



**HAL**  
open science

## Parametric Graph for Unimodal Ranking Bandit

Camille-Sovanneary Gauthier, Romaric Gaudel, Elisa Fromont, Boammani  
Aser Lompo

► **To cite this version:**

Camille-Sovanneary Gauthier, Romaric Gaudel, Elisa Fromont, Boammani Aser Lompo. Parametric Graph for Unimodal Ranking Bandit. ICML 2021 - International Conference on Machine Learning, Jul 2021, Virtual, Canada. pp.3630–3639. hal-03256621v2

**HAL Id: hal-03256621**

**<https://hal.science/hal-03256621v2>**

Submitted on 23 Jun 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# Parametric Graph for Unimodal Ranking Bandit

---

Camille-Sovanneary Gauthier<sup>\*12</sup> Romaric Gaudel<sup>\*3</sup> Elisa Fromont<sup>452</sup> Boammani Aser Lompo<sup>6</sup>

## Abstract

We tackle the online ranking problem of assigning  $L$  items to  $K$  positions on a web page in order to maximize the number of user clicks. We propose an original algorithm, easy to implement and with strong theoretical guarantees to tackle this problem in the Position-Based Model (PBM) setting, well suited for applications where items are displayed on a grid. Besides learning to rank, our algorithm, GRAB (for *parametric Graph for unimodal RAnking Bandit*), also learns the parameter of a compact graph over permutations of  $K$  items among  $L$ . The logarithmic regret bound of this algorithm is a direct consequence of the unimodality property of the bandit setting with respect to the learned graph. Experiments against state-of-the-art learning algorithms which also tackle the PBM setting, show that our method is more efficient while giving regret performance on par with the best known algorithms on simulated and real life datasets.

## 1. Introduction

*Online Recommendation Systems* are used to choose relevant items, such as songs, adds or movies on a website. At each call, they select  $K$  items among  $L$  potential ones,  $K \leq L$ . User feedbacks are collected for each displayed items, reflecting how relevant these choices are: listening time, clicks, rates, etc. These feedbacks are available only for the items which were presented to the user. This corresponds to an instance of the *multi-armed bandit problem with semi-bandit feedback* (Gai et al., 2012; Chen et al., 2013). Another problem, related to ranking, is to display the  $K$  chosen items at the right positions to maximize the

user attention. Typical examples of such displays are (i) a list of news, visible one by one by scrolling; (ii) a list of products, arranged by rows; or (iii) advertisements spread in a web page. Numerous approaches have been proposed to jointly learn how to choose the best positions for the corresponding best items (Radlinski et al., 2008; Combes et al., 2015; Li et al., 2019) referred as *multiple-play bandit* or *online learning to rank* (OLR). To take into account the user behaviour while facing such a list of items, several models exist (Richardson et al., 2007; Craswell et al., 2008) and have been transposed to the bandit framework (Kveton et al., 2015a; Komiyama et al., 2017) such as the *Position-Based Model* (PBM) (Richardson et al., 2007). PBM allows to take into account displays where the best position is *a priori* unknown. This is typically the case when items are displayed on a grid and not in an ordered list. PBM assumes that the click probability on an item  $i$  at position  $k$  results from the product of two independent factors: the item relevance and its position's visibility. Items displayed at other positions do not impact the probability to consider the item  $i$  at position  $k$ . According to PBM, a user may give more than one feedbacks: she may click on all items relevant for her, e.g. when looking for product on commercial websites. PBM is also particularly interesting when the display is dynamic, as often on modern web pages, and may depend on the reading direction of the user (which varies from one country to another) and on the ever-changing layout of the page.

In this paper, we tackle an online learning to rank bandit setting, which mainly covers PBM click model, with an *unimodal bandits* point of view (Combes & Proutière, 2014). First, we expose a family of parametric graphs of degree  $L - 1$  over permutations, such that the PBM setting is unimodal w.r.t. one graph in this family. While the corresponding graph is unknown from the learner, graphs of this family enable an efficient exploration strategy of the set of potential recommendations. Secondly, we introduce a new bandit algorithm, GRAB, which learns online the appropriate graph in this family and bases its recommendations on the learned graph. From an application point of view, this algorithm has several interesting features: it is simple to implement and efficient in terms of computation time; it handles the PBM bandit setting without any knowledge on the impact of positions (contrarily to many competitors); and it empirically exhibits a regret on par with other theoretically proven

---

<sup>\*</sup>Equal contribution <sup>1</sup>Louis Vuitton, F-75001 Paris, France <sup>2</sup>IRISA UMR 6074 / INRIA rba, F-35000 Rennes, France <sup>3</sup>Univ Rennes, Ensai, CNRS, CREST - UMR 9194, F-35000 Rennes, France <sup>4</sup>Univ. Rennes 1, F-35000 Rennes, France <sup>5</sup> Institut Universitaire de France, M.E.S.R.I., F-75231 Paris <sup>6</sup>ENS Rennes, F-35000 Rennes, France. Correspondence to: Camille-Sovanneary Gauthier <camille-sovanneary.gauthier@louisvuitton.com>.

Table 1. Settings and upper-bound on cumulative regret for state of the art algorithms. The main notations for the assumptions are given in Section 3.  $\mathcal{N}_{\pi^*}(\mathbf{a}^*)$  is a set of recommendations in the neighborhood of the best recommendation.  $K_{max}$  is the maximum number of differences between two arms; see. Theorem 2 for a specific definition.

ALGORITHM	HANDLED SETTINGS	REGRET	$\Delta$ , ASSUMING $\theta_1 \geq \theta_2 \geq \dots \geq \theta_L$
GRAB (OUR ALGORITHM)	PBM	$\mathcal{O}\left(\frac{L}{\Delta} \log T\right)$	$\min_{\mathbf{a} \in \mathcal{N}_{\pi^*}(\mathbf{a}^*)} \mu^* - \mu_{\mathbf{a}}$
COMBUCB1 (Kveton et al., 2015b)	PBM	$\mathcal{O}\left(\frac{LK^2}{\Delta} \log T\right)$	$\min_{\mathbf{a} \in \mathcal{P}_K^L} \mu^* - \mu_{\mathbf{a}}$
PBM-PIE (Lagrée et al., 2016)	PBM WITH $\kappa$ KNOWN	$\mathcal{O}\left(\frac{L-K}{\Delta} \log T\right)$	$\min_{i \in \{K+1, \dots, L\}} \mu^* - \mu_{\mathbf{a}[K:=i]}$
PMED-HINGE (Komiya et al., 2017)	PBM WITH $\kappa_1 \geq \dots \geq \kappa_K$	$\mathcal{O}(c^*(\boldsymbol{\theta}, \boldsymbol{\kappa}) \log T)$	$\emptyset$
TOPRANK (Lattimore et al., 2018)	PBM WITH $\kappa_1 \geq \dots \geq \kappa_K$ , CM, ...	$\mathcal{O}\left(\frac{LK}{\Delta} \log T\right)$	$\min_{(j,i) \in [L] \times [K]: j > i} \frac{\theta_i - \theta_j}{\theta_i}$
OSUB (Combes & Proutière, 2014)	UNIMODAL	$\mathcal{O}\left(\frac{\gamma}{\Delta} \log T\right)$	$\min_{\mathbf{a} \in \mathcal{N}_G(\mathbf{a}^*)} \mu^* - \mu_{\mathbf{a}}$
KL-COMBUCB (THEOREM 2)	COMBINATORIAL	$\mathcal{O}\left(\frac{ A K_{max}^2}{\Delta} \log T\right)$	$\min_{\mathbf{a} \in \mathcal{A}} \mu^* - \mu_{\mathbf{a}}$

algorithms on both artificial and real datasets. In particular, we prove a  $\mathcal{O}(L/\Delta \log T)$  regret upper-bound for GRAB. The corresponding proof extends OSUB’s proof (Combes & Proutière, 2014) both (i) to the context of a graph learned online, and (ii) to the combinatorial semi-bandit setting.

This paper is organized as follows: Section 2 presents the related work and Section 3 defines our target setting. We introduce GRAB and the hypotheses needed in Section 4. Theoretical guarantees and empirical performance are presented respectively in Section 5 and 6. We conclude in Section 7.

## 2. Related Work

A comparison of the assumptions and the regret upper-bounds of the related algorithms is shown in Table 1.

The Position-Based Model (PBM) (Richardson et al., 2007; Craswell et al., 2008) relies on two vectors of parameters:  $\boldsymbol{\theta} \in [0, 1]^L$  and  $\boldsymbol{\kappa} \in [0, 1]^K$ , where  $\theta_i$  is the probability for the user to click on item  $i$  when she observes this item, and  $\kappa_k$  is the probability for the user to observe position  $k$ . Several bandit algorithms are designed to handle PBM (Komiya et al., 2015; Lagrée et al., 2016; Komiya et al., 2017). However, each of them assumes some knowledge about the ranking of positions. (Komiya et al., 2015) and (Lagrée et al., 2016) assume  $\boldsymbol{\kappa}$  known beforehand. Thanks to this very strong assumption (that we do not make in this paper), the theoretical results from (Lagrée et al., 2016) depend on the  $L - K$  worst items and their regret is expressed as  $\mathcal{O}((L - K)/\Delta \log T)$ . (Komiya et al., 2017) and (Gauthier et al., 2021) propose respectively PMED and PB-MHB, the only approaches learning both  $\boldsymbol{\theta}$  and  $\boldsymbol{\kappa}$  while recommending. However, PMED still requires

the  $\kappa_k$  values to be organized in decreasing order. It derives a bound on the regret in  $\mathcal{O}(c^*(\boldsymbol{\theta}, \boldsymbol{\kappa}) \log T)$ , where  $c^*(\boldsymbol{\theta}, \boldsymbol{\kappa})$  only depends on  $\boldsymbol{\theta}$  and  $\boldsymbol{\kappa}$  and is asymptotically optimal in this setting. Unfortunately, to the best of our knowledge, there is no known closed-form for  $c^*(\boldsymbol{\theta}, \boldsymbol{\kappa})$ , which hinders the comparison to other algorithms, including ours. (Gauthier et al., 2021) has shown very good performances on PBM, but PB-MHB, based on Thomson sampling, does not have any theoretical guarantees. Other learning to rank algorithms such as TopRank (Lattimore et al., 2018) and BubbleRank (Li et al., 2019) cover many click models (including PBM). They exhibit a regret upper-bound for  $T$  iterations of  $\mathcal{O}(LK/\Delta \log T)$ , where  $\Delta$  depends on the attraction probability  $\boldsymbol{\theta}$  of items. They also assume that the best recommendation is the one displaying the items from the most attractive to the  $K$ -ith most attractive, which implies that the first position is the most-observed one, the second position is the second most-observed one, and so on.

Although the hypotheses taken by PMED and TopRank are often assumed by the approaches handling PBM setting. In this paper we tackle a full PBM setting where there is no a priori hypothesis on the ordering of positions. Our algorithm, GRAB, suffers a  $\mathcal{O}(L/\Delta \log T)$  regret that we conjecture to be on par with the best theoretical results provided by PMED ( $\mathcal{O}(c^*(\boldsymbol{\theta}, \boldsymbol{\kappa}) \log T)$ ).

Combinatorial algorithms (Gai et al., 2012; Chen et al., 2013) can also handle the PBM bandit setting. Typically, CombUCB1 (Kveton et al., 2015b) applied to PBM leads to an algorithm which suffers a  $\mathcal{O}(LK^2/\Delta \log T)$  regret (see the appendix for more details), which is higher than the upper-bound on the regret of GRAB by a factor  $K^2$ . Note that the proof of the upper-bound on the regret of GRAB is based on the same reduction of the PBM bandit

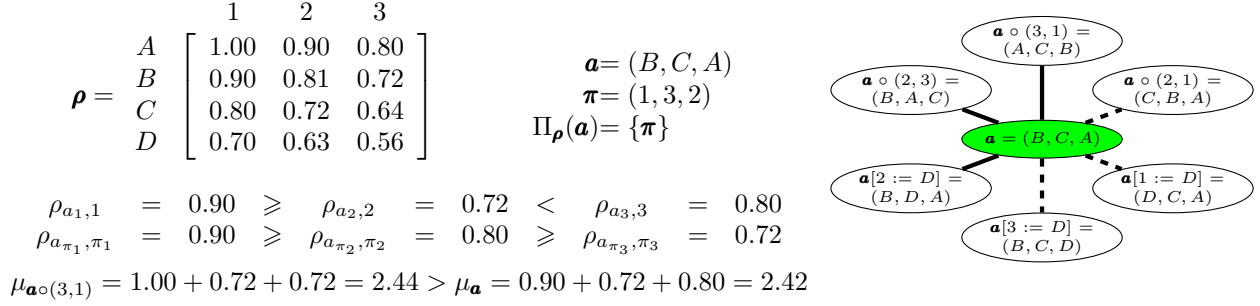


Figure 1. Assumption 1 in practice. To distinguish between items and positions, the 4 items are denoted  $A, B, C$ , and  $D$ . **On the left:** parameters and considered recommendation  $\mathbf{a}$ . We consider a matrix of probabilities of clicks  $\rho$  which corresponds to a PBM click model, and a sub-optimal recommendation  $\mathbf{a}$ . The corresponding set  $\Pi_{\rho}(\mathbf{a})$  of appropriate rankings of positions is composed of a unique permutation  $\boldsymbol{\pi}$ . **On the right:** corresponding neighborhoods. Solid lines identify the neighborhood  $\mathcal{N}_{\boldsymbol{\pi}}(\mathbf{a})$  used by GRAB, and both solid and dashed lines correspond to the neighborhood  $\mathcal{N}_G(\mathbf{a})$  used by S-GRAB. Note that there is a recommendation better than  $\mathbf{a}$  in both neighborhoods:  $\mathbf{a} \circ (3, 1) = (A, C, B)$ .

to a combinatorial semi-bandit, but with two additional properties derived from the design of GRAB.

