.$$

Hence, we can describe a BTL model entirely using the collection of its n skill score parameters $\{\alpha_i : i \in [n]\}$.

We next describe how a pairwise comparison model characterizes the likelihood of a tournament between n agents. To this end, fix any pairwise comparison model $\{p_{ij} \in (0, 1) : i, j \in [n], i \neq j\}$. For any $i \neq j$, define the outcome of the m th i vs j pairwise comparison between them as the Bernoulli random variable

$$Z_{m_{ij}} \triangleq \begin{cases} 1, & \text{if } j \text{ beats } i \text{ (with probability } p_{ij}), \\ 0, & \text{if } i \text{ beats } j \text{ (with probability } 1 - p_{ij}), \end{cases} \quad (3)$$

for $m \in [k_{ij}]$, where k_{ij} denotes the number of i vs j comparisons. We assume throughout that the observation random variables $\mathcal{Z} \triangleq \{Z_{mij} : i, j \in [n], i \neq j, m \in [k_{ij}]\}$ are mutually independent. Let $Z_{ij} \triangleq \sum_{m=1}^{k_{ij}} Z_{mij}$. Clearly, it follows that for any $i \neq j$, Z_{ij} is a binomial random variable, i.e., $Z_{ij} \sim \text{Bin}(k_{ij}, p_{ij})$, and for simplicity, we set $Z_{ii} = 0$. We also make the following assumption on the pairwise comparison model.

Assumption 1 (Dynamic Range): We assume that there is a constant $\delta \in (0, 1)$ such that for all $i, j \in [n]$,

$$\frac{\delta}{1+\delta} \leq p_{ij} \leq \frac{1}{1+\delta}. \quad (4)$$

Goal. Given the observations \mathcal{Z} of a tournament as defined above, our objective is to determine whether the underlying pairwise comparison model is a BTL model. This corresponds to solving a *composite hypothesis testing* problem:

$$\begin{aligned} H_0 : \mathcal{Z} &\sim \text{BTL model for some } \alpha_1, \dots, \alpha_n > 0, \\ H_1 : \mathcal{Z} &\sim \text{pairwise comparison model that is not BTL,} \end{aligned} \quad (5)$$

where the null hypothesis H_0 states that \mathcal{Z} is distributed according to a BTL model, and the alternative hypothesis H_1 states that \mathcal{Z} is distributed according to a general non-BTL pairwise comparison model. To pose this hypothesis testing problem more rigorously, we demonstrate an interesting relation between a BTL model and its canonical Markov chain.

Recall that a Markov chain on the state space $[n]$, defined by the row stochastic matrix $W \in \mathbb{R}^{n \times n}$, is said to be *reversible* if it satisfies the *detailed balance conditions* [38, Section 1.6]:

$$\forall i, j \in [n], i \neq j, \quad \pi_i W_{ij} = \pi_j W_{ji}, \quad (6)$$

where W_{ij} denotes the probability of transitioning from state i to state j , and $\pi = (\pi(1), \dots, \pi(n))$ denotes the invariant distribution of the Markov chain (which always exists). Equivalently, the Markov chain W is reversible if and only if

$$\text{diag}(\pi)W = W^T \text{diag}(\pi). \quad (7)$$

It turns out that there is a tight connection between reversible Markov chains and the BTL model. This is elucidated in the ensuing proposition, cf. [39, Lemma 6], [10].

Proposition 1 (BTL Model and Reversibility): A pairwise comparison model $\{p_{ij} \in (0, 1) : i, j \in [n], i \neq j\}$ is a BTL model if and only if its canonical Markov chain $S \in \mathbb{R}^{n \times n}$ is reversible and $p_{ij} + p_{ji} = 1$ for all $i, j \in [n]$.

