

Dueling Bandits with Adversarial Sleeping

Aadirupa Saha*

Pierre Gaillard †

January 18, 2022

Abstract

We introduce the problem of sleeping dueling bandits with stochastic preferences and adversarial availabilities (DB-SPAA). In almost all dueling bandit applications, the decision space often changes over time; eg, retail store management, online shopping, restaurant recommendation, search engine optimization, etc. Surprisingly, this ‘sleeping aspect’ of dueling bandits has never been studied in the literature. Like dueling bandits, the goal is to compete with the best arm by sequentially querying the preference feedback of item pairs. The non-triviality however results due to the non-stationary item spaces that allow any arbitrary subsets items to go unavailable every round. The goal is to find an optimal ‘no-regret’ policy that can identify the best available item at each round, as opposed to the standard ‘fixed best-arm regret objective’ of dueling bandits. We first derive an instance-specific lower bound for DB-SPAA $\Omega(\sum_{i=1}^{K-1} \sum_{j=i+1}^K \frac{\log T}{\Delta(i,j)})$, where K is the number of items and $\Delta(i, j)$ is the gap between items i and j . This indicates that the sleeping problem with preference feedback is inherently more difficult than that for classical multi-armed bandits (MAB). We then propose two algorithms, with near optimal regret guarantees. Our results are corroborated empirically.

1 Introduction

The problem of *Dueling-Bandits* has gained much attention in the machine learning community [34, 38, 36], which is an online learning framework that generalizes the standard multiarmed bandit (MAB) [6] setting for identifying a set of ‘good’ arms from a fixed decision-space (set of arms/items) by querying preference feedback of actively chosen item-pairs. More formally, in dueling bandits, the learning proceeds in rounds: At each round, the learner selects a pair of arms and observes stochastic preference feedback of the winner of the comparison (duel) between the selected arms; the objective of the learner is to minimize the regret with respect to a (or set of) ‘best’ arm(s) in hindsight. Towards this several algorithms have been proposed [2, 37, 21, 14]. Due to the inherent exploration-vs-exploitation tradeoff of the learning framework and several advantages of preference feedback [9, 35], many real-world applications can be modeled as dueling bandits, including movie recommendations, retail management, search engine optimization, job scheduling, etc.

However, in reality, the decision spaces might often change over time due to the non-availability of some items, which are considered to be ‘sleeping’. This ‘sleeping-aspect’ of online decision making problems has been widely studied in the standard multiarmed bandit (MAB) literature [17, 24, 15, 19, 18, 11]. There the goal is to learn a ‘no-regret’ policy that maps to the ‘best awake item’ of any available (non-sleeping) subset of items, and the learner’s performance is measured with respect to the optimal policy in hindsight. This setting is famously known as *Sleeping Bandits* in MAB [17, 24, 15, 11]. More discussions are given in Related Works.

Surprisingly, however, the ‘*sleeping problem*’ is completely unaddressed in the preference bandits literature, even for the special case of pairwise preference feedback, which is famously studied as *Dueling Bandits* [37, 34], even though the setup of changing decision spaces are quite relevant in almost every practical applications: Be that in retail stores where some items might go out of production over time, for search engine optimization

*Indian Institute of Science, Bangalore, India; aadirupa@iisc.ac.in.

†Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, 38000 Grenoble, France. pierre.gaillard@inria.fr

Parameters. Item set: $[K]$ (known), Preference: \mathbf{P} (unknown), Available item sets: \mathcal{S}_T (observed sequentially)

For $t = 1, 2, \dots, T$, the learner:

- Observes $S_t \subseteq [K]$ the set of available items
- Chooses $(x_t, y_t) \in S_t^2$
- Observes $o_t := \mathbf{1}(x_t \succ y_t) \sim \text{Ber}(\mathbf{P}(x_t, y_t))$
- Incurs $r_t := 1/2(\mathbf{P}(i_t^*, x_t) + \mathbf{P}(i_t^*, y_t) - 1)$;
where i_t^* is such that $\min_{j \in S_t} \mathbf{P}(i_t^*, j) \geq 1/2$

Figure 1: Setting of DB-SPAA($\mathbf{P}, \mathcal{S}_T$)

some websites could be down on certain days, in recommender systems some restaurants might be closed or movies could be outdated, in clinical trials certain drugs could be out of stock, and many more. This work is the first to consider the problem of *Sleeping Dueling Bandits*, where we formulated the stochastic K -armed dueling bandit problem with adversarial item availabilities. Here at each round $t \in \{1, 2, \dots, T\}$ the item preferences are considered to be generated from a fixed underlying (and of course unknown) preference matrix $\mathbf{P} \in [0, 1]^{K \times K}$, however, the set of available actions $S_t \subseteq \{1, 2, \dots, K\}$ is assumed to be adversarially chosen by the environment. We call the problem as *Sleeping-Dueling Bandit with Stochastic Preferences and Adversarial Availabilities* or in brief DB-SPAA($\mathbf{P}, \mathcal{S}_T$), where $\mathcal{S}_T = \{S_1, S_2, \dots, S_T\}$ denotes the sequence of available subsets over T rounds. We also assume the preference \mathbf{P} follows a ‘total-ordering assumption to ensure the existence of a best-item per available subset S_t . We describe the setting in Fig. 1 with a formal description in Sec. 2.

Our specific contributions are as follows:

1. We first analyze the fundamental performance limit for the DB-SPAA($\mathbf{P}, \mathcal{S}_T$) problem in Sec. 3: Thm. 1 gives an instance-specific regret lower bound of

$$\Omega\left(\sum_{i=1}^{K-1} \sum_{j=i+1}^K \frac{\log T}{\Delta(i,j)}\right),$$

with $\Delta(i, j)$ being the ‘preference gap’ of item i -vs- j (see Eqn. (2)). Our lower bound, which can be of order $\Omega(K^2 \log T / \Delta)$, $\Delta := \min_{i,j} \Delta(i, j)$ being the worst case gap, indicates that the *problem of sleeping dueling bandits is inherently more difficult than standard sleeping bandits (MAB)*, unlike the ‘non-sleeping’ case where both dueling bandits (with ‘total-ordering’ assumption on \mathbf{P}) and MAB are known to have the same fundamental performance limit of $\Omega(K \log T / \Delta)$ (Rem. 1).

2. We next design a ‘fixed confidence regret’ algorithm S1DB-UCB (Alg. 1), inspired from the pairwise upper confidence bound (UCB) based algorithm [37]. However due to the fixed confidence and ‘adversarial-sleeping’ nature of the problem, we need to differently maintain pairwise confidence bounds per item (based on the availability sequence $\{S_t\}$), which makes the resulting algorithm and its subsequent analysis significantly different than standard UCB based dueling bandit algorithms: Precisely given any $\delta > 0$, S1DB-UCB achieves a regret of $O\left(\frac{K^3 \log(1/\delta)}{\Delta^2}\right)$ with probability at least $1 - \delta$ over any problem instance of DB-SPAA($\mathbf{P}, \mathcal{S}_T$) (Sec. 4).

3. In Sec. 5, we design another computationally efficient algorithm, S1DB-ED (Alg. 2), for ‘expected regret’ guarantee. Unlike the previous algorithm (S1DB-UCB), S1DB-ED uses empirical divergence (ED) based measures to filter out the ‘good’ set of arms, inspired from the idea of RMED algorithm of [21] for standard dueling bandits; however, due to sleeping nature of the items, it requires a different maintenance of ‘good’ arms and the regret analysis of the algorithm requires derivation of new results (as described in Sec. 5). The algorithm is shown to perform near optimally with an expected *non-asymptotic* regret upper-bound of (Thm. 6). Note that for any problem instance with constant suboptimality gaps $\Delta(i, j) = \Delta$ for all $i < j$, regret bound of S1DB-ED is tight and matches the lower-bound ensuring the near optimality of S1DB-ED in the worst case. Furthermore, a novelty of our finite time regret analysis lies in showing a cleaner tradeoff between regret vs. availability sequence \mathcal{S}_T which automatically adapts to the inherent ‘hardness’ of the

sequence of available subsets S_T , compared to existing sleeping bandits work for adversarial availabilities in the MAB setting [19] which only gives a worst-case regret bound over all possible availability sequences (Rem. 3).

4. Finally we corroborate our theoretical results with extensive empirical evaluations. (Sec. 6).

Related Works. The problem of regret minimization for stochastic multiarmed bandits (MAB) is extremely well studied in the online learning literature [6, 1, 22, 5, 16], where the learner gets to see a noisy draw of absolute reward feedback of an arm upon playing a single arm per round.

A well motivated generalization of MAB framework is *Sleeping Bandits* [17, 24, 18, 15], much studied in the online learning community, where at any round the set of available actions could vary stochastically based on some unknown distributions over the decision space of K items [24, 11] or adversarially [15, 19, 18]. Besides the reward model, the set of available actions could also vary stochastically or adversarially [17, 24]. The problem is NP-hard when both rewards and availabilities are adversarial [19, 18, 15]. In case of stochastic reward and adversarial availabilities [19] proposed an UCB based no-regret algorithm, which was also shown to be provably optimal. The case of adversarial reward and stochastic availabilities has also been studied where the achievable regret lower bound is known to be $\Omega(\sqrt{KT})$ by the inefficient EXP4 algorithm [19, 15].

On the other hand over the last decade, the relative feedback variants of stochastic MAB problem has seen a widespread resurgence in the form of the Dueling Bandit problem, where, instead of getting noisy feedback of the reward of the chosen arm, the learner only gets to see a noisy feedback on the pairwise preference of two arms selected by the learner. The objective of the learner is to minimize the regret with respect to ‘best arm in the stochastic model’. Several algorithms have been proposed to address this dueling bandits problem, for different notions of ‘best arms’ or preference models [10, 31, 38, 37, 36, 21, 33, 13], or even extending the pairwise preference to subsetwise preferences [29, 8, 26, 27, 25]. However, surprisingly, unlike the ‘sleeping bandits generalization’ of MAB, no parallel has been drawn for dueling bandits, which remains our main focus.

2 Problem Formulation

Notations. Decision space (or item/arm set) $[K] := \{1, 2, \dots, K\}$. The available set of items at round t is denoted by $S_t \subseteq [K]$. For any matrix $\mathbf{M} \in \mathbb{R}^{K \times K}$, we define $m_{ij} := M(i, j)$, $\forall i, j \in [K]$. We write $S_{\setminus i} = S \setminus \{i\}$, for any $S \subseteq [K]$ and $i \in S$. $\mathbf{1}(\cdot)$ denotes the indicator random variable which takes value 1 if the predicate is true and 0 otherwise and \lesssim a rough inequality which holds up to universal constants. For any two items $x, y \in [K]$, we use the symbol $x \succ y$ to denote x is preferred over y . Σ_K denotes the set of all permutations of the items in set $[K]$. The KL-divergence of two Bernoullis with biases p and q respectively is written $\text{kl}(p, q) := p \log(p/q) + (1-p) \log((1-p)/(1-q))$. We assume $\frac{0}{0} := 0.5$ (in Alg. 1 and 2).

Setup. We consider the problem of stochastic K -armed dueling bandits with adversarial availabilities: At every iteration $t = 1, \dots, T$, a set of available items (actions) $S_t \subseteq [K]$ is revealed, and the learner is asked to choose two items $x_t, y_t \in S_t$. Then, the learner receives a preference feedback $o_t = \mathbf{1}(x_t \succ y_t) \sim \text{Ber}(\mathbf{P}(x_t, y_t))$, where $\mathbf{P} \in [0, 1]^{K \times K}$ is an underlying pairwise preference matrix, unknown to the learner. The setting is described in Figure 1. We assume that \mathbf{P} respects a ‘total ordering’, say $\sigma^* \in \Sigma_K$. Without loss of generality, we set $\sigma^* = (1, 2, \dots, K)$ throughout the paper. This implies $\mathbf{P}(i, j) \geq 0.5$ for $i \leq j$. One possible pairwise probability model which respects ‘total ordering’ is Plackett-Luce [7], where it is assumed that the K items are associated to positive score parameters $\theta_1, \dots, \theta_K$, and $\mathbf{P}(i, j) = \theta_i / (\theta_i + \theta_j)$ for all $i, j \in [K]$. In fact any well random utility (RUM) based preference model would have the above property, like [7, 28]. Note also that our assumption corresponds to assuming the existence of a Condorcet winner for every subset $S_t \subseteq [K]$.