Finally, GRAB extends the *unimodal* bandit setting (Combes & Proutière, 2014) which assumes the existence of a known graph  $G$  carrying a partial order on the set of bandit arms denoted  $\mathcal{A}$ . The unimodal bandit algorithms are aware of  $G$ , but ignore the partial order induced by the edges of  $G$ . However, they rely on  $G$  to efficiently browse the arms up to the best one. Typically, the algorithm OSUB (Combes & Proutière, 2014) selects at each iteration  $t$ , an arm  $\mathbf{a}(t)$  in the neighborhood  $\mathcal{N}_G(\tilde{\mathbf{a}}(t))$  given  $G$  of the current best arm  $\tilde{\mathbf{a}}(t)$  (a.k.a. the *leader*). By restricting the exploration to this neighborhood, the regret suffered by OSUB scales as  $\mathcal{O}(\gamma/\Delta \log T)$ , where  $\gamma$  is the maximum degree of  $G$ , to be compared with  $\mathcal{O}(|\mathcal{A}|/\Delta \log T)$  if the arms were independent.

### 3. Learning to Rank in a Semi-Bandit Setting

We consider the following *online learning to rank (OLR) problem with clicks feedback* which encompasses the PBM setting. For any integer  $n$ , let  $[n]$  denote the set  $\{1, \dots, n\}$ . An instance of our OLR problem is a tuple  $(L, K, (\rho_{i,k})_{(i,k) \in [L] \times [K]})$ , where  $L$  is the number of items to be displayed,  $K \leq L$  is the number of positions to display the items, and for any item  $i$  and position  $k$ ,  $\rho_{i,k}$  is the probability for a user to click on item  $i$  when displayed at position  $k$ , independently of the items displayed at other positions. Under PBM click-model, there exists two vectors  $\boldsymbol{\theta} \in \mathbb{R}^L$  and  $\boldsymbol{\kappa} \in \mathbb{R}^K$ , such that  $\rho_{i,k} = \theta_i \kappa_k$  (i.e.  $\rho$  is of rank 1).

A recommendation algorithm is only aware of  $L$  and  $K$  and has to deliver  $T$  consecutive recommendations. At each iteration  $t \in [T]$ , the algorithm recommends a permutation  $\mathbf{a}(t) = (a_1(t), \dots, a_K(t))$  of  $K$  distinct items among  $L$ ,

where  $a_k(t)$  is the item displayed at position  $k$ . We denote  $\mathcal{A} = \mathcal{P}_K^L$  the set of such permutations, which corresponds to the set of arms of the bandit setting. Throughout the paper, we will use the terms *permutation* and *recommendation* interchangeably to denote an element of  $\mathcal{P}_K^L$ . Thereafter, the algorithm observes the clicks vector  $\mathbf{c}(t) \in \{0, 1\}^K$ , where for any position  $k$ ,  $c_k(t)$  equals 1 if the user clicks on item  $a_k(t)$  displayed at position  $k$ , and 0 otherwise. Note that the recommendation at time  $t$  is only based on previous recommendations and observations.

While the individual clicks are observed, the reward of the algorithm is their total number  $r(t) \stackrel{\text{def}}{=} \sum_{k=1}^K c_k(t)$ . Let  $\mu_{\mathbf{a}}$  denote the expectation of  $r(t)$  when the recommendation is  $\mathbf{a}(t) = \mathbf{a}$ , and  $\mu^* \stackrel{\text{def}}{=} \max_{\mathbf{a} \in \mathcal{P}_K^L} \mu_{\mathbf{a}}$  the highest expected reward. The aim of the algorithm is to minimize the *cumulative (pseudo-) regret*

$$R(T) = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T \mu_{\mathbf{a}(t)} \right], \quad (1)$$

where the expectation is taken w.r.t. the recommendations from the algorithm and the clicks. Note that for any recommendation  $\mathbf{a} \in \mathcal{P}_K^L$ ,  $\mu_{\mathbf{a}} = \sum_{k=1}^K \rho_{a_k, k}$ .

#### 3.1. Modeling Assumption

Apart from the independency of the clicks, the proposed algorithm assumes a relaxed version of unimodality. Here we present this assumption and state its relation with PBM. We first define the set of *appropriate rankings of positions*: for each recommendation  $\mathbf{a} \in \mathcal{P}_K^L$ , we denote  $\Pi_{\rho}(\mathbf{a}) \subseteq \mathcal{P}_K^K$  the set of permutations  $\boldsymbol{\pi}$  of the  $K$  positions such that  $\rho_{a_{\pi_1}, \pi_1} \geq \rho_{a_{\pi_2}, \pi_2} \geq \dots \geq \rho_{a_{\pi_K}, \pi_K}$ . Therefore, an appropriate ranking of positions orders the positions from the one with the highest probability of click to the one with the

lowest probability of click. See Figure 1 for an example.

With this notation, our assumption is the following:

**Assumption 1** (Relaxed Unimodality). *For any recommendation  $\mathbf{a} \in \mathcal{P}_K^L$  and any ranking of positions  $\pi \in \Pi_\rho(\mathbf{a})$ , if  $\mu_{\mathbf{a}} \neq \mu^*$ , then either there exists  $k \in [K-1]$  such that*

$$\mu_{\mathbf{a}} < \mu_{\mathbf{a} \circ (\pi_k, \pi_{k+1})} \quad (2)$$

or there exists  $i \in [L] \setminus \mathbf{a}([K])$  such that

$$\mu_{\mathbf{a}} < \mu_{\mathbf{a}[\pi_K := i]}, \quad (3)$$

where

- $\mathbf{a} \circ (\pi_k, \pi_{k+1})$  is the permutation for which the items at positions  $\pi_k$  and  $\pi_{k+1}$  are swapped,
- $\mathbf{a}[\pi_K := i]$  is the permutation which is the same as  $\mathbf{a}$  for any position  $k \neq \pi_K$ , and such that  $\mathbf{a}[\pi_K := i]_{\pi_K} = i$ ,
- and  $\mathbf{a}([K])$  is the set of items recommended by  $\mathbf{a}$ , namely  $\mathbf{a}([K]) \stackrel{\text{def}}{=} \{a_1, \dots, a_K\}$ .

Assumption 1 relates to a natural property of standard click models: (i) for the optimal recommendation, the position with the  $k$ -th highest probability to be observed is the one displaying the  $k$ -th most attractive item, (ii) for a sub-optimal recommendation, swapping two consecutive items, given this order, leads to an increase of the expected reward. However, Assumption 1 considers the order based on the click probabilities  $\rho_{a_k, k}$ , not on the observation probabilities  $\kappa_k$ . Figure 1 gives an example of both orders and of the neighborhood associated to the ranking  $\pi$  defined after the order on click probabilities  $\rho_{a_k, k}$ .

While the existence of a better recommendation in the neighborhood defined given this order is less intuitive, it remains true for state of the art click models (PBM, the cascading model, and the dependent click model) and paves the way to an algorithm based on observed random variables. Note also that while there exists a better recommendation both in the neighborhood based on the order on observation probability and in the neighborhood based on the order on click probability, this is not true for any neighborhood based on any arbitrary order (as soon as  $K \geq 4$ ).

Hereafter, Lemma 1, states that Assumption 1 is weaker than the PBM one. The proof of this Lemma is deferred to the appendix.

**Lemma 1.** *Let  $(L, K, (\theta_i \kappa_k)_{(i,k) \in [L] \times [K]})$  be an online learning to rank problem with users following PBM, with positive probabilities of looking at a given position. Then Assumption 1 is fulfilled.*

### 3.2. Relation with Unimodality

Assumption 1 relates to the unimodality of the set of expected rewards  $(\mu_{\mathbf{a}})_{\mathbf{a} \in \mathcal{P}_K^L}$ . Let us first recall the definition of unimodality in (Combes & Proutière, 2014) and then express this relation.

**Definition 1** (Unimodality). *Let  $\mathcal{A}$  be a set of arms, and  $(\nu_{\mathbf{a}})_{\mathbf{a} \in \mathcal{A}}$  a set of reward distributions of respective expectations  $(\mu_{\mathbf{a}})_{\mathbf{a} \in \mathcal{A}}$ . Let  $G = (V, E)$  be an undirected graph with vertices  $V = \mathcal{A}$  and edges  $E \subseteq V^2$ . The set of expected rewards  $(\mu_{\mathbf{a}})_{\mathbf{a} \in \mathcal{A}}$  is unimodal w.r.t.  $G$ , if and only if:*

1. *the set of expected rewards admits a unique best arm:  $\operatorname{argmax}_{\mathbf{a} \in \mathcal{A}} \mu_{\mathbf{a}} = \{\mathbf{a}^*\}$ ;*
2. *and from any arm  $\mathbf{a} \neq \mathbf{a}^*$ , there exists a path  $(\mathbf{a}^0, \mathbf{a}^1, \dots, \mathbf{a}^n)$  in  $G$  such that  $\mathbf{a}^0 = \mathbf{a}$ ,  $\mathbf{a}^n = \mathbf{a}^*$ , and  $\forall i \in [n], \mu_{\mathbf{a}^i} > \mu_{\mathbf{a}^{i-1}}$ .*

Note that the second property of unimodal sets of expected rewards is equivalent to the property stating that from any sub-optimal arm  $\mathbf{a}$ , there exists an arm  $\mathbf{a}' \in \mathcal{N}_G(\mathbf{a})$  such that  $\mu_{\mathbf{a}'} > \mu_{\mathbf{a}}$ , where  $\mathcal{N}_G(\mathbf{a})$  is the neighborhood of  $\mathbf{a}$  in  $G$ .

Let's assume that there exists a unique recommendation  $\mathbf{a}^*$  with maximum expected reward, and denote  $\mathcal{F} = (\pi_{\mathbf{a}})_{\mathbf{a} \in \mathcal{P}_K^L}$  a set of rankings of positions such that for any recommendation  $\mathbf{a}$ ,  $\pi_{\mathbf{a}} \in \Pi_\rho(\mathbf{a})$ . Then, by denoting  $G_{\mathcal{F}} = (V, E_{\mathcal{F}})$  the directed graph with vertices  $V = \mathcal{P}_K^L$  and edges

$$E_{\mathcal{F}} \stackrel{\text{def}}{=} \left\{ (\mathbf{a}, \mathbf{a} \circ (\pi_{\mathbf{a}k}, \pi_{\mathbf{a}(k+1)})) : k \in [K-1] \right\} \cup \left\{ (\mathbf{a}, \mathbf{a}[\pi_{\mathbf{a}K} := i]) : i \in [L] \setminus \mathbf{a}([K]) \right\},$$

$(\mu_{\mathbf{a}})_{\mathbf{a} \in \mathcal{P}_K^L}$  is unimodal<sup>1</sup> with respect to  $G_{\mathcal{F}}$ . Note that this graph is unknown from the algorithm as it builds upon the unknown mapping  $\mathcal{F}$ . However, this mapping may be learned online, paving the way to an OSUB-like algorithm to explore the space of recommendations.

## 4. GRAB Algorithm

Our algorithm, GRAB, takes inspiration from the unimodal bandit algorithm OSUB (Combes & Proutière, 2014) by selecting at each iteration  $t$  an arm  $\mathbf{a}(t)$  in the neighborhood of the current best arm (a.k.a. the leader). While in OSUB the neighborhood is known beforehand, here we learn it online. GRAB is described in Algorithm 1. This algorithm uses the following notations:

<sup>1</sup>While the definition of unimodality in (Combes & Proutière, 2014) involves an **undirected** graph, OSUB only requires a **directed** graph and the existence of a strictly increasing path from any sub-optimal arm to the optimal one.