Proof: We provide a proof for completeness. If the pairwise comparison model is BTL, it implies that for some weight vector $\alpha \in \mathbb{R}_+^n$, the pairwise comparison matrix P is given by

$$p_{ij} = \frac{\alpha_j}{\alpha_i + \alpha_j}$$

for $i \neq j$. It is easy to verify that $\pi \triangleq (\sum_{i=1}^n \alpha_i)^{-1} [\alpha_1 \dots \alpha_n]^T$ is the stationary distribution of canonical Markov chain matrix S corresponding to P . Moreover, S is reversible as

$$\pi_i S_{ij} = \frac{\alpha_i}{\sum_{i=1}^n \alpha_i} \times \frac{\alpha_j}{n(\alpha_i + \alpha_j)} = \pi_j S_{ji}.$$

For the converse, since $p_{ij} > 0$ for all $i, j \in [n]$, S is irreducible and has a unique stationary distribution π . By reversibility of S , we have for all $i \neq j$,

$$\pi_i S_{ij} = \pi_j S_{ji} \implies \pi_i p_{ij} = \pi_j p_{ji} \implies p_{ij} = \frac{\pi_j}{\pi_i + \pi_j},$$

where last step follows from the fact that $p_{ij} + p_{ji} = 1$. Thus, P corresponds to a BTL model with weight vector π . \square

Let π be the stationary distribution of the canonical Markov chain matrix S corresponding to a valid pairwise comparison matrix P . By Proposition 1, any pairwise comparison matrix P is BTL if and only if it satisfies the reversibility condition $\Pi P = P^T \Pi$, where $\Pi = \text{diag}(\pi)$, and translated skew-symmetry $P + P^T = \mathbf{1}_n \mathbf{1}_n^T$. It turns out that both conditions are elegantly captured by the matrix $\Pi P + P^T \Pi - \mathbf{1}_n \pi^T$ as illustrated in the Proposition 2, and we will later use the norm of this matrix to quantify the deviation of a pairwise comparison matrix from being BTL.

Proposition 2 (Orthogonal Decomposition): For any pairwise comparison matrix $P \in \mathbb{R}^{n \times n}$ and vector $\pi \in \mathbb{R}^n$ with strictly positive entries, we have

$$\begin{aligned} \|\Pi P + P^T \Pi - \mathbf{1}_n \pi^T\|_{\pi^{-1}, F}^2 &= \|\Pi P - P^T \Pi\|_{\pi^{-1}, F}^2 \\ &\quad + \|P + P^T - \mathbf{1}_n \mathbf{1}_n^T\|_{\pi, F}^2, \end{aligned}$$

where $\Pi = \text{diag}(\pi)$.

The proof can be found in the extended version of this paper [40]. It is important to note here that Assumption 1 and the *Perron-Frobenius theorem* [41, Chapter 8] imply that $\pi_i > 0$ for all $i \in [n]$. Hence, the norm $\|\cdot\|_{\pi^{-1}, F}$ is always well-defined.

Hypothesis testing problem. For a given tolerance parameter $\epsilon > 0$, we can formulate the hypothesis testing problem in (5) as:

$$\begin{aligned} H_0 : \quad &\Pi P + P^T \Pi = \mathbf{1}_n \pi^T, \\ H_1 : \quad &\frac{1}{n} \|\Pi P + P^T \Pi - \mathbf{1}_n \pi^T\|_F \geq \epsilon \|\pi\|_\infty, \end{aligned} \quad (8)$$

where $\Pi = \text{diag}(\pi)$ and π is the stationary distribution of the canonical Markov chain matrix S corresponding to P . Proposition 3, which is proved in [40], verifies that the null hypothesis indeed captures BTL models.

Proposition 3 (BTL Model Characterization): The pairwise comparison matrix P defined in Section II-B corresponds to a BTL model if and only if the hypothesis H_0 in (8) is true.