Objective. The objective of the learner is to minimize his regret over T rounds with respect to the best policy in the policy class $\Pi = \{\pi : 2^K \mapsto [K] \mid \forall t \in [T], \pi(S_t) \in S_t\}$, i.e. any $\pi \in \Pi$ is such that for any $t \in [T]$, $\pi(S_t) \in S_t$. More formally we define the regret as follows:

$$R_T = \max_{\pi \in \Pi} \sum_{t=1}^T \frac{\mathbf{P}(\pi(S_t), x_t) + \mathbf{P}(\pi(S_t), y_t) - 1}{2}. \quad (1)$$

We analyze both *fixed-confidence* and *expected* regret guarantees in this paper respectively in Sec. 4 (see Thm. 3) and Sec. 5 (see Thm. 6). It is easy to note that under our preference modelling assumptions, the best policy, say π^* , turns out to be $\pi^*(S) = \min\{S\}$ for any $S \subseteq [K]$. We henceforth denote by $i_t^* = \pi^*(S_t)$. We define the above problem to be *Sleeping-Dueling Bandit with Stochastic Preferences and Adversarial Availabilities* over the stochastic preference matrix $\mathbf{P} \in [0, 1]^{K \times K}$ and the sequence of available subsets $\mathcal{S}_T = \{S_1, \dots, S_T\}$, or in short **DB-SPAA**($\mathbf{P}, \mathcal{S}_T$). For ease of notation we respectively define the gaps and the non-zeros gaps as $\Delta(i, j) := \mathbf{P}(i, j) - 1/2$, and

$$\Delta(i, j)_+ := \begin{cases} \Delta(i, j) & \text{if } \Delta(i, j) \neq 0 \\ +\infty & \text{if } \Delta(i, j) = 0 \end{cases} \quad (2)$$

The regret thus can be rewritten as $R_T := \sum_{t=1}^T r_t$, where $r_t := (\Delta(i_t^*, x_t) + \Delta(i_t^*, y_t))/2$ denotes the instantaneous regret. We also denote by $n_{ij}(t) := \sum_{\tau=1}^t \mathbf{1}(\{x_\tau, y_\tau\} = \{i, j\})$ the number of times the pair (i, j) is played until time t and by $w_{ij}(t)$ the number of times i beats j in t rounds.

3 Lower Bound

We first derive a worst case regret lower bound over all possible sequences of \mathcal{S}_T . The proof idea essentially lies in constructing *hard enough* availability sequences \mathcal{S}_T , where no learner can escape learning the preferences of every distinct pair of items (i, j) . This leads to a potential lower bound of $\Omega(K^2 \log(T)/\Delta)$. For this section we denote \mathbf{P} by \mathbf{P}_K to make the dependency on K more precise.

Theorem 1 (Lower Bound for **DB-SPAA**($\mathbf{P}_K, \mathcal{S}_T$)). *For any No-regret learning algorithm \mathcal{A} , there exists a problem instance **DB-SPAA**($\mathbf{P}_K, \mathcal{S}_T$) with $T \geq K^4$, such that its expected regret is lower-bounded as:*

$$\mathbf{E}[R_T(\mathcal{A})] \geq \Omega\left(\sum_{i=1}^{K-1} \sum_{j=i+1}^K \frac{\log T}{\Delta(i, j)_+}\right).$$

The *No-regret learning algorithm* refers to the class of ‘consistent algorithms’ which do not pull any suboptimal pair more than $O(T^\alpha)$, $\alpha \in [0, 1]$ (see Def. 7, Appendix 8).

Proof (sketch) The main argument lies behind the fact that in the worst case the adversary can force the algorithm to learn the preference of every distinct pair (i, j) as for the ‘worst-case’ sequences \mathcal{S}_T , a knowledge of the already ‘learnt’ pairwise preferences would not disclose any information on the remaining pairs; e.g. assuming $\sigma^* = (1, 2, \dots, K)$, revealing the available subsets in the following sequence $(1, 2), (1, 3), \dots, (1, K), (2, 3), (2, 4), \dots, (K-1, K)$ would force the learner to explore (learn the preferences) all $\binom{K}{2}$ distinct pairs. The remaining proof establishes this formally, towards which we first show a $\Omega(\ln(T)/\Delta(1, 2))$ regret lower bound for a **DB-SPAA** instance with just two items (i.e. $K = 2$) as shown in Lem. 2. The lower bound for any general K can now be derived applying the above bound on independent $\binom{K}{2}$ subintervals, with the availability sequence $(1, 2), (1, 3), \dots, (1, K), (2, 3), (2, 4), \dots, (K-1, K)$. The full proof is given in Appendix 8. \square

Lemma 2 (Lower Bound of **DB-SPAA**($\mathbf{P}_K, \mathcal{S}_T$) for 2 items). *For any No-regret learning algorithm \mathcal{A} , there exists a problem instance **DB-SPAA**($\mathbf{P}_2, \mathcal{S}_T$) such that the expected regret incurred by \mathcal{A} on that can be lower bounded as: $\mathbf{E}[R_T(\mathcal{A})] \geq \Delta^{-1} \log(T)$, Δ being the ‘preference-gap’ between the two items (i.e. $\Delta = \mathbf{P}_{12} - 1/2$, assuming $P_{12} > 1/2$ or equivalently $\Delta > 0$).*

Remark 1 (Implication of the lower bound). The above result indicates that in the preference based learning setup, the fundamental problem complexity lies in distinguishing every pair of items $1 \leq i < j \leq K$. If the

Algorithm 1 S1DB-UCB

1: **input:** Arm set: $[K]$, parameters $\alpha > 0.5$, Confidence parameter $\delta \in [0, 1]$
2: **init:** $w_{ij}(1) \leftarrow 0, \forall i, j \in [K]$.
3: **define:** $n_{ij}(t) := w_{ij}(t) + w_{ji}(t), \forall t \in [T]$
4: **for** $t = 1, 2, \dots, T$ **do**
5: Receive $S_t \subseteq [K]$
6: $\hat{p}_{ij}(t) = \frac{w_{ij}(t)}{n_{ij}(t)}, c_{ij}(t) \leftarrow \sqrt{\frac{\alpha \log a_{ij}(t)}{n_{ij}(t)}}, \forall i, j \in S_t$, (assume $\frac{x}{0} := 0.5, \forall x \in \mathbb{R}$)
7: $u_{ij}(t) \leftarrow \hat{p}_{ij}(t) + c_{ij}(t), u_{ii}(t) \leftarrow 1/2, \forall i, j \in S_t$ ▷ UCB of empirical preferences
8: **for** $k \in S_t$ **do**
9: $\mathcal{C}_k(t) := \{j \in S_t \mid u_{kj}(t) > 1/2\}$ ▷ Potential losers to k
10: **end for**
11: $\mathcal{C}_t = \{i \in S_t \mid |\mathcal{C}_i(t)| = \max_{j \in S_t} |\mathcal{C}_j(t)|\}$ ▷ Potential best items
12: Select a random x_t from \mathcal{C}_t . Choose $y_t \leftarrow \arg \max_{i \in \mathcal{C}_t} u_{ix_t}(t)$
13: Play (x_t, y_t) . Receive preference o_t
14: Update: $\forall i, j \in [K], w_{x_t y_t}(t+1) \leftarrow w_{x_t y_t}(t) + o_t, w_{y_t x_t}(t+1) \leftarrow w_{y_t x_t}(t) + (1 - o_t), a_{ij}(t+1) \leftarrow \max\{n_{ij}(t+1), C(K, \delta)\}, \forall i, j \in [K]$
15: **end for**

learner fails to learn the preference of any pair (i, j) , the adversary can make the learner suffer $O(T)$ regret by setting $S_t = \{i, j\}$ henceforth at all round. It is worth noting that in the ‘no sleeping’ case both dueling bandits and MAB are known to have the same fundamental performance limit of $\Omega(K \log(T)/\Delta)$ (assuming \mathbf{P} respects a *condorcet winner* [6, 21]). Thus Thm. 1 shows that the ‘sleeping-aspect’ of dueling bandits makes the problem K -times harder than ‘sleeping-MAB’ for which the regret lower bound is known to be only $\Omega(\sum_{i=1}^K \log(T)/\Delta(i, i+1))$ [19].

4 S1DB-UCB: A Fixed-Confidence Algorithm

In this section, we design an efficient algorithm for the DB-SPAA(K, T) problem with instance-dependent regret guarantee.

Main ideas. Our algorithm, described in Alg. 1, depends on an hyper-parameter $\alpha > 0.5$ and a confidence parameter $\delta > 0$. It maintains, for each item $k \in [K]$, its own record of empirical pairwise estimates of the duels, $(i, j) \in [K] \times [K]$ and their respective upper confidence bounds defined as:

$$\hat{p}_{ij}(t) := \frac{w_{ij}(t)}{n_{ij}(t)} \quad \text{and} \quad u_{ij}(t) := \hat{p}_{ij}(t) + c_{ij}(t), \quad \text{with} \quad c_{ij}(t) := \sqrt{\frac{\alpha \log a_{ij}(t)}{n_{ij}(t)}},$$

where $w_{ij}(t)$ denotes the total number of times item i beats j up to round t , $n_{ij}(t) := w_{ij}(t) + w_{ji}(t)$, and for all $i, j \in [K]$ and $t \in [T]$

$$a_{ij}(t) := \max\{C(K, \delta), n_{ij}(t)\} \quad \text{and} \quad C(K, \delta) := \left(\frac{(4\alpha-1)K^2}{(2\alpha-1)\delta}\right)^{\frac{1}{2\alpha-1}}.$$

A key observation is that our careful choice of the confidence bounds $c_{ij}(t)$ ensures that with high probability $p_{ij}(t) \in [\hat{p}_{ij}(t) - c_{ij}(t), \hat{p}_{ij}(t) + c_{ij}(t)]$ for any duel $i, j \in [K]$ and any $t \in [T]$ (Lem. 4). Now at any round $t \geq 1$, the algorithm first computes a set of potential winners of S_t as $\mathcal{C}_t = \{k \in S_t \mid |\mathcal{C}_k(t)| = \max_{j \in S_t} |\mathcal{C}_j(t)|\}$, where $\mathcal{C}_k(t) := \{j \in S_t \mid u_{kj}(t) > \frac{1}{2}\}$ denotes the set of items that item k dominates (optimistically). At each round, we play a random item from the set potential winners \mathcal{C}_t as the left arm x_t . Finally the right-arm y_t is chosen to be the most competitive opponent of x_t as $y_t \leftarrow \arg \max_{i \in \mathcal{C}_t} u_{ji}(t)$ from the potential winners. Our arm selection strategy ensures that eventually for all t , algorithm plays the optimal pair (i_t^*, i_t^*) frequently enough as desired.

Theorem 3 (Fixed-confidence regret analysis: S1DB-UCB). *Given any $\delta > 0$ and $\alpha > 1/2$, with probability at least $1 - \delta$, the regret incurred by S1DB-UCB (Alg. 1) is upper-bounded as:*

$$R_T \leq 2 \sum_{i=1}^{K-1} \sum_{j=i+1}^K M_{ij} \log(2C(K, \delta)M_{ij})$$

where

$$C(K, \delta) := \left(\frac{(4\alpha - 1)K^2}{(2\alpha - 1)\delta} \right)^{\frac{1}{2\alpha - 1}} \quad \text{and} \quad M_{ij} = \sum_{k=1}^i \frac{4\alpha}{\min\{\Delta(k, i)_+, \Delta(k, j)_+\}^2}.$$

The complete proof with a precise dependencies on the model parameters is deferred to Appendix 9.

Remark 2. *The dependency on $\Delta = \min_{i,j} \Delta_+(i, j)$ does not match the lower-bound of Thm. 1, which is of order $O(\log(T)/\Delta)$. Instead, Thm. 3 proves $O(\log(1/\delta)/\Delta^2)$. Yet, the bounds are not directly comparable because the lower-bound is on the expected regret while the upper-bound considers fixed-confidence δ and is hence independent of T . All existing dueling bandit algorithms, that minimize the expected regret, suffer an additional constant term of order $O(1/\Delta^2)$ –see for instance [37, 21]. Achieving an order $O(1/\Delta)$ dependence is an interesting question for future work.*

Proof sketch of Thm. 3. The key steps lie in proving the following four lemmas. The first lemma follows along the line of Lem. 1 of RUCB algorithm [37] and shows that all the pairwise estimates are contained within their respective confidence intervals with high probability.

Lemma 4. *Let $\alpha > 0.5$ and $\delta > 0$. Then, for any $i, j \in [K]$, with probability at least $1 - \delta/K^2$,*

$$\hat{p}_{ij}(t) - c_{ij}(t) \leq p_{ij} \leq u_{ij}(t) := \hat{p}_{ij}(t) + c_{ij}(t), \quad \forall t \in [T].$$

The lemma below shows that once the algorithm can not play any suboptimal pair ‘too many times’.

Lemma 5. *Let $\alpha > 0.5$. Under the notations and the high-probability event of Lem. 4, for all $i, j, k \in [K]$ such that $\{i, j\} \neq \{k, k\}$, and for any $\tau \geq 1$*

$$\sum_{t=1}^{\tau} \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) \leq \frac{4\alpha \log a_{i,j}(\tau)}{\min\{\Delta(k, i)_+, \Delta(k, j)_+\}^2},$$

where recall $a_{ij}(\tau) = \max(C(K, \delta), n_{ij}(\tau))$.

With probability of at least $1 - \delta$, the event of Lem. 4 holds and thus so do the ones of Lem. 5. The regret of Alg. 1 then follows from applying the above lemmas with the following careful decomposition of the regret:

$$R_T \leq \sum_{i=1}^{K-1} \sum_{j=i+1}^K \sum_{t=1}^T \sum_{k=1}^i \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) = \sum_{i=1}^{K-1} \sum_{j=i+1}^K n_{ij}(T)$$

and the proof is concluded by using Lemma 5 to upper-bound $n_{ij}(T)$. The complete proof given in Appendix 9. \square

5 S1DB-ED: An Expected Regret Algorithm

In this section, we propose another computationally efficient algorithm, S1DB-ED (Alg. 2), which achieves near-optimal expected-regret for DB-SPAA problem, and also performs competitively against S1DB-UCB empirically (see Sec. 6). Furthermore, a novelty of our finite time regret analysis of S1DB-ED lies in showing a cleaner trade-off between regret vs availability sequence \mathcal{S}_T which automatically adapts to the inherent ‘hardness’ of the sequence of available subsets \mathcal{S}_T , unlike the previous attempts made in standard sleeping bandits for adversarial availabilities [19] (Rem. 3).

