# **Choice Bandits**

#### **Arpit Agarwal**

University of Pennsylvania Philadelphia, PA 19104, USA aarpit@seas.upenn.edu

#### Nicholas Johnson\*

University of Minnesota Minneapolis, MN 55455, USA njohnson@cs.umn.edu

#### Shivani Agarwal

University of Pennsylvania Philadelphia, PA 19104, USA ashivani@seas.upenn.edu

## **Abstract**

There has been much interest in recent years in the problem of dueling bandits, where on each round the learner plays a pair of arms and receives as feedback the outcome of a relative pairwise comparison between them. Here we study a natural generalization, that we term *choice bandits*, where the learner plays a set of up to k > 2 arms, and receives limited relative feedback in the form of a single multiway choice among the pulled arms, drawn from an underlying multiway choice model. We study choice bandits under a very general class of choice models that is characterized by the existence of a unique 'best' arm (which we term generalized Condorcet winner), and includes as special cases the well-studied multinomial logit (MNL) and multinomial probit (MNP) choice models, and more generally, the class of random utility models with i.i.d. noise (IID-RUMs). We propose an algorithm for choice bandits, termed Winner Beats All (WBA), with a distribution dependent  $O(\log T)$  regret bound under all these choice models. The challenge in our setting is that the decision space is  $\Theta(n^k)$ , which is large for even moderate k. Our algorithm addresses this challenge by extracting just  $O(n^2)$ statistics from multiway choices and exploiting the existence of a unique 'best' arm to find arms that are competitive to this arm in order to construct sets with low regret. Since these statistics are extracted from the same choice observations, one needs a careful martingale analysis in order to show that these statistics are concentrated. We complement our upper bound result with a lower bound result, which shows that our upper bound is order-wise optimal. Our experiments demonstrate that for the special case of k=2, our algorithm is competitive with previous dueling bandit algorithms, and for the more general case k > 2, outperforms the recently proposed MaxMinUCB algorithm designed for the MNL model.

#### 1 Introduction

The dueling bandit problem has received a lot of interest in recent years [1] [2] [3] [4] [5] [6] [7] [8] [9] [10] [11] [12] [13] [14]. Here there are n arms  $\{1,\ldots,n\}$ ; on each trial t, the learner pulls a pair of arms  $(i_t,j_t)$ , and receives relative feedback indicating which of the two arms has a better quality/reward. In the regret minimization setting, the goal is to identify the 'best' arm(s) while also minimizing the regret due to playing sub-optimal arms in the learning (exploration) phase.

<sup>\*</sup>Work done while at the University of Pennsylvania.

In many applications, however, it can be natural for the learner to pull more than two arms at a time, and seek relative feedback among them. For example, in recommender systems, it is natural to display several items or products to a user, and seek feedback on the most preferred item among those shown. In online advertising, it is natural to display several ads at a time, and observe which of them is clicked (preferred). In online ranker evaluation for information retrieval, one can easily imagine a generalization of the setting studied by Yue & Joachims [15], where one may want to "multi-leave" several rankers at a time to help identify the best ranking system while also presenting good/acceptable results to users using the system during the exploration phase. In general, there is also support in the marketing literature for showing customers more than two items at a time [16].

Motivated by such applications, we consider a framework that generalizes the dueling bandit problem to allow the learner to pull more than two arms at a time. Here, on each trial t, the learner pulls a set  $S_t$  of up to k arms (for fixed  $k \in \{2, \ldots, n\}$ ), and receives relative feedback in the form of a multiway choice  $y_t \in S_t$  indicating which arm in the set has the highest quality/reward. The goal of the learner is again to identify a 'best' arm (to be formalized below) while minimizing a suitable notion of regret that penalizes the learner for playing sub-optimal arms during the exploration phase. We term the resulting framework *choice bandits*.

In the (stochastic) dueling bandits framework, the underlying probabilistic model from which feedback is observed is a *pairwise comparison model*, which for each pair of arms (i,j), defines a probability  $P_{ij}$  that arm i has higher reward/quality than arm j. In our choice bandits framework, the underlying probabilistic model is a *multiway choice model*, which for each set of arms  $S \subseteq [n]$  with  $|S| \le k$  and each arm  $i \in S$ , defines a probability  $P_{i|S}$  that arm i has the highest reward/quality in the set S.

We study choice bandits under a new class of choice models, that are characterized by the existence of a unique *generalized Condorcet winner* (GCW), which we define to be an arm that has larger probability of being chosen than any other arm in any choice set. This class includes as special cases the well-studied multinomial logit (MNL) [17] [18] [19] and multinomial probit (MNP) [20] choice models, and more generally, the class of random utility models with i.i.d. noise (IID-RUMs) [21] [22].

Our main contribution is a computationally efficient algorithm, termed Winner Beats All (WBA), that achieves a distribution dependent  $O(n^2 \log n + n \log T)$  regret bound under any choice model that exhibits a unique GCW, where T is the time-horizon. We complement our upper bound result with an order-wise lower bound of  $\Omega(n \log T)$  for any no-regret algorithm, showing that our algorithm has asymptotically order optimal regret under our general class of choice models. If the underlying model is MNL, then WBA achieves an instance-wise asymptotically optimal regret bound, which is better than the regret bound for the recent MaxMinUCB algorithm under MNL [23].

The main challenge in designing an algorithm under our framework is that the space of exploration (number of possible sets the learner can play) is  $\Theta(n^k)$  which is large even for moderate k. Therefore, it can be challenging to simultaneously *explore/learn* the choice sets with low regret out of the possible  $\Theta(n^k)$  sets and *exploit* these low regret sets. We overcome these challenges by extracting just  $O(n^2)$  pairwise statistics from the observed multiway choices under different sets, and using these statistics to find choice sets with low regret. Since these pairwise statistics are extracted from multiway choices under different sets, a technical challenge is to show that these statistics are concentrated. We resolve this challenge by using a novel coupling argument that couples the stochastic process generating choices with another stochastic process, and showing that pairwise estimates according to this other process are concentrated. We believe that our results for efficient learning under this large class of choice models that is considerably more general than the MNL class are of independent interest.

We also run experiments on several synthetic and real-world datasets. Our experiments on these datasets show that our algorithm for the special case of k=2 is competitive as compared to previous dueling bandit algorithms, even though it is designed for a more general setting. For the case of k>2, we compare our algorithm with the MaxMinUCB algorithm of [23] which was designed for the MNL model. We observe that our algorithm performs better in terms of regret than MaxMinUCB under all datasets (even under synthetic MNL datasets). We further observe that under several datasets the regret achieved by our algorithm for k>2 is better than the regret for k=2.

**Related Work.** There has been some recent interest in bandit settings where more than two arms are pulled at a time, although no work that we are aware of considers the types of general choice models that we do. (1) A related setting to ours is that of *multi-dueling bandits* [24, 25, 26], where the learner also pulls a set  $S_t$  of k items; however, the feedback received by the learner is assumed

to be drawn from a pairwise comparison model (in particular, the learner observes some subset of the  $\binom{k}{2}$  possible pairwise comparisons among items in  $S_t$ ). In contrast, in our choice bandits setting, the learner receives the outcome of a direct multiway choice among the items in  $S_t$ , generated from a multiway choice model. (2) In *combinatorial bandit with relative feedback* [23], the learner pulls a set  $S_t$  of up to k arms, and observes top-m ordered feedback drawn according to the MNL model, for some m < k. In contrast, we only observe the (top-1) choice feedback from the set  $S_t$ that is played, but, we study a much more general class of choice models than the MNL model. (3) Another related setting is that of battling bandits [27], where the learner pulls a multiset of k arms and receives feedback indicating which arm was chosen. However, their setting considers a specific pairwise-subset (PS) choice model that is defined in terms of a pairwise comparison model, whereas we consider more general choice models. (4) In stochastic click bandits [28], the learner pulls an *ordered* set of k arms/documents, and observes *clicks* on a subset of these documents, drawn according to an underlying *click model* which is a probabilistic model for click generation over an ordered set. However, click models in their setting are different than choice models in our setting, and neither can be cast as a special case of the other. (5) Another related setting is that of best-of-k bandits [29], where again the learner pulls a set  $S_t$  of k arms. Of the various types of feedback considered in [29], the marked bandit feedback corresponds to the type of feedback that we study, however, the choice models studied in [29] correspond only to a subclass of random utility choice models, and moreover, the analysis in [29] is in the PAC/pure exploration setting, while ours is in the regret minimization setting. (6) Other recent work has specifically considered active learning problems, either in the context of dynamic assortment optimization under MNL where the goal is to maximize expected revenue [30, 31, 32, 33, 34]; or in the context of best arm(s) identification under MNL or IID-RUMs [35] 36 in a PAC/pure exploration setting. (7) Finally, we also mention *combinatorial* bandits, which have a different goal but also involve pulling subsets of arms [37] 38, 39, 40]. See the supplementary material for more detailed discussion.

Organization. We set up the choice bandits problem in Section 2. We give our lower bound result in Section 3. We present our algorithm in Section 4. and its regret analysis in Section 5. We give experimental results in Section 6. We finally conclude with a brief discussion in Section 7. All the proofs can be found in the supplementary material.

## 2 Problem Setup and Preliminaries

In the choice bandits problem, there are n arms  $[n]:=\{1,\ldots,n\}$ , and a set size parameter  $2\leq k\leq n$ . On each trial t, the learner pulls (selects/plays) a choice set  $S_t\subseteq [n]$  of up to k arms, i.e. with  $|S_t|\leq k$ , and receives as feedback  $y_t\in S_t$ , indicating the arm that is most preferred in  $S_t$ . We assume the feedback  $y_t$  is generated probabilistically from an underlying  $\mathit{multiway choice model}$ , which defines for each  $S\subseteq [n]$  such that  $|S|\leq k$ , and arm  $i\in S$ , a  $\mathit{choice probability } P_{i|S}$  which corresponds to the probability that arm i is the most preferred arm in S. Before defining appropriate notions of 'best' arm and regret for the learner we will give some examples of multiway choice models.

Random utility models with i.i.d. noise (IID-RUMs). IID-RUMs are a well-known class of choice models that have origins in the econometrics and marketing literature [21] [41]. Under an IID-RUM, the (random) utility associated with arm  $i \in [n]$  is given by  $U_i = v_i + \epsilon_i$  where  $v_i \in \mathbb{R}$  is a deterministic utility and  $\epsilon_i \in \mathbb{R}$  is the noise drawn i.i.d. from a distribution  $\mathcal{D}$  over reals. For a set S, the probability of choosing  $i \in S$  is given by  $P_{i|S} = \Pr\left(U_i > U_j, \forall j \in S \setminus \{i\}\right)$ . We will sometimes also refer to  $v_i$  as the weight of item i. Under any IID-RUM if  $v_i > v_j$  for some  $i, j \in [n]$  then arm i will be more likely to be chosen than arm j in any set. The IID-RUM class contains some popular models, such as the multinomial logit (MNL) [17] [18] [19], and multinomial probit (MNP) [20].

**Example 1** (MNL). Under MNL, the noise distribution  $\mathcal{D}$  is a Gumbel(0,1) and the probability  $P_{i|S}$  of choosing an item i from a set S has the following closed form expression:  $P_{i|S} := e^{v_i}/(\sum_{j \in S} e^{v_j})$ . It is clear from this expression that arms with higher weights are more likely to be chosen.

**Example 2** (MNP). Under the MNP model, the noise distribution  $\mathcal{D}$  is the standard Normal distribution  $\mathcal{N}(0,1)$ , however, unlike the MNL there is no closed form expression for the choice probabilities.

Under IID-RUMs there is a clear notion of 'best' arm: an arm that has the highest weight  $\max_{i \in [n]} v_i$ . We now define a strictly more general class of models where there is a clear notion of 'best' arm.

<sup>&</sup>lt;sup>2</sup>Note that for the special case of k=2, our framework reduces to dueling bandits; the pairwise comparison probabilities  $P_{ij} := \Pr(i \succ j)$  in dueling bandits can be viewed as pairwise choice probabilities  $P_{i|\{i,j\}}$ .

A New Class of Choice Models. We introduce a new class of multiway choice models that are characterized by the following condition that requires the existence of a unique 'best' arm.

**Definition 1** (Generalized Condorcet Condition (GCC)). A choice model is said to satisfy the GCC condition if there exists a unique arm  $i^* \in [n]$  such that for every choice set  $S \subseteq [n]$  that contains  $i^*$ , we have  $P_{i^*|S} > P_{j|S}$  for all  $j \in S \setminus \{i^*\}$ .

Intuitively, the above condition requires the existence of a unique arm that is always (stochastically) preferred to all other arms, no matter what other arms are shown with it. This condition is a generalization of the Condorcet condition studied for pairwise comparison models [6] [1]]. Just as the Condorcet condition need not be satisfied for all pairwise comparison models, similarly, GCC need not be satisfied by all multiway choice models. Below we show that the GCC condition is satisfied for all IID-RUMs subject to a minor technical condition.

**Lemma 1** (IID-RUMs satisfy GCC). For any IID-RUM choice model with utility for arm  $i \in [n]$  given by  $U_i = v_i + \epsilon_i$ , the GCC condition is satisfied if  $|\operatorname{argmax}_{i \in [n]} v_i| = 1$ .

In this paper, we study the class of all choice models where the GCC is satisfied. Under GCC, we will refer to this unique 'best' arm as the generalized Condorcet winner (GCW) and denote it by  $i^*$ . Note that for any set S containing the GCW  $i^*$ , we must have  $P_{i^*|S} \geq \frac{1}{|S|}$ .

**Regret Notion.** Similar to dueling bandits, the goal of the learner in our setting is to identify the best arm while also playing good/competitive sets with respect to this arm during the exploration phase Hence, our notion of regret measures the sub-optimality of a choice set S relative to  $i^*$ , and is a generalization of the regret defined by [23] for the special case of MNL choice model. Moreover, under our notion of regret it is optimal to play  $S^* = \{i^*\}$ , i.e. regret of playing  $S^*$  is 0. The regret of a set is defined to be the sum of regret due to individual arms in the set, and the regret for an arm corresponds to the 'margin' by which the best arm  $i^*$  beats this arm. In other words, the regret of an arm corresponds to the *shortfall in preference probability* due to pulling this arm over the 'best' arm.

**Definition 2.** The regret 
$$r(S)$$
 for  $S \subseteq [n]$  is defined as:  $r(S) := \sum_{i \in S} (P_{i^*|S \cup \{i^*\}} - P_{i|S \cup \{i^*\}})$ .

This notion of regret can be interpreted as: r(S) is the sum over all arms  $i \in S$ , the fraction of consumers that will choose  $i^*$  minus the fraction of consumers that will choose i when  $i^*$  is played together with S. It is easy to see that  $r(\{i^*\}) = 0$ , and  $0 \le r(S) \le |S|$  for any set  $S \subseteq [n]$ .

**Example 3.** Consider a choice model where arm 1 is the GCW, and for each set 
$$S$$
 containing arm 1, we have  $P_{1|S} = 0.51$  and  $P_{i|S} = \frac{0.49}{|S|-1} \ \forall i \in S \setminus \{1\}$ . Then  $r(\{1,\ldots,m\}) = 0.51 \times (m-1) - 0.49$ .

In the above example, the regret increases linearly as we increase m. The following gives an example where the arms are much more 'competitive' and regret is smaller.

**Example 4.** Consider the MNL choice model with weights 
$$v_1 = \log(1 + \epsilon)$$
, for  $\epsilon > 0$ , and  $v_2 = \cdots = v_n = 0$ . Then  $r(\{1, \cdots, m\}) = \sum_{i \in S} \frac{e^{v_1} - e^{v_i}}{\sum_{j \in S} e^{v_j}} = \frac{\epsilon(m-1)}{m+\epsilon}$ .

The regret here increases much more slowly in terms of m. Note that our regret is not necessarily well-defined in the dueling bandits setting, due to the need to consider choice probabilities for sets of size 3 even when one plays only sets of size 2. In the supplementary material, we give results for an additional notion of regret that is a direct generalization of the dueling bandit regret, and allows for a more direct comparison between our framework and the dueling bandits framework.

Under the above notion of regret, the goal of an algorithm A is to minimize its cumulative regret over T trials defined as:  $R(T) = \sum_{t=1}^{T} r(S_t)$ .

## 3 A Fundamental Lower Bound

In this section we present a regret lower bound for our choice bandits problem. We say that an algorithm is *strongly consistent* under GCC if its expected regret over T trials is  $o(T^a)$  for any a>0 under any model in this class. Before presenting the lower bound let us define the following distribution dependent quantities.

$$\Delta_{i^*i|S} = \frac{P_{i^*|S} - P_{i|S}}{P_{i^*|S} + P_{i|S}}, \quad \Delta_{\max}^{GCC} := \max_{S:|S| \le k} \max_{i \in S} \Delta_{i^*i|S}, \quad \Delta_{\min}^{GCC} := \min_{S:|S| \le k} \min_{i \in S} \Delta_{i^*i|S}. \quad (3.1)$$

<sup>&</sup>lt;sup>3</sup>Note that we are *not* working in the pure exploration setting, where all sets incur equal cost during exploration.