C. Minimax Risk and Decision Rule

Let ϕ denote a hypothesis test (or decision rule) that maps the consolidated observations $\{Z_{ij}, k_{ij}\}_{i,j \in [n]}$ to $\{0, 1\}$, where 0 represents the null hypothesis and 1 represents the alternative hypothesis. Let \mathbb{P}_{H_0} and \mathbb{P}_{H_1} denote the probability distributions of the input variables under H_0 and H_1 , respectively. Let \mathcal{M}_0 and $\mathcal{M}(\epsilon)$ denote the sets of matrices P that satisfy the null and alternative hypotheses in (8), respectively. Now, define the *minimax risk* as

$$\mathcal{R}_m \triangleq \inf_{\phi} \left\{ \sup_{P \in \mathcal{M}_0} \mathbb{P}_{H_0}(\phi = 1) + \sup_{P \in \mathcal{M}(\epsilon)} \mathbb{P}_{H_1}(\phi = 0) \right\}, \quad (9)$$

where the infimum is taken over all $\{0, 1\}$ -valued tests ϕ . Finally, we can define the *critical threshold* of the hypothesis testing problem in (8) as the smallest value of ϵ for which the minimax risk is bounded by $\frac{1}{3}$:

$$\epsilon_c = \inf \left\{ \epsilon > 0 : \mathcal{R}_m \leq \frac{1}{3} \right\}. \quad (10)$$

The constant $\frac{1}{3}$ is arbitrary and can be replaced by any constant in $(0, 1)$.

Formally, our objective is to characterize the scaling of the critical radius with respect to n . To this end, we consider a hypothesis test which takes the consolidated observations $\{Z_{ij}, k_{ij}\}_{i,j \in [n]}$ as input and evaluates the following expression:

$$T = \sum_{i=1}^n \sum_{j=1}^n (\hat{\pi}_i + \hat{\pi}_j)^2 \frac{Z_{ij}(Z_{ij} - 1)}{k_{ij}(k_{ij} - 1)} + \hat{\pi}_j^2 - 2\hat{\pi}_j(\hat{\pi}_i + \hat{\pi}_j) \frac{Z_{ij}}{k_{ij}} \quad (11)$$

where $\hat{\pi}$ denotes the stationary distribution (choosing one arbitrarily if there are several) of the empirical Markov chain matrix $\hat{S} \in \mathbb{R}^{n \times n}$ defined via

$$\hat{S}_{ij} \triangleq \begin{cases} \frac{Z_{ij}}{k_{ij}n}, & i \neq j, \\ 1 - \frac{1}{n} \sum_{u: u \neq i} \frac{Z_{iu}}{k_{iu}}, & i = j. \end{cases} \quad (12)$$

The alternative hypothesis is selected if $T > \gamma/n$ for some appropriately chosen constant γ independent of n . The test has constructed such that if $\hat{\pi} = \pi$ then $\mathbb{E}[T] = \|\Pi P + P\Pi - \mathbf{1}_n \pi^\top\|_F^2$, i.e., we “plug-in” $\hat{\pi}$ in an unbiased estimator of $\|\Pi P + P\Pi - \mathbf{1}_n \pi^\top\|_F^2$.

III. UPPER BOUND ON CRITICAL THRESHOLD

The ensuing theorem establishes an upper bound on the critical radius of the hypothesis testing problem for the BTL model. For simplicity of analysis, we will assume that $k_{ij} = k_{ji} = k$ for all $i, j \in [n]$ throughout the sequel.

Theorem 1 (Upper Bound on ϵ_c): Consider the hypothesis testing problem in (8), and suppose the number of comparisons per pair of agents satisfies $k \geq \max\{2, \frac{36C^2 \log(n)}{n\delta^4}\}$ for some constant $C > 0$. Then, there exists another constant $c > 0$ such that for any $\epsilon > 0$ with $\epsilon^2 \geq \frac{c}{n}$, we have $\mathcal{R}_m \leq \frac{1}{3}$. Hence, we obtain the bound

$$\epsilon_c^2 \leq \frac{c}{n}.$$

The proof is provided in [40].

IV. LOWER BOUND ON CRITICAL THRESHOLD

We now prove an information theoretic lower bound on the critical radius for the BTL hypothesis testing problem, thus proving the minimax optimality of the scaling provided in the upper bound (up to constant factors).