Algorithm 2 S1DB-ED

1: **input:** Arm set: $[K]$, exploration parameter $t_0 > 0$, parameter $\alpha > 0$
 2: **for** $t = 1, 2, \dots, T$ **do**
 3: $\hat{p}_{ij}(t) := \frac{w_{ij}(t)}{n_{ij}(t)}$, $\hat{p}(i, i) \leftarrow 1/2$, $\forall i, j \in [K]$ (assume $\frac{x}{0} := 0.5$, $\forall x \in \mathbb{R}$)
 4: Receive $S_t \subseteq [K]$
 5: **if** $|S_t| \geq 2$ and $\exists i, j \in S_t$ s.t. $n_{ij}(t) < t_0$, $i \neq j$ **then**
 6: Set $x_t \leftarrow i$, $y_t \leftarrow j$ ▷ Exploration rounds
 7: **else**
 8: $\hat{\mathcal{B}}_i(t) := \{j \in [K] \setminus \{i\} \mid \hat{p}_{ij}(t) \leq 1/2\} \cap S_t$, $\forall i \in [K]$ ▷ Empirical winners over i
 9: $\mathcal{I}_i(t) := \sum_{j \in \hat{\mathcal{B}}_i(t)} n_{ij}(t) \text{kl}(\hat{p}_{ij}(t), 1/2)$, $\forall i \in [K]$ and $\hat{i}_t^* \leftarrow \arg \min_{i \in S_t} \mathcal{I}_i(t)$
 10: $\mathcal{C}_t := \{i \in S_t \mid \mathcal{I}_i(t) - \mathcal{I}_{\hat{i}_t^*}(t) \leq \alpha \log t\}$ ▷ Potential good arms
 11: Select any x_t from \mathcal{C}_t uniformly at random
 12: **if** ($\hat{i}_t^* \in \hat{\mathcal{B}}_{x_t}(t)$ or $\hat{\mathcal{B}}_{x_t}(t) = \emptyset$): set $y_t \leftarrow \hat{i}_t^*$, **else:** $y_t \leftarrow \arg \max_{i \in S_t \setminus \{x_t\}} \hat{p}_{ix_t}(t)$
 13: **end if**
 14: Play (x_t, y_t) Receive preference feedback o_t
 15: **end for**

Main ideas. We again use the notations $w_{ij}(t), n_{ij}(t)$ as used for S1DB-UCB (Alg. 1), with the same initializations. Same as S1DB-UCB, this algorithm also maintains the empirical pairwise preferences $\hat{p}_{ij}(t)$ for each item pair $i, j \in [K]$. However, unlike the earlier case here we need to ensure an initial t_0 rounds of exploration ($t_0 = 1$ in the theorem) for every distinct pairs (i, j) , and instead of maintaining pairwise UCBs, in this case the set of ‘good-items’ is defined in terms of *empirical divergences* for all $i \in S_t$

$$\mathcal{I}_i(t) := \sum_{j \in \hat{\mathcal{B}}_i(t)} n_{ij}(t) \text{kl}(\hat{p}_{ij}(t), 1/2), \quad \hat{\mathcal{B}}_i(t) := \left\{ j \in [K] \setminus \{i\} \mid \hat{p}_{ij}(t) \leq 1/2 \right\} \cap S_t$$

denotes the empirical winners of item i in set S_t . Now intuitively since $\exp(-\mathcal{I}_i(t))$ can be interpreted as the likelihood of i being the *best-item* of S_t , we denote by $\hat{i}_t^* \leftarrow \arg \min_{i \in S_t} \mathcal{I}_i(t)$ the *empirical-best* item of round t and define the set of ‘near-best’ items $\mathcal{C}_t := \{i \in S_t \mid \mathcal{I}_i(t) - \mathcal{I}_{\hat{i}_t^*}(t) \leq \alpha \log t\}$, whose likelihood is close enough to that of \hat{i}_t^* . Finally the algorithm selects an arm pair (x_t, y_t) such that x_t is a potential candidate of good arm (which ensures the required exploration) and y_t being the strongest challenger of x_t w.r.t the empirical preferences. The algorithm is given in Alg. 2.

Theorem 6 (Expected regret analysis S1DB-ED). *Let $t_0 = 1$ and $\alpha = 4K$. Then as $T \rightarrow \infty$, the expected regret incurred by S1DB-ED (Alg. 2) can be upper bounded as: For all $\varepsilon_2, \dots, \varepsilon_K \geq 0$*

$$\begin{aligned} \mathbf{E}[R_T] &\lesssim K^2 + \sum_{1 \leq i < j \leq K} \left(\frac{K \mathbf{1}_{\{\Delta(i,j) > \varepsilon_j\}}}{\Delta(i,j)^2} + n_{ij}(T) \min\{\varepsilon_j, \Delta(i,j)\} \right) + \sum_{j=2}^K \frac{K \log T}{\max\{\varepsilon_j, \Delta(j-1, j)_+\}} \\ &\leq O\left(\min \left\{ \sum_{j=2}^K \frac{K \log T}{\Delta(j-1, j)_+}, KT^{2/3} \right\} \right). \end{aligned}$$

The proof is deferred to Appendix 10. Although it borrows some high-level ideas from [19] for sleeping bandits and from [21] for RMED in standard dueling bandits, our analysis needed new ingredients in order to obtain $O(K^2(\log T)/\Delta)$. This is especially the case for the proofs of the technical Lemmas 8 and 9 which significantly differ from ‘‘corresponding’’ technical lemmas of [21]. Specifically, both regret bounds of RMED and ours need to control the length of an initial exploration t_0 after which pairwise preferences are well estimated by $\hat{p}_{ij}(t)$. This is done respectively by our Lemma 8 and Lemma 5 of [21]. Yet, RMED’s original analysis is based on a union bound over all possible subsets $S \subset \{1, \dots, K\}$ of items (see Equation (19) in [21]), whose number is exponential in K . Despite our efforts, we could not follow the proof of Lemma 5

of [21], which to the best of our understanding, should yield to an exploration t_0 exponentially large in K contrary to $O(1)$ claimed in [21]. Instead, in our proof of Lem. 8, we carefully apply concentration inequalities to run union bounds over the items directly instead of sets of items.

The upper-bound of Thm. 6 is close to optimal. It suffers at most a suboptimal factor K and exactly matches the lower-bound for some problems. The distribution-free upper-bound of order $O(T^{2/3})$ matches obtainable standard dueling bandit problems [21, 37], since the latter algorithms also suffer constant terms of order Δ^{-2} . Yet, it is unclear whether it is optimal or if $O(\sqrt{T})$ can be obtained.

Remark 3 (Sequence \mathcal{S}_T adaptivity of Alg. 2). It is worth pointing out that the regret bound of Thm. 6 is finite time and automatically adapts to the sequence of available sets \mathcal{S}_T . In the worst-case, the complexity lies in identifying for all items j the gap with the earlier item $j-1$. Yet, our regret-bound, which holds for any $\varepsilon_j \geq 0$, will automatically perform a trade-off for each j between the gap $\Delta(j-1, j)_{+1}^{-1}$ and $\varepsilon_j \sum_{i=1}^{j-1} n_{ij}(T)$ the number of times j is played together with a better item $i < j$. In particular, $\sum_{i=1}^{j-1} n_{ij}(T)$ can be small if j is rarely available in S_t while not optimal. Notably, this adaptivity to \mathcal{S}_T item per item improves the regret guarantee of Thm. (10), [19], which also addresses the problem of sleeping bandits with ‘adversarial availabilities’ but for the stochastic multi-armed bandit setup and only provides worst-case guarantees over all \mathcal{S}_T and a trade-off ε independent of j .

Remark 4 (S1DB-ED in standard dueling bandits). Even in the dueling bandit setting (without the sleeping component), S1DB-ED and Thm. 6 have advantages compared to the RMED algorithm and analysis of [21]. Our regret bound is valid for all number of items K , while the one of Thm. 3 of [21] is only asymptotic when $K \rightarrow \infty$. This is due to the fact that the algorithm of [21] depends on a hyper-parameter $f(K)$ which needs to be larger than AK , where A is a constant in K and T but which depends on the unknown sub-optimality gaps $\Delta(i, j)$. Thus, [21] chooses $f(K) \approx K^{1+\varepsilon}$ so that eventually the bound is satisfied when $K \rightarrow \infty$. Instead, our algorithm only depends on easily tunable hyper-parameters t_0 and α , whose optimal values are independent of unknown parameters.

6 Experiments

In this section, we compare the empirical performances of our two proposed algorithms (Alg. 1 and 2). Note that there are no other existing algorithms for our problem (see Sec. 1).

Constructing Preference Matrices (P). We use the following three different utility based Plackett-Luce(θ) preference models (see Sec. 2) that ensures a *total-ordering*. We now construct three types of problem instances 1. *Easy* 2. *Medium* 3. *Hard*, for any given K , such that items with their respective θ parameters are assigned as follows: 1. *Easy*: $\theta(1 : \lfloor K/2 \rfloor) = 1$, $\theta(\lfloor K/2 \rfloor + 1 : K) = 0.5$. 2. *Medium*: $\theta(1 : \lfloor K/3 \rfloor) = 1$, $\theta(\lfloor K/3 \rfloor + 1 : \lfloor 2K/3 \rfloor) = 0.7$, $\theta(\lfloor 2K/3 \rfloor + 1 : K) = 0.4$. 3. *Hard*: $\theta(i) = 1 - (i-1)/K$, $\forall i \in [K]$. Note for each $\sigma^* = (1 > 2 > \dots > K)$.

In every experiment, we set the learning parameters $\alpha = 0.51$, $\delta = 1/T$ for S1DB-UCB (Alg. 1) and as per Thm. 6 for S1DB-ED (Alg. 2). All results are averaged over 50 runs.

Regret over Varying Preference Matrices. We first plot the cumulative regret of our two algorithms (Alg. 1 and 2) over time on the above three Plackett-Luce datasets for $K = 10$. We generate availability sequence \mathcal{S}_T randomly by sampling every item $i \in [K]$ independently with probability 0.5. Fig 2 shows that, as their names suggest too, *instance-Easy* is easiest to learn as the best-vs-worst item preferences are well separated and the diversity of the item preferences across different groups are least. Consequently the algorithms yield slightly more regret on *instance-Medium* due to higher preference diversity, and the hardest instance being *Hard* where the learner really needs to differentiate the ranking of every item for any arbitrary set sequences \mathcal{S}_T . Empirically S1DB-UCB is seen to slightly outperform S1DB-ED, though otherwise they perform competitively.

Regret over Varying Set Availabilities. In these set of experiments, the idea is to understand how the regret improves over completely random subset availabilities as now the learner may not have to distinguish all item preferences as some of the item combinations occurs rarely. We choose $K = 10$ and to enforce

item dependencies we generate each set S_t by drawing a random sample from Gaussian(μ, Σ) such that $\mu_i = 0, \forall i \in [10]$, and Σ is a fixed 10×10 positive definite matrix which controls the set dependencies: Precisely we use two different block diagonal matrices for *Low-Correlation* and *High-Correlation* with the following correlations: 1. *Low-Correlation*: Σ is a separated block diagonal matrix on item partitions $\{1, 2, 3\}, \{4, 5, 6\}, \{7, 8, 9, 10\}$. 2. *High-Correlation*: Σ is constructed by merging three all-1 matrices on partitions $\{1 \dots 5\}, \{2, \dots 8\}$, and $\{6, \dots 10\}$, however as the resulting matrix is positive semi-definite, so we further take its SVD and reconstruct the matrix back eliminating the negative eigenvalues. At every round we sample a random vector from Gaussian(μ, Σ), and S_t is considered to be the set of items whose value exceeds 0.5. Both experiments are run on *instance-Hard*. Fig. 3 shows, as expected, on *Low-Correlation* both algorithms converge to σ^* relatively faster and at lower regret compared to *High-Correlation* (as the latter induces higher variability of the available subsets).

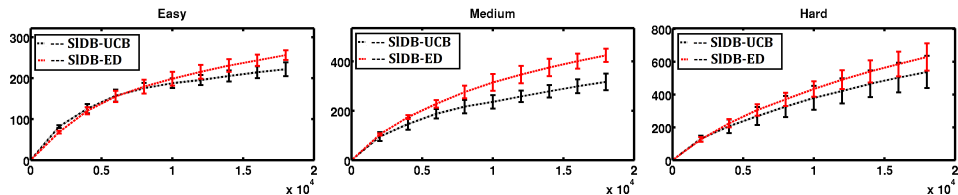


Figure 2: Regret (R_t) vs time (t) over three preference instances (\mathbf{P})

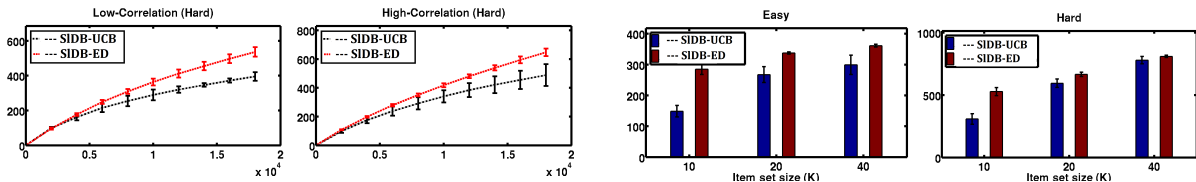


Figure 3: Regret (R_t) vs time (t) over availability sequences S_T

Figure 4: Final regret (R_T) at $T = 10^4$ with varying sizes (K)

Final Regret vs Setsize(K). We also compared the (averaged) final regret of the two algorithms over varying item sizes K . We additionally constructed two larger Plackett-Luce (θ) *Easy* and *Hard* instances for $K = 20$ and 40 , using similar θ assignments explained before. We set $T = 10000$ and use itemwise independent set generation idea, as described for Fig. 2. As expected, Fig. 4 shows the regret of both algorithms scales up with increasing K with effect on S1DB-ED being slightly worse than S1DB-UCB, though the latter generally exhibits a higher variance.