The following theorem presents a lower bound for any strongly consistent algorithm.

**Theorem 1.** Given a set of arms [n], choice set size bound  $k \leq n$ , parameter  $\Delta \in (0,1)$ , and any strongly consistent algorithm  $\mathcal{A}$  under GCC, there exists a GCC choice model with  $\Delta_{\min}^{\text{GCC}} = \Delta$  such that when choice outcomes are drawn from this model we have

$$\liminf_{T \to \infty} \frac{\mathbf{E}\left[R(T)\right]}{\log T} = \Omega\left(\frac{n-1}{\Delta}\right) \,,$$

where T is the time-horizon. If the underlying model is MNL with parameters  $v_1, v_2, \cdots v_n \in \mathbb{R}$ , then:  $\liminf_{T \to \infty} \frac{\mathbf{E}[R(T)]}{\log T} = \Omega\left(\sum_{i \in [n] \setminus \{i^*\}} \frac{1}{\Delta_{i^*i}^{\text{MNL}}}\right)$  where  $\Delta_{i^*i}^{\text{MNL}} = \frac{e^{v_i *} - e^{v_i}}{e^{v_i *} + e^{v_i}}$ , for  $i \in [n] \setminus \{i^*\}$ .

**Discussion.** The above bound shows that any algorithm for the choice bandits problem needs to incur  $\Omega(n \log T)$  regret in the worst case. Note that the above lower bound does not depend on the choice set size parameter k. If the choices are generated from an underlying MNL model, then the above theorem gives an instance-dependent lower bound for the regret of any algorithm. Note that [23] also provided a lower bound under MNL for our notion of regret, however, their bound depends on the worst-case gap between  $i^*$  and any other arm  $i \neq i^*$ , while we provide a more fine-grained bound under MNL which depends on gaps between  $i^*$  and each individual arm  $i \in [n]$ .

In order to prove the above bound we construct a pair of instances that have different GCW arms, and use the information divergence lemma of [42] in order to characterize the minimum number of samples needed in order to collect the 'information' needed to separate these two instances.

**Remark 1.** In order to prove a lower bound for our choice bandits problem one may also be able to use the lower bound given in [43], by casting our problem as a structured bandit problem. However, the lower bound of [43] is in terms of the solution of a linear program, and one will then need to design a distribution over hard instances in order to quantify the solution of this linear program in terms of the gap parameter  $\Delta_{\min}^{GCC}$ . One of the main novelty of our bound is this construction of hard instances that allows us to quantify the lower bound in terms of  $\Delta_{\min}^{GCC}$ .

## 4 Algorithm

In this section we will present our algorithm for the choice bandits problem, termed Winner Beats All (WBA). WBA divides its execution into rounds and each round can contain multiple trials. We will use r as an global index for a round, and t as an global index for a trial. For each round r, WBA maintains a set  $A_r$  of active arms, which are a set of arms for which the algorithm is still not confident enough that these are 'bad' arms. Note that an arm that is inactive in a particular round, can become active in a later round. We also maintain a set Q that is initialized to being empty at the beginning of each round and keeps track of the arms in  $A_r$  that have been played so far in the round.

Given a trial t belonging to round r, WBA selects a special arm termed the 'anchor' arm, and a set  $S \subseteq A_r \setminus Q$  (arbitrarily) of up to k-1 arms in  $A_r$  that have not been played so far in round r. The set S and  $a_t$  are selected such that  $a_t$  empirically performs better than each arm in S. The set S is then played together with arm  $a_t$  (if |S| < k-1, then other arbitrary arms from  $A_r$  are added to the played set). The anchor arm is updated in every trial and is chosen so that one can *quickly* find evidence that arms in S are not good.

Let  $y_t$  be the feedback received in trial t when  $S_t$  was played including anchor  $a_t$ . For all  $i, j \in [n]$ , let  $N_{ij}(t)$  denote the number of times (up to round t) that either arm i or j was chosen when arm j is the anchor, i.e.  $N_{ij}(t) := \sum_{t'=1}^t \mathbb{1}(a_{t'} = j, \{i, j\} \subseteq S_{t'}, y_{t'} \in \{i, j\})$ .

For each  $i, j \in [n]$  and trial t, such that  $N_{ij}(t) > 0$ , the algorithm maintains an estimate of the marginal probability of arm i beating the arm j as

$$\hat{P}_{ij}(t) := \frac{1}{N_{ij}(t)} \sum_{t=1}^{t} \mathbb{1}(a_{t'} = j, \{i, j\} \subseteq S_{t'}, y_{t'} = i), \qquad (4.1)$$

which is the fraction of times i was selected (compared to j) when both i and j were played together and j was the anchor. (When  $N_{ij}(t)=0$ , we can simply take  $\hat{P}_{ia}(t)$  to be 1/2.) Similar to [44], let us define an *empirical divergence*  $I_i(t,S)$  which provides a certificate that an arm i is worse than (some) arms in S, as  $I_i(t,S)=\sum_{j\in S}\mathbb{1}[\hat{P}_{ij}(t)\leq \frac{1}{2}]\cdot N_{ij}(t)\cdot d(\hat{P}_{ij}(t),\frac{1}{2})$ , where  $d(\hat{P}_{ij},\frac{1}{2})$  is the KL-divergence defined as  $d(P,Q)=P\log(\frac{P}{Q})+(1-P)\log(\frac{1-P}{1-Q})$ , for  $P,Q\in[0,1]$ . If  $I_i(t,S)$  is

0, it means that arm i is empirically at least as good as all other arms in S, and a higher  $I_i(t,S)$  would suggest that arm i is most likely 'bad'. For a constant C, we define the condition  $\mathcal{J}_i(t,C)$  for arm  $i \in [n]$  and round t as  $\mathcal{J}_i(t,C) = \mathbbm{1} \Big\{ \exists S \subseteq [n] : I_i(t,S) \geq |S| \log(nC) + \log(t) \Big\}$ . If  $\mathcal{J}_i(t,C) = 1$  for some i, it means that there exists a certificate S to show that i is not likely the best arm as it loses to some arms in S by a large 'margin' f The larger the set S the larger the margin needs to be. This condition can be evaluated in polynomial time by computing  $\arg\max_{S\subseteq [n]} I_i(t,S) - |S| \cdot \log(nC)$  and checking if it is greater than  $\log(t)$  (details in supplementary material).

Finally, let t be the final round in a round r. In order to decide which arms should be included in the next set of active arms  $A_{r+1}$  we simply check the condition  $\mathcal{J}_i(t,C)$  for each  $i\in[n]$  and include all arms for which  $\mathcal{J}_i(t,C)=0$  holds. Note that the set of active arms  $A_{r+1}$  can be empty, in which case we will simply play the anchor arm until it becomes non-empty in the future. The anchor arm in each trial is the arm which empirically beats the maximum number of unplayed arms in the current round. Detailed pseudo-code for WBA is given in Algorithm  $\P$ 

## 5 Regret Analysis

In this section we will prove a regret upper bound for our WBA algorithm. The following theorem gives the upper bound.

**Theorem 2.** Let n be the number of arms,  $k \leq n$  be the choice set size parameter, and  $i^*$  be the GCW arm . If the multiway choices are drawn according to a GCC choice model with  $\Delta_{\min}^{GCC}$  and  $\Delta_{\max}^{GCC}$  defined in Equation 3.1 then for any  $C \geq 1/(\Delta_{\min}^{GCC})^4$ , the expected regret incurred by WBA is upper bounded by

$$\mathbf{E}\left[R(T)\right] \leq O\left(\frac{n^2 \log n}{(\Delta_{\min}^{\text{GCC}})^2}\right) + O\left(n \log(TC) \cdot \frac{\Delta_{\max}^{\text{GCC}}}{(\Delta_{\min}^{\text{GCC}})^2}\right),\,$$

where T is the (unknown) time-horizon. If the underlying model is MNL with weights  $v_1, \dots, v_n \in \mathbb{R}$ , then for any  $C \geq 1/(\Delta_{\min}^{\text{MNL}})^4$ , we have

$$\mathbf{E}\left[R(T)\right] \le O\left(\frac{n^2 \log n}{(\Delta_{\min}^{\text{MNL}})^2}\right) + O\left(\sum_{i \in [n] \setminus i^*} \frac{\log(TC)}{\Delta_{i^*i}^{\text{MNL}}}\right) ,$$

where  $\Delta^{\mathrm{MNL}}_{i^*i} = \frac{e^{v_{i^*}} - e^{v_{i}}}{e^{v_{i^*}} + e^{v_{i}}}$  and  $\Delta^{\mathrm{MNL}}_{\min} := \min_{i \neq i^*} \Delta^{\mathrm{MNL}}_{i^*i}$ .

Remark 2 (Selecting C). A value of  $T^4$  for the parameter C suffices for Theorem 2 to hold, giving a regret upper bound of  $O(\log(TC)) = O(\log(T^5)) = O(\log(T))$ . (If T is not known, one can use the doubling trick.) To see this note that in order to obtain any non-trivial upper bound for our algorithm,  $\Delta_{\min}$  has to be larger than 1/T. Hence, either  $\Delta_{\min}$  is upper bounded by 1/T, or the instance is too hard to allow any non-trivial upper bound. Therefore,  $C \geq T^4$  would suffice whenever the instance is not already too hard. We actually believe setting  $C = T^4$  may be somewhat pessimistic (it arises from taking a union bound over all possible states of the algorithm in our regret analysis – indeed, in our experiments, we set C = 1 for all datasets, and our algorithm still demonstrates sublinear regret with this choice – but it certainly suffices, and the regret bound with  $C = T^4$  is at most a constant factor 5 times what one might get with C = 1 if the regret bound holds in that case.

**Discussion.** The above theorem yields a  $O(n^2 \log n + n \log T)$  upper bound on regret. Comparing this bound with the lower bound given in Section one can observe this upper bound is *asymptotically* order-optimal. This upper bound is similar (in order-wise sense) to the upper bounds obtained for some popular dueling bandit algorithms such as RUCB , RMED . DTS . Test is also important to note that our regret bound does not depend directly on the choice set size k. However, the behavior of this bound is more subtle and depends on the specific multiway choice model through the gap parameters  $\Delta_{\max}^{GCC}$  and  $\Delta_{\min}^{GCC}$ . We also note that while in general the regret can behave differently for different models, in our experiments, we find that there are choice models (including some in real data) where our algorithm empirically achieves smaller regret when allowed to play sets of size k > 2 as compared to k = 2. If the underlying model is MNL, then our algorithm achieves asymptotically

<sup>&</sup>lt;sup>4</sup>Note that the above condition is similar to condition used in [44], except that they only use the set [n] as a certificate instead of all possible subsets  $S \subseteq [n]$ . In our analysis and experiments will show that this condition is an improvement over the condition used in [44] for the case of dueling bandits.

## Algorithm 1 Winner Beats All (WBA)

```
1: Input: set of arms [n], size of choice set k, parameter C
 2: t \leftarrow 1, r \leftarrow 1, A_r \leftarrow [n], a_t \leftarrow \text{Unif}([n]), Q \leftarrow \emptyset
 3: \hat{P}_{ij} \leftarrow \frac{1}{2}, \forall i, j \in [n]
4: while t \leq T do
            Select largest S \subseteq A_r \setminus \{Q \cup a_t\} with |S| \le k - 1 and
            \hat{P}_{ia_t} \leq \frac{1}{2}, \forall i \in S
 6:
            Let S_t \leftarrow S \cup \{a_t\}; while |S_t| < k and A_r \setminus S_t \neq \emptyset: add
            an (arbitrary) arm from A_r \setminus S_t to S_t
            Play set S_t and receive y_t \in S_t as feedback; Q \leftarrow Q \cup S
 7:
            For all i \in S_t, calculate \hat{P}_{ia_t}(t) and \mathcal{J}_i(t,C)
 8:
            if \not\exists i \in A_r \setminus \{Q \cup a_t\} such that \hat{P}_{ia_t}(t) \leq \frac{1}{2} then
 9:
                 a_{t+1} \leftarrow \operatorname{argmax}_{i \in [n]} \sum_{i \in [n] \setminus O} \mathbb{1}[\hat{P}_{ji}(t) \leq \frac{1}{2}]
10:
11:
12:
                 a_{t+1} \leftarrow a_t
13:
            if Q = A_r or S = \emptyset then
                 A_{r+1} \leftarrow \emptyset, r \leftarrow r+1

for i \in [n] do
14:
15:
                     if \mathcal{J}_i(t,C) = 0, then A_r \leftarrow A_r \cup \{i\}
16:
                 a_{t+1} \leftarrow \operatorname{argmax}_{i \in [n]} \sum_{j \in [n]} \mathbb{1}[\hat{P}_{ji}(t) \leq \frac{1}{2}], Q \leftarrow \emptyset
17:
            t \leftarrow t + 1
18:
```

Figure 1: Datasets used in our experiments

- 1. **MNL-Exp:** MNL weights drawn i.i.d. from  $\text{Exp}(\lambda = 3.5)$ ;
- 2. **MNL-Geom:** Geometrically decreasing MNL weights:  $1, \frac{1}{2}, \dots, \frac{1}{2^{n-1}}$ ;
- GCC-One: GCC model defined in Example 3;
- 4. **GCC-Two:** GCC model similar to Example 3 but with different choice probabilities:
- GCC-Three: GCC model similar to Example 3 but with different choice probabilities;
- 6. **Sushi:** Choice model extracted from the Sushi dataset [47];
- 7. **Irish-Dublin:** Choice model extracted from Irish-Dublin election dataset;
- 8. **Irish-Meath:** Choice model extracted from Irish-Meath election dataset.

optimal instance-wise regret that does not depend on k. This instance-wise bound under MNL is an improvement over the upper bound for the MaxMinUCB algorithm under MNL for (top-1) choice feedback which depends on worst-case gap parameters [23]. An important point to note is that we do not need a specialized algorithm for MNL in order to achieve an instance-wise bound under MNL.

**Proof Ideas.** Our algorithm is based on the idea of isolating a 'good' anchor arm and playing arms that are competitive against this anchor. Hence, in order to prove a regret upper bound we need to show that the GCW  $i^*$  would eventually beat every other arm i, i.e.  $\hat{P}_{i^*i}(t)$  (Equation 4.1) would eventually become larger than 1/2. In this case  $i^*$  would become the anchor arm. However, an important technical challenge here is to bound the deviation in these pairwise estimates  $\hat{P}_{i^*i}(t)$  obtained from multiway choices. In the past, [46] have shown that if one uses rank breaking to extract pairwise estimates under the MNL model, then these pairwise estimates will be concentrated. However, this concentration result relies crucially on the independence from irrelevant attributes (IIA) property of MNL which states that for any two arms, the odds of choosing one over the other in any set remains the same *regardless of which set is shown*. This concentration result does not apply to our setting as the IIA property does not hold for general GCC models beyond the MNL.

Below we outline a novel coupling argument that allows us to prove concentration for the extracted pairwise estimates between the GCW arm  $i^*$  and any other arm  $i \in [n]$ 

**Lemma 2** (Concentration). Consider a GCC choice model with GCW  $i^*$ . Fix  $i \in [n]$ . Let  $S_1, \dots, S_T$  be a sequence of subsets of [n] and  $y_1, \dots, y_T$  be a sequence of choices according to this model, let  $\mathcal{F}_t = \{S_1, y_1, \dots, S_t, y_t\}$  be a filtration such that  $S_{t+1}$  is a measurable function of  $\mathcal{F}_t$ . We have

$$\Pr(\hat{P}_{i^*i}(t) \le P_{i^*i}^{GCC} - \epsilon \text{ and } N_{i^*i}(t) \ge m) \le e^{-d(P_{i^*i}^{GCC} - \epsilon, P_{i^*i}^{GCC}) \cdot m}$$
(5.1)

where  $P_{i^*i}^{\text{GCC}} = \min_{S:|S| \leq k, \{i^*,i\} \subseteq S} \frac{P_{i^*|S}}{P_{i^*|S} + P_{i|S}}$ , and  $d(\cdot, \cdot)$  is the KL-divergence.