Theorem 2 (Lower Bound on ϵ_c): Consider the hypothesis testing problem in (8). Then, there exists a constant $c > 0$ such that the critical radius ϵ_c is lower bounded as

$$\epsilon_c^2 \geq \frac{c}{kn}.$$

Proof: We will use the *Ingster-Suslina method* for constructing a lower bound on the critical radius [42]. The method is similar to the well-known *Le Cam’s method*, but it establishes a minimax lower bound by considering a point and a mixture on the parameter space instead of just two points. (Although Le Cam’s method could also be used for this proof in principle, the Ingster-Suslina method greatly simplifies the calculations to bound total variation distance in our setting.)

Under the null hypothesis, we assume that the pairwise comparison matrix P is fixed to be an all $1/2$ matrix, i.e.,

$$H_0 : P = P_0 \triangleq \frac{1}{2} \mathbf{1}_n \mathbf{1}_n^\top. \quad (13)$$

We will denote the distribution corresponding to the pairwise comparison matrix P_0 by \mathbb{P}_0 . Moreover, note that under H_0 , the stationary distribution of the canonical Markov chain matrix S is uniform, i.e., $\pi = \frac{1}{n} \mathbf{1}_n$. Under the alternative hypothesis, we assume that the pairwise comparison matrix P_θ is generated by sampling the parameter θ uniformly from the set Θ , i.e.,

$$H_1 : P = P_\theta \text{ and } \theta \sim \text{Unif}(\Theta), \quad (14)$$

and for some $\theta \in \Theta$, P_θ is given by

$$P_\theta = \begin{bmatrix} \frac{1}{2} \mathbf{1}_{n/2} \mathbf{1}_{n/2}^\top & \frac{1}{2} \mathbf{1}_{n/2} \mathbf{1}_{n/2}^\top + \eta Q_\theta \\ \frac{1}{2} \mathbf{1}_{n/2} \mathbf{1}_{n/2}^\top - \eta Q_\theta & \frac{1}{2} \mathbf{1}_{n/2} \mathbf{1}_{n/2}^\top \end{bmatrix}, \quad (15)$$

where Θ is set of all permutation matrices, and Q_θ is the $n/2 \times n/2$ permutation matrix corresponding to the permutation θ . Let \mathbb{P}_Θ denote the distribution corresponding to the pairwise comparison matrix P_θ . The construction of this mixture was inspired by [36]. However, there are two notable differences. Firstly, the problem in [36] is distinguishing whether two sets of data samples consisting of pairwise comparisons are coming from the same underlying distribution or two different distributions described by a pairwise comparison model. In contrast, our work tests whether or not a single dataset is sampled from a BTL model. Secondly, the manner in which a notion of distance is used to define the deviation of the given data from the null hypothesis is very different in the two works.

Let S_θ denote the canonical Markov chain matrix corresponding to P_θ . It is straightforward to verify that the stationary distribution of S_θ is independent of the permutation θ . Let π denote the stationary distribution of S_θ . By the symmetry of S_θ , the set of first $n/2$ elements, and respectively, last $n/2$ elements, of π are equal, i.e., $\pi_1 = \dots = \pi_{n/2} \triangleq x$ and $\pi_{(n/2)+1} = \dots = \pi_n \triangleq y$. Now x and y can be determined by solving the set of linear equations:

$$\pi^\top = \pi^\top S_\theta \text{ and } \sum_{i=1}^n \pi_i = 1.$$

Solving these equations gives

$$x = \frac{1}{n} \left(1 - \frac{4\eta}{n} \right) \text{ and } y = \frac{1}{n} \left(1 + \frac{4\eta}{n} \right).$$