Worst Case Regret vs Time. We run an experiment to analyse the regret of our two algorithms on the worst case problem instances. Towards this we use preference matrices \mathbf{P}_Δ of the form: $\mathbf{P}_\Delta(i, j) = 0.5 + \Delta, \forall 1 \leq i < j \leq K$, i.e. all items are spaced with equidistant gap $\Delta \in (0, 0.5]$. As before, we choose $T = 20,000$ and $K = 10$, and run the algorithms on above problem instances varying Δ in the range $[10/T, \dots, 0.5]$ with uniform grid-size of 0.005 (i.e. total 100 values of Δ , each corresponds to a separate problem instance \mathbf{P}_Δ with different ‘gap-complexity’). At the end we plot the worst case regret of both the algorithms over time, by plotting $\max_\Delta R_t(\mathbf{P}_\Delta)$ vs t . We run the experiments over three availability sequences: 1. Independent (as used in Fig. 2), 2. Low-Correlation, and 3. High-Correlation (as used in Fig. 3). As a consequence the resulting plots reflect the worst case (w.r.t. Δ) performances of the algorithms, which seem to be scaling as $O(T^{2/3})$ for S1DB-ED, as conjectured to be its distribution free upper bound (see discussion after Thm. 6), and with a slightly lower rate for S1DB-UCB. Fig. 5 shows the comparative performances.

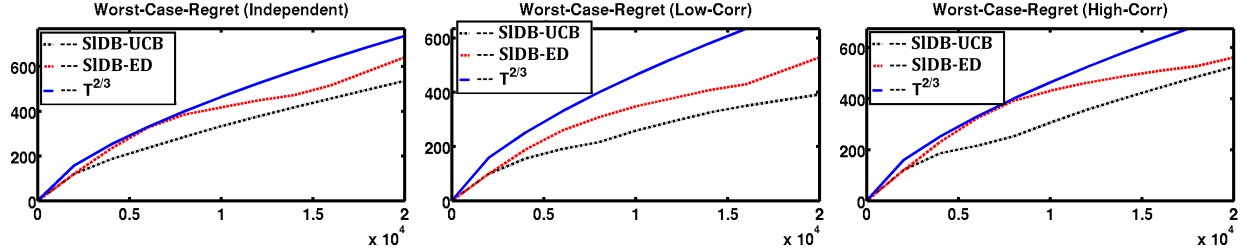


Figure 5: “Worst Case Regret” ($\max_{\Delta} R_t(P_{\Delta})$) vs time (t) over three availability sequences \mathcal{S}_T

7 Conclusion and Perspective

We introduce the problem of sleeping dueling bandits with stochastic preferences and adversarial availabilities, which, despite of great practical relevance, was left unaddressed till date. Towards this we adapt two dueling bandit algorithms for the problem and give regret analysis for both. We also derive an instance dependent regret lower bound for our problem setup which shows that our second algorithm is asymptotically near-optimal (up to the problem dependent constants). Finally, we compare both our algorithms empirically where usually the first algorithm is shown to outperform the second, although having a relatively weaker regret.

Future Works. Moving forward, one can address many open questions along this direction, including relaxing the *total-ordering* assumption on the stochastic preferences assuming more general ranking objective based on borda [32] or copeland scores [36], or extending the framework to a general contextual scenario with subsetwise feedback. Another direction worth understanding is to analyze the connection of this problem with other bandit setups, e.g., learning with feedback graphs [3, 4] or other side information [23, 20]. It would also be interesting to consider the dueling bandit problem for adversarial preference and stochastic availabilities [24, 17], and also analyzing these class of problems for general subsetwise preferences [25, 30, 8].

References

- [1] Shipra Agrawal and Navin Goyal. Analysis of Thompson sampling for the multi-armed bandit problem. In *Conference on Learning Theory*, pages 39–1, 2012.
- [2] Nir Ailon, Zohar Shay Karnin, and Thorsten Joachims. Reducing dueling bandits to cardinal bandits. In *ICML*, volume 32, pages 856–864, 2014.
- [3] Noga Alon, Nicolo Cesa-Bianchi, Ofer Dekel, and Tomer Koren. Online learning with feedback graphs: Beyond bandits. In *JMLR Workshop and Conference Proceedings*, volume 40. Microtome Publishing, 2015.
- [4] Noga Alon, Nicolo Cesa-Bianchi, Claudio Gentile, Shie Mannor, Yishay Mansour, and Ohad Shamir. Nonstochastic multi-armed bandits with graph-structured feedback. *SIAM Journal on Computing*, 46(6):1785–1826, 2017.
- [5] Jean-Yves Audibert and Sébastien Bubeck. Best arm identification in multi-armed bandits. In *COLT-23th Conference on Learning Theory-2010*, pages 13–p, 2010.
- [6] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- [7] Hossein Azari, David Parkes, and Lirong Xia. Random utility theory for social choice. In *Advances in Neural Information Processing Systems*, pages 126–134, 2012.

- [8] Brian Brost, Yevgeny Seldin, Ingemar J. Cox, and Christina Lioma. Multi-dueling bandits and their application to online ranker evaluation. *CoRR*, abs/1608.06253, 2016.
- [9] Róbert Busa-Fekete and Eyke Hüllermeier. A survey of preference-based online learning with bandit algorithms. In *International Conference on Algorithmic Learning Theory*, pages 18–39. Springer, 2014.
- [10] Róbert Busa-Fekete, Eyke Hüllermeier, and Balázs Szörényi. Preference-based rank elicitation using statistical models: The case of mallows. In *Proceedings of The 31st International Conference on Machine Learning*, volume 32, 2014.
- [11] Corinna Cortes, Giulia Desalvo, Claudio Gentile, Mehryar Mohri, and Scott Yang. Online learning with sleeping experts and feedback graphs. In *International Conference on Machine Learning*, pages 1370–1378, 2019.
- [12] Imre Csiszár. The method of types. *IEEE Transactions on Information Theory*, 44(6):2505–2523, 1998.
- [13] Miroslav Dudík, Katja Hofmann, Robert E Schapire, Aleksandrs Slivkins, and Masrour Zoghi. Contextual dueling bandits. In *Conference on Learning Theory*, pages 563–587, 2015.
- [14] Pratik Gajane, Tanguy Urvoy, and Fabrice Clérot. A relative exponential weighing algorithm for adversarial utility-based dueling bandits. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 218–227, 2015.
- [15] Satyen Kale, Chansoo Lee, and Dávid Pál. Hardness of online sleeping combinatorial optimization problems. In *Advances in Neural Information Processing Systems*, pages 2181–2189, 2016.
- [16] Shivaram Kalyanakrishnan, Ambuj Tewari, Peter Auer, and Peter Stone. PAC subset selection in stochastic multi-armed bandits. In *ICML*, volume 12, pages 655–662, 2012.
- [17] Varun Kanade, H Brendan McMahan, and Brent Bryan. Sleeping experts and bandits with stochastic action availability and adversarial rewards. 2009.
- [18] Varun Kanade and Thomas Steinke. Learning hurdles for sleeping experts. *ACM Transactions on Computation Theory (TOCT)*, 6(3):11, 2014.
- [19] Robert Kleinberg, Alexandru Niculescu-Mizil, and Yogeshwer Sharma. Regret bounds for sleeping experts and bandits. *Machine learning*, 80(2-3):245–272, 2010.
- [20] Tomas Kocak, Gergely Neu, Michal Valko, and Rémi Munos. Efficient learning by implicit exploration in bandit problems with side observations. In *Advances in Neural Information Processing Systems*, pages 613–621, 2014.
- [21] Junpei Komiyama, Junya Honda, Hisashi Kashima, and Hiroshi Nakagawa. Regret lower bound and optimal algorithm in dueling bandit problem. In *COLT*, pages 1141–1154, 2015.
- [22] Tor Lattimore and Csaba Szepesvári. Bandit algorithms. *preprint*, 2018.
- [23] Shie Mannor and Ohad Shamir. From bandits to experts: On the value of side-observations. In *Advances in Neural Information Processing Systems*, pages 684–692, 2011.
- [24] Gergely Neu and Michal Valko. Online combinatorial optimization with stochastic decision sets and adversarial losses. In *Advances in Neural Information Processing Systems*, pages 2780–2788, 2014.
- [25] Wenbo Ren, Jia Liu, and Ness B Shroff. PAC ranking from pairwise and listwise queries: Lower bounds and upper bounds. *arXiv preprint arXiv:1806.02970*, 2018.
- [26] Aadirupa Saha and Aditya Gopalan. Battle of bandits. In *Uncertainty in Artificial Intelligence*, 2018.

- [27] Aadirupa Saha and Aditya Gopalan. PAC Battling Bandits in the Plackett-Luce Model. In *Algorithmic Learning Theory*, pages 700–737, 2019.
- [28] Hossein Azari Soufiani, David C Parkes, and Lirong Xia. Preference elicitation for general random utility models. In *Uncertainty in Artificial Intelligence*, page 596. Citeseer, 2013.
- [29] Yanan Sui, Vincent Zhuang, Joel Burdick, and Yisong Yue. Multi-dueling bandits with dependent arms. In *Conference on Uncertainty in Artificial Intelligence*, UAI’17, 2017.
- [30] Yanan Sui, Masrour Zoghi, Katja Hofmann, and Yisong Yue. Advancements in dueling bandits. In *IJCAI*, pages 5502–5510, 2018.
- [31] Balázs Szörényi, Róbert Busa-Fekete, Adil Paul, and Eyke Hüllermeier. Online rank elicitation for plackett-luce: A dueling bandits approach. In *Advances in Neural Information Processing Systems*, pages 604–612, 2015.
- [32] Tanguy Urvoy, Fabrice Clerot, Raphael Féraud, and Sami Naamane. Generic exploration and k-armed voting bandits. In *International Conference on Machine Learning*, pages 91–99, 2013.
- [33] Huasen Wu and Xin Liu. Double Thompson sampling for dueling bandits. In *Advances in Neural Information Processing Systems*, pages 649–657, 2016.
- [34] Yisong Yue, Josef Broder, Robert Kleinberg, and Thorsten Joachims. The k -armed dueling bandits problem. *Journal of Computer and System Sciences*, 78(5):1538–1556, 2012.
- [35] Yisong Yue and Thorsten Joachims. Interactively optimizing information retrieval systems as a dueling bandits problem. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 1201–1208. ACM, 2009.
- [36] Masrour Zoghi, Zohar S Karnin, Shimon Whiteson, and Maarten De Rijke. Copeland dueling bandits. In *Advances in Neural Information Processing Systems*, pages 307–315, 2015.
- [37] Masrour Zoghi, Shimon Whiteson, Remi Munos, Maarten de Rijke, et al. Relative upper confidence bound for the k -armed dueling bandit problem. In *JMLR Workshop and Conference Proceedings*, number 32, pages 10–18. JMLR, 2014.
- [38] Masrour Zoghi, Shimon A Whiteson, Maarten De Rijke, and Remi Munos. Relative confidence sampling for efficient on-line ranker evaluation. In *Proceedings of the 7th ACM international conference on Web search and data mining*, pages 73–82. ACM, 2014.

Supplementary: Dueling Bandits with Adversarial Sleeping

8 Appendix for Sec. 3

Definition 7 (*No-regret algorithm*). An algorithm \mathcal{A} for *Sleeping-Dueling Bandit with Stochastic Preferences and Adversarial Availabilities* problem is defined to be a *No-regret algorithm*, if for each problem instance *DB-SPAA*($\mathbf{P}, \mathcal{S}_T$) model, the expected number of times \mathcal{A} plays any suboptimal duel $(i, j) \in [K] \times [K]$ is sublinear in T , or more precisely, $\forall (i, j) \neq (i_t^*, i_t^*)$, $\mathbf{E}[n_{ij}(T)] = o(T^\alpha)$, for some $\alpha \in (0, 1)$, where recall that we define $n_{ij}(t) := \sum_{\tau=1}^t \mathbf{1}(\{x_\tau, y_\tau\} = \{i, j\})$ denotes the number of times the pair (i, j) is played by \mathcal{A} in T rounds. ($\mathbf{E}[\cdot]$ denotes expectation under the randomization of \mathcal{A} and the *DB-SPAA*($\mathbf{P}, \mathcal{S}_T$) model.)