At each iteration  $t$ , we denote

$$\hat{\rho}_{i,k}(t) \stackrel{def}{=} \frac{1}{T_{i,k}(t)} \sum_{s=1}^{t-1} \mathbb{1}\{a_k(s) = i\} c_k(s)$$

the average number of clicks obtained at position  $k$  when displaying item  $i$  at this position, where

$$T_{i,k}(t) \stackrel{def}{=} \sum_{s=1}^{t-1} \mathbb{1}\{a_k(s) = i\}$$

is the number of time item  $i$  has been displayed at position  $k$ ;  $\hat{\rho}_{i,k}(t) \stackrel{def}{=} 0$  when  $T_{i,k}(t) = 0$ .

We also denote  $\tilde{\mathbf{a}}(t)$  the *leader*, meaning the recommendation with the best *pseudo average reward*  $\bar{\mu}_{\mathbf{a}}(t) \stackrel{def}{=} \sum_{k=1}^K \hat{\rho}_{a_k,k}(t)$ , and we note

$$\tilde{T}_{\mathbf{a}}(t) \stackrel{def}{=} \sum_{s=1}^{t-1} \mathbb{1}\{\tilde{\mathbf{a}}(s) = \mathbf{a}\}$$

the number of times the leader is  $\mathbf{a}$  for iterations 1 to  $t-1$ .

Finally, the statistics  $\hat{\rho}_{i,k}(t)$  are paired with their respective *indices*

$$b_{i,k}(t) \stackrel{def}{=} f\left(\hat{\rho}_{i,k}(t), T_{i,k}(t), \tilde{T}_{\tilde{\mathbf{a}}(t)}(t) + 1\right),$$

where  $f(\hat{\rho}, s, t)$  stands for

$$\sup\{p \in [\hat{\rho}, 1] : s \times \text{kl}(\hat{\rho}, p) \leq \log(t) + 3 \log(\log(t))\},$$

with

$$\text{kl}(p, q) \stackrel{def}{=} p \log\left(\frac{p}{q}\right) + (1-p) \log\left(\frac{1-p}{1-q}\right)$$

the *Kullback-Leibler divergence* from a Bernoulli distribution of mean  $p$  to a Bernoulli distribution of mean  $q$ ;  $f(\hat{\rho}, s, t) \stackrel{def}{=} 1$  when  $\hat{\rho} = 1$ ,  $s = 0$ , or  $t = 0$ .

At each iteration  $t$ , GRAB first identifies the leader  $\tilde{\mathbf{a}}(t)$ , and then recommends either  $\tilde{\mathbf{a}}(t)$  every  $L$ -th iteration, or the best permutation in the inferred neighborhood, given the sum of indices  $\sum_{k=1}^K b_{a_k,k}(t)$  (see Figure 1 for an example of a neighborhood). Each time an argmax is computed, the ties are randomly broken.

To finish the presentation of GRAB, let us now discuss its initialisation and its time-complexity.

**Remark 1** (Initialisation). The initialisation of the algorithm is handled through the default value of indices  $b_{i,k}$ : 1. This value ensures that any permutation is recommended at least once, as soon as it belongs to the neighborhood of an arm which is often the leader. If a permutation is not in such neighborhood, the theoretical analysis in Section 5 proves that this permutation is sub-optimal, hence it does not matter whether this permutation is explored at least once or not.

**Algorithm 1** GRAB: parametric Graph for unimodal Ranking Bandit

**Input:** number of items  $L$ , number of positions  $K$

1: **for**  $t = 1, 2, \dots$  **do**

2:  $\tilde{\mathbf{a}}(t) \leftarrow \operatorname{argmax}_{\mathbf{a} \in \mathcal{P}_K^L} \sum_{k=1}^K \hat{\rho}_{a_k,k}(t)$

3: **find**  $\tilde{\pi}(t)$  **s.t.**  $\hat{\rho}_{\tilde{a}_{\tilde{\pi}_1(t)}(t), \tilde{\pi}_1(t)}(t) \geq \hat{\rho}_{\tilde{a}_{\tilde{\pi}_2(t)}(t), \tilde{\pi}_2(t)}(t) \geq \dots \geq \hat{\rho}_{\tilde{a}_{\tilde{\pi}_K(t)}(t), \tilde{\pi}_K(t)}(t)$   $\geq$

4: **recommend**

$$\mathbf{a}(t) = \begin{cases} \tilde{\mathbf{a}}(t) & , \text{ if } \frac{\tilde{T}_{\tilde{\mathbf{a}}(t)}(t)}{L} \in \mathbb{N}, \\ \operatorname{argmax}_{\substack{\mathbf{a} \in \{\tilde{\mathbf{a}}(t)\} \\ \cup \mathcal{N}_{\tilde{\pi}}(\tilde{\mathbf{a}}(t))}} \sum_{k=1}^K b_{a_k,k}(t) & , \text{ otherwise} \end{cases}$$

where  $\mathcal{N}_{\tilde{\pi}}(\mathbf{a}) = \{\mathbf{a} \circ (\pi_k, \pi_{k+1}) : k \in [K-1]\} \cup \{\mathbf{a}[\pi_K := i] : i \in [L] \setminus \mathbf{a}([K])\}$

5: **observe** the clicks vector  $\mathbf{c}(t)$

6: **end for**

**Remark 2** (Algorithmic Complexity). Even though the two optimization steps might seem costly, at each iteration the choice of a recommendation is done in a polynomial time w.r.t.  $L$  and  $K$ : first, the maximization at Line 2 is a *linear sum assignment problem* which is solvable in  $\mathcal{O}(K^2(L + \log K))$  time (Ramshaw & Tarjan, 2012); it is not required to scan the  $L!/(L-K)!$  permutations of  $K$  distinct items among  $L$ . Secondly, the maximization at Line 4 is over a set of  $L-1$  recommendations and is equivalent to the maximization of

$$B_{\mathbf{a}}(t) = \sum_{k=1}^K b_{a_k,k}(t) - \sum_{k=1}^K b_{\tilde{a}_k(t),k}(t)$$

which reduces to the sum of up to four  $b_{a_k,k}(t)$  terms as we are looking at recommendations  $\mathbf{a}$  in the neighborhood of the leader. Specifically, either

- $\mathbf{a} = \tilde{\mathbf{a}}(t)$  and  $B_{\mathbf{a}}(t) = 0$ ,
- or  $\mathbf{a} = \tilde{\mathbf{a}}(t) \circ (k, k')$  and  $B_{\mathbf{a}}(t) = b_{\tilde{a}_{k'}(t),k}(t) + b_{\tilde{a}_k(t),k'}(t) - b_{\tilde{a}_k(t),k}(t) - b_{\tilde{a}_{k'}(t),k'}(t)$ ,
- or  $\mathbf{a} = \tilde{\mathbf{a}}(t)[k := i]$  and  $B_{\mathbf{a}}(t) = b_{i,k}(t) - b_{\tilde{a}_k(t),k}(t)$ .

Hence, this maximization requires  $\mathcal{O}(L)$  computation time. Overall, the computation time per iteration is a  $\mathcal{O}(K^2(L + \log K))$ .

## 5. Theoretical Analysis

As already discussed in Section 2, the proof of the upper-bound on the regret of GRAB follows a similar path as the

proof of OSUB (Combes & Proutière, 2014): (1) apply standard bandit analysis to control the regret under the condition that the leader  $\tilde{\mathbf{a}}(t)$  is the best arm  $\mathbf{a}^*$ , and (2) upper-bound the expected number of iterations such that  $\tilde{\mathbf{a}}(t) \neq \mathbf{a}^*$  by a  $\mathcal{O}(\log \log T)$ . The inference of the rankings on positions adds up a third step (3) upper-bounding the expected number of iterations such that  $\tilde{\pi}(t) \notin \Pi_{\rho}(\tilde{\mathbf{a}}(t))$ .

The first step differs from (Combes & Proutière, 2014), as we have to account for the semi-bandit feedback. We note that when the leader is the best arm, GRAB behaves as a Kullback-Leibler variation of CombUCB1 (Kveton et al., 2015b) that we call KL-CombUCB in the following (see the appendix for a complete definition of KL-CombUCB). We derive an upper-bound specific to KL-CombUCB which accounts for the fact that the maximization at Line 4 of Algorithm 1 can be reduced to the maximization over sums of at most 4 terms (see Remark 2). In the context of GRAB, this new result, expressed by Theorem 2, reduces the regret-bound by a factor  $K$  w.r.t. the standard upper-bound for CombUCB1.

The second part of the analysis is based on the fact that with high probability  $\bar{\mu}_{\mathbf{a}}(t) > \bar{\mu}_{\mathbf{a}'}(t)$  if  $\mu_{\mathbf{a}} > \mu_{\mathbf{a}'}$ , which derives from the control of the deviation of each  $\hat{\rho}_{i,k}(t)$ . Here lies the second main difference with Combes & Proutière's analysis: we control the deviation of each individual  $\hat{\rho}_{i,k}(t)$  while they control the deviation of  $\hat{\mu}_{\mathbf{a}}(t) \stackrel{def}{=} (\sum_{s=1}^{t-1} \mathbb{1}\{\mathbf{a}(s) = \mathbf{a}\})^{-1} \sum_{s=1}^{t-1} \mathbb{1}\{\mathbf{a}(s) = \mathbf{a}\} r(s)$ . Again, the analysis benefits from the small number of differences between recommendations in the neighborhood of the leader. Moreover, the analysis handles the fact that the neighborhoods may change from an iteration to another, while the neighborhoods are constant in Combes & Proutière's analysis. The corresponding result is expressed, in the following, by Lemma 2.

Finally, the number of iterations at which the inferred ranking on the positions is inappropriate is controlled by Lemma 3. The proof of this lemma is eased by the fact that the number of times the leader is played is at least proportional to the number of times it is the leader.

We now propose and prove the main theorem that upper-bounds the regret of GRAB. Its proof is given after the presentation of all the necessary theorems and lemmas.

**Theorem 1** (Upper-Bound on the Regret of GRAB). *Let  $(L, K, (\rho_{i,k})_{(i,k) \in [L] \times [K]})$  be an online learning to rank problem satisfying Assumption 1 and such that there exists a unique recommendation  $\mathbf{a}^*$  with maximum expected reward.*

When facing this problem, GRAB fulfills:

$$\forall \mathbf{a} \in \mathcal{N}_{\pi^*}(\mathbf{a}^*), \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\tilde{\mathbf{a}}(t) = \mathbf{a}^*, \tilde{\pi}(t) = \pi^*, \mathbf{a}(t) = \mathbf{a}\} \right] \leq \frac{8}{\Delta_{\mathbf{a}}^2} \log T + \mathcal{O}(\log \log T), \quad (4)$$

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\tilde{\mathbf{a}}(t) \neq \mathbf{a}^*\} \right] = \mathcal{O}(\log \log T), \quad (5)$$

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\tilde{\pi}(t) \notin \Pi_{\rho}(\tilde{\mathbf{a}}(t))\} \right] = \mathcal{O}(1), \quad (6)$$

and hence

$$\begin{aligned} R(T) &\leq \sum_{\mathbf{a} \in \mathcal{N}_{\pi^*}(\mathbf{a}^*)} \frac{8}{\Delta_{\mathbf{a}}^2} \log T + \mathcal{O}(\log \log T) \quad (7) \\ &= \mathcal{O} \left( \frac{L}{\Delta_{\min}} \log(T) \right), \end{aligned}$$

where  $\pi^*$  is the unique ranking of positions in  $\Pi_{\rho}(\mathbf{a}^*)$ ,  $\Delta_{\mathbf{a}} \stackrel{def}{=} \mu^* - \mu_{\mathbf{a}}$ , and  $\Delta_{\min} \stackrel{def}{=} \min_{\mathbf{a} \in \mathcal{N}_{\pi^*}(\mathbf{a}^*)} \Delta_{\mathbf{a}}$ .

The first upper-bound (Equation (4)) deals with the expected number of iterations at which GRAB recommends a sub-optimal permutation while the leader is the best permutation. It derives from Theorem 2 hereafter, which detailed proof is in the appendix.