It is also easy to verify that the deviation from BTL $\|\Pi P_\theta + P_\theta \Pi - \mathbf{1}_n \pi^T\|_F$ is independent of the permutation θ and is given by

$$\begin{aligned} \|\Pi P_\theta + P_\theta^T \Pi - \mathbf{1}_n \pi^T\|_F^2 &= \frac{n}{2} \left((x+y) \left(\frac{1}{2} + \eta \right) - y \right)^2 \\ &\quad + \frac{n}{2} \left((x+y) \left(\frac{1}{2} - \eta \right) - x \right)^2 \\ &\quad + \frac{n}{2} \left(\frac{n}{2} - 1 \right) \left(\frac{x+y}{2} - y \right)^2 + \frac{n}{2} \left(\frac{n}{2} - 1 \right) \left(\frac{x+y}{2} - x \right)^2 \\ &= \frac{2\eta^2}{n} \left(1 - \frac{2}{n} \right)^2 + \frac{2\eta^2}{n^2} \left(1 - \frac{2}{n} \right). \end{aligned} \quad (16)$$

Let $\epsilon = \|\Pi P_\theta + P_\theta \Pi - \mathbf{1}_n \pi^T\|_F / (n \|\pi\|_\infty)$ to ensure that P_θ 's satisfy the condition of the alternative hypothesis in (8). Substituting the values of $\|\pi\|_\infty = y$ and $\|\Pi P_\theta + P_\theta^T \Pi - \mathbf{1}_n \pi^T\|_F$ implies that $\epsilon \leq C\eta/\sqrt{n}$ for some constant $C > 0$.

Now, the Ingster-Suslina method [42] states that

$$\mathcal{R}_m \geq \frac{1}{2} \left(1 - \sqrt{\chi^2(\mathbb{P}_0 \| \mathbb{P}_\Theta)} \right),$$

where $\chi^2(\cdot \| \cdot)$ denotes the χ^2 -divergence. Similar to the analysis in [36, Theorem 3], for the distributions \mathbb{P}_0 and \mathbb{P}_Θ , if we have $\eta^2 \leq \frac{c}{k}$ for some constant c independent of n, k , then we can upper bound the $\chi^2(\mathbb{P}_0 \| \mathbb{P}_\Theta)$ term as $\chi^2(\mathbb{P}_0 \| \mathbb{P}_\Theta) \leq \frac{1}{9}$. Hence, we have shown that there exists a constant $c > 0$, such that if $n\epsilon^2 \leq c/k$, then $\chi^2(\mathbb{P}_0 \| \mathbb{P}_\Theta) \leq \frac{1}{9}$, which implies that the minimax risk $\mathcal{R}_m \geq \frac{1}{3}$. Hence, $\epsilon_c^2 \geq c/(kn)$ as desired. \square

We defer readers to [40] for a more detailed proof.

V. STABILITY OF THE BTL ASSUMPTION

In this section, we analyze the stability of the BTL assumption in the context of rankings. The BTL ranking orders agents based on the stationary distribution π of the canonical Markov-chain matrix. Meanwhile, the Borda count ranking is more general, as it doesn't rely on the BTL assumption, and instead is based on Borda scores [27] (defined below). If the BTL assumption holds, then Borda ranking equals the BTL ranking. Our goal is to determine the size of the deviation from the BTL condition in a pairwise comparison model that is sufficient to produce a discrepancy between the BTL and Borda count rankings. For this section, we will consider the symmetric setting in which the pairwise comparison matrix P satisfies, $p_{ij} + p_{ji} = 1$. Define the *Borda count* $\tau_i(P)$ of an agent $i \in [n]$ as the (scaled) probability that i beats any other agent selected uniformly at random [27]:

$$\tau_i(P) \triangleq \sum_{j=1}^n (1 - p_{ij}). \quad (17)$$

Now we show that the stability of the BTL assumption decreases as n grows meaning that even smaller deviations from the BTL condition can lead to a reversed ranking.

Proposition 4 (Stability of BTL Assumption): There exists a pairwise comparison matrices $P \in \mathbb{R}^{n \times n}$ such that two rows $i, j \in [n]$ having constant $\Delta\tau \triangleq \tau_i(P) - \tau_j(P) > 0$, but has an

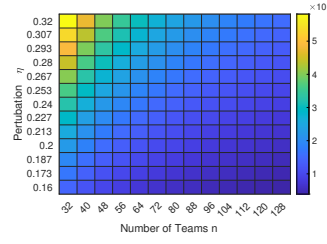


Fig. 1. Plot of empirical average of T under hypothesis H_1 .