Proof Sketch. Let us consider an alternate process for generating multiway choices  $y'_t$  from sets  $S_t$ . In this process, given any t and a set  $S_t$  such that  $i^*, i \in S_t$  with  $a_t = i$ , we first generate a Bernoulli random variable  $X_t$  with probability  $P_{i^*|S} + P_{i|S}$ . If  $X_t = 0$  we set  $y'_t = j$  with probability  $\frac{P_{j|S}}{1 - P_{i^*|S} - P_{i|S}}$ , for  $j \in S \setminus \{i, i^*\}$ . If  $X_t = 1$  then we sample another Bernoulli random variable  $Z_t$  with probability  $P_{i^*|S}^{GCC}$ . If  $Z_t = 1$  then we let  $y'_t = i^*$ , otherwise if  $Z_t = 0$  we set  $y'_t = i$ . Let  $P_{i^*i|S_t} = P_{i^*|S_t}/(P_{i^*|S_t} + P_{i|S_t})$ . Now, we couple  $y'_t$  and  $y_t$  as follows: if  $y'_t \in S_t \setminus \{i\}$  then we let  $y_t = y'_t$ , otherwise if  $y'_t = i$  then we let  $y_t = i^*$  with probability  $(P_{i^*i|S_t} - P_{i^*i}^{GCC})/(1 - P_{i^*i}^{GCC})$  and let  $y_t = i$  with probability  $(1 - P_{i^*i|S_t})/(1 - P_{i^*i}^{GCC})$ . One can verify that  $y_t$  is distributed according to the correct underlying choice distribution. It is now easy to observe that the estimates  $\hat{P}_{i^*i}(t)$  under  $y_t$  will always be larger than the estimates  $\hat{P}_{i^*i}(t)$  under  $y'_t$ , hence, we will have that

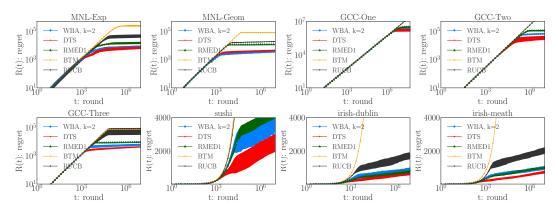


Figure 2: Regret v/s trials for our algorithm WBA (for k=2) against dueling bandit algorithms (DTS, BTM, RUCB and RMED1) (the shaded region corresponds to std. deviation). As can be observed, our algorithm is competitive against these algorithms.

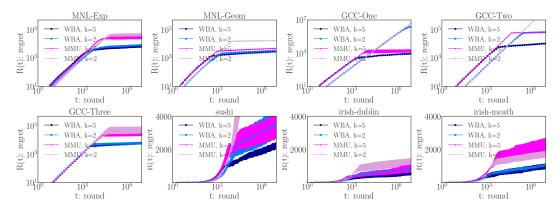


Figure 3: Regret v/s trials for our algorithm WBA against the MaxMinUCB (MMU) algorithm for k=2 and k=5 (the shaded region corresponds to std. deviation). We observe that our algorithm is better than MaxMinUCB on all datasets for both values of k. We further observe that under several datasets the regret achieved by our algorithm for k>2 is better than the regret of our algorithm for k=2.

 $\Pr(\hat{P}_{i^*i}(t) \leq x) \leq \Pr(\hat{P}'_{i^*i}(t) \leq x)$  for any x > 0. One can then show concentration for the coupled estimates  $\hat{P}'_{i^*i}(t)$ , and use it to bound the deviation in  $\hat{P}_{i^*i}(t)$ .

Note that the above lemma only shows concentration for the pairwise estimates  $\hat{P}_{i^*i}(t)$  between  $i^*$  and any other arm  $i \in [n]$ , but not for estimates  $\hat{P}_{ij}(t)$  between two arbitrary arms  $i \in [n]$  and  $j \in [n]$ . However, in order to prove our result we only need concentration of estimates between  $i^*$  and any other arm  $i \in [n]$ . We believe that the above concentration lemma is of independent interest, and might be useful in other learning from multiway choice settings beyond MNL.

Once we have bounded the deviation for the pairwise estimates, we bound the number of rounds r in which  $i^*$  is not a part of the active set  $A_r$ . We then bound the expected number of times that there exists an arm i such that  $\hat{P}_{i^*i}(t) < \frac{1}{2}$ , thus bounding the number of trials until  $i^*$  becomes the anchor. Finally, once  $i^*$  is the anchor arm, we bound the regret incurred due to sub-optimal arms.

## 6 Experiments

We compared the performance of our WBA algorithm and other existing algorithms on our choice bandit problem under different choice models. The first two choice models were MNL models, the next three were from the GCC class, and the last three we extracted from real-world datasets. Details of these models are in Figure [1] (additional details can be found in the supplementary material).

Below we describe the different sets of experiments that were performed. Each experiment was repeated 10 times. The value of n was 100 for all synthetic datasets, 16 for Sushi, 8 for Irish-Dublin, and 12 for Irish-Meath. The parameter C in our algorithm was set to 1.

Comparison with Dueling Bandit Algorithms (k=2). For the special case of k=2, we compared our algorithms with a representative set of dueling bandit algorithms (RMED1 [11], DTS [45], RUCB [6], BTM [2]). Note that the purpose of these experiments is merely to perform a sanity check and ensure that our algorithm performs reasonably well compared with dueling bandit baselines when k=2; the goal is not to argue that our choice bandit algorithm beats the state-of-the-art for the specialized dueling bandit (k=2) setting. We set  $\alpha=0.51$  for RUCB and DTS, and  $f(K)=0.3K^{1.01}$  for RMED, and  $\gamma=1.3$  for BTM. Figure 2 contain plots for these comparisons. Our algorithm either performs better or similar to RMED1, RUCB, and BTM on all datasets; and is competitive with DTS on most of the datasets.

Comparison with MaxMinUCB Algorithm [23] (k > 2). We compared the performance of our algorithm with the recent MaxMinUCB algorithm [23] that was designed and analyzed primarily for MNL choice models under the same notion of regret as ours. We set the parameter  $\alpha$  to be 0.51 for MaxMinUCB. Figure contain plots for these experiments for k = 2 and k = 5. We observe that our algorithm is much better in terms of regret than MaxMinUCB under all datasets for both values of k. One should note that WBA performs better than MaxMinUCB even under the MNL datasets, even though MaxMinUCB is specialized to MNL while our algorithm works under more general models. We further observe that under several datasets (GCC-One, GCC-Two, Sushi, Irish-Dublin) the regret achieved by our algorithm for k > 2 is better than for k = 2. Note that even though our study of more general choice feedback is motivated by applications where it might be desirable to pull sets of size larger than 2 due to reasons other than improving regret, these experimental results show that there exist settings of choice models (including some in real data) where our algorithm empirically achieves a smaller regret when allowed to play sets of size k > 2 as compared to k = 2.

## 7 Conclusion

We have introduced a new framework for bandit learning from choice feedback that generalizes the dueling bandit framework. Our main result is to show that computationally efficient learning is possible in this more general framework under a wide class of choice models that is considerably more general than the previously studied class of MNL models. Our algorithm for this general setting, termed Winner Beats All (WBA), achieves order-wise optimal regret for the general class of GCC models. For the special case k=2, WBA is competitive with previous dueling bandit algorithms; for k>2, WBA outperforms the recently proposed MaxMinUCB (MMU) algorithm even on MNL models for which MMU was designed.

## Acknowledgement

Thanks to Aadirupa Saha for early discussions related to this work, and to the anonymous referees for helpful comments. This material is based upon work supported in part by the US National Science Foundation (NSF) under Grant Nos. 1717290 and 1934876. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

## **Broader Impact**

The purpose of this paper is to understand whether efficient learning is possible in a bandit setting where one does not receive quantitative feedback for an individual arm but rather relative feedback in the form of a multiway choice. It is well-known that quantitative judgments of humans can have biases; our algorithm, which learns from relative multiway choices, can help alleviate these biases. Moreover, by receiving larger choice sets from our algorithm, humans can have a better sense of the quality distribution of arms, and can make more informed choices.

 $<sup>^5</sup>$  We also considered the SelfSparring algorithm of [26] and the battling bandit algorithms of [27], which are applicable to choice models defined in terms of an underlying pairwise comparison model P. However, these algorithms all return *multisets*  $S_t$ , and any simple reduction of such multisets to strict sets as considered in our setting (as well as the setting of [23]) can end up throwing away important information learned by the algorithms, resulting in a comparison that could be unfair to those algorithms. We did explore such reductions and our algorithm easily outperformed them, but we chose not to include the results here due to this issue of fairness. (Moreover, under the MNL model, [23] already established that MaxMinUCB outperforms those algorithms – presumably under similar reductions – so in the end, we decided such a comparison would provide little additional value here.)

Another advantage of our setting is that we do not rely on historic data as our data collection is online. Hence, one does not need to worry about past biases being reflected in the choice datasets. However, one has to be cautious about the use of our algorithm in applications where arms represent individuals/entities such as job applicants, property renters etc. In these applications, the choices of people can be biased against certain individuals/groups, thereby hurting the chances of these individuals/groups to be selected by our algorithm. Here, depending on the application, one might need to consider imposing some form of fairness constraints on the choice sets output by our algorithm in order to prevent any discrimination against such individuals/groups.

#### References

- [1] Yisong Yue, Josef Broder, Robert Kleinberg, and Thorsten Joachims. The K-armed Dueling Bandits Problem. In *Proceedings of the 22nd Conference on Learning Theory*, 2009.
- [2] Yisong Yue and Thorsten Joachims. Beat the mean bandit. In *Proceedings of the 28th International Conference on Machine Learning*, 2011.
- [3] Yisong Yue, Josef Broder, Robert Kleinberg, and Thorsten Joachims. The k-armed dueling bandits problem. J. Comput. Syst. Sci., 78(5):1538–1556, 2012.
- [4] Tanguy Urvoy, Fabrice Clerot, Raphael Feraud, and Sami Naamane. Generic Exploration and K-armed Voting Bandits. In *Proceedings of the 30th International Conference on Machine Learning*, 2013.
- [5] Nir Ailon, Zohar Karnin, and Thorsten Joachims. Reducing Dueling Bandits to Cardinal Bandits. In *Proceedings of the 31st International Conference on Machine Learning*, 2014.
- [6] Masrour Zoghi, Shimon Whiteson, Remi Munos, and Maarten de Rijke. Relative Upper Confidence Bound for the K-Armed Dueling Bandit Problem. In *Proceedings of the 31st International Conference on Machine Learning*, 2014.
- [7] Masrour Zoghi, Zohar Karnin, Shimon Whiteson, and Maarten de Rijke. Copeland Dueling Bandits. In *Advances in Neural Information Processing Systems* 28, 2015.
- [8] Masrour Zoghi, Shimon Whiteson, and Maarten de Rijke. MergeRUCB: A method for large-scale online ranker evaluation. In Proceedings of the 8th ACM International Conference on Web Search and Data Mining, 2015.
- [9] Miroslav Dudik, Katja Hofmann, Robert E. Schapire, Aleksandrs Slivkins, and Masrour Zoghi. Contextual Dueling Bandits. In *Proceedings of the 28th Conference on Learning Theory*, 2015.
- [10] Kevin Jamieson, Sumeet Katariya, Atul Deshpande, and Robert Nowak. Sparse Dueling Bandits. In *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics*, 2015.
- [11] Junpei Komiyama, Junya Honda, Hisashi Kashima, and Hiroshi Nakagawa. Regret Lower Bound and Optimal Algorithm in Dueling Bandit Problem. In *Proceedings of the 28th Conference on Learning Theory*, 2015.
- [12] Junpei Komiyama, Junya Honda, and Hiroshi Nakagawa. Copeland Dueling Bandit Problem: Regret Lower Bound, Optimal Algorithm, and Computationally Efficient Algorithm. In *Proceedings of the 33rd International Conference on Machine Learning*, 2016.
- [13] Siddartha Ramamohan, Arun Rajkumar, and Shivani Agarwal. Dueling Bandits: Beyond Condorcet Winners to General Tournament Solutions. In Advances in Neural Information Processing Systems 29, 2016
- [14] Bangrui Chen and Peter I. Frazier. Dueling Bandits with Weak Regret. In *Proceedings of the 34th International Conference on Machine Learning*, 2017.
- [15] Yisong Yue and Thorsten Joachims. Interactively Optimizing Information Retrieval Systems as a Dueling Bandits Problem. In *Proceedings of the 26th International Conference on Machine Learning*, 2009.
- [16] Eric J Johnson, Suzanne B Shu, Benedict GC Dellaert, Craig Fox, Daniel G Goldstein, Gerald Häubl, Richard P Larrick, John W Payne, Ellen Peters, David Schkade, et al. Beyond nudges: Tools of a choice architecture. *Marketing Letters*, 23(2):487–504, 2012.
- [17] Robert Duncan Luce. Individual choice behavior: A theoretical analysis. Wiley, 1959.

- [18] Robin L. Plackett. The Analysis of Permutations. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 24(2):193–202, 1975.
- [19] Daniel McFadden. Conditional Logit Analysis of Qualitative Choice Analysis. New York: Academic Press, 1974.
- [20] Louis L Thurstone. A law of comparative judgment. Psychological review, 34(4):273, 1927.
- [21] Jacob Marschak. Binary choice constraints and random utility indicators. In Stanford Symposium on Mathematical Methods in the Social Sciences, page 312–329, 1960.
- [22] TA Domencich and D McFadden. Urban travel demand; a behavioural analysis. North-Holland, 1975.
- [23] Aadirupa Saha and Aditya Gopalan. Combinatorial bandits with relative feedback. In *Advances in Neural Information Processing Systems*, pages 983–993, 2019.
- [24] Brian Brost, Yevgeny Seldin, Ingemar J. Cox, and Christina Lioma. Multi-Dueling Bandits and Their Application to Online Ranker Evaluation. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management*, 2016.
- [25] Anne Schuth, Harrie Oosterhuis, Shimon Whiteson, and Maarten de Rijke. Multileave Gradient Descent for Fast Online Learning to Rank. In Proceedings of the 9th ACM International Conference on Web Search and Data Mining, 2016.
- [26] Yanan Sui, Vincent Zhuang, Joel W. Burdick, and Yisong Yue. Multi-dueling Bandits with Dependent Arms. In Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence, 2017.
- [27] Aadirupa Saha and Aditya Gopalan. Battle of bandits. In UAI, pages 805–814, 2018.
- [28] Masrour Zoghi, Tomas Tunys, Mohammad Ghavamzadeh, Branislav Kveton, Csaba Szepesvari, and Zheng Wen. Online learning to rank in stochastic click models. In *ICML*, pages 4199–4208, 2017.
- [29] Max Simchowitz, Kevin Jamieson, and Benjamin Recht. Best-of-K Bandits. In Proceedings of the 29th Annual Conference on Learning Theory, 2016.
- [30] Paat Rusmevichientong, Zuo-Jun Max Shen, and David B. Shmoys. Dynamic Assortment Optimization with a Multinomial Logit Choice Model and Capacity Constraint. *Operations Research*, 58(6):1666–1680, 2010.
- [31] Denis Sauré and Assaf Zeevi. Optimal Dynamic Assortment Planning with Demand Learning. Manufacturing & Service Operations Management, 15(3):387–404, 2013.
- [32] Shipra Agrawal, Vashist Avadhanula, Vineet Goyal, and Assaf Zeevi. A Near-Optimal Exploration-Exploitation Approach for Assortment Selection. In Proceedings of the 17th ACM Conference on Economics and Computation, 2016.
- [33] Shipra Agrawal, Vashist Avadhanula, Vineet Goyal, and Assaf Zeevi. Thompson Sampling for the MNL-Bandit. In *Proceedings of the 30th Conference on Computational Learning Theory*, 2017.
- [34] Xi Chen and Yining Wang. A Note on Tight Lower Bound for MNL-Bandit Assortment Selection Models. Technical report, arXiv:1709.06109v2, 2017.
- [35] Xi Chen, Yuanzhi Li, and Jieming Mao. A Nearly Instance Optimal Algorithm for Top-k Ranking under the Multinomial Logit Model. In *Proceedings of the 29th Annual ACM-SIAM Symposium on Discrete* Algorithms, 2018.
- [36] Aadirupa Saha and Aditya Gopalan. Best-item learning in random utility models with subset choices. In *AISTATS*, 2020.
- [37] Yi Gai, Bhaskar Krishnamachari, and Rahul Jain. Combinatorial Network Optimization With Unknown Variables: Multi-Armed Bandits With Linear Rewards and Individual Observations. *IEEE/ACM Transac*tions on Networking, 20(5):1466–1478, 2012.
- [38] Wei Chen, Yajun Wang, and Yang Yuan. Combinatorial Multi-Armed Bandit: General Framework, Results and Applications. In *Proceedings of the 30th International Conference on Machine Learning*, 2013.
- [39] Branislav Kveton, Zheng Wen, Azin Ashkan, and Csaba Szepesvari. Tight Regret Bounds for Stochastic Combinatorial Semi-Bandits. In Proceedings of the 18th International Conference on Artificial Intelligence and Statistics, 2015.