**Theorem 2** (New Upper-Bound on the Regret of KL-CombUCB). *We consider a combinatorial semi-bandit setting. Let  $E$  be a set of elements and  $\mathcal{A} \subseteq \{0, 1\}^E$  be a set of arms, where each arm  $\mathbf{a}$  is a subset of  $E$ . Let's assume that the reward when drawing the arm  $\mathbf{a} \in \mathcal{A}$  is  $\sum_{e \in \mathbf{a}} c_e$ , where for each element  $e \in E$ ,  $c_e$  is an independent draw of a Bernoulli distribution of mean  $\rho_e \in [0, 1]$ . Therefore, the expected reward when drawing the arm  $\mathbf{a} \in \mathcal{A}$  is  $\mu_{\mathbf{a}} = \sum_{e \in \mathbf{a}} \rho_e$ .*

When facing this bandit setting, KL-CombUCB (CombUCB1 equipped with Kullback-Leibler indices, see the appendix) fulfills

$$\begin{aligned} \forall \mathbf{a} \in \mathcal{A} \text{ s.t. } \mu_{\mathbf{a}} \neq \mu^*, \\ \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\mathbf{a}(t) = \mathbf{a}\} \right] \leq \frac{2K_{\mathbf{a}}^2}{\Delta_{\mathbf{a}}^2} \log T + \mathcal{O}(\log \log T), \end{aligned}$$

and hence

$$\begin{aligned} R(T) &\leq \sum_{\mathbf{a} \in \mathcal{A}: \mu_{\mathbf{a}} \neq \mu^*} \frac{2K_{\mathbf{a}}^2}{\Delta_{\mathbf{a}}^2} \log T + \mathcal{O}(\log \log T) \\ &= \mathcal{O} \left( \frac{|\mathcal{A}| K_{\max}^2}{\Delta_{\min}} \log(T) \right), \end{aligned}$$

where  $\mu^* \stackrel{\text{def}}{=} \max_{\mathbf{a} \in \mathcal{A}} \mu_{\mathbf{a}}$ ,  $\Delta_{\mathbf{a}} \stackrel{\text{def}}{=} \mu^* - \mu_{\mathbf{a}}$ ,  $\Delta_{\min} \stackrel{\text{def}}{=} \min_{\mathbf{a} \in \mathcal{A}: \Delta_{\mathbf{a}} > 0} \Delta_{\mathbf{a}}$ ,  $K_{\mathbf{a}} \stackrel{\text{def}}{=} \min_{\mathbf{a}^* \in \mathcal{A}: \mu_{\mathbf{a}^*} = \mu^*} |\mathbf{a} \setminus \mathbf{a}^*|$  is the smallest number of elements to remove from  $\mathbf{a}$  to get an optimal arm, and  $K_{\max} \stackrel{\text{def}}{=} \max_{\mathbf{a} \in \mathcal{A}: \mu_{\mathbf{a}} \neq \mu^*} K_{\mathbf{a}}$ .

Secondly, the expected number of iterations at which the leader is not the optimal arm (Equation (5)) is controlled by Lemma 2, which detailed proof is in the appendix.

**Lemma 2** (Upper-Bound on the Number of Iterations of GRAB for which  $\tilde{\mathbf{a}}(t) \neq \mathbf{a}^*$ ). *Under the hypotheses of Theorem 1 and using its notations,*

$$\forall \tilde{\mathbf{a}} \in \mathcal{P}_K^L \setminus \{\mathbf{a}^*\}, \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}\} \right] = \mathcal{O}(\log \log T).$$

Finally, the number of iterations at which the inferred ranking on the positions is inappropriate (Equation (6)) is controlled by Lemma 3, which detailed proof is in the appendix.

**Lemma 3** (Upper-Bound on the Number of Iterations of GRAB for which  $\boldsymbol{\pi}(t) \notin \Pi_{\boldsymbol{\rho}}(\tilde{\mathbf{a}})$ ). *Under the hypotheses of Theorem 1 and using its notations,*

$$\forall \tilde{\mathbf{a}} \in \mathcal{P}_K^L, \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{\boldsymbol{\pi}}(t) \notin \Pi_{\boldsymbol{\rho}}(\tilde{\mathbf{a}})\} \right] = \mathcal{O}(1).$$

We assemble these results to get the proof of Theorem 1.

*Proof of Theorem 1.* First note that, since there is a unique optimal permutation, there is a unique appropriate ranking  $\boldsymbol{\pi}^*$  of positions w.r.t.  $\mathbf{a}^*$ :  $\Pi_{\boldsymbol{\rho}}(\mathbf{a}^*) = \{\boldsymbol{\pi}^*\}$ . Then, the proof is based on the following decomposition of the set  $[T]$  of iterations:

$$\begin{aligned} [T] &= \bigcup_{\substack{\mathbf{a} \in \{\mathbf{a}^*\} \\ \cup \mathcal{N}_{\boldsymbol{\pi}^*}(\mathbf{a}^*)}} \{t \in [T] : \tilde{\mathbf{a}}(t) = \mathbf{a}^*, \tilde{\boldsymbol{\pi}}(t) = \boldsymbol{\pi}^*, \mathbf{a}(t) = \mathbf{a}\} \\ &\cup \{t \in [T] : \tilde{\mathbf{a}}(t) \neq \mathbf{a}^*\} \cup \{t \in [T] : \tilde{\boldsymbol{\pi}}(t) \notin \Pi_{\boldsymbol{\rho}}(\tilde{\mathbf{a}}(t))\}. \end{aligned}$$

As for any recommendation  $\mathbf{a}$ ,  $\Delta_{\mathbf{a}} \leq K$ , this decomposition leads to the inequality  $R(T) \leq \sum_{\mathbf{a} \in \mathcal{N}_{\boldsymbol{\pi}^*}(\mathbf{a}^*)} \Delta_{\mathbf{a}} A_{\mathbf{a}} + KB + KC$ , with

$$\begin{aligned} A_{\mathbf{a}} &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\tilde{\mathbf{a}}(t) = \mathbf{a}^*, \tilde{\boldsymbol{\pi}}(t) = \boldsymbol{\pi}^*, \mathbf{a}(t) = \mathbf{a}\} \right], \\ B &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\tilde{\mathbf{a}}(t) \neq \mathbf{a}^*\} \right], \\ C &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\tilde{\boldsymbol{\pi}}(t) \notin \Pi_{\boldsymbol{\rho}}(\tilde{\mathbf{a}}(t))\} \right]. \end{aligned}$$

The term  $A_{\mathbf{a}}$  is smaller than the expected number of times the arm  $\mathbf{a}$  is chosen by KL-CombUCB when it plays on the set of arms  $\{\mathbf{a}^*\} \cup \mathcal{N}_{\boldsymbol{\pi}^*}(\mathbf{a}^*)$ . As any of these arms differs with  $\mathbf{a}^*$  at at most two positions, Theorem 2 upper-bounds  $A_{\mathbf{a}}$  by

$$\frac{8}{\Delta_{\mathbf{a}}} \log T + \mathcal{O}(\log \log T)$$

and hence  $\sum_{\mathbf{a} \in \mathcal{N}_{\boldsymbol{\pi}^*}(\mathbf{a}^*)} \Delta_{\mathbf{a}} A_{\mathbf{a}} = \mathcal{O}(L/\Delta_{\min} \log T)$  as  $|\mathcal{N}_{\boldsymbol{\pi}^*}(\mathbf{a}^*)| = L - 1$ .

Note that Theorem 5 of (Kveton et al., 2015b), upper-bounding the regret of CombUCB1, leads to a  $\mathcal{O}(LK/\Delta \log T)$  bound<sup>2</sup> for  $\sum_{\mathbf{a} \in \mathcal{N}_{\boldsymbol{\pi}^*}(\mathbf{a}^*)} \Delta_{\mathbf{a}} A_{\mathbf{a}}$ , which we reduce by a factor  $K$  by using Theorem 2.

From Lemma 2, the term  $B$  is upper-bounded by

$$B = \sum_{\tilde{\mathbf{a}} \in \mathcal{P}_K^L \setminus \{\mathbf{a}^*\}} \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}\} \right] = \mathcal{O}(\log \log T),$$

and we upper-bound the term  $C$  with Lemma 3:

$$C = \sum_{\tilde{\mathbf{a}} \in \mathcal{P}_K^L} \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{\boldsymbol{\pi}}(t) \notin \Pi_{\boldsymbol{\rho}}(\tilde{\mathbf{a}})\} \right] = \mathcal{O}(1).$$

Finally, the regret of GRAB is upper-bounded by summing these three terms, which concludes the proof.  $\square$

## 5.1. Discussion

Assuming  $\theta_1 \geq \dots \geq \theta_L$  and  $\kappa_1 \geq \dots \geq \kappa_K$ , the detailed formula for the regret upper-bound (7) is  $\sum_{k=1}^{K-1} \frac{8 \log T}{(\kappa_k - \kappa_{k+1})(\theta_k - \theta_{k+1})} + \sum_{k=K+1}^L \frac{8 \log T}{\kappa_K(\theta_K - \theta_k)}$ , where the first sum corresponds to the set of neighbors of  $\mathbf{a}^*$  which recommend the same items as  $\mathbf{a}^*$ , and the second sum relates to the set of neighbors of  $\mathbf{a}^*$  which replace the ‘last’ item in  $\mathbf{a}^*$ . Hence, the number of displayed items does not impact the total number of terms, but the gaps  $\Delta_{\mathbf{a}}$ .

Note also that GRAB is, by design, robust to miss-specifications. Typically, GRAB would properly handle a matrix  $\boldsymbol{\rho} = \boldsymbol{\theta}^T \boldsymbol{\kappa} + \mathcal{E}$ , if  $\max_{i,j} |\mathcal{E}_{i,j}|$  is smaller than half of the minimum gap between two entries of the matrix  $\boldsymbol{\theta}^T \boldsymbol{\kappa}$ .

However, if there is a set of optimal recommendations  $\mathcal{A}^*$  (instead of a unique one), after convergence, the leader will be picked in that set at each iteration. So the neighborhood of each optimal recommendation will be explored, and we will get a regret bound in  $\mathcal{O}(|\mathcal{A}^*|L)$ . This behavior

<sup>2</sup>In this setting, the ground set is  $E \stackrel{\text{def}}{=} \bigcup_{k \in [K]} \{(a_{\max(k-1,1)}, k), (a_k, k), (a_{\min(k+1, K)}, k)\} \cup \bigcup_{k \in [L] \setminus [K]} \{(a_k, K)\}$  and is of size  $L + 2K - 2$ , and any arm is composed of exactly  $K$  elements in  $E$ .



questions the applicability of unimodality to the *Cascading Model* (CM), as with this model there is at least  $K!$  optimal recommendations. Moreover, while Assumption 1 is valid for CM and the *Dependent Click Model* (DCM), our setting also assumes the existence of the matrix  $\rho$ , which is false for CM and DCM: in both settings the probability of clicking on item  $i$  in position  $\ell$  depends on other displayed items.

## 6. Experiments

In this section, we compare GRAB to PMED (Komiyama et al., 2017), to TopRank (Lattimore et al., 2018), to PB-MHB (Gauthier et al., 2021), to  $\epsilon_n$ -Greedy, to Static Graph for unimodal Ranking Bandit (S-GRAB), a simplified version of GRAB, and to KL-CombUCB, an adaptation of CombUCB1 (Kveton et al., 2015b) (see the appendix for details regarding S-GRAB and KL-CombUCB). The experiments are conducted on the Yandex dataset (Yandex, 2013) and on purely simulated data. We use the cumulative regret to evaluate the performance of each algorithm, where the cumulative regret is averaged over 20 independent runs of  $T = 10^7$  iterations each. Code and data for replicating our experiments are available at [https://github.com/gaudel/ranking\\_bandits](https://github.com/gaudel/ranking_bandits).

### 6.1. Experimental Setting

We use two types of datasets: a *simulated* one for which we set the values for  $\kappa$  and  $\theta$  and a *real* one, where parameters are inferred from real life logs of Yandex search engine (Yandex, 2013). Let’s remind that  $\theta_i$  is the probability for the user to click on item  $i$  when it observes this item, and  $\kappa_k$  is the probability for the user to observe position  $k$ .

Simulated data allow us to test GRAB in extreme situations. We consider  $L = 10$  items,  $K = 5$  positions, and  $\kappa = [1, 0.75, 0.6, 0.3, 0.1]$ . The range of values for  $\theta$  is either close to zero ( $\theta^- = [10^{-3}, 5 \cdot 10^{-4}, 10^{-4}, 5 \cdot 10^{-5}, 10^{-5}, 10^{-6}, \dots, 10^{-6}]$ ), or close to one ( $\theta^+ = [0.99, 0.95, 0.9, 0.85, 0.8, 0.75, \dots, 0.75]$ ).

Real data contain the logs of actions toward the Yandex search engine: 65 million search queries and 167 million hits (clicks). Common use of this database in the bandit setting consists first in extracting from these logs the parameters of the chosen real model, and then in simulating users’ interactions given these parameters (Lattimore et al., 2018). We use Pyclick library (Chuklin et al., 2015) to infer the PBM parameters of each query with the *expectation maximization* algorithm. This leads to  $\theta_i$  values ranging from 0.070 to 0.936, depending on the query. Similarly to (Lattimore et al., 2018), we look at the results averaged on the 10 most frequent queries, while displaying  $K = 5$  items among the  $L = 10$  most attractive ones.