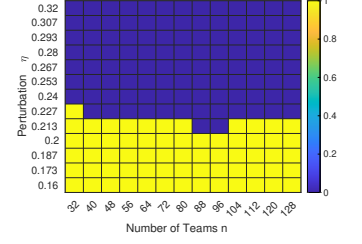


Fig. 2. Plot of $\mathbb{1}_{\mathcal{R}_m > 4/5}$.

opposite ranking BTL ranking, i.e., $\pi_i < \pi_j$. Let $\Pi = \text{diag}(\pi)$. Moreover, the deviation of P from BTL condition decays as

$$\|\Pi P + P \Pi - \mathbf{1}_n \pi^T\|_F \leq \frac{c}{\sqrt{n}}$$

for some constant $c > 0$.

The proof is provided in [40]. Proposition 4 highlights that the BTL assumption may potentially give a wrong ranking when the underlying pairwise matrix is $O(1/\sqrt{n})$ “distance” away from the BTL condition. Interestingly, this $O(1/\sqrt{n})$ deviation coincides with the critical threshold for the BTL testing problem (up to constant factors). It would be interesting to further explore the stability of the BTL assumption in the context of rankings.

VI. NUMERICAL SIMULATIONS

In this section, we will empirically analyze the behavior of the minimax risk and the deviation $\|\Pi P + P \Pi - \mathbf{1}_n \pi^T\|_F$ via a synthetic experiment. We will use the same construction for the pairwise comparison matrix P that we utilized in (2) under the null and alternate hypothesis which are presented in (13) and (14). We set the number of pairwise comparisons per pair of agents $k = 12$, the number of agents n is linearly increased from 32 to 128, and the perturbation η in (15) is increased from 0.16 to 0.32. Simulations are performed for each value of η and n , and the corresponding value of expected values of test T under hypothesis H_1 and *minimax risk* \mathcal{R}_m is estimated. The threshold used for the decision rule is set to η^2/n . Fig. 2 plots the empirical average of T under H_1 and $\mathbb{1}_{\mathcal{R}_m > 4/5}$ for different values of η and n . Note that the behavior of empirical average of T (under H_1) is consistent with (16) and moreover for a fixed value of η the behavior of \mathcal{R}_m is independent of n .

VII. CONCLUSION

In this work, we studied the problem of testing whether a BTL model accurately represents the data generated from an underlying pairwise comparison model. We developed a rigorous hypothesis test (8) for this purpose and established that the (minimax) critical threshold for this test is $\epsilon_c = \Theta(1/\sqrt{n})$. We also took the first steps towards rigorously characterizing the stability of the BTL assumption for rankings. Our results indicated that the BTL assumption may lead to different rankings if the pairwise comparison matrix deviates from the BTL condition by $O(1/\sqrt{n})$.