- [40] Richard Combes, M. Sadegh Talebi, Alexandre Proutiere, and Marc Lelarge. Combinatorial Bandits Revisited. In Advances in Neural Information Processing Systems 28, 2015.
- [41] Kenneth E. Train. Discrete Choice Methods with Simulation. Cambridge University Press, 2003.
- [42] Emilie Kaufmann, Olivier Cappe, and Aurelien Garivier. On the Complexity of Best-Arm Identification in Multi-Armed Bandit Models. *Journal of Machine Learning Research*, 17(1):1–42, 2016.
- [43] Richard Combes, Stefan Magureanu, and Alexandre Proutière. Minimal exploration in structured stochastic bandits. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA, pages 1763–1771, 2017
- [44] Junpei Komiyama, Junya Honda, and Hiroshi Nakagawa. Optimal Regret Analysis of Thompson Sampling in Stochastic Multi-armed Bandit Problem with Multiple Plays. In *Proceedings of the 32nd International Conference on Machine Learning*, 2015.
- [45] Huasen Wu and Xin Liu. Double thompson sampling for dueling bandits. In *Advances in Neural Information Processing Systems* 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain, pages 649–657, 2016.
- [46] Aadirupa Saha and Aditya Gopalan. PAC battling bandits in the plackett-luce model. In *Algorithmic Learning Theory, ALT 2019, 22-24 March 2019, Chicago, Illinois, USA*, pages 700–737, 2019.
- [47] Toshihiro Kamishima. Nantonac collaborative filtering: recommendation based on order responses. In Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 24 - 27, 2003, pages 583–588, 2003.

# Choice Bandits Supplementary Material

# A Organization

We provide additional discussion about the related work in Appendix B. We provide the proof of our regret lower bound (Theorem I) in Appendix C. We prove a concentration inequality for pairwise estimates in Appendix D. We then provide the proof of our regret upper bound (Theorem 2) in Appendix E. In Appendix F we provide additional details about our experimental setup. In Appendix G we provide experimental results for an alternate notion of regret. Appendix H contains some technical lemmas used in the proof of the upper bound result in Theorem 2.

## **B** Related Work

There has been some recent interest in bandit settings where more than two arms are played at once (although no previous work considers choice models at the level of generality we do). We review related work here and provide a summary in Table 1.

**Multi-dueling bandits:** In *multi-dueling bandits* [1] [2] [3], the learner pulls a set  $S_t$  of k items; however, the feedback received by the learner is assumed to be drawn from a pairwise comparison model (in particular, the learner observes some subset of the  $\binom{k}{2}$  possible pairwise comparisons among items in  $S_t$ ). In contrast, in our choice bandits setting, the learner receives the outcome of a direct multiway choice among the items in  $S_t$ , generated from a multiway choice model.

Combinatorial bandits: In combinatorial (semi) bandits [4, 5, 6, 7], each arm i is associated with an unknown random variable (stochastic reward)  $Y_i$ ; the learner pulls a set  $S_t$  of up to k arms, and observes the realized rewards  $y_t(i)$  for all arms i in  $S_t$ . In contrast, we only observe the arm that is chosen from the set  $S_t$  that is played.

Combinatorial bandits with relative feedback: In this very recent framework [8], the learner pulls a set  $S_t$  of up to k arms, and observes top-m ordered feedback drawn according to the MNL model, for some  $m \leq k$ . In contrast, we only observe the (top-1) choice feedback from the set  $S_t$  that is played. Moreover, we study a much more general class of choice models than the MNL model studied by them.

**Stochastic Click Bandits:** In *stochastic click bandits* [D], the learner pulls an *ordered* set of k arms/documents, and observes *clicks* on a subset of these documents, drawn according to an underlying *click model* which is a probabilistic model for click generation over an ordered set. However, click models in their setting are different than choice models in our setting, and neither can be cast as a special case of the other.

**Battling Bandits:** Another related setting is that of *battling bandits* [10], where the learner pulls a set  $S_t$  of *exactly k* arms and receives a feedback indicating which arm was chosen. However, their setting considers a specific pairwise-subset (PS) choice model that is defined in terms of a pairwise comparison model, whereas we consider much more general choice models.

**Preselection Bandits:** There has been a recent framework called *preselection bandits*  $[\![\![\!]\!]\!]$  where two settings are considered: (1) where the learner pulls a set  $S_t$  of size exactly k, (2) where the learner pulls a set  $S_t$  of any size less than n. In both settings the learner receives feedback drawn from the MNL model. Firstly, the two settings considered by this paper are different than our setting where the learner plays a set of size up to k. Secondly, we study a much

	Rep.	Arms Pulled	Feedback in	
Problem	Paper	in Round $t$	Round $t$	Goal
Dueling Bandits	[19]	$(i_t, j_t) \in [n]^2$	$y_t \in \{i_t, j_t\}$	Min. regret w.r.t. best arm
Multi-dueling Bandits	[3]	$S_t \in [n]^k$	$Y_t = \{0, 1, \emptyset\}^{k \times k}$	Min. regret w.r.t. best arm
Combinatorial Bandits	<u>6</u>	$S_t \in \mathcal{S} \subseteq 2^{[n]} :  S_t  \le k$	$y_t(i) \in \mathbb{R} \ \forall i \in S_t$	Min. regret w.r.t. top- $k$ arms
Com. Ban. Relative Feed.	[8]	$S_t \subseteq [n]: S_t  \le k$	$O_t \subseteq S_t,  O_t  \le m$	Min. reg. w.r.t. best arm (MNL)
Battling Bandits	[10]	$S_t \in [n]^k$	$y_t \in S_t$	Min. reg. w.r.t. best arm (PS)
Stochastic Click Bandits	9	$O_t \subseteq [n]:  O_t  = k,$	$y_t \subseteq O_t$	Max. expected clicks
Dynamic Assortment	[13]	$\{0\} \cup S_t \subseteq [n]:  S_t  \le k$	$y_t \in S_t$	Max. expected revenue
Choice Bandits	This paper	$S_t \subseteq [n]: S_t  \le k$	$y_t \in S_t$	Min. regret w.r.t. best arm

Table 1: Overview of related work in regret minimization settings. There are several definitions of 'best' arm; the reader is encouraged to refer to the relevant papers and to our problem setting for details. (Note: in multi-dueling bandits,  $\emptyset$  denotes no feedback; in stochastic click bandits,  $O_t$  denotes an ordered set; in combinatorial bandits, S denotes a set of allowed subsets; in dynamic assortment optimization, S denotes the "no-purchase" option.)

more general class of choice models than the MNL model studied by them.

**Dynamic assortment optimization:** In dynamic assortment optimization [12] [13] [14] [15] [16], there are n products and each product is associated with a revenue. The learner plays an assortment  $S_t$  of up to k products, and observes a feedback indicating which (if any) of the products was purchased; the goal of the learner is to maximize the expected revenue.

**Best-of-**k bandits (PAC setting). [17] consider a best-of-k bandits setting, where again the learner pulls a set  $S_t$  of k arms; however here each arm i is associated with an unknown random variable (stochastic reward)  $Y_i$ . Of the various types of feedback that are considered, the marked bandit feedback corresponds to a setting that is similar to our choice bandits framework, however, the analysis in [17] is in the PAC/pure exploration setting, while ours is in the regret minimization setting.

**Top-**k identification under MNL model (PAC setting). Recently, there has also been work on identifying the top-k items under an MNL model from actively selected sets  $S_t$  in the PAC/pure exploration setting [18].

# C Proof of Lower Bound (Theorem 1)

We say that an algorithm is *strongly consistent* under GCC if its expected regret over T trials is  $o(T^a)$  for a constant a < 1 under any model in this class.

**Theorem 1.** Given a set of arms [n], choice set size bound  $k \le n$ , parameter  $\Delta \in (0,1)$ , and any strongly consistent algorithm A under GCC, there exists a GCC choice model with  $\Delta_{\min}^{GCC} = \Delta$  such that when choice outcomes are drawn from this model we have

$$\liminf_{T \to \infty} \frac{\mathbf{E}\left[R(T)\right]}{\log T} = \Omega\left(\frac{n-1}{\Delta}\right) \,,$$

where T is the time-horizon. If the underlying model is MNL with parameters  $v_1, v_2, \cdots v_n \in \mathbb{R}$ , then:

$$\liminf_{T \to \infty} \frac{\mathbf{E}\left[R(T)\right]}{\log T} = \Omega\left(\sum_{i \in [n] \setminus \{i^*\}} \frac{1}{\Delta_{i^*i}^{\text{MNL}}}\right) \,,$$

where 
$$\Delta^{\text{MNL}}_{i^*i} = \frac{e^{v_{i^*}} - e^{v_i}}{e^{v_{i^*}} + e^{v_i}}$$
, for  $i \in [n] \setminus \{i^*\}$ .

The above theorem states that there exists a model in the GCC class where any strongly consistent algorithm needs to incur  $\Omega(n \log T)$  regret. If the underlying model is MNL, then such an algorithm will again incur  $\Omega(n \log T)$  regret,

however, we provide a more refined instance-wise bound in this case. Also note the difference in quantifiers 'there exists' for GCC and 'for any' for MNL.

We will prove this theorem using the following change of measure lemma of [20].

**Lemma 3** ([20]). Consider two multi-armed bandit instances where A is the set of arms, and the two different collections of reward distributions are  $\mu = \{\mu_i : \forall i \in A\}$  and  $\mu' = \{\mu'_i : \forall i \in A\}$ , let  $i_t$  be the arm played at trial t by an algorithm and  $X_t$  be the reward at time t, and let  $\mathcal{F}_t = \sigma(i_1, X_1, \dots, i_t, X_t)$  be the sigma algebra upto time t. Consider a  $\mathcal{F}_T$  measurable random variable  $Z \in [0, 1]$ , then

$$\sum_{i \in A} \mathbf{E}_{\mu}[N_i(T)] KL(\mu_i, \mu_i') \ge d(\mathbf{E}_{\mu}[Z], \mathbf{E}_{\mu'}[Z]),$$

where  $N_i(T)$  denotes the number of pulls of arm i in T trials and KL is the Kullback-Leibler divergence between two distributions, and d(p;q) is the Kullback-Leibler divergence between Bernoulli distributions with parameters p and q.

In the proof of the lower bound we first bound the number of times an arm is played using the above lemma, and then bound the total regret due to this arm. Let us first define the regret per arm  $i \in [n]$  as

$$R(T,i) = \sum_{t=1}^{T} \mathbb{1}[i \in S_t] \cdot (P_{i^*|S_t \cup i^*} - P_{i|S_t \cup i^*}).$$

We will now provide the proof of the lower bound.

Proof of Theorem I Given a  $\Delta \in (0,1)$ , we will construct instance I of the choice bandits problem with n arms such that the GCW arm  $i^*$  is arm 1. Under this instance, given any set S such that  $i^* \in S$ , we have  $P_{i^*|S} = \frac{1+\Delta}{|S|(1-\Delta)+2\Delta}$  and for any  $i \in S \setminus \{i^*\}$ ,  $P_{i|S} = \frac{1-\Delta}{|S|(1-\Delta)+2\Delta}$ . Given any set S such that  $i^* \notin S$ , we will let an arbitrary chosen arm  $i^*_S \in S$  be the arm with the highest choice probability in S. We have  $P_{i^*_S|S} = \frac{1+\Delta}{|S|(1-\Delta)+2\Delta}$ , and for any  $i \in S \setminus \{i^*_S\}$ ,  $P_{i|S} = \frac{1-\Delta}{|S|(1-\Delta)+2\Delta}$ . Note that  $i^*_S$  will be equal to  $i^*$  when  $i^* \in S$ . For any set S with  $|S| \ge 2$  and  $i \in S$ , the instance I also satisfies that I and I are I are I and I are I are I and I are I and I are I are I and I are I and I are I and I are I are I and I are I are I and I are

For  $i \in [n] \setminus \{1\}$ , we will now modify this instance to create a new instance  $\mathbf{P}'$  where the GCW arm is i. Now, in the new instance, for any set S, we will have that  $P'_{i_S^*|S} := P_{i|S}$  and  $P'_{i|S} := P_{i_S^*|S}$  and for all  $j \in S \setminus \{i_S^*, i\}$  we will have  $P'_{j|S} := P_{j|S}$ . Clearly, the best arm in this new instance is the arm i as it has the highest choice probability in any choice set. It is also easy to verify that both instances belong to the GCC class.

Now, given any set S, the probability distributions  $P_S$  and  $P_S'$  associated with this set are categorical distributions where the feedback is j with probability  $P_{j|S}$  and  $P_{j'|S}$ , respectively. Now, let  $A := \{S \subseteq [n] : |S| \le k\}$  be the set of choice sets of size at most k. We can then use Lemma  $\mathfrak{Z}$  with arms corresponding to sets in A and the reward for set S being drawn from categorical distributions  $P_S$  and  $P_S'$ . We then have the following bound—

$$\sum_{S \in A} \mathbf{E}_{\mathbf{P}}[N_S(T)] KL(P_S, P_S') \ge d(\mathbf{E}_{\mathbf{P}}[Z], \mathbf{E}_{\mathbf{P}'}[Z]).$$

where  $N_S(T)$  is the number of times set S is played in T rounds, and Z is any  $\mathcal{F}_T$  measurable random variable. Also, let  $A^i = \{S \in A \setminus \{i\} : i \in S\}$  be all sets that contain i except the singleton set  $\{i\}$ . Since, we have that for any  $S \in A \setminus A^i$  the KL divergence  $KL(P_S, P_S') = 0$ , then the above bound becomes:

$$\sum_{S \in A^i} \mathbf{E}_{\mathbf{P}}[N_S(T)] KL(P_S, P_S') \geq d(\mathbf{E}_{\mathbf{P}}[Z], \mathbf{E}_{\mathbf{P}'}[Z]) \,.$$

Given any set  $S \in A^i$  we can now calculate the KL divergence between the two categorical distributions using the

inequality  $KL(p,q) \leq \sum_{x \in \mathcal{X}} \frac{(p(x) - q(x))^2}{q(x)}$ , where  $\mathcal{X}$  is the support of the two distributions.

$$\begin{split} KL(P_S, P_S') &\leq \sum_{j \in S} \frac{(P_{j|S} - P_{j|S}')^2}{P_{j|S}'} \\ &= \frac{(P_{i|S} - P_{i|S}')^2}{P_{i|S}'} + \frac{(P_{i_S^*|S} - P_{i_S^*|S}')^2}{P_{i_S^*|S}'} \\ &= \frac{(P_{i|S} - P_{i_S^*|S})^2}{P_{i_S^*|S}} + \frac{(P_{i|S} - P_{i_S^*|S})^2}{P_{i|S}} \end{split}$$

Now, similar to [3], let Z be the fraction of times out of T the singleton set  $\{i\}$  is played, i.e.  $Z = N_i(T)/T$  where  $N_i(T)$  counts the number of times set  $\{i\}$  is played. We will then have

$$d(\mathbf{E}_{\mathbf{P}}[Z], \mathbf{E}_{\mathbf{P}'}[Z]) \ge \left(1 - \frac{\mathbf{E}_{\mathbf{P}}[N_i(T)]}{T}\right) \ln \frac{T}{T - \mathbf{E}_{\mathbf{P}'}[N_i(T)]} - \ln 2.$$

Since, the algorithm is strongly consistent it can only play a suboptimal arm  $\{i\}$  only a sublinear number of times, i.e.  $\mathbf{E}_{\mathbf{P}}[N_i(T)] = o(T^{\alpha})$  and  $T - \mathbf{E}_{\mathbf{P}'}[N_i(T)] = o(T^{\alpha})$  for some  $\alpha < 1$ . Hence, we have that

$$\lim_{T \to \infty} \frac{1}{\ln T} d(\mathbf{E}_{\mathbf{P}}[Z], \mathbf{E}_{\mathbf{P}'}[Z]) \ge \lim_{T \to \infty} \frac{1}{\ln T} \left( 1 - \frac{o(T^{\alpha})}{T} \right) \ln \frac{T}{o(T^{\alpha})} - \ln 2 \ge (1 - \alpha). \tag{C.1}$$