8.1 Proof of Thm. 1

Proof. The main argument lies behind the fact that in the worst case the adversary can force the algorithm to learn the preference of every distinct pair (i, j) as the in the ‘worst-case’ sequence \mathcal{S}_T knowledge of the already ‘learnt’ pairwise preferences would not disclose any information on the remaining pairs; e.g. assuming $\sigma^* = (1, 2, \dots, K)$, revealing the available subsets in the following sequence $(1, 2), (1, 3), \dots (1, K), (2, 3), (2, 4), \dots (K-1, K)$ would force the learner to explore (learn the preferences) all $\binom{K}{2}$ distinct pairs.

The remaining proof establishes this formally, towards which we first show a $\Omega\left(\frac{\ln T}{\Delta(1,2)}\right)$ regret lower bound for a *DB-SPAA* instance with just two items (i.e. $K = 2$) as shown in Lem. 2. The lower bound for any general K can now be derived applying the above bound on independent $\binom{K}{2}$ subintervals, with the availability sequence $(1, 2), (1, 3), \dots (1, K), (2, 3), (2, 4), \dots (K-1, K)$.

For the interest of the problem instance construction to prove the lower bound, we would assume $\Delta(i, i+1) > 0$, $\forall i \in [K-1]$ and thus we use $\Delta(i, i+1)_+ = \Delta(i, i+1)$ for the rest of this proof (as also assumed for Lem. 2). Note that this is without loss of generality since otherwise the regret lower bound in Lem. 2 is trivially 0.

We add the details below for completeness.

Let $K' = \binom{K}{2}$ and suppose we divide the time horizon into sub-intervals $1, 2, \dots, K'$ each of length $T' := T/K'$, where the available subsets are fixed inside every subinterval, and follows the sequence $(1, 2), (1, 3), \dots (1, K), (2, 3), (2, 4), \dots (K-1, K)$ across subintervals. Note that with above construction, the regret minimization problem within each sub-interval boils down to the standard stochastic dueling bandit problem over 2 arms.

Further since the preferences of available set S_t s are independent across different sub-intervals, applying the lower bound of Lem. 2 individually to every K' subintervals the total cumulative regret of \mathcal{A} in T rounds can be lower bounded as:

$$\mathbf{E}[R_T(\mathcal{A})] = \sum_{i=1}^{K-1} \sum_{j=i+1}^K \mathbf{E}[R_{T'}(\mathcal{A})] \geq \Omega\left(\sum_{i=1}^{K-1} \sum_{j=i+1}^K \frac{\log T'}{\Delta(i, j)}\right) = \Omega\left(\sum_{i=1}^{K-1} \sum_{j=i+1}^K \frac{\log T}{\Delta(i, j)}\right)$$

where the last inequality holds since $T \geq (K')^2$, which implies $\log \frac{T}{K'} \geq \log T - \log \sqrt{T} = 1/2 \log T$, and this concludes the proof. □

Lemma 2 (Lower Bound of *DB-SPAA*($\mathbf{P}_K, \mathcal{S}_T$) for 2 items). For any *No-regret learning algorithm* \mathcal{A} , there exists a problem instance *DB-SPAA*($\mathbf{P}_2, \mathcal{S}_T$) such that the expected regret incurred by \mathcal{A} on that can be lower bounded as: $\mathbf{E}[R_T(\mathcal{A})] \geq \Delta^{-1} \log(T)$, Δ being the ‘preference-gap’ between the two items (i.e. $\Delta = \mathbf{P}_{12} - 1/2$, assuming $P_{12} > 1/2$ or equivalently $\Delta > 0$).

Proof. Note that for $K = 2$, the only non-trivial available set is $\{1, 2\}$, therefore we assume $S_t = \{1, 2\}$, $\forall t \in [T]$. The proof now simply follows by applying the existing lower bound (Thm. 2) of [34] for standard stochastic dueling bandit problem for only 2 arms. \square

9 Appendix for Sec. 4

Notations. Let us start with defining useful notation for the analysis. We write for any pair $1 \leq i < j \leq K$

$$M_{ij} = \sum_{k=1}^i \frac{4\alpha}{\min \{\Delta(k, i)_+, \Delta(k, j)_+\}^2}.$$

We also denote $S_{\setminus i} = S \setminus \{i\}$, $i \in S$, for any $S \subseteq [K]$.

9.1 Complete proof of Thm. 3

Theorem 3. *Given any $\delta > 0$ and $\alpha \geq 1$, with probability at least $1 - \delta$, the regret incurred by S1DB-UCB (Alg. 1) is upper-bounded as:*

$$R_T \leq 2 \sum_{i=1}^{K-1} \sum_{j=i+1}^K M_{ij} \log(2C(K, \delta)M_{ij})$$

where $C(K, \delta) := ((4\alpha - 1)K^2 / ((2\alpha - 1)\delta))^{\frac{1}{2\alpha - 1}}$.

Proof of Thm. 3. The key steps lie in proving the following four lemmas. The first lemma follows along the line of Lem. 1 of RUCB algorithm [37]. It shows after $C(K, \delta)$ rounds all the pairwise estimates are contained within their respective confidence intervals:

Lemma 4. *Let $\alpha > 0.5$ and $\delta > 0$. Then, with probability at least $1 - \delta$, for any $i, j \in [K]$*

$$\widehat{p}_{ij}(t) - c_{ij}(t) \leq p_{ij} \leq u_{ij}(t) := \widehat{p}_{ij}(t) + c_{ij}(t), \quad \forall t \in [T].$$

The lemma below is adapted from Proposition 2, [37]. It basically states that once the algorithm has explored enough (i.e., more than $C(K, \delta)$) the algorithm will not play a suboptimal pair too many times.

Lemma 5. *Let $\alpha > 0.5$. Under the notations and the high-probability event of Lem. 4, for all $i, j, k \in [K]$ such that $\{i, j\} \neq \{k, k\}$, and for any $\tau \geq 1$*

$$\sum_{t=1}^{\tau} \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) \leq \frac{4\alpha \log a_{i,j}(\tau)}{\min \{\Delta(k, i)_+, \Delta(k, j)_+\}^2},$$

where recall $a_{ij}(\tau) = \max(C(K, \delta), n_{ij}(\tau))$.

Given the above results, we are ready to analyze the regret guarantee of S1DB-UCB. For ease on notation we denote $\mathcal{X}_t = \{x_t, y_t\}$. Let us assume the ‘good event’ of Lem. 4 holds good for all $t \in [T]$, which is true with probability of at least $1 - \delta$. Conditioned on that, note that Lem. 5 is satisfied. Based on this we now analyze the regret of Alg. 1:

$$\begin{aligned} R_T &= \sum_{t=1}^T \sum_{k=1}^{K-1} \mathbf{1}(i_t^* = k) r_t \\ &= \sum_{t=1}^T \sum_{k=1}^K \sum_{i=k}^K \sum_{j=i}^K \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) r_t \quad \leftarrow \text{because } x_t \geq k \text{ and } y_t \geq k \end{aligned}$$

$$\begin{aligned}
&= \sum_{t=1}^T \sum_{k=1}^{K-1} \sum_{i=k}^{K-1} \sum_{j=i+1}^K \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) r_t \leftarrow \text{because } i = j = k \text{ implies } r_t = 0 \\
&\leq \sum_{t=1}^T \sum_{k=1}^{K-1} \sum_{i=k}^{K-1} \sum_{j=i+1}^K \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) \leftarrow \text{because } r_t \leq 1 \\
&= \sum_{i=1}^{K-1} \sum_{j=i+1}^K \sum_{t=1}^T \sum_{k=1}^i \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) \\
&= \sum_{i=1}^{K-1} \sum_{j=i+1}^K n_{ij}(T). \tag{3}
\end{aligned}$$

Now, fix $1 \leq i < j \leq K$ and let us upper-bound $n_{ij}(T)$ the number of times such a pair is played. Summing the upper-bound of Lemma 5 over $k \leq i$, we get

$$n_{ij}(T) = \sum_{k=1}^i \sum_{t=1}^T \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) \leq \sum_{k=1}^i \frac{4\alpha \log(\max\{C(K, \delta), n_{ij}(T)\})}{\min\{\Delta(k, i)_+, \Delta(k, j)_+\}^2}.$$

Therefore, since $C(K, \delta) \geq 1$,

$$n_{ij}(T) \leq M_{ij} (\log(C(K, \delta)) + \log(n_{ij}(T))), \quad \text{where } M_{ij} = \sum_{k=1}^i \frac{4\alpha}{\min\{\Delta(k, i)_+, \Delta(k, j)_+\}^2}.$$

which implies

$$n_{ij}(T) \leq 2M_{ij} (\log C(K, \delta) + \log(2M_{ij})).$$

Substituting into Inequality (3) entails

$$R_T \leq 2 \sum_{i=1}^{K-1} \sum_{j=i+1}^K M_{ij} \log(2C(K, \delta)M_{ij}),$$

which concludes the proof. \square

9.2 Technical lemmas for Thm. 3

Lemma 4. *Let $\alpha > 0.5$ and $\delta > 0$. Then, with probability at least $1 - \delta$, for any $i, j \in [K]$*

$$\widehat{p}_{ij}(t) - c_{ij}(t) \leq p_{ij} \leq u_{ij}(t) := \widehat{p}_{ij}(t) + c_{ij}(t), \quad \forall t \in [T].$$

Proof. The proof of this lemma is adapted from a similar result (Lemma 1) of [37]. Suppose $\mathcal{G}_{ij}(t)$ denotes the event that at time $t \in [T]$ and item-pair $i, j \in [K]$, $p_{ij} \in [l_{ij}(t), u_{ij}(t)]$. We also define $\mathcal{G}_{ij}^c(t)$ its complement. Let $i, j \in [K]$.

Note that for any such that pair (i, i) , $\mathcal{G}_{ii}(t)$ always holds true for any $t \in [T]$ and $i \in [n]$, as $p_{ii} = u_{ii} = l_{ii} = \frac{1}{2}$. We can thus assume $i \neq j$. Moreover, for any t and i, j , $\mathcal{G}_{ij}(t)$ holds if and only if $\mathcal{G}_{ij}(t)$ as $|\widehat{p}_{ji}(t) - p_{ji}| = |(1 - \widehat{p}_{ij}(t)) - (1 - p_{ij})| = |\widehat{p}_{ij}(t) - p_{ij}|$. Thus we will restrict our focus only to pairs $i < j$ for the rest of the proof. Hence, to prove the lemma it suffices to show

$$\mathbf{P}\left(\exists t \in [T], i < j, \text{ such that } \mathcal{G}_{ij}^c(t)\right) \leq \delta,$$

which we do now. Recall from the definition of $c_{ij}(t)$ that $\mathcal{G}_{ij}(t)$ can be rewritten as:

$$|\widehat{p}_{ij}(t) - p_{ij}| \leq \sqrt{\frac{\alpha \ln(a_{ij}(t))}{n_{ij}(t)}}.$$

Let $\tau_{ij}(n)$ the time step $t \in [T]$ when the pair (i, j) was updated (i.e. i and j was compared) for the n^{th} time. We now bound the probability of the confidence bound ($\mathcal{G}_{ij}(t)$) getting violated at any round $t \in [T]$ for some duel (i, j) as follows:

$$\begin{aligned} \mathbf{P}\left(\exists t \in [T], i < j, \text{ such that } \mathcal{G}_{ij}^c(t)\right) &\leq \sum_{i < j} \mathbf{P}\left(\exists n \geq 0, |p_{ij} - \widehat{p}_{ij}(\tau_{ij}(n))| > \sqrt{\frac{\alpha a_{ij}(\tau_{ij}(n))}{n_{ij}(\tau_{ij}(n))}}\right) \\ &= \sum_{i < j} \left[\mathbf{P}\left(\exists n \leq C(K, \delta), |p_{ij} - \widehat{p}_{ij}(n)| > \sqrt{\frac{\alpha \ln(C(K, \delta))}{n}}\right) \right. \\ &\quad \left. + \mathbf{P}\left(\exists n > C(K, \delta), |p_{ij} - \widehat{p}_{ij}(\tau_{ij}(n))| > \sqrt{\frac{\alpha \ln(n_{ij}(\tau_{ij}(n)))}{n}}\right) \right], \end{aligned}$$

where $\widehat{p}_{ij}(t) = \frac{w_{ij}(t)}{w_{ij}(t) + w_{ij}(t)}$ is the frequentist estimate of p_{ij} at round t (after $n = n_{ij}(t)$ comparisons between arm i and j). To ease the notation, denote $F = C(K, \delta)$. Noting $n_{ij}(\tau_{ij}(n)) = n$, and using Hoeffding's inequality, we further get

$$\begin{aligned} \mathbf{P}\left(\exists t \in [T], i < j, \text{ such that } \mathcal{G}_{ij}^c(t)\right) &\leq \sum_{i < j} \left[\sum_{n=1}^F 2e^{-2n \frac{\alpha \ln F}{n}} + \sum_{n=F+1}^{\infty} 2e^{-2n \frac{\alpha \ln n}{n}} \right] \\ &= \frac{n(n-1)}{2} \left[2 \sum_{n=1}^F \frac{1}{F^{2\alpha}} + \sum_{n=F+1}^{\infty} \frac{2}{n^{2\alpha}} \right] \\ &\leq \frac{n^2}{F^{2\alpha-1}} + n^2 \int_F^{\infty} \frac{dx}{x^{2\alpha}} \leq \frac{n^2}{F^{2\alpha-1}} - \frac{n^2}{(1-2\alpha)F^{2\alpha-1}} = \frac{(2\alpha)n^2}{(2\alpha-1)F^{2\alpha-1}} = \delta. \end{aligned}$$

where the last inequality is because $F = C(K, \delta) = \left[\frac{2\alpha n^2}{(2\alpha-1)\delta} \right]^{\frac{1}{2\alpha-1}}$. This concludes the claim. \square

Lemma 5. *Let $\alpha > 0.5$. Under the notations and the high-probability event of Lem. 4, for all $i, j, k \in [K]$ such that $\{i, j\} \neq \{k, k\}$, and for any $\tau \geq 1$*

$$\sum_{t=1}^{\tau} \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) \leq \frac{4\alpha \log a_{i,j}(\tau)}{\min\{\Delta(k, i)_+, \Delta(k, j)_+\}^2},$$

where $a_{ij}(\tau) = \max(C(K, \delta), n_{ij}(\tau))$.