REFERENCES

- [1] R. A. Bradley and M. E. Terry, "Rank analysis of incomplete block designs. I. The method of paired comparisons," *Biometrika*, vol. 39, no. 3/4, pp. 324–345, December 1952.
- [2] R. D. Luce, *Individual Choice Behavior: A Theoretical Analysis*. New York, NY, USA: John Wiley & Sons Inc., 1959.
- [3] L. L. Thurstone, "A law of comparative judgment," *Psychological Review*, vol. 34, no. 4, pp. 273–286, 1927.
- [4] R. L. Plackett, "The analysis of permutations," *Journal of the Royal Statistical Society, Series C (Applied Statistics)*, vol. 24, no. 2, pp. 193–202, 1975.
- [5] D. McFadden, "Conditional logit analysis of qualitative choice behavior," in *Frontiers in Econometrics*, ser. Economic Theory and Mathematical Economics, P. Zarembka, Ed. New York, NY, USA: Academic Press, 1973, pp. 105–142.
- [6] E. Zermelo, "Die berechnung der turnier-ergebnisse als ein maximumproblem der wahrscheinlichkeitsrechnung," *Mathematische Zeitschrift*, vol. 29, no. 1, pp. 436–460, December 1929.
- [7] D. R. Hunter, "MM algorithms for generalized Bradley-Terry models," *The Annals of Statistics*, vol. 32, no. 1, pp. 384–406, February 2004.
- [8] J. Guiver and E. Snelson, "Bayesian inference for Plackett-Luce ranking models," in *Proceedings of the 26th Annual International Conference on Machine Learning (ICML)*, Montreal, QC, Canada, June 14–18 2009, pp. 377–384.
- [9] F. Caron and A. Doucet, "Efficient Bayesian inference for generalized Bradley-Terry models," *Journal of Computational and Graphical Statistics*, vol. 21, no. 1, pp. 174–196, March 2012.
- [10] Sahand Negahban and Sewoong Oh and Devavrat Shah, "Rank centrality: Ranking from pairwise comparisons," *Operations Research, INFORMS*, vol. 65, no. 1, pp. 266–287, January-February 2017.
- [11] Y. Chen, J. Fan, C. Ma, and K. Wang, "Spectral method and regularized MLE are both optimal for top- K ranking," *The Annals of Statistics*, vol. 47, no. 4, pp. 2204–2235, 2019.
- [12] H. Bong, W. Li, S. Shrotriya, and A. Rinaldo, "Nonparametric estimation in the dynamic Bradley-Terry model," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2 2020, pp. 3317–3326.
- [13] J. Hendrickx, A. Olshevsky, and V. Saligrama, "Minimax rate for learning from pairwise comparisons in the BTL model," in *Proceedings of the 37th Annual International Conference on Machine Learning (ICML), Proceedings of Machine Learning Research (PMLR)*, vol. 119, Vienna, Austria, July 13–18 2020, pp. 4193–4202.
- [14] A. Jadbabaie, A. Makur, and D. Shah, "Estimation of skill distribution from a tournament," in *Proceedings of the Advances in Neural Information Processing Systems 33 (NeurIPS)*, Vancouver, BC, Canada, December 6–12 2020, pp. 8418–8429.
- [15] A. Jadbabaie, A. Makur, and D. Shah, "Estimation of skill distributions," June 2020, arXiv:2006.08189 [stat.ML]. [Online]. Available: <https://arxiv.org/abs/2006.08189>
- [16] D. Davidson and J. Marschak, *Experimental tests of a stochastic decision theory*. Springer Netherlands, 1959, vol. 17, p. 274.
- [17] D. H. McLaughlin and R. D. Luce, "Stochastic transitivity and cancellation of preferences between bitter-sweet solutions," *Psychonomic Science*, vol. 2, pp. 89–90, 1 1965.
- [18] A. Tversky, "Elimination by aspects: A theory of choice," *Psychological Review*, vol. 79, pp. 281–299, 7 1972. [Online]. Available: [/record/1973-00249-001](https://doi.org/10.1037/0033-2909.79.4.281)
- [19] B. Morley and D. Thomas, "An investigation of home advantage and other factors affecting outcomes in english one-day cricket matches," *Journal of Sports Sciences*, vol. 23, pp. 261–268, 3 2005.
- [20] S. R. Clarke and J. M. Norman, "Home ground advantage of individual clubs in english soccer," *Journal of the Royal Statistical Society. Series D (The Statistician)*, vol. 44, pp. 509–521, 1995.