Combining this with the previous inequality, we have that

$$\lim_{T \to \infty} \frac{1}{\ln T} \sum_{S \in A^{i}} \mathbf{E}_{\mathbf{P}}[N_{S}(T)] \left( \frac{(P_{i|S} - P_{i_{S}^{*}|S})^{2}}{P_{i_{S}^{*}|S}} + \frac{(P_{i|S} - P_{i_{S}^{*}|S})^{2}}{P_{i|S}} \right) \ge (1 - \alpha)$$

$$\implies \lim_{T \to \infty} \frac{1}{\ln T} \sum_{S \in A^{i}} \mathbf{E}_{\mathbf{P}}[N_{S}(T)] \cdot (P_{i|S} - P_{i_{S}^{*}|S}) \left( \frac{(P_{i|S} - P_{i_{S}^{*}|S})}{P_{i_{S}^{*}|S}} + \frac{(P_{i|S} - P_{i_{S}^{*}|S})}{P_{i|S}} \right) \ge (1 - \alpha)$$

$$\implies \lim_{T \to \infty} \frac{1}{\ln T} \sum_{S \in A^{i}} \mathbf{E}_{\mathbf{P}}[N_{S}(T)] \cdot \frac{3}{2} \cdot (P_{i^{*}|S \cup i^{*}} - P_{i|S \cup i^{*}}) \left( \frac{(P_{i_{S}^{*}|S} - P_{i|S})}{P_{i_{S}^{*}|S}} + \frac{(P_{i_{S}^{*}|S} - P_{i|S})}{P_{i|S}} \right) \ge (1 - \alpha)$$

$$\lim_{T \to \infty} \frac{1}{\ln T} \mathbf{E}[R(T, i)] \cdot \frac{3}{2} \left( \frac{(P_{i_{S}^{*}|S} - P_{i|S})}{P_{i_{S}^{*}|S}} + \frac{(P_{i_{S}^{*}|S} - P_{i|S})}{P_{i|S}} \right) \ge (1 - \alpha),$$

where the second last equation follows from the properties of the underlying instance, and the last equation follows from the definition of regret per arm. We will now argue that

$$\left(\frac{(P_{i_S^*|S} - P_{i|S})}{P_{i_S^*|S}} + \frac{(P_{i_S^*|S} - P_{i|S})}{P_{i|S}}\right) = \frac{2\Delta}{1 + \Delta} + \frac{2\Delta}{1 - \Delta} = \frac{4\Delta}{(1 + \Delta)(1 - \Delta)}.$$

Using this we will have that

$$\lim_{T \to \infty} \frac{1}{\ln T} \mathbf{E}[R(T, i)] \cdot \frac{3}{2} \left( \frac{(P_{i_S^*|S} - P_{i|S})}{P_{i_S^*|S}} + \frac{(P_{i_S^*|S} - P_{i|S})}{P_{i|S}} \right) \ge (1 - \alpha)$$

$$\implies \lim_{T \to \infty} \frac{1}{\ln T} \mathbf{E}[R(T, i)] \cdot \frac{4\Delta}{(1 + \Delta)(1 - \Delta)} \ge (1 - \alpha) \cdot \frac{2}{3}$$

$$\implies \lim_{T \to \infty} \frac{1}{\ln T} \mathbf{E}[R(T, i)] \ge (1 - \alpha) \cdot \frac{(1 + \Delta)(1 - \Delta)}{4\Delta} \cdot \frac{2}{3}.$$

We also have that  $\frac{(1+\Delta)(1-\Delta)}{4\Delta} = \Omega(\frac{1}{\Delta})$  for any  $\Delta$  bounded away from 1. Since, we have that  $R(T) = \sum_{i \in [n]} R(T,i)$  we get that

$$\lim_{T \to \infty} \frac{1}{\ln T} \mathbf{E}[R(T)] = \Omega\left(\frac{n-1}{\Delta}\right) ,$$

which concludes the proof of the lower bound for the general GCC class.

Now, given any MNL instance, we also derive a regret lower bound which gives the minimum instance-wise regret any strongly-consistent algorithm for the GCC class needs to incur under this MNL instance.

Consider an instance  $\mathbf P$  with an underlying MNL model with weights  $v_1,\cdots,v_n$ . We will assume that all these weights are distinct for simplicity, otherwise we can add a small perturbation to these weights to break ties. We will re-parameterize this instance, and let  $w_i := \log v_i$  for any  $i \in [n]$ . Given any set S, let  $w_S = \sum_{j \in [n]} w_j$ . We have that  $P_{i|S} = w_i/w_S$  for any  $i \in S$ . Given S, we will again let  $i_S^*$  to be the arm that has the highest choice probability in S, i.e.  $i_S^* = \operatorname{argmax}_{i \in S} w_i$ . We will denote by  $\kappa$  the ratio of the maximum weight to minimum weight, i.e.  $\kappa = \max_i w_i / \min_j w_j$ .

For  $i \in [n] \setminus \{1\}$ , we will now modify this instance to create a new instance  $\mathbf{P}'$  where the GCW arm is i. In the new instance, for any set S, we will have that  $P'_{i_S^*|S} := P_{i|S}$  and  $P'_{i|S} := P_{i_S^*|S}$  and for all  $j \in S \setminus \{i_S^*, i\}$  we will have  $P'_{j|S} := P_{j|S}$ . Clearly, the best arm in this new instance is the arm i as it has the highest choice probability in any choice set. It is also easy to verify that this new instance  $\mathbf{P}'$  belongs to the GCC class. Note that  $\mathbf{P}'$  might not belong to the MNL class. Under the instance  $\mathbf{P}$  we have that  $(1 + \kappa)(P_{i^*|S \cup i^*} - P_{i|S \cup i^*}) \geq (P_{i_S^*|S} - P_{i|S})$ .

Given these two instances, we can follow steps analogous to the proof of the GCC case, to derive the following bound

$$\lim_{T \to \infty} \frac{1}{\ln T} \mathbf{E}[R(T, i)] \cdot (1 + \kappa) \left( \frac{(P_{i_S^*|S} - P_{i|S})}{P_{i_S^*|S}} + \frac{(P_{i_S^*|S} - P_{i|S})}{P_{i|S}} \right) \ge (1 - \alpha).$$

We now have that

$$\begin{split} \left(\frac{(P_{i_S^*|S} - P_{i|S})}{P_{i_S^*|S}} + \frac{(P_{i_S^*|S} - P_{i|S})}{P_{i|S}}\right) &= \frac{w_{i_S^*} - w_i}{w_i} + \frac{w_{i_S^*} - w_i}{w_{i_S^*}} = \frac{w_{i_S^*} - w_i}{w_{i_S^*} + w_i} \left(\frac{w_{i_S^*} + w_i}{w_i} + \frac{w_{i_S^*} + w_i}{w_{i_S^*}}\right) \\ &\leq \frac{w_{i^*} - w_i}{w_{i^*} + w_i} \left(3 + \kappa\right) = \Delta_{i^*i}^{\text{MNL}} \left(3 + \kappa\right) \end{split}$$

Using the same steps as above we have that

$$\lim_{T \to \infty} \frac{1}{\ln T} \mathbf{E}[R(T, i)] \ge (1 - \alpha) \cdot \frac{1}{\Delta_{i*i}^{\text{MNL}}} \cdot \frac{1}{(3 + \kappa)(1 + \kappa)}.$$

Since, we have that  $R(T) = \sum_{i \in [n]} R(T,i)$  we get that

$$\lim_{T \to \infty} \frac{1}{\ln T} \mathbf{E}[R(T)] = \Omega \left( \sum_{i \in [n] \setminus \{i^*\}} \frac{1}{\Delta_{i^*i}^{\text{MNL}}} \right) ,$$

which concludes the proof of the lower bound for the MNL case.

Note that the lower bound for the MNL model also implies a lower bound for the general GCC class. However, we chose to construct an instance outside MNL for the GCC lower bound in order to show that such a lower bound also holds beyond the MNL. Also, note that the lower bound in [8] for MNL under MNL consistent algorithms is worst-case while our lower bound for MNL under GCC consistent algorithms applies to all MNL instances.

# **D** A Concentration Inequality for Pairwise Estimates

In this section we will prove our concentration inequality that would be needed to bound the deviation in the pairwise preference estimates extracted from multiway comparisons.

**Lemma 2.** Consider a GCC choice model with GCW  $i^*$ . Fix  $i \in [n]$ . Let  $S_1, \dots, S_T$  be a sequence of subsets of [n] and  $y_1, \dots, y_T$  be a sequence of choices according to this model, let  $\mathcal{F}_t = \{S_1, y_1, \dots, S_t, y_t\}$  be a filtration containing the history of execution of the algorithm such that  $S_{t+1}$  is a measurable function of  $\mathcal{F}_t$ . Let  $\hat{P}_{i^*i}(t)$  be the empirical probability estimate of  $i^*$  beating i calculated according to Equation 4.1 then for any given  $t \in [T]$  we have that

$$\Pr\left(\hat{P}_{i^*i}(t) \le P_{i^*i}^{\text{GCC}} - \epsilon \text{ and } N_{i^*i}(t) \ge m\right) \le e^{-d(P_{i^*i}^{\text{GCC}} - \epsilon, P_{i^*i}^{\text{GCC}}) \cdot m} \tag{D.1}$$

where

$$P_{i^*i}^{\text{GCC}} = \min_{S:|S| \le k, \{i^*, i\} \subseteq S} \frac{P_{i^*|S}}{P_{i^*|S} + P_{i|S}},$$
(D.2)

and  $d(\cdot, \cdot)$  is the KL-divergence between two Bernoulli distributions, and  $N_{i^*i}(t) := \sum_{t'=1}^t \mathbb{1}(a_{t'} = i, \{i^*, i\} \subseteq S_{t'}, y_{t'} \in \{i^*, i\})$ . The above bound implies the following bound

$$\Pr\left(\hat{P}_{i^*i}(t) \le \frac{1}{2}; N_{i^*i}(t) \ge m\right) \le e^{-d(\frac{1}{2}, P_{i^*i}^{GCC})m}$$
(D.3)

We also have the following bound—

$$\Pr(\hat{P}_{ii^*}(t) \ge P_{ii^*}^{GCC} + \epsilon; N_{i^*i}(t) \ge m) \le e^{-d(P_{i^*i^*}^{GCC} - \epsilon, P_{i^*i}^{GCC}) \cdot m}$$
(D.4)

where  $P_{ii^*}^{GCC} = 1 - P_{i^*i}^{GCC}$ .

Proof. We will first prove inequality  $\boxed{\text{D.1}}$ . Let  $Z_1, Z_2, \cdots$  be a sequence of i.i.d. Bernoulli random variables with probability of success  $P_{i*i}^{\text{GCC}}$ . We will initialize a counter C to 0. Let us consider an alternate process for generating multiway choices  $y'_t$  from sets  $S_t$ . In this process, given any t and a set  $S_t$  such that  $i^*, i \in S_t$  with  $a_t = i$ , we first generate a Bernoulli random variable  $X_t$  with probability  $P_{i*|S} + P_{i|S}$ . If  $X_t = 0$  we sample a multinomial random variable  $Y_t$  such that  $Y_t = j$  with probability  $\frac{P_{j|S}}{1 - P_{i*|S} - P_{i|S}}$ , for  $j \in S \setminus \{i, i^*\}$ , and let  $y'_t = Y_t$ . If  $X_t = 1$ , then we increase the counter C by 1, and sample the Bernoulli random variable  $Z_C$  with probability  $P_{i*i}^{\text{GCC}}$ . If  $Z_C = 1$  we declare  $i^*$  as the choice, i.e.  $y'_t = i^*$ , otherwise if  $Z_C = 0$  we declare i to be the choice. Let  $P_{i*i|S} = P_{i*|S}/(P_{i*|S} + P_{i|S})$ . Now, we couple the process generating  $y'_t$  and the process generating  $y_t$  as follows: if  $y'_t \in S_t \setminus \{i\}$  then we let  $y_t = y'_t$ , otherwise if  $y'_t = i$  then we let  $y_t = i^*$  with probability  $(P_{i*i|S_t} - P_{i*i}^{\text{GCC}})/(1 - P_{i*i}^{\text{GCC}})$  and let  $y_t = i$  with probability  $(1 - P_{i*i|S_t})/(1 - P_{i*i}^{\text{GCC}})$ . The first thing to check is that  $y_t$  is drawn from the correct probabilities  $P_{y_t|S_t}$  according to the underlying choice model. We have, for any  $j \in S_t \setminus \{i^*, i\}$ 

$$\begin{aligned} \Pr\{y_t = j | S_t\} &= \Pr\{X_t = 0, Y_t = j | S_t\} \\ &= \Pr\{X_t = 0 | S_t\} \Pr\{Y_t = j | X_t = 0, S_t\} \\ &= \left(1 - P_{i^*|S_t} - P_{i|S_t}\right) \cdot \frac{P_{j|S_t}}{1 - P_{i^*|S_t} - P_{i|S_t}} \\ &= P_{j|S_t} \end{aligned}$$

We also have that

$$\Pr\{y_{t} = i^{*}|S_{t}\} = \Pr\{X_{t} = 1, Y_{t} = i^{*}|S_{t}\} + \frac{P_{i^{*}i|S_{t}} - P_{i^{*}i}^{GCC}}{1 - P_{i^{*}i}^{GCC}} \cdot \Pr\{X_{t} = 1, Y_{t} = i|S_{t}\}$$

$$= \left(P_{i^{*}|S_{t}} + P_{i|S_{t}}\right) \cdot \left(P_{i^{*}i}^{GCC} + \left(1 - P_{*i}^{GCC}\right) \cdot \frac{P_{i^{*}i|S_{t}} - P_{i^{*}i}^{GCC}}{1 - P_{i^{*}i}^{GCC}}\right)$$

$$= \left(P_{i^{*}|S_{t}} + P_{i|S_{t}}\right) \cdot \left(P_{i^{*}i|S_{t}}\right)$$

$$= P_{i^{*}|S_{t}}$$

where the last inequality follows from definition of  $P_{i^*i|S}$ . The fact that  $\Pr\{y_t = i|S_t\} = P_{i|S}$  follows from the fact that the choice probabilities sum to 1.

Let  $W_{i^*i}(t) = \sum_{t'=1}^t \mathbb{1}(a_{t'} = i, \{i^*, i\} \subseteq S_{t'}, y_{t'} = i^*)$  and  $W'_{i^*i}(t) = \sum_{t'=1}^t \mathbb{1}(a_{t'} = i, \{i^*, i\} \subseteq S_{t'}, y'_{t'} = i^*)$ . Due to the above coupling, we immediately have that  $\Pr(W_{i^*i}(t)) \ge \Pr(W'_{i^*i}(t))$  for any  $t \in [T]$ . Then

$$\Pr(W_{i^*i}(t) \le r) \le \Pr(W'_{i^*i}(t) \le r)$$

for any  $r \geq 0$ , and any  $t \in [T]$ . Using this, we have that

$$\Pr\left(\hat{P}_{i^*i}(t) \le P_{i^*i}^{GCC} - \epsilon; N_{i^*i}(t) \ge m\right) = \Pr\left(W_{i^*i}(t) \le N_{i^*i}(t) \cdot (P_{i^*i}^{GCC} - \epsilon); N_{i^*i}(t) \ge m\right)$$

$$\le \Pr\left(W'_{i^*i}(t) \le N_{i^*i}(t) \cdot (P_{i^*i}^{GCC} - \epsilon); N_{i^*i}(t) \ge m\right)$$

Now, using techniques similar to [21], we have the following bound

$$\Pr\left(\frac{W'_{i*i}(t)}{N_{i*i}(t)} \le P_{i*i}^{GCC} - \epsilon; N_{i*i}(t) \ge m\right) = \Pr\left(\frac{\sum_{s=1}^{N_{i*i}(t)} Z_s}{N_{i*i}(t)} \le P_{i*i}^{GCC} - \epsilon; N_{i*i}(t) \ge m\right)$$

$$= \sum_{r=m}^{t} \Pr\left(\frac{\sum_{s=1}^{r} Z_s}{r} \le P_{i*i}^{GCC} - \epsilon; N_{i*i}(t) = r\right)$$

$$= \sum_{r=m}^{t} \Pr\left(\frac{\sum_{s=1}^{r} Z_s}{r} \le P_{i*i}^{GCC} - \epsilon\right) \Pr(N_{i*i}(t) = r)$$

where the last equality holds because of the fact that  $Z_1, Z_2, \cdots$  is an independent sequence of random variables that do not lie in the sigma algebra of  $S_1, \cdots, S_t, X_1, \cdots, X_t$ . Using the KL-divergence based concentration inequality from [22] we have that

$$\Pr\left(\frac{\sum_{s=1}^{r} Z_s}{r} \le P_{i^*i}^{\text{GCC}} - \epsilon\right) \le e^{-d(P_{i^*i}^{\text{GCC}} - \epsilon, P_{i^*i}^{\text{GCC}})r}.$$

We then have that

$$\sum_{r=m}^{t} \Pr \bigg( \frac{\sum_{s=1}^{r} Z_{s}}{r} \leq P_{i^{*}i}^{\text{GCC}} - \epsilon \bigg) \Pr (N_{i^{*}i}(t) = r) \leq \sum_{r=m}^{t} e^{d(P_{i^{*}i}^{\text{GCC}} - \epsilon, P_{i^{*}i}^{\text{GCC}})r} \Pr (N_{i^{*}i}(t) = r) \\ \leq e^{-d(P_{i^{*}i}^{\text{GCC}} - \epsilon, P_{i^{*}i}^{\text{GCC}})m}$$

The proof of reverse direction follows from a similar coupling argument followed by the above concentration inequality.  $\Box$ 

Note that the above coupling technique has similarity to the coupling used in [21] in order to show concentration of pairwise estimates under the MNL model. However, this argument relies on the IIA property of MNL, which does not hold under general GCC models.