Proof. We assume the confidence bound of Lem. 4 is holds good for all pair $(i, j) \in [K]^2$, at all round $t \in [T]$, which we know happens with probability at least $(1 - \delta)$. Let us define $l_{ij}(t) := 1 - u_{ji}(t)$. Let $t \geq 1$. Let $i, j, k \in [K]$ such that $i_t^* = k$, $x_t = i$, and $y_t = j$ and $\{i, j\} \neq \{k, k\}$. Since $i_t^* = k$, this implies both $i \geq k$ and $j \geq k$. Furthermore, we recall that $i_t^* = k$ is unique by definition, $i_t^* = \min\{S_t\}$. We consider the following cases.

- **Case 1** ($i = j > k$). Then, $x_t = y_t = i = j$. By the arm selection strategy (Step 14. of Algorithm 1)

$$y_t \leftarrow \arg \max_{m \in \mathcal{C}_t} u_{mx_t}(t)$$

which implies $1/2 = u_{jj}(t) > u_{kj}(t)$. But, on the other hand, since $k < j$, by Lemma 4, $u_{kj}(t) \geq p_{kj} > \frac{1}{2}$. This causes a contraction and this case is not possible.

- **Case 2** ($j > i = k$). Then, $y_t = j$ and $x_t = i_t^* = k$. We again proceed by contradiction. Assume that $n_{kj}(t) > \frac{4\alpha \ln a_{kj}(t)}{\Delta(k, j)_+^2}$. Then, by definition of $c_{kj}(t)$, it implies

$$2c_{kj}(t) = 2\sqrt{\frac{\alpha \log a_{kj}(t)}{n_{ij}(t)}} < \Delta(k, j)_+,$$

which by Lem. 4 entails

$$u_{jk}(t) = \widehat{p}_{jk}(t) + c_{jk}(t) < p_{jk} + 2c_{jk}(t) < \frac{1}{2} - \Delta(k, j)_+ + \Delta(k, j)_+ < \frac{1}{2}.$$

Again since our arm selection strategy enforces $y_t \leftarrow \arg \max_{i \in C_t} u_{ix_t}(t)$, clearly $\frac{1}{2} = u_{kk}(t) > u_{jk}(t)$, so that j can not be selected as y_t . Therefore, recalling that $k = i$,

$$n_{ij}(t) \leq \frac{4\alpha \ln a_{kj}(t)}{\Delta(i, j)_+^2}. \quad (4)$$

- **Case 3** ($i > j = k$). Then, $x_t = i$ and $y_t = i_t^* = k$. This can be proved similarly as the previous case. Assuming $n_{ik}(t) > 4\alpha \ln a_{ik}(t) \Delta(i, k)_+^{-2}$ yields $u_{ik}(t) < 1/2$. Therefore, since $u_{ii}(t) = 1/2$, it entails

$$|C_i(t)| = |\{m \in S_t | u_{im}(t) > 1/2\}| \leq |S_t \setminus \{i, k\}| \leq |S_t| - 2.$$

But by Lemma 4, for all $m > k$, $u_{km}(t) \geq p_{km} = 1/2 + \Delta(k, m) > 1/2$. Thus, since $k = i_t^*$, we also have $|C_k(t)| = |S_t| - 1$ and thus

$$|C_k(t)| > |C_i(t)|.$$

By Step 12 of Algorithm 1, this implies that $i \notin C_t$ and thus $x_t \neq i$ as x_t is selected from C_t , which causes a contradiction. Therefore, recalling $j = k$,

$$n_{ij}(t) \leq \frac{4\alpha \ln a_{i,j}(t)}{\Delta(i, j)_+^2}. \quad (5)$$

- **Case 4.** ($i \neq j > k$). Then, assuming $n_{ij}(t) > 4\alpha \log a_{i,j}(t) \min\{\Delta(k, i)_+, \Delta(k, j)_+\}^{-2}$, note that

$$u_{ij}(t) - l_{ij}(t) = 2c_{ij}(t) = 2\sqrt{\frac{\alpha \log a_{i,j}(t)}{n_{ij}(t)}} < \min(\Delta(k, i)_+, \Delta(k, j)_+).$$

But, on the other hand, $x_t = i$ implies $u_{ij}(t) > 1/2$, and $y_t = j$ implies $u_{ji}(t) > u_{jk}(t) > p_{jk}$, and then $l_{ij}(t) = 1 - u_{ji}(t) < 1 - p_{jk}$. So we have $u_{ij}(t) - l_{ij}(t) > 1/2 - p_{jk} = \Delta(k, i)_+$ which gives a contradiction. Thus,

$$n_{i,j}(t) \leq \frac{4\alpha \log a_{i,j}(t)}{\min\{\Delta(k, i)_+, \Delta(k, j)_+\}^2}. \quad (6)$$

Note that the case $x_t = j$, $y_t = i$, and $i_t^* = k$ is symmetric with the above cases and can be considered similarly. Denote by τ' the last time before $\tau \geq 1$ such that a pair $\{i, j\} \neq \{k, k\}$ is pulled when $k = i_t^*$, that is

$$\tau' = \operatorname{argmax}_{1 \leq t \leq \tau} \{k = i_t^*, \{x_t, y_t\} = \{i, j\}\}.$$

Then,

$$\sum_{t=1}^{\tau} \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) \leq \sum_{t=1}^{\tau'} \mathbf{1}(i_t^* = k) \mathbf{1}(\{x_t, y_t\} = \{i, j\}) \leq n_{ij}(\tau').$$

But, at time τ' , the suboptimal pair $\{i, j\}$ got pulled, thus one of the above four cases is true, which implies from Inequalities (4), (5), and (6) that

$$n_{ij}(\tau') \leq \frac{4\alpha \log a_{i,j}(\tau)}{\min\{\Delta(k, i)_+, \Delta(k, j)_+\}^2}.$$

Substituting into the previous inequality concludes the proof. \square

10 Appendix for Sec. 5

10.1 Technical lemmas

Before proving the regret guarantee of SLDB-ED (Alg. 2) in Thm. 6, we would like to introduce three lemmas which are crucially used towards bounding Alg. 2's regret. Lemma 8 below states that after some exploration, the algorithm estimates well all p_{ij} with $\hat{p}_{ij}(t)$.

Lemma 8. *Let $t_0 \geq 1$ and $\alpha \geq 4K$. Let $\varepsilon_i > 0$ for all $i \in [K]$. For any $t \in [T]$, let us define $\mathcal{E}_t := \{n_{ij}(t) > t_0, \forall i, j \in S_t, i \neq j\}$ to be the event when all distinct pairs $i, j \in [K]$ is played for at least for t_0 times. Let us also denote the event $\mathcal{G}(t) := \{\forall j > i_t^*, \Delta(i_t^*, j) > \varepsilon_i, \hat{p}_{i_t^* j}(t) > 1/2\}$. $\mathcal{G}^c(t)$ denotes the complement event of $\mathcal{G}(t)$. Then SLDB-ED satisfies:*

$$\mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{E}_t) \mathbf{1}(\mathcal{G}^c(t)) \right] = 2K + K \sum_{i=1}^{K-1} \sum_{j=i+1}^K \frac{e^{-(t_0-1)\Delta(i,j)^2}}{\Delta(i,j)^2}.$$

Proof. First, we show that with high probability for all $t = 1, \dots, T$, $i_t^* \in \mathcal{C}_t$, i_t^* belongs to the set of potential winners. Let $t \geq 1$. By definition of \mathcal{C}_t , we have

$$\begin{aligned} \mathbf{P}(i_t^* \notin \mathcal{C}_t) &\leq \mathbf{P} \left(\sum_{j \in \tilde{\mathcal{B}}_{i_t^*}(t)} n_{i_t^* j}(t) \text{kl}(\hat{p}_{i_t^* j}(t), 0.5) \geq \alpha \log t \right) \\ &\leq \mathbf{P} \left(\exists j > i_t^* \quad \text{s.t.} \quad n_{i_t^* j}(t) \text{kl}(\hat{p}_{i_t^* j}(t), 0.5) \geq \frac{\alpha \log t}{K} \right) \\ &\leq \sum_{j=i_t^*+1}^K \sum_{n=1}^t \mathbf{P} \left(\text{kl}(\tilde{p}_{i_t^* j}(n), 0.5) \geq \frac{\alpha \log t}{nK} \right), \end{aligned}$$

where $\tilde{p}_{ij}(n)$ denotes the frequentist empirical estimate of $P(i, j)$ after n pairwise comparisons between i and j (i.e., $\tilde{p}_{ij}(n) = \hat{p}_{ij}(t)$ with $n_{ij}(t) = n$). From Lemma II.1 of [12], this yields

$$\mathbf{P}(i_t^* \notin \mathcal{C}_t) \leq \sum_{j=i_t^*+1}^K \sum_{n=1}^t n \exp \left(-\frac{\alpha \log t}{K} \right) \leq K t^2 \exp \left(-\frac{\alpha \log t}{K} \right) \leq \frac{K}{t^2},$$

since $\alpha \geq 4K$. Therefore,

$$\sum_{t=1}^T \mathbf{P}(i_t^* \notin \mathcal{C}_t) \leq K \sum_{t=1}^T \frac{1}{t^2} \leq 2K. \quad (7)$$

Then,

$$\begin{aligned} \mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{E}_t) \mathbf{1}(\mathcal{G}^c(t)) \right] &\leq \mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{E}_t) \mathbf{1}(\mathcal{G}^c(t)) \mathbf{1}(i_t^* \in \mathcal{C}_t) \right] + \mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(i_t^* \notin \mathcal{C}_t) \right] \\ &\stackrel{(7)}{\leq} \mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{E}_t) \mathbf{1}(\mathcal{G}^c(t)) \mathbf{1}(i_t^* \in \mathcal{C}_t) \right] + 2K \\ &= \mathbf{E} \left[\sum_{t=1}^T \sum_{i=1}^K \mathbf{1}(\mathcal{E}_t) \mathbf{1}(\mathcal{G}^c(t)) \mathbf{1}(i_t^* = i) \mathbf{1}(i \in \mathcal{C}_t) \right] + 2K. \end{aligned}$$

Now, since x_t is uniformly sampled from \mathcal{C}_t from Line 11 of Algorithm 2, given that $i \in \mathcal{C}_t$, the probability that $x_t = i$ is at least $1/|\mathcal{C}_t| \geq 1/K$. Thus,

$$\mathbf{E} \left[\mathbf{1}(\mathcal{E}_t) \mathbf{1}(\mathcal{G}^c(t)) \mathbf{1}(i_t^* = i) \mathbf{1}(i \in \mathcal{C}_t) \right] \leq K \mathbf{E} \left[\mathbf{1}(\mathcal{E}_t) \mathbf{1}(\mathcal{G}^c(t)) \mathbf{1}(i \in \mathcal{C}_t) \mathbf{1}(i_t^* = x_t = i) \right],$$

which yields

$$\begin{aligned}
\mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{E}_t) \mathbf{1}(\mathcal{G}^c(t)) \right] &\leq K \sum_{i=1}^K \mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{E}_t) \mathbf{1}(\mathcal{G}^c(t)) \mathbf{1}(i \in \mathcal{C}_t) \mathbf{1}(i_t^* = x_t = i) \right] + 2K \\
&\stackrel{(*)}{=} K \sum_{i=1}^K \sum_{j: \Delta(i,j) > \varepsilon_i} \mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{E}_t) \mathbf{1}(i \in \mathcal{C}_t) \mathbf{1}(i_t^* = i) \mathbf{1}((x_t, y_t) = (i, j)) \mathbf{1}(\widehat{p}_{ij}(t) < 1/2) \right] + 2K \\
&\leq K \sum_{i=1}^K \sum_{j: \Delta(i,j) > \varepsilon_i} \mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{E}_t) \mathbf{1}((x_t, y_t) = (i, j)) \mathbf{1}(\widehat{p}_{ij}(t) < 1/2) \right] + 2K \tag{8}
\end{aligned}$$

where $(*)$ is because $\mathcal{G}^c(t) := \{\exists j > i, \Delta(i, j) > \varepsilon_{i_t^*}, \widehat{p}_{ij}(t) < 1/2\}$ when $i = i_t^*$, and since y_t is chosen such that $\widehat{p}_{i y_t}(t) < 1/2$ (see Line 12. of Alg. 2). Recall that \mathcal{E}_t ensures that (i, j) was pulled at least t_0 times during the exploration phase. Recalling that $\tilde{p}_{ij}(n)$ equals $\widehat{p}_{ij}(t)$ where t is such that $n = n_{ij}(t)$, we have

$$\sum_{t=1}^T \mathbf{1}(\mathcal{E}_t) \mathbf{1}((x_t, y_t) = (i, j)) \mathbf{1}(\widehat{p}_{ij}(t) < 1/2) \leq \sum_{n=t_0}^{\infty} \mathbf{1}(\tilde{p}_{ij}(n) < 1/2).$$