Among our opponents, TopRank and PMED require de-

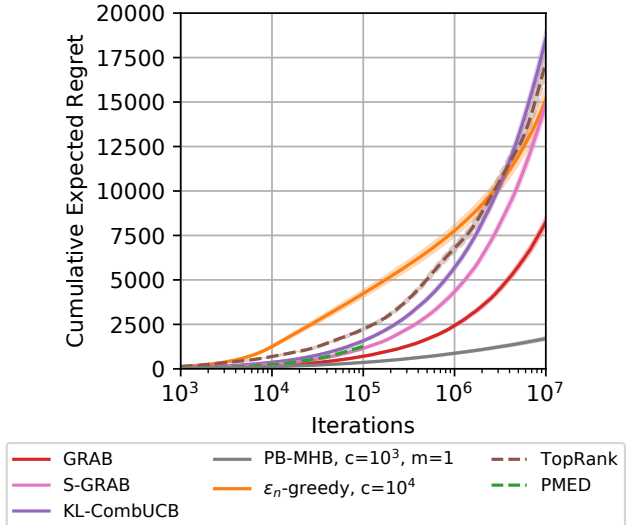


Figure 2. Cumulative regret w.r.t. iterations on Yandex dataset. The plotted curves correspond to the average over 200 independent sequences of recommendations (20 sequences per query). The shaded area depicts the standard error of our regret estimates.

creasing values of  $\kappa$  which may not be fulfilled by PBM. We pre-order them to fulfill these algorithms’ requirements. Otherwise,  $\kappa$  is shuffled at the beginning of each sequence of recommendations. We also carefully tune the exploration hyper-parameter  $c$  of  $\epsilon_n$ -greedy taking values ranging exponentially from  $10^0$  to  $10^6$ . For PB-MHB, we use the hyper-parameters recommended in (Gauthier et al., 2021).

### 6.2. Results

Figure 2 shows the results for the algorithms on Yandex and Figure 3 on the simulated data. We measure the performance of each algorithm according to the cumulative regret (see Equation 1). It is the sum, over  $T$  consecutive recommendations, of the difference between the expected reward of the best answer and of the answer of a given recommender system. The best algorithm is the one with the lowest regret. We average the results of each algorithm over 20 independent sequences of recommendations per query or simulated setting. Although PMED theoretically yields an asymptotically optimal regret, we stop it at iteration  $t = 10^5$  due to its heavy computation-time.

**Ablation Study** The two main ingredients of GRAB are the use of a graph to explore the set of recommendations, and the online inference of this graph. Without these ingredients, GRAB boils down to KL-CombUCB which recommends at each iteration the best permutation given the sum of indices  $b_{i,k}$  and has a  $\mathcal{O}(LK^2/\Delta \log T)$  regret. With only the first ingredient (namely a static graph of degree

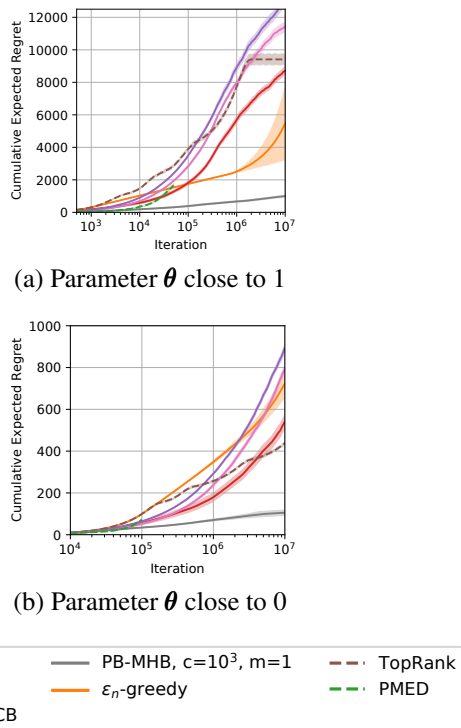


Figure 3. Cumulative regret w.r.t. iterations on simulated data. The plotted curves correspond to the average over 20 independent sequences of recommendations. The shaded area depicts the standard error of our regret estimates. For  $\epsilon_n$ -Greedy,  $c$  is set to  $10^5$  when  $\theta$  is close to 0, and to  $10^3$  when  $\theta$  is close to 1.

$\Theta(LK)$ ), we get S-GRAB which regret is upper-bounded by  $\mathcal{O}(LK/\Delta \log T)$ , while GRAB’s regret is upper-bounded by  $\mathcal{O}(L/\Delta \log T)$  thanks to a set of graphs of degree  $L - 1$ .

We want to assert the empirical impact of these ingredients. On Figures 2 and 3, we see that GRAB has a better regret than S-GRAB and KL-CombUCB in every settings. This confirms that the proposed graphs are relevant to explore the set of recommendations, and that GRAB quickly infer the appropriate graph in the family of potential ones.

**Results Analysis** Figure 2 compares the empirical regret of all algorithms on Yandex dataset. GRAB is the second best with a regret at  $T = 10^7$  about two times smaller than the rest of the algorithms. Only PB-MHB yields a smaller regret, but PB-MHB is more than ten times slower to deliver a recommendation than GRAB and it does not have any theoretical guarantees.

Figure 3 shows our results on purely simulated data illustrating extreme settings even though these settings are less realistic. In both settings, GRAB is in the top-3 algorithms. PB-MHB is still the algorithm yielding the best regret. However, while TopRank provides better or similar result as

Table 2. Average computation time for sequences of  $10^7$  recommendations vs. all queries of Yandex dataset

ALGORITHM	(HOUR/MIN)	TRIAL (MS)
GRAB	2H24	0.9
S-GRAB	9H56	3.6
$\epsilon_n$ -GREEDY $c = 10^4$	1H13	0.4
PB-MHB $c = 10^3, m = 1$	44H50	16
KL-COMBUCB	2H03	0.7
PMED	474H13*	170
TOPRANK	9H29	3

\* EXTRAPOLATION FROM  $10^5$  RECOMMENDATIONS.

GRAB at iteration  $10^7$ , its regret is higher than the one of GRAB up to iteration  $t = 4 \times 10^6$ . TopRank only catches-up GRAB at the end of the sequences of recommendations. We note that in the setting close to 1, TopRank manages to find the perfect order after  $10^6$  iterations. In this setting too,  $\epsilon_n$ -greedy has better performance during the  $10^6$  first iterations, but suffers from its greedy behaviour during the last steps with a large variance.

**Computation Time** As shown in Table 2, the fastest algorithm is  $\epsilon_n$ -greedy. KL-CombUCB and GRAB are two times slower. The exploration of S-GRAB multiplies its computation time by 4 compared to GRAB. TopRank is about three times slower than GRAB, and PB-MHB, despite its good regret is one of the slowest algorithm with PMED.

## 7. Conclusion

Our work targets the full PBM setting, which aims at recommending a ranking of  $K$  items among  $L$  and display them, without prior knowledge on the attractiveness of the positions. We learn online both the user preferences and their gaze habits. To solve this problem, we define a graph parametrized by rankings of positions, and we extend the unimodal bandit setting to this family of graphs. We also design GRAB, an algorithm that learns online the proper parametrization of the graph, and we prove a regret upper-bound in  $\mathcal{O}(L/\Delta \log T)$  for this algorithm which reduces by a factor  $K^2$  (respectively  $K$ ) the bound which would be obtained without the unimodal setting (resp. with the standard unimodal setting). On real and simulated data, GRAB quickly delivers good recommendations.

The extension of the unimodal setting is a promising tool which may benefit to recommendations to users with a more general behavior, or to other combinatorial semi-bandit scenarios. The integration of unimodal bandit algorithms working on parametric spaces (Combes et al., 2020) may also pave the way to efficient contextual recommendation systems handling larger sets of items and positions.

## References

- Chen, W., Wang, Y., and Yuan, Y. Combinatorial multi-armed bandit: General framework and applications. In *proc. of the 30th Int. Conf. on Machine Learning, ICML'13*, 2013.
- Chuklin, A., Markov, I., and de Rijke, M. *Click Models for Web Search*. Morgan & Claypool Publishers, 2015.
- Combes, R. and Proutière, A. Unimodal bandits: Regret lower bounds and optimal algorithms. In *proc. of the 31st Int. Conf. on Machine Learning, ICML'14*, 2014.
- Combes, R., Magureanu, S., Proutière, A., and Laroche, C. Learning to rank: Regret lower bounds and efficient algorithms. In *proc. of the ACM SIGMETRICS Int. Conf. on Measurement and Modeling of Computer Systems*, 2015.
- Combes, R., Proutière, A., and Fauquette, A. Unimodal bandits with continuous arms: Order-optimal regret without smoothness. *Proc. ACM Meas. Anal. Comput. Syst.*, 4(1), May 2020.
- Craswell, N., Zoeter, O., Taylor, M., and Ramsey, B. An experimental comparison of click position-bias models. In *proc. of the Int. Conf. on Web Search and Data Mining, WSDM '08*, 2008.
- Gai, Y., Krishnamachari, B., and Jain, R. Combinatorial network optimization with unknown variables: Multi-armed bandits with linear rewards and individual observations. *IEEE/ACM Trans. Netw.*, 20(5):1466–1478, October 2012.
- Garivier, A. and Cappé, O. The kl-ucb algorithm for bounded stochastic bandits and beyond. In *proc. of the 24th Annual Conf. on Learning Theory, COLT'11*, 2011.
- Gauthier, C.-S., Gaudel, R., Fromont, E., and Lompo, B. A. Position-based multiple-play bandits with thompson sampling. In *proc. of the 19th Symposium on Intelligent Data Analysis, IDA'21*, 2021.
- Komiyama, J., Honda, J., and Nakagawa, H. Optimal regret analysis of thompson sampling in stochastic multi-armed bandit problem with multiple plays. In *proc. of the 32nd Int. Conf. on Machine Learning, ICML'15*, 2015.
- Komiyama, J., Honda, J., and Takeda, A. Position-based multiple-play bandit problem with unknown position bias. In *Advances in Neural Information Processing Systems 30, NIPS'17*, 2017.
- Kveton, B., Szepesvári, C., Wen, Z., and Ashkan, A. Cascading bandits: Learning to rank in the cascade model. In *proc. of the 32nd Int. Conf. on Machine Learning, ICML'15*, 2015a.
- Kveton, B., Wen, Z., Ashkan, A., and Szepesvari, C. Tight Regret Bounds for Stochastic Combinatorial Semi-Bandits. In *proc. of the 18th Int. Conf. on Artificial Intelligence and Statistics, AISTATS'15*, 2015b.
- Lagrée, P., Vernade, C., and Cappé, O. Multiple-play bandits in the position-based model. In *Advances in Neural Information Processing Systems 30, NIPS'16*, 2016.
- Lattimore, T., Kveton, B., Li, S., and Szepesvari, C. Toprank: A practical algorithm for online stochastic ranking. In *Advances in Neural Information Processing Systems 31, NIPS'18*, 2018.
- Li, C., Kveton, B., Lattimore, T., Markov, I., de Rijke, M., Szepesvári, C., and Zoghi, M. Bubblerank: Safe online learning to re-rank via implicit click feedback. In *proc. of the 35th Uncertainty in Artificial Intelligence Conference, UAI'19*, 2019.
- Radlinski, F., Kleinberg, R., and Thorsten, J. Learning diverse rankings with multi-armed bandits. In *proc. of the 25th Int. Conf. on Machine Learning, ICML'08*, 2008.
- Ramshaw, L. and Tarjan, R. E. On minimum-cost assignments in unbalanced bipartite graphs. Technical report, HP research labs, 2012.
- Richardson, M., Dominowska, E., and Ragno, R. Predicting clicks: Estimating the click-through rate for new ads. In *proc. of the 16th International World Wide Web Conference, WWW '07*, 2007.
- Yandex. Yandex personalized web search challenge. 2013. URL <https://www.kaggle.com/c/yandex-personalized-web-search-challenge>.

The appendix is organized as follows. We first list most of the notations used in the paper in Appendix A. Lemma 1 is proved in Appendix B. In Appendix C, we recall a Lemma from (Combes & Proutière, 2014) used by our own Lemmas and Theorems, and then in Appendices D to F we respectively prove Theorem 2, Lemma 2, and Lemma 3. In Appendix G we define KL-CombUCB and discuss its regret and its relation to GRAB. Finally in Appendix H we introduce and discuss S-GRAB.