# E Proof of Regret Bound for WBA

In this section we will prove the regret bound for our WBA algorithm. The following theorem presents the bound.

**Theorem 2.** Let n be the number of arms,  $k \leq n$  be the choice set size parameter, and  $i^*$  be the GCW arm. If the multiway choices are drawn according to a GCC choice model with  $\Delta_{\min}^{GCC}$  and  $\Delta_{\max}^{GCC}$  defined in Equation 3.1 then for any  $C \geq 1/(\Delta_{\min}^{GCC})^4$ , the expected regret incurred by WBA is upper bounded by

$$\mathbf{E}\left[R(T)\right] \leq O\left(\frac{n^2 \log n}{(\Delta_{\min}^{\mathrm{GCC}})^2}\right) + O\left(n \log(TC) \cdot \frac{\Delta_{\max}^{\mathrm{GCC}}}{(\Delta_{\min}^{\mathrm{GCC}})^2}\right) \,,$$

where T is the (unknown) time-horizon. If the underlying model is MNL with weights  $v_1, \dots, v_n \in \mathbb{R}$ , then for any  $C \geq 1/(\Delta_{\min}^{\text{MNL}})^4$ , we have

$$\mathbf{E}\left[R(T)\right] \le O\left(\frac{n^2 \log n}{(\Delta_{\min}^{\text{MNL}})^2}\right) + O\left(\sum_{i \in [n] \setminus i^*} \frac{\log(TC)}{\Delta_{i^*i}^{\text{MNL}}}\right) ,$$

where  $\Delta^{\mathrm{MNL}}_{i^*i} = \frac{e^{v_{i^*}} - e^{v_i}}{e^{v_{i^*}} + e^{v_i}}$  and  $\Delta^{\mathrm{MNL}}_{\min} := \min_{i \neq i^*} \Delta^{\mathrm{MNL}}_{i^*i}$ .

The proof of the above theorem hinges on three main lemmas given below. Before stating these lemmas, we would like to remind the reader that the execution of our algorithm is divided in rounds and each round contain up to n trials. The first lemma bounds the number of rounds arm  $i^*$  is not in the active set.

**Lemma 4** (Number of rounds where  $i^*$  is not active). Fix an anchor arm  $a \in [n] \setminus \{i^*\}$ . The expected number of rounds arm  $i^*$  will not be a part of the active set is bounded as

$$\mathbf{E}\left[\sum_{r=1}^T \mathbb{1}[i^* \not\in A_r]\right] \le 2.$$

We will define  $a_r$  to be the arm that empirically beats all other arms at the end of round r-1 if such an arm exists, i.e.  $\sum_{j\in[n]}\mathbb{1}[\hat{P}_{ja_r}(t)\leq \frac{1}{2}]=n-1$ , where t is the last trial in round r-1. If there is no arm that empirically beats all other arms then we will let  $a_r=0$ . If there are multiple such arms, then we will choose one arbitrarily. The following lemma will now bound the number of rounds arm  $i^*$  does not empirically beat every other arm.

**Lemma 5** (Time when  $i^*$  is not the empirically best arm). The total number of rounds when the best arm  $i^*$  will not be the empirically best arm, even when it is in the active set, is upper bounded as

$$\mathbf{E}\left[\sum_{r=1}^{T} \mathbb{1}[a_r \neq i^*, i^* \in A_r]\right] \leq \sum_{i \in [n] \setminus \{i^*\}} \frac{1}{\exp\{d(1/2, P_{i^*i}^{GCC})\} - 1},$$

where  $P_{i*i}^{GCC}$  is defined in Equation D.2

Note that if  $a_r = i^*$  then the anchor arm in all the trials in that round becomes  $i^*$ . Let us define the regret per arm  $i \in [n]$  for a set S as

$$r(S,i) = \mathbb{1}[i \in S] \cdot (P_{i^*|S \cup i^*} - P_{i|S \cup i^*}).$$

The following lemma now bounds the regret incurred due to each suboptimal arm when played against the anchor  $i^*$ .

**Lemma 6** (Regret due to a bad arm). Given an arm  $i \in [n] \setminus \{i^*\}$  the expected regret incurred due to arm i when arm  $i^*$  is the anchor is upper bounded as

$$\mathbf{E}\left[\sum_{t=1}^{T} r(S_t, i) \cdot \mathbb{1}[a_t = i^*, i \in S_t]\right] \leq \Delta_{i^*i}^{\text{GCC}} \cdot \frac{2e}{e-1} \cdot \left(\frac{(1+\delta)\log(TnC)}{d(P_{i^*i}^{\text{GCC}}, \frac{1}{2})} + \frac{1}{\Omega(\delta^2)}\right),$$

where  $\delta > 0$  is some constant, and  $\Delta_{i^*i}^{GCC} = \max_{S:|S| \leq k} \Delta_{i^*i|S}$ .

We will now prove the above theorem using the three lemmas above.

Proof of Theorem 2 The execution of the algorithm can roughly be divided into three intermittent phases—(1) when the GCW arm  $i^*$  is not in the active set, (2) when  $i^*$  is in the active set but does not beat all other arms empirically, i.e.  $a_r \neq i^*$ , (3) when  $i^*$  is in the active set and also beats all other arms empirically. The three lemmas above bound the number of rounds spent in these three phases.

However, in order to prove a regret upper bound we will also have to bound the total regret incurred due to a single round. The first thing to observe is that each arm is played at most once in each round except a few arms that might be played multiple times due to step 6 of the algorithm. Hence, the regret for all steps except step 6 is upper bounded by n as the regret for each arm is at most 1. Now, in order to bound the regret for step 6, we need to observe that the number of times the anchor arm is changed in a single round can be at most  $\log n$ . This is due to the fact that  $A_r \setminus Q$  reduces by a factor of at least 2 each time a new anchor arm is selected by the algorithm. Now, we can bound the regret incurred due to step 6 of the algorithm by  $k \log n \le n \log n$  as the regret for each arm is upper bounded by 1 and there can be at most k arms added in step 6 per anchor arm.

Hence, we now have that

$$\begin{split} \mathbf{E}[R(T)] &\leq n \log n \cdot \left( \mathbf{E}\left[\sum_{r=1}^{T} \mathbb{1}[i^* \not\in A_r]\right] + \mathbf{E}\left[\sum_{r=1}^{T} \mathbb{1}[a_r \neq i^*, i^* \in A_r]\right] \right) + \sum_{i \in [n] \setminus \{i^*\}} \mathbf{E}\left[\sum_{t=1}^{T} r(S_t, i) \cdot \mathbb{1}[a_t = i^*, i \in S_t]\right] \\ &\leq n \log n \cdot \left( 2 + \sum_{i \in [n] \setminus \{i^*\}} \frac{1}{\exp\{d(1/2, P_{i^*i}^{\text{GCC}})\} - 1} \right) + \sum_{i \in [n] \setminus \{i^*\}} \Delta_{i^*i}^{\text{GCC}} \cdot \frac{2e}{e - 1} \cdot \left(\frac{(1 + \delta) \log(TnC)}{d(P_{i^*i}^{\text{GCC}}, \frac{1}{2})} + \frac{1}{\Omega(\delta^2)}\right) \\ &\leq O\left(\frac{n^2 \log n}{(\Delta_{\min}^{\text{GCC}})^2}\right) + n \cdot \Delta_{\max}^{\text{GCC}} \cdot \frac{2e}{e - 1} \cdot \frac{1}{\Omega(\delta^2)} + \sum_{i \in [n] \setminus \{i^*\}} \Delta_{i^*i}^{\text{GCC}} \cdot \frac{2e}{e - 1} \cdot \frac{(1 + \delta) \log(TnC)}{d(P_{i^*i}^{\text{GCC}}, \frac{1}{2})} \\ &= O\left(\frac{n^2 \log n}{(\Delta_{\min}^{\text{GCC}})^2}\right) + \sum_{i \in [n] \setminus \{i^*\}} \Delta_{i^*i}^{\text{GCC}} \cdot \frac{2e}{e - 1} \cdot \frac{(1 + \delta) \log(TC)}{d(P_{i^*i}^{\text{GCC}}, \frac{1}{2})} \end{split}$$

where the third inequality follows from the well-known Pinsker's inequality  $d(P,Q) \geq 2(P-Q)^2$  and the last inequality holds for any constant  $\delta$ . Now again using the Pinsker's inequality we have that  $d(P_{i^*i}^{\text{GCC}}, \frac{1}{2}) \geq (\Delta_{\min}^{\text{GCC}})^2/2$ . For a general GCC model, we then have that

$$\begin{split} \mathbf{E}[R(T)] &\leq O\left(\frac{n^2 \log n}{(\Delta_{\min}^{\text{GCC}})^2}\right) + \sum_{i \in [n] \backslash \{i^*\}} \Delta_{i^*i}^{\text{GCC}} \cdot \frac{4e}{e-1} \cdot \frac{(1+\delta) \log(TC)}{\Delta_{\min}^{\text{GCC}}} \\ &\leq O\left(\frac{n^2 \log n}{(\Delta_{\min}^{\text{GCC}})^2}\right) + n \cdot \Delta_{\max}^{\text{GCC}} \cdot \frac{4e}{e-1} \cdot \frac{(1+\delta) \log(TC)}{(\Delta_{\min}^{\text{GCC}})^2} \end{split}$$

which gives the desired bound under any GCC model.

Now, if the underlying GCC model is MNL, then we have  $d(P_{i^*i}^{\text{GCC}}, \frac{1}{2}) \geq (\Delta_{i^*i}^{\text{MNL}})^2/2$  and  $\Delta_{i^*i}^{\text{GCC}} = \Delta_{i^*i}^{\text{MNL}}$ . We then have that

$$\begin{split} \mathbf{E}[R(T)] &\leq O\left(\frac{n^2 \log n}{(\Delta_{\min}^{\text{MNL}})^2}\right) + \sum_{i \in [n] \backslash \{i^*\}} \Delta_{i^*i}^{\text{MNL}} \cdot \frac{4e}{e-1} \cdot \frac{(1+\delta) \log(TC)}{(\Delta_{i^*i}^{\text{MNL}})^2} \\ &= O\left(\frac{n^2 \log n}{(\Delta_{\min}^{\text{MNL}})^2}\right) + \sum_{i \in [n] \backslash \{i^*\}} \frac{4e}{e-1} \cdot \frac{(1+\delta) \log(TC)}{\Delta_{i^*i}^{\text{MNL}}} \,. \end{split}$$

## E.1 Proof of Lemma 4

The following lemma calculates the expected number of rounds arm  $i^*$  will not be played.

**Lemma 4** (Number of rounds where  $i^*$  is not active). Fix an anchor arm  $a \in [n] \setminus \{i^*\}$ . The expected number of rounds arm  $i^*$  will not be a part of the active set is bounded as

$$\mathbf{E}\left[\sum_{r=1}^{T} \mathbb{1}[i^* \not\in A_r]\right] \le 2.$$

Proof. We have that

$$\mathbf{E}\left[\sum_{r=1}^{T}\mathbb{1}[i^* \notin A_r]\right] = \mathbf{E}\left[\sum_{r=2}^{T}\mathbb{1}[i^* \notin A_r]\right] \leq \mathbf{E}\left[\sum_{t=2}^{T}\mathbb{1}[\neg \mathcal{J}_{i^*}(t,C)]\right].$$

The first equality above follows due to the fact that  $A_1$  will always include  $i^*$ . Using the union bound we have the following inequality-

$$\begin{split} \mathbb{1}[\neg \mathcal{J}_{i^*}(t,C)] &\leq \sum_{S \subseteq [n] \backslash \{i^*\}} \sum_{\{n_a\} \in [T]^S} \cdots \sum_{\{n_a\} \in [T]^S} \\ \mathbb{1}[\bigcap_{a \in S} \{N_{i^*a}(t) = n_a, \hat{P}_{i^*a}(t) < \frac{1}{2}\} \cap \bigcap_{a \notin S} \{\hat{P}_{i^*a}(t) \geq \frac{1}{2}\} \cap \{\neg \mathcal{J}_{i^*}(t,C)\}] \,. \end{split}$$

Fix some set  $S \subseteq [n] \setminus \{i^*\}$ . Also, let s := |S|. Fix some  $n_a \in [T]$  for all  $a \in S$ . Let  $\hat{P}^{n_a}_{i^*a}$  be the empirical probability of  $i^*$  beating a after being pulled together  $n_a$  times. We will analyze the number of rounds that  $i^*$  is excluded from the active set due to the above configuration of S,  $\{n_a\}$ . The conditions  $\mathcal{J}_{i^*}(t,C)$  will hold when

$$\sum_{a \in S} n_a d(\hat{P}_{i^*a}^{n_a}, \frac{1}{2}) \le \log(t) + s \log(nC) \implies t \ge \exp\left(\sum_{a \in S} n_a d(\hat{P}_{i^*a}^{n_a}, \frac{1}{2}) - s \log(nC)\right).$$

Hence, we have that

$$\sum_{t=2}^{\infty} \mathbb{1}\left[\bigcap_{a \in S} \{N_{i^*a}(t) = n_a, \hat{P}_{i^*a}(t) < \frac{1}{2}\} \cap \bigcap_{a \notin S} \{\hat{P}_{i^*a}(t) \ge \frac{1}{2}\} \cap \{\neg \mathcal{J}_{i^*}(t, C)\}\right]$$

$$\leq \exp\left(\sum_{a \in S} n_a d(\hat{P}_{i^*a}^{n_a}, \frac{1}{2}) - s\log(nC)\right).$$

Now, we will use the method similar to the one used in Lemma 5 of [23], to bound the expectation of the above quantity. Fix  $x_a \in [0, \log 2]$  for all  $a \in S$ . Let  $P_a(x_a) = \Pr\left(\hat{P}_{i^*a}^{n_a} \leq \frac{1}{2}, d^+(\hat{P}_{i^*a}^{n_a}, \frac{1}{2}) \geq x_a\right)$ , where  $d^+(P, Q) = \Pr\left(\hat{P}_{i^*a}^{n_a} \leq \frac{1}{2}, d^+(\hat{P}_{i^*a}^{n_a}, \frac{1}{2}) \geq x_a\right)$ 

 $\mathbb{1}[P \leq Q] \cdot d(P,Q)$ . We then have

$$\begin{split} \mathbf{E}\left[\sum_{t=2}^{T}\mathbb{1}[\bigcap_{a\in S}\{N_{i^*a}(t)=n_a,\hat{P}_{i^*a}(t)<\frac{1}{2}\}\cap\bigcap_{a\notin S}\{\hat{P}_{i^*a}(t)\geq\frac{1}{2}\}\cap\{\neg\mathcal{J}_{i^*}(t,C)\}]\right] \\ &\leq \int_{\{x_a\}\in[0,\log(2)]^{|S|}}\exp\left(\sum_{a\in S}n_ax_a-s\log(nC)\right)\prod_{a\in S}\mathrm{d}(-P_a(x_a)) \\ &=\exp\left(-s\log(nC)\right)\cdot\prod_{a\in S}\int_{x_a\in[0,\log(2)]}\exp\left(n_ax_a\right)\mathrm{d}(-P_a(x_a)) \\ &\qquad \qquad (\mathrm{due\ to\ the\ independence\ of\ comparisons\ with\ respect\ to\ different\ anchors)} \\ &=\exp\left(-s\log(nC)\right)\cdot\prod_{a\in S}\left(\left[-\exp(n_ax_a)P_a(x_a)\right]_0^{\log(2)}+\int_{x_a\in[0,\log(2)]}n_a\exp\left(n_ax_a\right)P_a(x_a)\mathrm{d}x_a\right) \\ &\qquad \qquad (\mathrm{integration\ by\ parts)} \\ &\leq\exp\left(-s\log(nC)\right)\cdot\prod_{a\in S}\left(P_a(0)+\int_{x_a\in[0,\log(2)]}n_a\exp\left(n_ax_a\right)\exp\left\{-n_a(x_a+C_1(P_{i^*a}^{\mathrm{GCC}},\frac{1}{2}))\right\}\mathrm{d}x_a\right) \\ &\qquad \qquad (\mathrm{Using\ concentration\ inequality\ (Lemma\ 2)\ and\ Fact\ 10\ in\ 23,\ with\ C_1(p,q)=(p-q)^2/2p(1-q))} \\ &=\exp\left(-s\log(nC)\right)\cdot\prod_{a\in S}\left(\exp\left\{-n_ad(\frac{1}{2},P_{i^*a}^{\mathrm{GCC}})\right\}+\int_{x_a\in[0,\log(2)]}n_a\exp\left\{-n_aC_1(P_{i^*a}^{\mathrm{GCC}},\frac{1}{2})\right\}\mathrm{d}x_a\right) \\ &=\exp\left(-s\log(nC)\right)\cdot\prod_{a\in S}\left(\exp\left\{-n_ad(\frac{1}{2},P_{i^*a}^{\mathrm{GCC}})\right\}+\log(2)n_a\exp\left\{-n_aC_1(P_{i^*a}^{\mathrm{GCC}},\frac{1}{2})\right\}\right). \end{split}$$