Therefore, plugging the latter inequality into the previous upper-bound (8), it yields

$$\begin{aligned}
\mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{E}_t) \mathbf{1}(\mathcal{G}^c(t)) \right] &\leq K \sum_{i=1}^K \sum_{j: \Delta(i,j) > \varepsilon_i} \sum_{n=t_0}^{\infty} \mathbf{P}(\tilde{p}_{ij}(n) < 1/2) + 2K \\
&\stackrel{(a)}{=} K \sum_{i=1}^K \sum_{j: \Delta(i,j) > \varepsilon_i} \sum_{n=t_0}^{\infty} \mathbf{P}(\tilde{p}_{ij}(n) < P(i, j) - \Delta(i, j)) + 2K \\
&\stackrel{(b)}{\leq} K \sum_{i=1}^K \sum_{j: \Delta(i,j) > \varepsilon_i} \sum_{n=t_0}^{\infty} \exp(-n\Delta(i, j)^2) + 2K \\
&\leq K \sum_{i=1}^K \sum_{j: \Delta(i,j) > \varepsilon_i} \frac{e^{-(t_0-1)\Delta(i,j)^2}}{e^{\Delta(i,j)^2} - 1} + 2K \\
&\leq K \sum_{i=1}^K \sum_{j: \Delta(i,j) > \varepsilon_i} \frac{e^{-(t_0-1)\Delta(i,j)^2}}{\Delta(i, j)^2} + 2K
\end{aligned}$$

where (a) follows by definition $\Delta(i, j) := \mathbf{P}(i, j) - 1/2$ and (b) is by Hoeffding's inequality. \square

The high-level idea of Lemma 9 below is that for any pair $1 \leq i < j \leq K$, j will not be played too much more than $M_{ij}(\delta)$ times together with items $k \leq i$. In other words, after sufficiently enough rounds j is detected as worse than all items $k < i$.

Lemma 9. *Let $1 \leq i < j \leq K$. Then, SLDB-ED (Alg. 2) satisfies:*

$$\mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \mathbf{1}(N_{ij}(t) > M_{ij}(\delta)) \right] \leq \frac{32}{\delta^2 \Delta(i, j)^2},$$

where $\mathcal{G}(t)$ is as defined in Lem. 8 and $N_{ij}(t) := \sum_{k=1}^i n_{kj}(t)$ is the number of times j was compared with some arm in $1, \dots, i$ and $M_{ij}(\delta) := (\alpha + \delta)(\log T) / \text{kl}(\mathbf{p}_{ji}, 0.5)$.

Proof. Let $1 \leq i \leq K - 1$. We start by recalling some useful notations:

$$\widehat{\mathcal{B}}_i(t) := \left\{ j \mid j \in [K], \widehat{p}_{i,j}(t) \leq 1/2 \right\}, \quad \mathcal{I}_i(t) := \sum_{j \in \widehat{\mathcal{B}}_i(t)} n_{ij}(t) \text{kl}(\widehat{p}_{ij}(t), 0.5),$$

where $\widehat{i}^*(t) := \arg \min_{i \in [K]} \mathcal{I}_i(t)$, and for simplicity we here denote $\mathcal{I}_{\widehat{i}^*}(t) = \mathcal{I}^*(t) := \min_{i \in [K]} \mathcal{I}_i(t)$. We also denote that event $\mathcal{J}_i(t) := \{\mathcal{I}_i(t) - \mathcal{I}^*(t) \leq \alpha \log t\}$. Then for any fixed $j > i$, we have

$$\begin{aligned} S_T(i, j) &:= \mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \mathbf{1}(N_{ij}(t) > M_{ij}(\delta)) \right] \\ &= \mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \mathbf{1}(N_{ij}(t) > M_{ij}(\delta)) \mathbf{1}(\mathcal{J}_j(t)) \right] \quad \leftarrow \text{as } x_t = j \text{ implies } \mathbf{1}(\mathcal{J}_j(t)) = 1 \\ &= \mathbf{E} \left[\sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \mathbf{1}(N_{ij}(t) > M_{ij}(\delta), \mathcal{J}_j(t)) \right] \end{aligned}$$

Substituting $\mathcal{J}_i(t) := \{\mathcal{I}_i(t) - \mathcal{I}^*(t) \leq \alpha \log t\}$, and using that $\mathcal{G}(t)$ implies $\mathcal{I}^*(t) = 0$, we get

$$\begin{aligned} S_T(i, j) &\leq \mathbf{E} \left[\sum_{t=1}^T \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \mathbf{1}(N_{ij}(t) > M_{ij}(\delta), \mathcal{I}_j(t) \leq \alpha \log t) \right] \\ &\leq \mathbf{E} \left[\sum_{t=1}^T \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \mathbf{1}(N_{ij}(t) > M_{ij}(\delta), \sum_{k \in \widehat{\mathcal{B}}_j(t)} n_{jk}(t) \text{kl}(\widehat{p}_{jk}(t), 0.5) \leq \alpha \log T) \right] \\ &\leq \mathbf{E} \left[\sum_{t=1}^T \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \mathbf{1}(N_{ij}(t) > M_{ij}(\delta), \sum_{k=1}^i n_{kj}(t) \text{kl}^+(\widehat{p}_{jk}(t), 0.5) \leq \alpha \log T) \right] \end{aligned}$$

where $\text{kl}^+(p, q) := \text{kl}(p, q) \mathbf{1}(p < q)$. But, from convexity of $\text{kl}^+(\cdot, 0.5)$ together with Jensen's inequality

$$\sum_{k=1}^i n_{kj}(t) \text{kl}^+(\widehat{p}_{jk}(t), 0.5) \geq N_{ij}(t) \text{kl}^+\left(\frac{1}{N_{ij}(t)} \sum_{k=1}^i n_{kj}(t) \widehat{p}_{jk}(t), 0.5\right).$$

Therefore, denoting

$$\tilde{p}_{1:ij}(N_{ij}(t)) := \frac{1}{N_{ij}(t)} \sum_{k=1}^i n_{kj}(t) \widehat{p}_{jk}(t) = \frac{1}{N_{ij}(t)} \sum_{k=1}^i w_{ki}(t)$$

the frequentist empirical estimate obtained after $N_{ij}(t)$ comparisons of j with any item better than i , we have

$$S_T(i, j) \leq \mathbf{E} \left[\sum_{t=1}^T \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \mathbf{1}(N_{ij}(t) > M_{ij}(\delta), N_{ij}(t) \text{kl}^+(\tilde{p}_{1:ij}(N_{ij}(t)), 0.5) \leq \alpha \log T) \right]$$

But, for each $n > M_{ij}(\delta)$, $N_{ij}(t) = n$ is only possible for one of the above rounds since $(x_t, y_t) = (i, k)$ with $k \leq i$, which increases $N_{ij}(t)$ by one. Thus,

$$\begin{aligned} S_T(i, j) &\leq \mathbf{E} \left[\sum_{n=M_{ij}(\delta)}^T \mathbf{1}(n \text{kl}^+(\tilde{p}_{1:ij}(n), 0.5) \leq \alpha \log T) \right] \\ &\leq \mathbf{E} \left[\sum_{n=\lceil M_{ij}(\delta) \rceil}^T \mathbf{1} \left(M_{ij}(\delta) \text{kl}^+(\tilde{p}_{1:ij}(n), 0.5) \leq \alpha \log T \right) \right] \\ &\leq \mathbf{E} \left[\sum_{n=\lceil M_{ij}(\delta) \rceil}^T \mathbf{1} \left(\text{kl}^+(\tilde{p}_{1:ij}(n), 0.5) \leq \frac{\text{kl}(p_{ji}, 0.5)}{1 + \delta} \right) \right] \end{aligned}$$

Now, let $\mu_i \in (p_{ji}, 0.5)$ such that $\text{kl}(\mu_i, 0.5) = \text{kl}(p_{ji}, 0.5)/(1 + \delta)$. By monotonicity of $\text{kl}^+(\cdot, 0.5)$,

$$\begin{aligned}
S_T(i, j) &= \sum_{n=\lceil M_{ij}(\delta) \rceil}^T \mathbf{P}\left(\text{kl}^+(\tilde{p}_{1:ij}(n), 0.5) \leq \text{kl}^+(\mu_i, 0.5)\right) \\
&\leq \sum_{n=\lceil M_{ij}(\delta) \rceil}^T \mathbf{P}\left(\tilde{p}_{1:ij}(n) \leq \mu_i\right) \\
&\leq \sum_{n=\lceil M_{ij}(\delta) \rceil}^T \mathbf{P}\left(\tilde{p}_{ij}(n) \leq \mu_i\right) \\
&\leq \sum_{n=\lceil M_{ij}(\delta) \rceil}^T e^{-\text{kl}(\mu_i, p_{ij})n},
\end{aligned}$$

where the last inequality is by Chernoff's inequality (e.g. see Fact 8 of [21]). Then,

$$S_T(i, j) \leq \sum_{n=1}^{\infty} e^{-\text{kl}(\mu_i, p_{ij})n} \leq \frac{1}{\text{kl}(\mu_i, p_{ij})}.$$

The proof is concluded using Pinsker's inequality followed by $4\Delta(i, j)$ -Lipschitzness of $\text{kl}(\cdot, 0.5)$ over $(0.5 - \Delta(i, j), 0.5)$:

$$\begin{aligned}
\text{kl}(\mu_i, p_{ij}) &\geq 2(\mu_i - p_{ij})^2 && \leftarrow \text{Pinsker's inequality} \\
&\geq \frac{2}{16\Delta(i, j)^2} (\text{kl}(\mu_i, 0.5) - \text{kl}(p_{ij}, 0.5))^2 && \leftarrow \text{Lipschitzness} \\
&= \frac{2\text{kl}(p_{ij}, 0.5)^2\delta^2}{16\Delta(i, j)^2(1 + \delta)^2} && \leftarrow \text{def of } \mu_i \\
&\geq \frac{\text{kl}(p_{ij}, 0.5)^2\delta^2}{32\Delta(i, j)^2} && \leftarrow \delta \in (0, 1) \\
&\geq \frac{\Delta(i, j)^2\delta^2}{32}. && \leftarrow \text{Pinsker's inequality}
\end{aligned}$$

Therefore,

$$S_T(i, j) \leq \frac{32}{\Delta(i, j)^2\delta^2}.$$

□

Lemma 10. For any $\varepsilon_2, \dots, \varepsilon_K \geq 0$,

$$\sum_{1 \leq i < j \leq K | \Delta(i, j) > \varepsilon_j} \frac{\Delta(i, j) - \Delta(i + 1, j)}{\Delta(i, j)^2} \leq \sum_{j=2}^n \frac{2}{\max\{\varepsilon_j, \Delta(j - 1, j)\}}.$$

Proof. The proof is adapted from similar techniques used for proving Lem. 5 of [19]. First note that

$$\sum_{1 \leq i < j \leq n | \Delta(i, j) > \varepsilon_j} \frac{\Delta(i, j) - \Delta(i + 1, j)}{\Delta(i, j)^2} = \sum_{i=1}^{K-1} \sum_{j \in [K] \setminus [i] | \Delta(i, j) > \varepsilon_j} \frac{\Delta(i, j) - \Delta(i + 1, j)}{\Delta(i, j)^2}$$

Let us fix any arm $i \in [K - 1]$, and denote by $\nabla_{i, j} := \Delta(i, j) - \Delta(i + 1, j)$. Then we note

$$\sum_{j \in [K] \setminus [i] | \Delta(i, j) > \varepsilon_j} \frac{\Delta(i, j) - \Delta(i + 1, j)}{\Delta(i, j)^2} = \sum_{j \in [K] \setminus [i] | \Delta(i, j) > \varepsilon_j} \nabla_{i, j} \int_0^{\infty} \mathbf{1}(\Delta(i, j)^{-2} \geq x) dx$$

$$\begin{aligned}
&= \sum_{j=i+1}^K \mathbf{1}(\Delta(i, j) > \varepsilon_j) \nabla_{i,j} \int_0^\infty \mathbf{1}(\Delta(i, j)^{-2} \geq x) dx \\
&= \sum_{j=i+1}^K \nabla_{i,j} \int_0^\infty \mathbf{1}(\Delta(i, j) > \varepsilon_j, \Delta(i, j)^{-2} \geq x) dx \\
&= 2 \sum_{j=i+1}^K \nabla_{i,j} \int_0^\infty y^{-3} \mathbf{1}(\varepsilon_j < \Delta(i, j) < y) dy \quad \leftarrow \text{change of variable } x = y^{-2}, dx = -2y^{-3} dy \\
&= 2 \sum_{j=i+1}^K \nabla_{i,j} \int_{\varepsilon_j}^\infty y^{-3} \mathbf{1}(\varepsilon_j < \Delta(i, j) < y) dy
\end{aligned}$$