## A. Notations

The following table summarize the notations used through the paper and the appendix.

SYMBOL	MEANING
$T$	TIME HORIZON
$t$	ITERATION
$L$	NUMBER OF ITEMS
$i$	INDEX OF AN ITEM
$K$	NUMBER OF POSITIONS IN A RECOMMENDATION
$k$	INDEX OF A POSITION
$[n]$	SET OF INTEGERS $\{1, \dots, n\}$
$\mathcal{P}_K^L$	SET OF PERMUTATIONS OF $K$ DISTINCT ITEMS AMONG $L$
$\boldsymbol{\theta}$	VECTORS OF PROBABILITIES OF CLICK
$\theta_i$	PROBABILITY OF CLICK ON ITEM $i$
$\boldsymbol{\kappa}$	VECTORS OF PROBABILITIES OF VIEW
$\kappa_k$	PROBABILITY OF VIEW AT POSITION $k$
$\mathcal{A}$	SET OF BANDIT ARMS
$\mathbf{a}$	AN ARM IN $\mathcal{A}$
$\mathbf{a}(t)$	THE ARM CHOSEN AT ITERATION $t$
$\tilde{\mathbf{a}}(t)$	BEST ARM AT ITERATION $t$ GIVEN THE PREVIOUS CHOICES AND FEEDBACKS (CALLED LEADER)
$\mathbf{a}^*$	BEST ARM
$G$	GRAPH CARRYING A PARTIAL ORDER ON $\mathcal{A}$
$\gamma$	MAXIMUM DEGREE OF $G$
$\mathcal{N}_G(\tilde{\mathbf{a}}(t))$	NEIGHBORHOOD OF $\tilde{\mathbf{a}}(t)$ GIVEN $G$
$\rho_{i,k}$	PROBABILITY OF CLICK ON ITEM $i$ DISPLAYED AT POSITION $k$
$\mathbf{c}(t)$	CLICKS VECTOR AT ITERATION $t$
$r(t)$	REWARD COLLECTED AT ITERATION $t$ , $r(t) = \sum_{k=1}^K c_k(t)$
$\mu_{\mathbf{a}}$	EXPECTATION OF $r(t)$ WHILE RECOMMENDING $\mathbf{a}$ , $\mu_{\mathbf{a}} = \sum_{k=1}^K \rho_{a_k,k}$
$\mu^*$	HIGHEST EXPECTED REWARD, $\mu^* = \max_{\mathbf{a} \in \mathcal{P}_K^L} \mu_{\mathbf{a}}$
$\Delta_{\mathbf{a}}$	GAP BETWEEN $\mu_{\mathbf{a}}$ AND $\mu^*$
$\Delta_{min}$	MINIMAL VALUE FOR $\Delta_{\mathbf{a}}$
$\Delta$	GENERIC REWARD GAP BETWEEN ONE OF THE SUB-OPTIMAL ARMS AND ONE OF THE BEST ARMS
$R(T)$	CUMULATIVE (PSEUDO-)REGRET, $R(T) = T\mu^* - \mathbb{E} \left[ \sum_{t=1}^T \mu_{\mathbf{a}(t)} \right]$
$\Pi_{\boldsymbol{\rho}}(\mathbf{a})$	SET OF PERMUTATIONS IN $\mathcal{P}_K^K$ ORDERING THE POSITIONS S.T. $\rho_{a_{\pi_1}, \pi_1} \geq \rho_{a_{\pi_2}, \pi_2} \geq \dots \geq \rho_{a_{\pi_K}, \pi_K}$
$\boldsymbol{\pi}$	ELEMENT OF $\Pi_{\boldsymbol{\rho}}(\mathbf{a})$
$\tilde{\boldsymbol{\pi}}$	ESTIMATION OF $\boldsymbol{\pi}$
$\mathbf{a} \circ (\pi_k, \pi_{k+1})$	PERMUTATION SWAPPING ITEMS IN POSITIONS $\pi_k$ AND $\pi_{k+1}$
$\mathbf{a}[\pi_K := i]$	PERMUTATION LEAVING $\mathbf{a}$ THE SAME FOR ANY POSITION EXCEPT $\pi_K$ FOR WHICH $\mathbf{a}[\pi_K := i]_{\pi_K} = i$
$\mathcal{F}$	RANKINGS OF POSITIONS RESPECTING $\Pi_{\boldsymbol{\rho}}$ , $\mathcal{F} = (\boldsymbol{\pi}_{\mathbf{a}})_{\mathbf{a} \in \mathcal{P}_K^L}$ S.T. $\forall \mathbf{a} \in \mathcal{P}_K^L, \boldsymbol{\pi}_{\mathbf{a}} \in \Pi_{\boldsymbol{\rho}}(\mathbf{a})$
$T_{i,k}(t)$	NUMBER OF ITERATIONS S.T. ITEM $i$ HAS BEEN DISPLAYED AT POSITION $k$ , $T_{i,k}(t) = \sum_{s=1}^{t-1} \mathbb{1}\{a_k(s) = i\}$
$\tilde{T}_{\mathbf{a}}(t)$	NUMBER OF ITERATIONS S.T. THE LEADER WAS $\mathbf{a}$ , $\tilde{T}_{\mathbf{a}}(t) \stackrel{def}{=} \sum_{s=1}^{t-1} \mathbb{1}\{\tilde{\mathbf{a}}(s) = \mathbf{a}\}$
$T_{\mathbf{a}}(t)$	NUMBER OF ITERATIONS S.T. THE CHOSEN ARM WAS $\mathbf{a}$ , $T_{\mathbf{a}}(t) = \sum_{s=1}^{t-1} \mathbb{1}\{\mathbf{a}(s) = \mathbf{a}\}$
$T_{\tilde{\mathbf{a}}}^{\mathbf{a}}(t)$	NUMBER OF ITERATIONS S.T. THE LEADER WAS $\tilde{\mathbf{a}}$ , THE CHOSEN ARM WAS $\mathbf{a}$ , AND $\mathbf{a}$ WAS CHOSEN BY THE ARGMAX ON $\sum_{k=1}^K b_{a_k,k}(t)$ : $T_{\tilde{\mathbf{a}}}^{\mathbf{a}}(t) = \sum_{s=1}^{t-1} \mathbb{1}\{\tilde{\mathbf{a}}(s) = \tilde{\mathbf{a}}, \mathbf{a}(s) = \mathbf{a}, \tilde{T}_{\tilde{\mathbf{a}}}(s)/L \notin \mathbb{N}\}$
$\hat{\rho}_{i,k}(t)$	ESTIMATION OF $\rho_{i,k}$ AT ITERATION $t$ , $\hat{\rho}_{i,k}(t) = \frac{1}{T_{i,k}(t)} \sum_{s=1}^{t-1} \mathbb{1}\{a_k(s) = i\} c_k(s)$
$b_{i,k}(t)$	KULLBACK-LEIBLER INDEX OF $\hat{\rho}_{i,k}(t)$ , $b_{i,k}(t) = f\left(\hat{\rho}_{i,k}(t), T_{i,k}(t), \tilde{T}_{\tilde{\mathbf{a}}}(t)(t) + 1\right)$
$f$	KULLBACK-LEIBLER INDEX FUNCTION, $f(\hat{\rho}, s, t) = \sup\{p \in [\hat{\rho}, 1] : s \times \text{kl}(\hat{\rho}, p) \leq \log(t) + 3 \log(\log(t))\}$ ,
$\text{kl}(p, q)$	KULLBACK-LEIBLER DIVERGENCE FROM A BERNOULLI DISTRIBUTION OF MEAN $p$ TO A BERNOULLI DISTRIBUTION OF MEAN $q$ , $\text{kl}(p, q) = p \log\left(\frac{p}{q}\right) + (1-p) \log\left(\frac{1-p}{1-q}\right)$
$B_{\mathbf{a}}(t)$	PSEUDO-SUM OF INDICES OF $\mathbf{a}$ AT ITERATION $t$ , $B_{\mathbf{a}}(t) = \sum_{k=1}^K b_{a_k,k}(t) - \sum_{k=1}^K b_{\tilde{\mathbf{a}}_k(t),k}(t)$

CONTINUED ON NEXT PAGE

## Parametric Graph for Unimodal Ranking Bandit

SYMBOL	MEANING
$\mathcal{N}_{\pi^*}(a^*)$	NEIGHBORHOOD OF THE BEST ARM
$K_{\mathbf{a}}$	(WITH COMBINATORIAL BANDIT SETTING) NUMBER OF ELEMENTS IN $\mathbf{a}$ BUT NOT IN $\mathbf{a}^*$ , $K_{\mathbf{a}} = \min_{\mathbf{a}^* \in \mathcal{A}: \mu_{\mathbf{a}^*} = \mu^*}  \mathbf{a} \setminus \mathbf{a}^* $
$K_{max}$	(WITH COMBINATORIAL BANDIT SETTING) MAXIMAL NUMBER OF ELEMENTS IN A SUB-OPTIMAL ARM $\mathbf{a}$ BUT NOT IN AN OPTIMAL ARM $\mathbf{a}^*$ , $K_{max} = \max_{\mathbf{a} \in \mathcal{A}: \mu_{\mathbf{a}} \neq \mu^*} K_{\mathbf{a}}$
$c^*(\boldsymbol{\theta}, \boldsymbol{\kappa})$	COEFFICIENT IN THE REGRET BOUND OF PMED
$c$	(IN $\varepsilon_n$ -GREEDY) PARAMETER CONTROLLING THE PROBABILITY OF EXPLORATION
$c$	(IN PB-MHB) PARAMETER CONTROLLING SIZE OF THE STEP IN THE METROPOLIS HASTING INFERENCE
$m$	(IN PB-MHB) NUMBER OF STEP IN THE METROPOLIS HASTING INFERENCE

### B. Proof of Lemma 1 (PBM Fulfills Assumption 1)

*Proof of Lemma 1.* Let  $(L, K, (\rho_{i,k})_{(i,k) \in [L] \times [K]})$  be an online learning to rank (OLR) problem with users following PBM, with positive probabilities of looking at a given position. Therefore, there exists  $\boldsymbol{\theta} \in [0, 1]^L$  and  $\boldsymbol{\kappa} \in (0, 1]^K$  such that for any item  $i$  and any position  $k$ ,  $\rho_{i,k} = \theta_i \kappa_k$ .

Let  $\mathbf{a} \in \mathcal{P}_K^L$  be a recommendation, and let  $\boldsymbol{\pi} \in \Pi_{\boldsymbol{\rho}}(\mathbf{a})$  be an appropriate ranking of positions. One of the four following properties is satisfied:

$$\exists k \in [K-1] \text{ s.t. } \theta_{a_{\pi_k}} < \theta_{a_{\pi_{k+1}}}, \quad (8)$$

$$\exists k \in [K-1] \text{ s.t. } \kappa_{\pi_k} < \kappa_{\pi_{k+1}}, \quad (9)$$

$$\exists i \in [L] \setminus \mathbf{a}([K]) \text{ s.t. } \theta_{a_{\pi_K}} < \theta_i, \quad (10)$$

$$\begin{cases} \forall k \in [K-1], \theta_{a_{\pi_k}} \geq \theta_{a_{\pi_{k+1}}} \\ \forall k \in [K-1], \kappa_{\pi_k} \geq \kappa_{\pi_{k+1}} \\ \forall i \in [L] \setminus \mathbf{a}([K]), \theta_{a_{\pi_K}} \geq \theta_i \end{cases}. \quad (11)$$

Let prove, by considering each of these properties one by one, that  $\mathbf{a}$  is either one of the best arms, or  $\mathbf{a}$  fulfills either Property (2) or Property (3) of Assumption 1.

If Property (8) is satisfied and  $\theta_{a_{\pi_k}} = 0$ , then by definition of  $\boldsymbol{\pi}$  and  $\Pi_{\boldsymbol{\rho}}(\mathbf{a})$ ,  $0 = \theta_{a_{\pi_k}} \kappa_{\pi_k} \geq \theta_{a_{\pi_{k+1}}} \kappa_{\pi_{k+1}} > 0$  which is absurd.

Therefore, If Property (8) is satisfied,  $\frac{\theta_{a_{\pi_{k+1}}}}{\theta_{a_{\pi_k}}} > 1$ .

Note that by definition of  $\boldsymbol{\pi}$  and  $\Pi_{\boldsymbol{\rho}}(\mathbf{a})$ , and as  $\rho_{i,k} = \theta_i \kappa_k$ ,  $\theta_{a_{\pi_k}} \kappa_{\pi_k} \geq \theta_{a_{\pi_{k+1}}} \kappa_{\pi_{k+1}}$ .