- [21] A. Agresti, *Categorical Data Analysis*, ser. Wiley Series in Probability and Statistics. John Wiley & Sons, Inc., 7 2002.
- [22] N. B. Shah and M. J. Wainwright, "Simple, robust and optimal ranking from pairwise comparisons," *Journal of Machine Learning Research*, vol. 18, pp. 1–38, 12 2018.
- [23] R. Heckel, N. B. Shah, K. Ramchandran, and M. J. Wainwright, "Active ranking from pairwise comparisons and when parametric assumptions do not help," *The Journal of Machine Learning Research*, vol. 47, pp. 3099–3126, 12 2019.
- [24] S. Chatterjee, "Matrix estimation by universal singular value thresholding," *The Annals of Statistics*, vol. 43, no. 1, pp. 177–214, February 2015.
- [25] N. B. Shah, S. Balakrishnan, J. Bradley, A. Parekh, K. Ramchandran, and M. J. Wainwright, "Estimation from pairwise comparisons: Sharp minimax bounds with topology dependence," *Journal of Machine Learning Research*, vol. 17, no. 58, pp. 1–47, 2016.
- [26] S. Negahban, S. Oh, and D. Shah, "Iterative ranking from pair-wise comparisons," in *Proceedings of the Advances in Neural Information Processing Systems 25 (NIPS)*, Lake Tahoe, NV, USA, December 3–8 2012, pp. 2474–2482.
- [27] J. Borda, "Memoire sur les elections au scrutin.(1781)," *Histoire de l'Académie Royale des Sciences*, 1781.
- [28] P. Diaconis, *Group Representations in Probability and Statistics*, ser. Lecture Notes-Monograph Series, S. S. Gupta, Ed. Hayward, CA, USA: Institute of Mathematical Statistics, 1988, vol. 11.
- [29] D. Shah, "Computing choice: Learning distribution over permutations," *Cambridge University Press Bulletin*, pp. 1–40, October 2019.
- [30] J. I. Yellott, Jr., "The relationship between Luce's choice axiom, Thurstone's theory of comparative judgment, and the double exponential distribution," *Journal of Mathematical Psychology*, vol. 15, no. 2, pp. 109–144, April 1977.
- [31] G. Simons and Y.-C. Yao, "Asymptotics when the number of parameters tends to infinity in the Bradley-Terry model for paired comparisons," *The Annals of Statistics*, vol. 27, no. 3, pp. 1041–1060, June 1999.
- [32] S. Negahban, S. Oh, K. K. Thekumparampil, and J. Xu, "Learning from comparisons and choices," *Journal of Machine Learning Research*, vol. 19, no. 40, pp. 1–95, August 2018.
- [33] E. L. Lehmann, J. P. Romano, S. N. York, and J. Steinebach, *E. L. Lehmann, J. P. Romano: Testing statistical hypotheses*. Springer, 8 2006, vol. 64.
- [34] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola, "A kernel two-sample test," *Journal of Machine Learning Research*, vol. 13, p. 723–773, Mar. 2012.
- [35] Y. I. Ingster, "Minimax detection of a signal in l_p metrics," *Journal of Mathematical Sciences*, vol. 68, pp. 503–515, 2 1994.
- [36] C. Rastogi, S. Balakrishnan, N. B. Shah, and A. Singh, "Two-sample testing on ranked preference data and the role of modeling assumptions," *Journal of Machine Learning Research*, vol. 23, no. 225, pp. 1–48, 2022. [Online]. Available: <http://jmlr.org/papers/v23/20-1304.html>
- [37] A. Seshadri and J. Ugander, "Fundamental limits of testing the independence of irrelevant alternatives in discrete choice," *ACM EC 2019 - Proceedings of the 2019 ACM Conference on Economics and Computation*, pp. 65–66, 1 2020.
- [38] D. A. Levin, Y. Peres, and E. L. Wilmer, *Markov Chains and Mixing Times*, 1st ed. Providence, RI, USA: American Mathematical Society, 2009.
- [39] A. Rajkumar and S. Agarwal, "A statistical convergence perspective of algorithms for rank aggregation from pairwise data," in *International conference on machine learning*. PMLR, 1 2014, pp. 118–126.
- [40] A. Makur and J. Singh, "Hypothesis testing for the Bradley-Terry-Luce model," 2023, in preparation.
- [41] C. D. Meyer, *Matrix Analysis and Linear Algebra*, 2000.
- [42] Y. Ingster and I. Suslina, *Nonparametric Goodness-of-Fit Testing Under Gaussian Models*. Springer Science & Business Media, 01 2003, vol. 169.