We will now take a union bound over  $\{n_a\}$ . We have that

$$\begin{split} \sum \sum_{\{n_a\} \in [T]^S} & \exp\left(-s\log(nC)\right) \cdot \prod_{a \in S} \left( \exp\left\{-n_a d(\frac{1}{2}, P_{i^*a}^{\text{GCC}})\right\} + \log(2) n_a \exp\left\{-n_a C_1(P_{i^*a}^{\text{GCC}}, \frac{1}{2})\right\} \right) \\ & = \exp\left(-s\log(nC)\right) \cdot \prod_{a \in S} \sum_{n_a} \left( \exp\left\{-n_a d(\frac{1}{2}, P_{i^*a}^{\text{GCC}})\right\} + \log(2) n_a \exp\left\{-n_a C_1(P_{i^*a}^{\text{GCC}}, \frac{1}{2})\right\} \right) \\ & \leq \exp\{-s\log(nC)\} \cdot \prod_{a \in S} \left( \frac{1}{\exp\{d(\frac{1}{2}, P_{i^*a}^{\text{GCC}})\} - 1} + \frac{\exp\{C_1(P_{i^*a}^{\text{GCC}}, \frac{1}{2})\} - 1)^2}{(\exp\{C_1(P_{i^*a}^{\text{GCC}}, \frac{1}{2})\} - 1)^2} \right) \\ & \leq \exp\{-s\log(nC) + s\log(C')\}, \end{split}$$

where the constant C' is defined as

$$C' := \max_{a \in [n] \setminus i^*} \left( \frac{1}{\exp\{d(\frac{1}{2}, P_{i^*a}^{\text{GCC}})\} - 1} + \frac{\exp\{C_1(P_{i^*a}^{\text{GCC}}, \frac{1}{2})\}}{(\exp\{C_1(P_{i^*a}^{\text{GCC}}, \frac{1}{2})\} - 1)^2} \right) \le \frac{1}{(\Delta_{\min}^{\text{GCC}})^4}.$$

We will now apply the union bound over all subsets  $S \subseteq [n] \setminus i^*$ . Now, if the parameter C is larger than C', then we

have

$$\sum_{S \subseteq [n] \setminus \{i^*\}} \exp\{-|S| \log(nC) + |S| \log(C')\} = \sum_{s=1}^{n-1} \sum_{S \subseteq [n] \setminus \{i^*\}, |S| = s} \exp\{-s \log(nC) + s \log(C')\}$$

$$\leq \sum_{s=1}^{n-1} \left(\frac{en}{s}\right)^s \exp\{-s \log(nC) + s \log(C')\}$$

$$= \sum_{s=1}^{n-1} \exp\{-s \log(nC) + s \log(C') + s \log(n) + s - s \log(s)\}$$

$$\leq \sum_{s=1}^{n-1} \exp\{s - s \log(s)\} \leq 2.$$

E.2 Proof of Lemma 5

The following lemma will now bound the number of times arm  $i^*$  will not be the empirically best arm.

**Lemma** [5] (Time when  $i^*$  is not the anchor). The total number of rounds when the best arm  $i^*$  will not be the empirically best arm, even when it is in the active set, is upper bounded as

$$\mathbf{E}\left[\sum_{r=1}^{T} \mathbb{1}[a_r \neq i^*, i^* \in A_r]\right] \leq \sum_{i \in [n] \setminus \{i^*\}} \frac{1}{\exp\{d(1/2, P_{i^*i}^{GCC})\} - 1},$$

where  $P_{i*i}^{GCC}$  is defined in Equation D.2

*Proof.* In the following we overload notation slightly and for a round r define  $N_{ii^*}(r)$  and  $\hat{P}_{ii^*}(r)$  to be the equal to

 $N_{ii^*}(t)$  and  $\hat{P}_{ii^*}(t)$ , where t is the last trial in round r. We have the following set of inequalities:

$$\begin{split} \mathbf{E}\left[\sum_{r=1}^{T}\mathbbm{1}[a_{r}\neq i^{*},i^{*}\in A_{r}]\right] &= \mathbf{E}\left[\sum_{r=1}^{T}\mathbbm{1}[\exists i\neq i^{*},i^{*}\in A_{r},N_{ii^{*}}(r)>N_{ii^{*}}(r-1),\hat{P}_{i^{*}i}(r-1)\leq \frac{1}{2}]\right] \\ &\leq \mathbf{E}\left[\sum_{r=1}^{T}\sum_{i\in[n]\backslash\{i^{*}\}}\mathbbm{1}[i^{*}\in A_{r},N_{ii^{*}}(r)>N_{ii^{*}}(r-1),\hat{P}_{i^{*}i}(r-1)\leq \frac{1}{2}]\right] \\ &\leq \mathbf{E}\left[\sum_{r=1}^{T}\sum_{i\in[n]\backslash\{i^{*}\}}\sum_{n_{i}=0}^{T}\mathbbm{1}[N_{ii^{*}}(r-1)=n_{i},N_{ii^{*}}(r)>n_{i},\hat{P}_{i^{*}i}^{n_{i}}\leq \frac{1}{2}]\right] \\ &= \mathbf{E}\left[\sum_{i\in[n]\backslash\{i^{*}\}}\sum_{n_{i}=0}^{T}\mathbbm{1}[\hat{P}_{i^{*}i}^{n_{i}}\leq \frac{1}{2}]\right] \\ &\leq \mathbf{E}\left[\sum_{i\in[n]\backslash\{i^{*}\}}\sum_{n_{i}=0}^{T}\mathbbm{1}[\hat{P}_{i^{*}i^{*}i}^{n_{i}}\leq \frac{1}{2}]\right] \\ &=\sum_{i\in[n]\backslash\{i^{*}\}}\sum_{n_{i}=0}^{T}\mathbbm{1}[\hat{P}_{i^{*}i^{*}i}^{n_{i}}\leq \frac{1}{2}]\right] \\ &=\sum_{i\in[n]\backslash\{i^{*}\}}\sum_{n_{i}=0}^{T}\exp\{-n_{i}d(1/2,P_{i^{*}i^{*}i}^{GCC})\} \qquad \text{(using concentration Lemma 2)} \\ &=\sum_{i\in[n]\backslash\{i^{*}\}}\frac{1}{\exp\{d(1/2,P_{i^{*}i^{*}}^{GCC})\}-1} \end{split}$$

E.3 Proof of Lemma 6

In the next lemma we will bound the regret for the number of times an arm other than the best arm will be played.

**Lemma** [6] (Regre due to a bad arm). Given an arm  $i \in [n] \setminus \{i^*\}$  the expected regret incurred due to arm i when arm  $i^*$  is the anchor is upper bounded as

$$\mathbf{E}\left[\sum_{t=1}^{T} r(S_t, i) \cdot \mathbb{1}[a_t = i^*, i \in S_t]\right] \leq \Delta_{i^*i}^{\text{GCC}} \cdot \frac{2e}{e-1} \cdot \left(\frac{(1+\delta)\log(TnC)}{d(P_{i^*i}^{\text{GCC}}, \frac{1}{2})} + \frac{1}{\Omega(\delta^2)}\right),$$

where  $\delta > 0$  is some constant, and  $\Delta^{\mathrm{GCC}}_{i^*i} = \max_{S:|S| \leq k} \Delta_{i^*i|S}$ .

*Proof.* Fix a value  $n_i \in \{0, \dots, T\}$ . We will first upper bound the following quantity

$$\mathbf{E}\left[\sum_{t'=1}^{T} r(S_{t'}, i) \cdot \mathbb{1}[a_{t'} = i^*, i \in S_{t'}, N_{ii^*}(t'-1) = n_i]\right].$$
(E.1)

This quantity bounds the total regret until the time  $N_{ii^*}$  remains equal to  $n_i$ . Now,  $N_{ii^*}$  is incremented in trial t' if either  $i^*$  or i. Hence,  $N_{ii^*}$  is incremented in trial t' with probability  $P_{i|S_{t'}} + P_{i^*|S_{t'}}$ . The total regret incurred due to the playing i in trial t' is given by  $P_{i^*|S_{t'}} - P_{i|S_{t'}}$ . Let us define  $c_{t'} := P_{i|S_{t'}} + P_{i^*|S_{t'}}$ , and  $p_{t'} := P_{i^*|S_{t'}} - P_{i|S_{t'}}$ .

The quantity in Equation E.1 is upper bounded by the cost of an experiment described in Fact  $\boxed{1}$  where the probability of success of coin t' is given by  $p_{t'}$  and its cost is given by  $c_{t'}$ . Using Fact  $\boxed{1}$  we have that

$$\mathbf{E}\left[\sum_{t'=1}^{T} r(S_{t'}, i) \cdot \mathbb{1}[a_{t'} = i^*, i \in S_{t'}, N_{ii^*}(t'-1) = n_i]\right] \leq \Delta_{i^*i}^{GCC} \cdot \frac{2e}{e-1}.$$

Also, let  $n_i^{\text{suf}} = \frac{(1+\delta)\log(TnC)}{d(P_i^{\text{GCC}}, \frac{1}{2})}$ . We can now upper bound the regret due to arm i as

$$\begin{split} \mathbf{E}\left[\sum_{t=1}^{T}r(S_{t},i)\cdot\mathbb{1}[a_{t}=i^{*},i\in S_{t}]\right] &= \mathbf{E}\left[\sum_{n_{i}=0}^{T}\sum_{t=1}^{T}r(S_{t},i)\cdot\mathbb{1}[a_{t}=i^{*},i\in S_{t},N_{ii^{*}}(t-1)=n_{i}]\right] \\ &\leq \sum_{n_{i}=0}^{n_{i}^{\text{suf}}}\mathbf{E}\left[\sum_{t=1}^{T}r(S_{t},i)\cdot\mathbb{1}[a_{t}=i^{*},i\in S_{t},N_{ii^{*}}(t-1)=n_{i}]\right] \\ &+ \sum_{n_{i}=n_{i}^{\text{suf}}+1}^{T}\mathbf{E}\left[\sum_{t=1}^{T}r(S_{t},i)\cdot\mathbb{1}[a_{t}=i^{*},i\in S_{t},N_{ii^{*}}(t-1)=n_{i}]\right] \\ &\leq n_{i}^{\text{suf}}\cdot\Delta_{i^{*}i}^{\text{GCC}}\cdot\frac{2e}{e-1} \\ &+ \sum_{n_{i}=n_{i}^{\text{suf}}+1}^{T}\mathbf{E}\left[\sum_{t=1}^{T}r(S_{t},i)\cdot\mathbb{1}[a_{t}=i^{*},i\in S_{t},N_{ii^{*}}(t-1)=n_{i}]\right] \end{split}$$

We will now bound the second quantity in the above equation. Fix  $n_i \in \{0, 1, \dots, T\}$ . Let  $t' \in [T]$  be such that the event  $\mathbb{1}[a_{t'} = i^*, i \in S_{t'}, N_{ii^*}(t'-1) = n_i - 1, N_{ii^*}(t') = n_i]$  holds if such a t' exists, otherwise let t' = T + 1. We have

$$\begin{split} \mathbf{E}\left[\sum_{t=1}^{T}r(S_{t},i)\cdot\mathbb{1}[a_{t}=i^{*},i\in S_{t},N_{ii^{*}}(t-1)=n_{i}]\right] \\ &=\mathbf{E}\left[\sum_{t=1}^{T}r(S_{t},i)\cdot\mathbb{1}[a_{t}=i^{*},i\in S_{t},N_{ii^{*}}(t-1)=n_{i},n_{i}\cdot d(\hat{P}_{ii^{*}}(t-1),\frac{1}{2})\leq\log(t-1)+\log(nC)]\right] \\ &=\mathbf{E}\left[\sum_{t=1}^{T}r(S_{t},i)\cdot\mathbb{1}[a_{t}=i^{*},i\in S_{t},N_{ii^{*}}(t-1)=n_{i},n_{i}\cdot d(\hat{P}_{ii^{*}}(t-1),\frac{1}{2})\leq\log(t-1)+\log(nC)]\right] \\ &=\mathbf{E}\left[\mathbb{1}[\exists t'\in[T]:N_{ii^{*}}(t')=n_{i},\;n_{i}\cdot d(\hat{P}_{ii^{*}}^{n_{i}},\frac{1}{2})\leq\log(t'nC)]\cdot\sum_{t=1}^{T}r(S_{t},i)\cdot\mathbb{1}[a_{t}=i^{*},i\in S_{t},N_{ii^{*}}(t-1)=n_{i}]\right] \\ &=\Pr\left[\exists t'\in[T]:N_{ii^{*}}(t')=n_{i},\;n_{i}\cdot d(\hat{P}_{ii^{*}}^{n_{i}},\frac{1}{2})\leq\log(t'nC)\right] \\ &\cdot\mathbf{E}\left[\sum_{t=t'+1}^{T}r(S_{t},i)\cdot\mathbb{1}[a_{t}=i^{*},i\in S_{t},N_{ii^{*}}(t-1)=n_{i}]\Big|\exists t'\in[T]:N_{ii^{*}}(t')=n_{i},\;n_{i}\cdot d(\hat{P}_{ii^{*}}^{n_{i}},\frac{1}{2})\leq\log(t'nC)\right] \end{split}$$

We will bound the quantities in the above equation one by one. Using a similar argument as above and Fact we have that

$$\mathbf{E}\left[\sum_{t=t'+1}^{T} r(S_t, i) \cdot \mathbb{1}[a_t = i^*, i \in S_t, N_{ii^*}(t-1) = n_i] \middle| \exists t' \in [T] : N_{ii^*}(t') = n_i, \ n_i \cdot d(\hat{P}_{ii^*}^{n_i}, \frac{1}{2}) \leq \log(t'nC) \right] \leq \Delta_{i^*i}^{GCC} \cdot \frac{2e}{e-1}.$$

This holds because of the fact that conditioning does not effect that events that happen after trial t' + 1. We finally have

$$\Pr\left(\exists t \in [T] : N_{ii^*}(t) = n_i, \ n_i d(\hat{P}_{ii^*}(t), \frac{1}{2}) \le \log(tnC)\right)$$

$$\le \Pr\left(\exists t \in [T] : N_{ii^*}(t) = n_i, \ n_i d(\hat{P}_{ii^*}(t), \frac{1}{2}) \le \log(TnC)\right) \qquad (n_i \ge \frac{(1+\delta) \log(TnC)}{d(P_{i^*i}^{GCC}, \frac{1}{2})})$$

$$\le \Pr\left(\exists t \in [T] : N_{ii^*}(t) = n_i, \ d(\hat{P}_{ii^*}(t), \frac{1}{2}) \le \frac{d(P_{i^*i}^{GCC}, \frac{1}{2})}{1+\delta}\right).$$

We will let  $P \in (\frac{1}{2}, P^{\text{GCC}}_{i^*i})$  to be a real number such that  $d(P, \frac{1}{2}) = \frac{d(P^{\text{GCC}}_{i^*i}, \frac{1}{2})}{1+\delta}$ , and use the concentration bound proved in Lemma 2, so that the above inequality can be written

$$\Pr\left(\exists t \in [T] : N_{ii^*}(t) = n_i, \ d(\hat{P}_{ii^*}(t), \frac{1}{2}) \le d(P, \frac{1}{2})\right) = \Pr\left(\exists t \in [T] : N_{ii^*}(t) = n_i, \ \hat{P}_{ii^*}(t) \ge 1 - P\right)$$

$$\le \exp\left(-d(P, P_{i^*, i}^{GCC}) \cdot n_i\right).$$

Hence, we will have that

$$\begin{split} \sum_{n_i = n_i^{\text{suf}} + 1}^T \mathbf{E} \left[ \sum_{t = 1}^T r(S_t, i) \mathbbm{1}[a_t = i, N_{ii^*}(t - 1) = n_i, i^* \in S_t] \right] &\leq \sum_{n_i = n_i^{\text{suf}} + 1} \Delta_{i^*i}^{\text{GCC}} \cdot \frac{2e}{e - 1} \cdot \exp\left(-d(P, P_{i^*, i}^{\text{GCC}}) \cdot n_i\right) \\ &\leq \Delta_{i^*i}^{\text{GCC}} \cdot \frac{2e}{e - 1} \cdot \frac{1}{\exp\left(d(P, P_{i^*i}^{\text{GCC}})\right) - 1} \\ &\leq \Delta_{i^*i}^{\text{GCC}} \cdot \frac{2e}{e - 1} \cdot \frac{1}{\Omega(\delta^2)} \,. \end{split}$$

Hence, we have proved an upper bound as

$$\sum_{n_i=0}^T \mathbf{E} \left[ \sum_{t=1}^T \mathbb{1}[a_t = i, N_{ii^*}(t) = n_i, i^* \in A_t] \right] \leq \Delta_{i^*i}^{\text{GCC}} \cdot \frac{2e}{e-1} \cdot \left( \frac{(1+\delta) \log(TnC)}{d(P_{i^*i}^{\text{GCC}}, \frac{1}{2})} + \frac{1}{\Omega(\delta^2)} \right)$$

# F Additional Information About Experimental Setup

# F.1 Synthetic Datasets

In this section we provide additional information about our synthetic datasets.