Hence  $\kappa_{\pi_k} \geq \frac{\theta_{a_{\pi_{k+1}}}}{\theta_{a_{\pi_k}}} \kappa_{\pi_{k+1}} > \kappa_{\pi_{k+1}}$ , and

$$\begin{aligned} \mu_{\mathbf{a}} - \mu_{\mathbf{a} \circ (\pi_k, \pi_{k+1})} &= \theta_{a_{\pi_k}} \kappa_{\pi_k} + \theta_{a_{\pi_{k+1}}} \kappa_{\pi_{k+1}} - \left( \theta_{a_{\pi_{k+1}}} \kappa_{\pi_k} + \theta_{a_{\pi_k}} \kappa_{\pi_{k+1}} \right) \\ &= \left( \theta_{a_{\pi_k}} - \theta_{a_{\pi_{k+1}}} \right) \left( \kappa_{\pi_k} - \kappa_{\pi_{k+1}} \right) \\ &< 0, \end{aligned}$$

meaning  $\mu_{\mathbf{a}} < \mu_{\mathbf{a} \circ (\pi_k, \pi_{k+1})}$ , which corresponds to Property (2) of Assumption 1.

Similarly, if Property (9) is satisfied, then Property (2) of Assumption 1 is fulfilled.

If Property (10) is satisfied,

$$\begin{aligned} \mu_{\mathbf{a}} - \mu_{\mathbf{a}[\pi_K := i]} &= \theta_{a_{\pi_K}} \kappa_{\pi_K} - \theta_i \kappa_{\pi_K} \\ &= \left( \theta_{a_{\pi_K}} - \theta_i \right) \kappa_{\pi_K} \\ &< 0. \end{aligned}$$

Hence  $\mu_{\mathbf{a}} < \mu_{\mathbf{a}[\pi_K := i]}$ , which corresponds to Property (3) of Assumption 1.

Finally, if Property (11) is satisfied,  $\mu_{\mathbf{a}} = \mu^*$ .

Overall, either  $\mathbf{a}$  is one of the best arms, or  $\mathbf{a}$  fulfills Property (2) of Assumption 1, or  $\mathbf{a}$  fulfills Property (3) of Assumption 1, which concludes the proof.  $\square$

### C. Preliminary to the Analysis of GRAB

The analysis of GRAB requires a control of the number of high deviations, as expressed by Lemma B.1 of (Combes & Proutière, 2014). Let us recall this lemma, which we denote Lemma 4 in current paper.

**Lemma 4** (Lemma B.1 of (Combes & Proutière, 2014)). *Let  $i \in [L]$ ,  $k \in [K]$ ,  $\epsilon > 0$ . Define  $\mathcal{F}(T)$  the  $\sigma$ -algebra generated by  $(\mathbf{c}(t))_{t \in [T]}$ . Let  $\Lambda \subseteq \mathbb{N}$  be a random set of instants. Assume that there exists a sequence of random sets  $(\Lambda(s))_{s \geq 1}$  such that (i)  $\Lambda \subseteq \bigcup_{s \geq 1} \Lambda(s)$ , (ii) for all  $s \geq 1$  and all  $t \in \Lambda(s)$ ,  $T_{i,k}(t) \geq \epsilon s$ , (iii)  $|\Lambda(s)| \leq 1$ , and (iv) the event  $t \in \Lambda(s)$  is  $\mathcal{F}_t$ -measurable. Then for all  $\delta > 0$ ,*

$$\mathbb{E} \left[ \sum_{t \geq 1} \mathbb{1}\{t \in \Lambda, |\hat{\rho}_{i,k}(t) - \rho_{i,k}| \geq \delta\} \right] \leq \frac{1}{\epsilon \delta^2}$$

### D. Proof of Theorem 2 (Upper-bound on the Regret of KL-CombUCB)

*Proof of Theorem 2.* Let  $\mathbf{a} \in \mathcal{A}$  be a sub-optimal arm. Let  $\mathbf{a}^* \in \mathcal{A}$  be an optimal arm such that  $|\mathbf{a} \setminus \mathbf{a}^*| = K_{\mathbf{a}}$ .

We denote  $\bar{K}_{\mathbf{a}} \stackrel{\text{def}}{=} |\mathbf{a}^* \setminus \mathbf{a}|$ ,  $T_{\mathbf{a}}(t) \stackrel{\text{def}}{=} \sum_{s=1}^{t-1} \mathbb{1}\{\mathbf{a}(s) = \mathbf{a}\}$  the number of time the arm  $\mathbf{a}$  has been drawn, and  $T_e(t) \stackrel{\text{def}}{=} \sum_{s=1}^{t-1} \mathbb{1}\{e \in \mathbf{a}(s)\}$  the number of time the element  $e$  was in the drawn arm.

Let decompose the expected number of iterations at which the permutation  $\mathbf{a}$  is recommended:

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{\mathbf{a}(t) = \mathbf{a}\} \right] &\leq \sum_{e \in \mathbf{a} \setminus \mathbf{a}^*} \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1} \left\{ \mathbf{a}(t) = \mathbf{a}, |\hat{\rho}_e(t) - \rho_e| \geq \frac{\Delta_{\mathbf{a}}}{2K_{\mathbf{a}}} \right\} \right] \\ &+ \sum_{e \in \mathbf{a}^* \setminus \mathbf{a}} \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}\{b_e(t) \leq \rho_e\} \right] \\ &+ \mathbb{E} \left[ \sum_{t=|E|}^T \mathbb{1} \left\{ \mathbf{a}(t) = \mathbf{a}, \forall e \in \mathbf{a} \setminus \mathbf{a}^*, |\hat{\rho}_e(t) - \rho_e| < \frac{\Delta_{\mathbf{a}}}{2K_{\mathbf{a}}}, \forall e \in \mathbf{a}^* \setminus \mathbf{a}, b_e(t) > \rho_e \right\} \right] \\ &+ |E|. \end{aligned}$$

The proof consists in upper-bounding each term on the right-hand side.

**First Term** Let  $e \in \mathbf{a} \setminus \mathbf{a}^*$ , and denote  $A_e = \left\{ t \in [T] : \mathbf{a}(t) = \mathbf{a}, |\hat{\rho}_e(t) - \rho_e| \geq \frac{\Delta_{\mathbf{a}}}{2K_{\mathbf{a}}} \right\}$ .

$A_e \subseteq \bigcup_{s \in \mathbb{N}} \Lambda_k(s)$ , where  $\Lambda_k(s) \stackrel{\text{def}}{=} \{t \in A_e : T_{\mathbf{a}}(t) = s\}$ . For any integer value  $s$ ,  $|\Lambda_k(s)| \leq 1$  as  $T_{\mathbf{a}}(t)$  increases for each  $t \in A_e$ . Note that for each  $s \in \mathbb{N}$  and  $n \in \Lambda_k(s)$ ,  $T_e(n) \geq T_{\mathbf{a}}(n) = s$ . Then, by Lemma 4

$$\begin{aligned}
 \mathbb{E}[|A_e|] &\leq \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\{t \in A_e\}\right] \\
 &= \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\left\{t \in A_e, |\hat{\rho}_e(t) - \rho_e| \geq \frac{\Delta_{\mathbf{a}}}{2K_{\mathbf{a}}}\right\}\right] \\
 &\leq \frac{4K_{\mathbf{a}}^2}{\Delta_{\mathbf{a}}^2}.
 \end{aligned}$$

Hence,  $\sum_{e \in \mathbf{a} \setminus \mathbf{a}^*} \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\left\{\mathbf{a}(t) = \mathbf{a}, |\hat{\rho}_e(t) - \rho_e| \geq \frac{\Delta_{\mathbf{a}}}{2K_{\mathbf{a}}}\right\}\right] = \sum_{e \in \mathbf{a} \setminus \mathbf{a}^*} \mathbb{E}[|A_e|] \leq \frac{4K_{\mathbf{a}}^3}{\Delta_{\mathbf{a}}^2}$ .

**Second Term** Let  $e \in \mathbf{a}^* \setminus \mathbf{a}$ , and denote  $B_e \stackrel{def}{=} \{t \in [T] : b_e(t) \leq \rho_e\}$ .

By Theorem 10 of (Garivier & Cappé, 2011),  $\mathbb{E}[|B_e|] = O(\log \log T)$ , so  $\sum_{e \in \mathbf{a}^* \setminus \mathbf{a}} \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\{b_e(t) \leq \rho_e\}\right] = \mathcal{O}(\bar{K}_{\mathbf{a}} \log \log T)$ .

**Third Term** Let note  $C \stackrel{def}{=} \left\{t \in [T] \setminus [|E|] : \mathbf{a}(t) = \mathbf{a}, \forall e \in \mathbf{a} \setminus \mathbf{a}^*, |\hat{\rho}_e(t) - \rho_e| < \frac{\Delta_{\mathbf{a}}}{2K_{\mathbf{a}}}, \forall e \in \mathbf{a}^* \setminus \mathbf{a}, b_e(t) > \rho_e\right\}$ .

Let  $t \in C$ .

At each step of the initialization phase, the algorithm removes at least one element  $e$  of the set  $\tilde{E}$  of unseen elements. Therefore, the initialization lasts at most  $|E|$  iterations. Hence, at iteration  $t$ ,  $\mathbf{a}(t) = \mathbf{a}$  is chosen as  $\sum_{e \in \mathbf{a}} b_e(t) = \max_{\mathbf{a}' \in \mathcal{A}} \sum_{e \in \mathbf{a}'} b_e(t)$ .

Then, by Pinsker's inequality and the fact that  $t \leq T$ , and  $T_e(t) \geq T_{\mathbf{a}}(t)$  for any  $e$  in  $\mathbf{a}$ ,

$$\begin{aligned}
 0 &\leq \sum_{e \in \mathbf{a}} b_e(t) - \sum_{e \in \mathbf{a}^*} b_e(t) \\
 &= \sum_{e \in \mathbf{a} \setminus \mathbf{a}^*} b_e(t) - \sum_{e \in \mathbf{a}^* \setminus \mathbf{a}} b_e(t) \\
 &\leq \sum_{e \in \mathbf{a} \setminus \mathbf{a}^*} \hat{\rho}_e(t) + \sqrt{\frac{\log(t) + 3 \log(\log(t))}{2T_e(t)}} - \sum_{e \in \mathbf{a}^* \setminus \mathbf{a}} b_e(t) \\
 &< \sum_{e \in \mathbf{a} \setminus \mathbf{a}^*} \rho_e + \frac{\Delta_{\mathbf{a}}}{2K_{\mathbf{a}}} + \sqrt{\frac{\log(T) + 3 \log(\log(T))}{2T_{\mathbf{a}}(t)}} - \sum_{e \in \mathbf{a}^* \setminus \mathbf{a}} \rho_e \\
 &\leq \sum_{e \in \mathbf{a}} \rho_e - \sum_{e \in \mathbf{a}^*} \rho_e + K_{\mathbf{a}} \frac{\Delta_{\mathbf{a}}}{2K_{\mathbf{a}}} + K_{\mathbf{a}} \sqrt{\frac{\log(T) + 3 \log(\log(T))}{2T_{\mathbf{a}}(t)}} \\
 &= -\Delta_{\mathbf{a}} + \frac{2\Delta_{\mathbf{a}}}{2} + K_{\mathbf{a}} \sqrt{\frac{\log(T) + 3 \log(\log(T))}{2T_{\mathbf{a}}(t)}} \\
 &= -\frac{\Delta_{\mathbf{a}}}{2} + K_{\mathbf{a}} \sqrt{\frac{\log(T) + 3 \log(\log(T))}{2T_{\mathbf{a}}(t)}}.
 \end{aligned}$$

Hence,  $T_{\mathbf{a}}(t) < K_{\mathbf{a}}^2 \frac{2 \log(T) + 6 \log(\log(T))}{\Delta_{\mathbf{a}}^2}$ . Therefore,  $C \subseteq \left\{t \in [T] \setminus [|E|] : \mathbf{a}(t) = \mathbf{a}, T_{\mathbf{a}}(t) < K_{\mathbf{a}}^2 \frac{2 \log(T) + 6 \log(\log(T))}{\Delta_{\mathbf{a}}^2}\right\}$ , and

$$\begin{aligned}
 & \mathbb{E} \left[ \sum_{t=|E|}^T \mathbb{1} \left\{ \mathbf{a}(t) = \mathbf{a}, \forall e \in \mathbf{a} \setminus \mathbf{a}^*, |\hat{\rho}_e(t) - \rho_e| < \frac{\Delta_{\mathbf{a}}}{2K_{\mathbf{a}}}, \forall e \in \mathbf{a}^* \setminus \mathbf{a}, b_e(t) > \rho_e \right\} \right] \\
 &= \mathbb{E} [|C|] \\
 &\leq \mathbb{E} \left[ \left| \left\{ t \in [T] \setminus [|E|] : \mathbf{a}(t) = \mathbf{a}, T_{\mathbf{a}}(t) < K_{\mathbf{a}}^2 \frac{2 \log(T) + 6 \log(\log(T))}{\Delta_{\mathbf{a}}^2} \right\} \right| \right] \\
 &\leq K_{\mathbf{a}}^2 \frac{2 \log(T) + 6 \log(\log(T))}{\Delta_{\mathbf{a}}^2}.
 \end{aligned}$$