Further summing over all $i \in [K-1]$, we get

$$\begin{aligned}
A_T &:= \sum_{1 \leq i < j \leq n | \Delta(i, j) > \varepsilon_j} \frac{\Delta(i, j) - \Delta(i+1, j)}{\Delta(i, j)^2} \\
&= \sum_{i=1}^{K-1} \left(2 \sum_{j=i+1}^K \nabla_{i,j} \int_{\varepsilon_j}^\infty y^{-3} \mathbf{1}(\varepsilon_j < \Delta(i, j) < y) dy \right) \\
&= 2 \sum_{i=1}^{K-1} \sum_{j=i+1}^K \int_{\varepsilon_j}^\infty \left(\nabla_{i,j} y^{-3} \mathbf{1}(\varepsilon_j < \Delta(i, j) < y) dy \right) \\
&= 2 \sum_{j=2}^K \sum_{i=1}^{j-1} \int_{\varepsilon_j}^\infty y^{-3} \left((\Delta(i, j) - \Delta(i+1, j)) \mathbf{1}(\varepsilon_j < \Delta(i, j) \leq y) \right) dy \\
&= 2 \sum_{j=2}^K \int_{\varepsilon_j}^\infty y^{-3} \sum_{i=i_\varepsilon(j)}^{i_\varepsilon(j)-1} (\Delta(i, j) - \Delta(i+1, j)) dy,
\end{aligned}$$

where $i_\varepsilon(j) := \arg \min\{i | i \leq j, \Delta(i, j) \leq \varepsilon\}$ (with the convention that the sum is empty if the arg min is empty) and because $\varepsilon_j < \Delta(i, j) \leq y$ is equivalent to $i_\varepsilon(\varepsilon) \leq i \leq i_\varepsilon - 1$. Using telescoping summation over i , we further get:

$$\begin{aligned}
A_T &\leq 2 \sum_{j=2}^K \int_{\varepsilon_j}^\infty y^{-3} (\Delta(i_y(j), j) - \Delta(i_{\varepsilon_j}(j), j)) dy \\
&\leq 2 \sum_{j=2}^K \int_{\varepsilon_j}^\infty y^{-3} \Delta(i_y(j), j) dy \quad \leftarrow \text{since } \Delta(i_{\varepsilon_j}(j), j) > 0
\end{aligned}$$

Then, since $\Delta(i_y(j), j) = 0$ if $y < \Delta(j-1, j)$, we have

$$\begin{aligned}
A_T &\leq 2 \sum_{j=2}^K \int_{\max\{\varepsilon_j, \Delta(j-1, j)\}}^\infty y^{-3} \Delta(i_y(j), j) dy \\
&\leq 2 \sum_{j=2}^K \int_{\max\{\varepsilon_j, \Delta(j-1, j)\}}^\infty y^{-2} dy \quad \leftarrow \text{since } \Delta(i_y(j), j) \leq y \\
&\leq 2 \sum_{j=2}^K \frac{1}{\max\{\varepsilon_j, \Delta(j-1, j)\}},
\end{aligned}$$

which concludes the proof. \square

10.2 Proof of Theorem 6

Theorem 6 (Expected regret analysis S1DB-ED). *Let $t_0 = 1$ and $\alpha = 4K$. Then as $T \rightarrow \infty$, the expected regret incurred by S1DB-ED (Alg. 2) can be upper bounded as: For all $\varepsilon_2, \dots, \varepsilon_K \geq 0$*

$$\begin{aligned} \mathbf{E}[R_T] &\lesssim K^2 + \sum_{1 \leq i < j \leq K} \left(\frac{K \mathbf{1}_{\{\Delta(i,j) > \varepsilon_j\}}}{\Delta(i,j)^2} + n_{ij}(T) \min\{\varepsilon_j, \Delta(i,j)\} \right) + \sum_{j=2}^K \frac{K \log T}{\max\{\varepsilon_j, \Delta(j-1, j)_+\}} \\ &\leq O\left(\min \left\{ \sum_{j=2}^K \frac{K \log T}{\Delta(j-1, j)_+}, KT^{2/3} \right\} \right). \end{aligned}$$

Proof. We analyse the expected regret S1DB-ED (Alg 2) for some fixed sequence \mathcal{S}_T . Recall that t_0 is the budget spent on exploration of each pair (i, j) and the notation

$$\mathcal{E}(t) := \{n_{ij}(t) > t_0, \forall i, j \in S_t, i \neq j\},$$

to be the event when all distinct pairs in S_t have been explored t_0 times and

$$\mathcal{G}(t) := \{\forall j > i_t^*, \Delta(i_t^*, j) > \varepsilon_i, \widehat{p}_{i_t^* j}(t) > 1/2\}$$

the event when the probabilities p_{ij} have been well estimated by the algorithm. Then, from Lemma 8, we have

$$\begin{aligned} \mathbf{E}[R_T] &= \mathbf{E}\left[\sum_{t=1}^T r_t \right] = \mathbf{E}\left[\sum_{t=1}^T \mathbf{1}(\mathcal{E}^c(t)) r_t + \sum_{t=1}^T \mathbf{1}(\mathcal{E}(t)) \mathbf{1}(\mathcal{G}^c(t)) r_t + \sum_{t=T_0+1}^T \mathbf{1}(\mathcal{E}(t)) \mathbf{1}(\mathcal{G}(t)) r_t \right] \\ &\leq K^2 t_0 + 2K + K \sum_{i=1}^{K-1} \sum_{j=i+1}^K \frac{e^{-(t_0-1)\Delta(i,j)^2}}{\Delta(i,j)^2} + \underbrace{\mathbf{E}\left[\sum_{t=T_0+1}^T \mathbf{1}(\mathcal{E}(t)) \mathbf{1}(\mathcal{G}(t)) r_t \right]}_{E_T}. \end{aligned} \quad (9)$$

We now upper-bound the third term of (9). Remark that under $\mathcal{G}(t)$ the algorithm chooses $y_t = i_t^*$. Therefore,

$$\begin{aligned} E_T &:= \sum_{t=1}^T \mathbf{1}(\mathcal{E}(t)) \mathbf{1}(\mathcal{G}(t)) r_t \\ &\leq \sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) r_t \\ &= \sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \sum_{1 \leq i < j \leq K} \mathbf{1}(x_t = j, y_t = i) r_t \\ &= \sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \sum_{1 \leq i < j \leq K} \mathbf{1}(x_t = j, y_t = i) \frac{\Delta(i,j)}{2} \quad \leftarrow \text{because } \mathcal{G}(t) \text{ implies } y_t = i_t^* \\ &\leq \underbrace{\sum_{1 \leq i < j \leq K: \Delta(i,j) < \varepsilon_j} n_{ij}(T) \frac{\Delta(i,j)}{2} + \sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \sum_{1 \leq i < j \leq K: \Delta(i,j) > \varepsilon_j} \mathbf{1}(x_t = j, y_t = i) \frac{\Delta(i,j)}{2}}_{=: D_T} \end{aligned} \quad (10)$$

Moreover, recalling the notations $n_{ij}(T) := \sum_{t=1}^T \mathbf{1}(\{x_t, y_t\} = \{i, j\})$ and defining

$$\tilde{N}_{ij}(T) := \sum_{k=1}^i \sum_{s=1}^t \mathbf{1}(\mathcal{G}(s)) \mathbf{1}(x_s = k, y_s = j) \leq N_{ij}(T) := \sum_{k=1}^i n_{kj}(T),$$

we have

$$\begin{aligned}
D_T &:= \sum_{1 \leq i < j \leq K: \Delta(i,j) > \varepsilon_j} \sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \mathbf{1}(x_t = j, y_t = i) \frac{\Delta(i,j)}{2} \\
&= \sum_{1 \leq i < j \leq K: \Delta(i,j) > \varepsilon_j} (\tilde{N}_{ij}(T) - \tilde{N}_{(i-1)j}(T)) \frac{\Delta(i,j)}{2} \\
&= \sum_{j=2}^K \sum_{i=1}^{i_{\varepsilon_j}(j)} (\tilde{N}_{ij}(T) - \tilde{N}_{(i-1)j}(T)) \frac{\Delta(i,j)}{2} \\
&= \sum_{j=2}^K \tilde{N}_{i_{\varepsilon_j}j}(T) \frac{\varepsilon_j}{2} + \sum_{j=2}^K \sum_{i=1}^{i_{\varepsilon_j}(j)-1} \tilde{N}_{ij}(T) \frac{\Delta(i,j) - \Delta(i+1,j)}{2} \\
&\leq \sum_{1 \leq i < j \leq K: \Delta(i,j) \geq \varepsilon_j} n_{ij}(T) \frac{\varepsilon_j}{2} + \sum_{1 \leq i < j \leq K: \Delta(i,j) > \varepsilon_j} \tilde{N}_{ij}(T) \frac{\Delta(i,j) - \Delta(i+1,j)}{2}. \tag{11}
\end{aligned}$$

Now, we need to upper-bound $\tilde{N}_{ij}(T)$. We have,

$$\begin{aligned}
\tilde{N}_{i,j}(T) &:= \sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \\
&\leq \sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \left[\mathbf{1}(N_{ij}(t) \leq M_{ij}(\delta)) + \mathbf{1}(N_{ij}(t) > M_{ij}(\delta)) \right] \\
&\leq M_{ij}(\delta) + \sum_{t=1}^T \mathbf{1}(\mathcal{G}(t)) \sum_{k=1}^i \mathbf{1}(x_t = j, y_t = k) \mathbf{1}(N_{ij}(t) > M_{ij}(\delta)) \\
&\leq M_{ij}(\delta) + \frac{32}{\delta^2 \Delta(i,j)^2} \quad \leftarrow \text{Lemma 9} \\
&= \frac{(\alpha + \delta) \log T}{\text{kl}(p_{ji}, 0.5)} + \frac{32}{\delta^2 \Delta(i,j)^2} \\
&\leq \left(2(\alpha + \delta) \log T + \frac{32}{\delta^2} \right) \frac{1}{\Delta(i,j)^2},
\end{aligned}$$

where the last inequality comes from Pinsker's inequality. This entails

$$\begin{aligned}
&\sum_{1 \leq i < j \leq K: \Delta(i,j) > \varepsilon_j} \tilde{N}_{ij}(T) \frac{\Delta(i,j) - \Delta(i+1,j)}{2} \\
&\leq \left((\alpha + \delta) \log T + \frac{16}{\delta^2} \right) \sum_{1 \leq i < j \leq K: \Delta(i,j) > \varepsilon_j} \frac{\Delta(i,j) - \Delta(i+1,j)}{\Delta(i,j)^2} \\
&\leq 2 \sum_{j=2}^K \frac{(\alpha + \delta) \log T + 16\delta^{-2}}{\max\{\varepsilon_j, \Delta(j-1,j)\}}.
\end{aligned}$$

Combining this inequality with (9), (10), and (11) and choosing $t_0 = 1$ and $\alpha = 4K$ concludes

$$\begin{aligned}
\mathbf{E}[R_T] &\leq K^2 t_0 + 2K + K \sum_{i=1}^{K-1} \sum_{j=i+1}^K \frac{e^{-(t_0-1)\Delta(i,j)^2}}{\Delta(i,j)^2} \\
&\quad + \sum_{1 \leq i < j \leq K} n_{ij}(T) \frac{\min\{\varepsilon_j, \Delta(i,j)\}}{2} + 2 \sum_{j=2}^n \frac{(\alpha + \delta) \log T + 16\delta^{-2}}{\max\{\varepsilon_j, \Delta(j-1,j)\}}
\end{aligned}$$

$$\begin{aligned} &\leq K(K+2) + \sum_{1 \leq i < j \leq K | \Delta(i,j) > \varepsilon_j} \frac{K}{\Delta(i,j)^2} \\ &\quad + \sum_{1 \leq i < j \leq K} n_{ij}(T) \frac{\min\{\varepsilon_j, \Delta(i,j)\}}{2} + 2 \sum_{j=2}^K \frac{(4K + \delta) \log T + 16\delta^{-2}}{\max\{\varepsilon_j, \Delta(j-1, j)\}}. \end{aligned}$$

□

Proof. Recall that the proof was done for any $\varepsilon_2, \dots, \varepsilon_K \geq 0$ that are independent of the algorithm. In particular, choosing $\varepsilon_2, \dots, \varepsilon_K = \varepsilon$ entails that for any $\varepsilon > 0$

$$\mathbf{E}[R_T] \lesssim K^2 + \varepsilon T + \sum_{1 \leq i < j \leq K | \Delta(i,j) > \varepsilon} \frac{K}{\Delta(i,j)^2} + \sum_{j=2}^K \frac{K \log T}{\max\{\varepsilon, \Delta(j-1, j)\}}$$

which yields making $\varepsilon \rightarrow 0$ the distribution-dependent asymptotic upper-bound

$$\mathbf{E}[R_T] \leq O\left(K \log(T) \sum_{j=2}^K \frac{\mathbf{1}\{\Delta(j-1, j) > 0\}}{\Delta(j-1, j)}\right)$$

as $T \rightarrow \infty$ and for any fix $\varepsilon \geq 0$ and choosing $\delta = 1$. Furthermore, optimizing $\varepsilon_1 = \varepsilon_2 = \dots = \varepsilon_K = \varepsilon = 2^{1/3}KT^{-1/3}$ yields the distribution-free upper-bound

$$\mathbf{E}[R_T] \leq K(K+2) + \frac{K^3}{\varepsilon^2} + \frac{T\varepsilon}{2} + \frac{(8K+1)K \log T + 16K}{\varepsilon} \leq 2KT^{2/3} + O(K^2 + KT^{1/3} \log T).$$

□