- MNL-Exp: A MNL model was generated by drawing random weights from the exponential distribution with parameter  $\lambda = 3.5$ , i.e. for arm  $i \in [n]$ ,  $\log v_i$  was sampled i.i.d. from  $\operatorname{Exp}(\lambda = 3.5)$ .
- MNL-Geom: A MNL model was generated with weights  $v_1=e, v_2=e^{\frac{1}{2}},...,v_n=e^{1/2^{n-1}}.$
- GCC-One: For this choice model, we selected arm 1 to be the GCW, and for each set S containing arm 1, we set  $p_{1|S} = 0.51$  and  $p_{i|S} = \frac{0.49}{|S|-1} \ \forall i \in S \setminus \{1\}$ ; for sets S not containing the GCW 1, we selected the smallest-index arm in S to be the highest-probability arm  $i_S^*$  in S, and set  $p_{i_S^*|S} = 0.51$  and  $p_{i|S} = \frac{0.49}{|S|-1} \ \forall i \in S \setminus \{i_S^*\}$ ).
- GCC-Two: For this choice model, we selected arm 1 to be the GCW, and for each set S we defined  $\Delta_S = \min\{\frac{|S|-1}{10}, 0.99\}$ . If  $i^* \notin S$  we selected the smallest-index arm in S to be the highest-probability arm  $i_S^*$  in S, otherwise we let  $i_S^* := i^*$ . We defined  $P_{i_S^*|S} = \frac{1+\Delta_S}{|S|(1-\Delta_S)+2\Delta_S}$  and for any  $i \in S \setminus \{i_S^*\}$ ,  $P_{i|S} = \frac{1-\Delta_S}{|S|(1-\Delta_S)+2\Delta_S}$ .

• GCC-Three: For this choice model, we selected arm 1 to be the GCW, and for each set S we defined  $\Delta_S = \max\{\frac{11-|S|}{11}, 0.01\}$ . Given this definition of  $\Delta_S$ , the choice probabilities we defined in a similar manner as GCC-Two.

Note that the above GCC choice models are similar to the instance constructed in the proof of the lower bound, except that the  $\Delta$  term now depends on the size of the set.

## F.2 Real-World Datasets

In this section we provide additional information for our real-world datasets.

Estimation of choice models from real-world datasets. We estimate choice probabilities from several real-world preference datasets, which contain multiple partial preference orders over items. The choice probability  $P_{i|S}$  of an item i over S, was taken to be the fraction of times in these partial order item i was the top ranked items in S. More formally, let there be m partial orders,  $\mathcal{P}_1, \cdots, \mathcal{P}_m$ , over n items. For any subset  $S \subseteq [n]$ , and  $i \in [n]$ , let  $N_{i|S}$  be defined as:

$$N_{i|S} := \sum_{j \in [m]} \mathbb{1}[\forall i' \in S \setminus \{i\} : i \succ_{\mathcal{P}_j} i'].$$

The choice probability  $P_{i|S}$  is then estimated as:

$$P_{i|S} := \frac{N_{i|S}}{\sum_{i' \in S} N_{i'|S}}.$$

We conducted experiments on three real-world datasets.

- Sushi: This is a dataset from [24] which contains 5000 partial preference orders given by humans over 100 different types of sushis. Similar to [25], we selected a subset of 16 sushi types, such that there exists a GCW among them.
- **IrishMeath:** This is a dataset downloaded from *preflib.org* and contains data about elections held in Dublin, Ireland. The dataset contains 64,081 partial preference orders given by humans over 14 candidates. We selected a subset of 12 candidates, such that there exists a GCW among them.
- **IrishDublin:** This dataset was also downloaded from *preflib.org* and also contains data about elections held in Dublin, Ireland. The dataset contains 29,988 partial preference orders given by humans over 9 candidates. We again selected a subset of 8 candidates, such that there exists a GCW among them.

## F.3 Runtime and Space Complexity of WBA

The space complexity of our algorithm is  $O(n^2)$  as it only stores the pairwise statistics extracted from multiway choices. Each trial of our algorithm runs in time polynomial in n. The most non-trivial step is computing  $\mathcal{J}_i(t,C)$  for each arm. This step requires polynomial time because we can compute the quantity  $\arg\max_{S\subseteq[n]}I_i(t,S)-|S|\cdot\log(nC)$  and check if it is greater than  $\log(t)$ . We compute  $\arg\max_{S\subseteq[n]}I_i(t,S)-|S|\cdot\log(nC)$  by first sorting arms j in the order of values  $\mathbbm{1}[\hat{P}_{ij}(t)\leq\frac{1}{2}]\cdot N_{ij}(t)\cdot d(\hat{P}_{ij}(t),\frac{1}{2})$ . We then start with  $S\leftarrow\emptyset$  and add one arm at a time from this sorted ordering to S. We stop adding arms to the set S once the value  $\mathbbm{1}[\hat{P}_{ij}(t)\leq\frac{1}{2}]\cdot N_{ij}(t)\cdot d(\hat{P}_{ij}(t),\frac{1}{2})$  of the current arm j is less than  $\log(nC)$ . It is easy to see that computing  $I_i(t,S)-|S|\cdot\log(nC)$  for this set S gives the value of  $\arg\max_{S\subseteq[n]}I_i(t,S)-|S|\cdot\log(nC)$ .

## F.4 Hardware Specifications

We ran all our experiments on a 32 core machine with Intel(R) Xeon(R) CPU E5-2683 v4 @ 2.10GHz processor cores. No GPUs were used in the experiments.

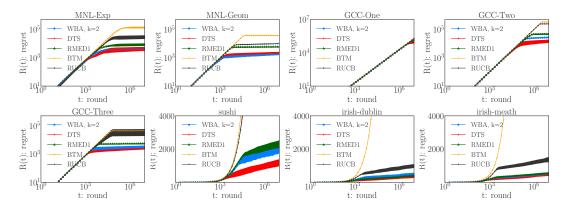


Figure 1: Dueling Bandit Regret ( $R_{DB}$ ) defined in Appendix  $\overline{G}$  v/s trials for our algorithm WBA (for k=2) against dueling bandit algorithms (DTS, BTM, RUCB and RMED1) (the shaded region corresponds to std. deviation). As can be observed, our algorithm is competitive against these algorithms.

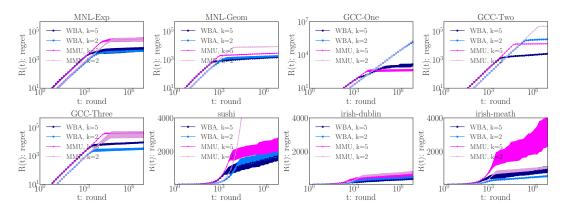


Figure 2: Dueling Bandit Regret  $(R_{\rm DB})$  defined in Appendix  $\overline{\bf G}$  v/s trials for our algorithm WBA against the MaxMinUCB (MMU) algorithm for k=2 and k=5 (the shaded region corresponds to std. deviation). We observe that our algorithm is better than MaxMinUCB on all datasets for both values of k. We further observe that under several datasets the regret achieved by our algorithm for k>2 is better than the regret of our algorithm for k=2.

# **G** Results for Additional Notion of Regret

In the section we define a simple generalization of dueling bandit regret. All our results can be be extended to this notion of regret. Under this notion the regret for an arm is measured as the shortfall in the preference probability in a direct pairwise comparison to the best arm  $i^*$ .

**Definition 1.** For a set  $S \subseteq [n]$ , we define the regret  $r_{DB}(S)$  to be

$$r_{\text{DB}}(S) = \sum_{i \in S} \left( P_{i^* | \{i, i^*\}} - P_{i | \{i, i^*\}} \right). \tag{G.1}$$

This notion of regret allows for a more direct comparison between the regret of a choice bandits algorithm and a dueling bandits algorithm, as the regret for pulling an arm i does not depend on the other arms pulled together with i. Using the definition of GCW  $i^*$ , it is easy to observe that  $r_{DB}(\{i^*\}) = 0$ , and  $0 \le r_{DB}(S) \le |S|$  for any set  $S \subseteq [n]$ .

We present additional experimental results for this notion of regret. Figure 1 contains plots for comparisons with the dueling bandit algorithms, and Figure 2 contains plots for the comparisons of our algorithm with the MaxMinUCB algorithm. The experimental setup was the same as the one described in Section 6 and Appendix F. The overall conclusion with these experiments match the conclusions drawn from the experiments given in Section 6.

## **H** Technical Fact

**Fact 1.** Consider the following experiment: we repeatedly toss (independent) coins from a finite set S of coins with different biases until we get a heads. Let the probability of heads for the i-th coin toss be given by  $p_i \ge 0$ , and the cost be given by  $c_i$ . The expected cost of this experiment is upper bounded as

$$\mathbf{E}\left[\sum_{i=1}^{|S|}\mathbb{1}[\text{no heads till }i-1]\cdot c_i\right] \leq \frac{2c}{p}\cdot \frac{e}{e-1}\,,$$

where  $\frac{c}{p} := \max_{i \in S} \frac{c_i}{p_i}$ 

*Proof.* We will group the sequence of coin tosses such that each group has a total probability mass of at least 1. Formally, group  $G_1$  will consist of the first  $l_1$  coins such that  $\sum_{i=1}^{l_1} p_i \ge 1$  and  $l_1$  is minimized, group  $G_2$  will consist of the next  $l_2$  coins such that  $\sum_{i=l_1+1}^{l_2} p_i \ge 1$  and  $l_2$  is minimized, and so on. The probability that we do not see a head in the first group  $G_1$  is upper bounded as

$$\prod_{i=1}^{l_1} (1 - p_i) \le \prod_{i=1}^{l_1} e^{-p_i} = e^{-\sum_{i=1}^{l_1} p_i} \le e^{-1}.$$

A similar calculation works for each group, showing that we will see a success in a particular group with probability at least 1 - 1/e.

Now, the amount of cost required for each group  $c_G := \sum_{i \in G} c_i$  is upper bounded by  $2 \max_{i \in S} \frac{c}{p}$ . This is due to the fact that each group contains a probability mass of at most 2; and the fact that the maximum cost per a probability mass of p is at most p, hence, the maximum cost per a probability mass of p can be at most p.

$$\begin{split} \mathbf{E}\left[\sum_{i=1}^{\infty}\mathbbm{1}[\text{no heads till }i-1]\cdot c_i\right] &= \mathbf{E}\left[\sum_{j=1}^{\infty}\mathbbm{1}[\text{no heads in group }G_{j-1}]\cdot c_{G_j}\right] \\ &\leq \frac{2c}{p}\cdot (1-\frac{1}{e}) + \frac{4c}{p}\cdot \frac{1}{e}\cdot (1-\frac{1}{e}) + \frac{6c}{p}\cdot \frac{1}{e^2}\cdot (1-\frac{1}{e})\cdots \\ &\leq \frac{2c}{p}\cdot \frac{e}{e-1}\,. \end{split}$$

## References

[1] Brian Brost, Yevgeny Seldin, Ingemar J. Cox, and Christina Lioma. Multi-Dueling Bandits and Their Application to Online Ranker Evaluation. In *Proceedings of the 25th ACM International Conference on Information and Knowledge Management*, 2016.

- [2] Anne Schuth, Harrie Oosterhuis, Shimon Whiteson, and Maarten de Rijke. Multileave Gradient Descent for Fast Online Learning to Rank. In *Proceedings of the 9th ACM International Conference on Web Search and Data Mining*, 2016.
- [3] Yanan Sui, Vincent Zhuang, Joel W. Burdick, and Yisong Yue. Multi-dueling Bandits with Dependent Arms. In *Proceedings* of the 33rd Conference on Uncertainty in Artificial Intelligence, 2017.
- [4] Yi Gai, Bhaskar Krishnamachari, and Rahul Jain. Combinatorial Network Optimization With Unknown Variables: Multi-Armed Bandits With Linear Rewards and Individual Observations. *IEEE/ACM Transactions on Networking*, 20(5):1466–1478, 2012.
- [5] Wei Chen, Yajun Wang, and Yang Yuan. Combinatorial Multi-Armed Bandit: General Framework, Results and Applications. In *Proceedings of the 30th International Conference on Machine Learning*, 2013.
- [6] Branislav Kveton, Zheng Wen, Azin Ashkan, and Csaba Szepesvari. Tight Regret Bounds for Stochastic Combinatorial Semi-Bandits. In *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics*, 2015.
- [7] Richard Combes, M. Sadegh Talebi, Alexandre Proutiere, and Marc Lelarge. Combinatorial Bandits Revisited. In *Advances in Neural Information Processing Systems* 28, 2015.
- [8] Aadirupa Saha and Aditya Gopalan. Combinatorial bandits with relative feedback. In *Advances in Neural Information Processing Systems*, pages 983–993, 2019.
- [9] Masrour Zoghi, Tomas Tunys, Mohammad Ghavamzadeh, Branislav Kveton, Csaba Szepesvari, and Zheng Wen. Online learning to rank in stochastic click models. In *ICML*, pages 4199–4208, 2017.
- [10] Aadirupa Saha and Aditya Gopalan. Battle of bandits. In UAI, pages 805-814, 2018.
- [11] Viktor Bengs and Eyke Hüllermeier. Preselection bandits under the plackett-luce model. CoRR, abs/1907.06123, 2019. URL http://arxiv.org/abs/1907.06123.
- [12] Paat Rusmevichientong, Zuo-Jun Max Shen, and David B. Shmoys. Dynamic Assortment Optimization with a Multinomial Logit Choice Model and Capacity Constraint. *Operations Research*, 58(6):1666–1680, 2010.
- [13] Denis Sauré and Assaf Zeevi. Optimal Dynamic Assortment Planning with Demand Learning. *Manufacturing & Service Operations Management*, 15(3):387–404, 2013.
- [14] Shipra Agrawal, Vashist Avadhanula, Vineet Goyal, and Assaf Zeevi. A Near-Optimal Exploration-Exploitation Approach for Assortment Selection. In *Proceedings of the 17th ACM Conference on Economics and Computation*, 2016.
- [15] Shipra Agrawal, Vashist Avadhanula, Vineet Goyal, and Assaf Zeevi. Thompson Sampling for the MNL-Bandit. In *Proceedings of the 30th Conference on Computational Learning Theory*, 2017.
- [16] Xi Chen and Yining Wang. A Note on Tight Lower Bound for MNL-Bandit Assortment Selection Models. Technical report, arXiv:1709.06109v2, 2017.
- [17] Max Simchowitz, Kevin Jamieson, and Benjamin Recht. Best-of-K Bandits. In Proceedings of the 29th Annual Conference on Learning Theory, 2016.
- [18] Xi Chen, Yuanzhi Li, and Jieming Mao. A Nearly Instance Optimal Algorithm for Top-k Ranking under the Multinomial Logit Model. In Proceedings of the 29th Annual ACM-SIAM Symposium on Discrete Algorithms, 2018.
- [19] Yisong Yue, Josef Broder, Robert Kleinberg, and Thorsten Joachims. The K-armed Dueling Bandits Problem. In *Proceedings* of the 22nd Conference on Learning Theory, 2009.

- [20] Emilie Kaufmann, Olivier Cappe, and Aurelien Garivier. On the Complexity of Best-Arm Identification in Multi-Armed Bandit Models. *Journal of Machine Learning Research*, 17(1):1–42, 2016.
- [21] Aadirupa Saha and Aditya Gopalan. PAC battling bandits in the plackett-luce model. In *Algorithmic Learning Theory, ALT* 2019, 22-24 March 2019, Chicago, Illinois, USA, pages 700–737, 2019.
- [22] Aurélien Garivier and Olivier Cappé. The KL-UCB algorithm for bounded stochastic bandits and beyond. In COLT 2011 The 24th Annual Conference on Learning Theory, June 9-11, 2011, Budapest, Hungary, pages 359–376, 2011.
- [23] Junpei Komiyama, Junya Honda, and Hiroshi Nakagawa. Optimal Regret Analysis of Thompson Sampling in Stochastic Multi-armed Bandit Problem with Multiple Plays. In Proceedings of the 32nd International Conference on Machine Learning, 2015.
- [24] Toshihiro Kamishima. Nantonac collaborative filtering: recommendation based on order responses. In *Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 24 27, 2003*, pages 583–588, 2003.
- [25] Junpei Komiyama, Junya Honda, Hisashi Kashima, and Hiroshi Nakagawa. Regret Lower Bound and Optimal Algorithm in Dueling Bandit Problem. In *Proceedings of the 28th Conference on Learning Theory*, 2015.