**Regret upper-bound** Overall,

$$\begin{aligned}
 \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1} \{ \mathbf{a}(t) = \mathbf{a} \} \right] &\leq \frac{4K_{\mathbf{a}}^3}{\Delta_{\mathbf{a}}^2} + \mathcal{O}(\bar{K}_{\mathbf{a}} \log \log T) + K_{\mathbf{a}}^2 \frac{2 \log(T) + 6 \log(\log(T))}{\Delta_{\mathbf{a}}^2} + |E| \\
 &= \frac{2K_{\mathbf{a}}^2}{\Delta_{\mathbf{a}}^2} \log(T) + \mathcal{O} \left( \left( \bar{K}_{\mathbf{a}} + \frac{K_{\mathbf{a}}^2}{\Delta_{\mathbf{a}}^2} \right) \log \log T \right)
 \end{aligned}$$

and

$$\begin{aligned}
 R(T) &= \sum_{\mathbf{a} \in \mathcal{A}: \mu_{\mathbf{a}} \neq \mu^*} \Delta_{\mathbf{a}} \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1} \{ \mathbf{a}(t) = \mathbf{a} \} \right] \\
 &\leq \sum_{\mathbf{a} \in \mathcal{A}: \mu_{\mathbf{a}} \neq \mu^*} \frac{2K_{\mathbf{a}}^2}{\Delta_{\mathbf{a}}} \log(T) + \mathcal{O} \left( \left( \bar{K}_{\mathbf{a}} \Delta_{\mathbf{a}} + \frac{K_{\mathbf{a}}^2}{\Delta_{\mathbf{a}}} \right) \log \log T \right) \\
 &= \mathcal{O} \left( \frac{|\mathcal{A}| K_{max}^2}{\Delta_{min}} \log T \right),
 \end{aligned}$$

which concludes the proof. □

## E. Proof of Lemma 2 (Upper-bound on the Number of Iterations of GRAB for which $\tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}} \neq \mathbf{a}^*$ )

*Proof of Lemma 2.* Let  $\tilde{\mathbf{a}} \in \mathcal{P}_K^L \setminus \{\mathbf{a}^*\}$  and prove that  $\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1} \{ \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}} \} \right] = \mathcal{O}(\log \log T)$ .

The proof requires notations related to the neighborhood of  $\tilde{\mathbf{a}}$ . Let  $\mathcal{N} \stackrel{def}{=} \bigcup_{\pi \in \mathcal{P}_K^K} \mathcal{N}_{\pi}(\tilde{\mathbf{a}})$  be the set of all the potential neighbors of  $\tilde{\mathbf{a}}$ . By definition of the neighborhoods,

$$\mathcal{N} = \{ \tilde{\mathbf{a}} \circ (k, k') : k, k' \in [K]^2, k > k' \} \cup \{ \tilde{\mathbf{a}}[k := i] : k \in [K], i \in [L] \setminus \tilde{\mathbf{a}}([K]) \},$$

and its size is  $N = K(2L - K - 1)/2$ . As  $\tilde{\mathbf{a}}$  is sub-optimal, and due to Assumption 1, for any appropriate ranking of positions  $\pi \in \Pi_{\rho}(\tilde{\mathbf{a}})$ , there exists a recommendation  $\mathbf{a}^+$  with a strictly better expected reward than  $\tilde{\mathbf{a}}$  in the neighborhood  $\mathcal{N}_{\pi}(\tilde{\mathbf{a}})$ . We denote

$$\mathcal{N}^+ \stackrel{def}{=} \bigcup_{\pi \in \Pi_{\rho}(\tilde{\mathbf{a}})} \left\{ \mathbf{a}^+ \in \mathcal{N}_{\pi}(\tilde{\mathbf{a}}) : \mu_{\mathbf{a}^+} = \max_{\mathbf{a} \in \mathcal{N}_{\pi}(\tilde{\mathbf{a}})} \mu_{\mathbf{a}} \right\}$$

the set of such recommendations. We also chose  $\epsilon < \min\{1/(2N), 1/L\}$  and note

$$\delta \stackrel{def}{=} \min_{\pi \in \Pi_{\rho}(\tilde{\mathbf{a}})} \min_{\mathbf{a} \in \mathcal{N}_{\pi}(\tilde{\mathbf{a}}) \cup \{\tilde{\mathbf{a}}\} \setminus \mathcal{N}^+} \left( \max_{\mathbf{a}' \in \mathcal{N}_{\pi}(\tilde{\mathbf{a}})} \mu_{\mathbf{a}'} - \mu_{\mathbf{a}} \right).$$



To bound  $\mathbb{E}[\mathbb{1}\{\tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}\}]$ , we use the decomposition  $\{t \in [T] : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}\} \subseteq \bigcup_{\mathbf{a}^+ \in \mathcal{N}^+} A_{\mathbf{a}^+} \cup B$  where for any permutation  $\mathbf{a}^+ \in \mathcal{N}^+$ ,

$$A_{\mathbf{a}^+} = \{t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, T_{\mathbf{a}^+}(t) \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t)\}$$

and

$$B = \{t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \forall \mathbf{a}^+ \in \mathcal{A}^+, T_{\mathbf{a}^+}(t) < \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t)\}.$$

Hence,

$$\mathbb{E}[\mathbb{1}\{\tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}\}] \leq \sum_{\mathbf{a}^+ \in \mathcal{A}^+} \mathbb{E}[|A_{\mathbf{a}^+}|] + \mathbb{E}[|B|].$$

**Bound on  $\mathbb{E}[|A_{\mathbf{a}^+}|]$**  Let  $\mathbf{a}^+$  be a permutation in  $\mathcal{N}^+$  and denote  $\mathcal{K}^+$  the set of positions for which  $\mathbf{a}^+$  and  $\tilde{\mathbf{a}}$  disagree:  $\mathcal{K}^+ = \{k \in [K] : a_k^+ \neq \tilde{a}_k\}$ . The permutation  $\mathbf{a}^+$  is in the neighborhood of  $\tilde{\mathbf{a}}$ , so either  $\mathbf{a}^+ = \tilde{\mathbf{a}} \circ (k, k')$  or  $\mathbf{a}^+ = \mathbf{a}[k := i]$ , with  $k$  and  $k'$  in  $[K]$ , and  $i$  in  $[L]$ . Overall,  $|\mathcal{K}^+| \leq 2$ .

By the design of the algorithm and by definition of  $\epsilon$ , we have that  $\forall t \in A_{\mathbf{a}^+}, T_{\tilde{\mathbf{a}}}(t) \geq \tilde{T}_{\tilde{\mathbf{a}}}(t)/L > \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t)$ . Moreover, at the considered iterations  $\tilde{\mathbf{a}}$  is the leader, so

$$\begin{aligned} A_{\mathbf{a}^+} &\subseteq \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{T}_{\tilde{\mathbf{a}}}(t) < \frac{1}{\epsilon} \right\} \cup \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \min\{T_{\tilde{\mathbf{a}}}(t), T_{\mathbf{a}^+}(t)\} \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t) \geq 1, \sum_{\ell} \hat{\rho}_{\tilde{a}_\ell, \ell}(t) \geq \sum_{\ell} \hat{\rho}_{a_\ell^+, \ell}(t) \right\} \\ &\subseteq \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{T}_{\tilde{\mathbf{a}}}(t) < \frac{1}{\epsilon} \right\} \cup \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \min\{T_{\tilde{\mathbf{a}}}(t), T_{\mathbf{a}^+}(t)\} \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t), \sum_{k \in \mathcal{K}^+} \hat{\rho}_{\tilde{a}_k, k}(t) \geq \sum_{k \in \mathcal{K}^+} \hat{\rho}_{a_k^+, k}(t) \right\} \\ &\subseteq \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{T}_{\tilde{\mathbf{a}}}(t) < \frac{1}{\epsilon} \right\} \\ &\quad \cup \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \min\{T_{\tilde{\mathbf{a}}}(t), T_{\mathbf{a}^+}(t)\} \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t), \exists k \in \mathcal{K}^+, |\hat{\rho}_{\tilde{a}_k, k}(t) - \rho_{\tilde{a}_k, k}| \geq \frac{\delta}{2|\mathcal{K}^+|} \text{ or } |\hat{\rho}_{a_k^+, k}(t) - \rho_{a_k^+, k}| \geq \frac{\delta}{2|\mathcal{K}^+|} \right\} \\ &\subseteq \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{T}_{\tilde{\mathbf{a}}}(t) < \frac{1}{\epsilon} \right\} \cup \bigcup_{k \in \mathcal{K}^+} \bigcup_{i \in \{\tilde{a}_k, a_k^+\}} \Lambda_{i, k}, \end{aligned}$$

with  $\Lambda_{i, k} \stackrel{def}{=} \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \min\{T_{\tilde{\mathbf{a}}}(t), T_{\mathbf{a}^+}(t)\} \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t), |\hat{\rho}_{i, k}(t) - \rho_{i, k}| \geq \frac{\delta}{2|\mathcal{K}^+|} \right\}$ .

Fix  $k$  in  $\mathcal{K}^+$  and  $i$  in  $\{\tilde{a}_k, a_k^+\}$ .  $\Lambda_{i, k} \subseteq \bigcup_{s \in \mathbb{N}} \Lambda_{i, k}(s)$ , with  $\Lambda_{i, k}(s) \stackrel{def}{=} \{t \in \Lambda_{i, k} : \tilde{T}_{\tilde{\mathbf{a}}}(t) = s\}$ .  $|\Lambda_{i, k}(s)| \leq 1$  as  $\tilde{T}_{\tilde{\mathbf{a}}}(t)$  increases for each  $t \in \Lambda_{i, k}$ . Note that for each  $s \in \mathbb{N}$  and  $n \in \Lambda_{i, k}(s)$ ,  $T_{i, k}(n) \geq \min\{T_{\tilde{\mathbf{a}}}(n), T_{\mathbf{a}^+}(n)\} \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(n) = \epsilon s$ . Then, by Lemma 4

$$\begin{aligned} \mathbb{E}[|\Lambda_{i, k}|] &= \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\{t \in \Lambda_{i, k}\}\right] \\ &= \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\left\{t \in \Lambda_{i, k}, |\hat{\rho}_{i, k}(t) - \rho_{i, k}| > \frac{\delta}{2|\mathcal{K}^+|}\right\}\right] \\ &\leq \frac{4|\mathcal{K}^+|^2}{\epsilon \delta^2} \end{aligned}$$

Hence,  $\mathbb{E}[|A_{\mathbf{a}^+}|] \leq \frac{1}{\epsilon} + \sum_{k \in \mathcal{K}^+} \sum_{i \in \{\tilde{a}_k, a_k^+\}} \mathbb{E}[|\Lambda_{i, k}|] \leq \frac{1}{\epsilon} + \frac{8|\mathcal{K}^+|^3}{\epsilon \delta^2}$ .

**Bound on  $\mathbb{E}[|B|]$**  We first split  $B$  in two parts:  $B = B^{t_0} \cup B_{t_0}^T$ , where  $B^{t_0} \stackrel{def}{=} \{t \in B : \tilde{T}_{\tilde{\mathbf{a}}}(t) \leq t_0\}$ ,  $B_{t_0}^T \stackrel{def}{=} \{t \in B : \tilde{T}_{\tilde{\mathbf{a}}}(t) > t_0\}$ , and  $t_0$  is chosen as small as possible to satisfy three constraints required in the rest of the proof.

Namely,  $t_0 = \max\left\{\frac{1}{\epsilon}, (1+N)(1 - \frac{1}{L} - \epsilon N)^{-1}, \inf\left\{t : 2\sqrt{\frac{\log(t+1)+3\log(\log(t+1))}{2\epsilon t}} < \frac{\delta}{8}\right\}\right\}$ . Note that  $t_0$  only depends on  $K, L$  and  $\delta$ , and that  $(1 - \frac{1}{L} - \epsilon N) > 0$  (assuming  $L \geq 2$ ) as  $\epsilon < 1/(2N)$ .

We also define

- $D \stackrel{def}{=} \bigcup_{(\mathbf{a},k) \in (\mathcal{N} \cup \{\tilde{\mathbf{a}}\}) \times [K]} D_{\mathbf{a},k}$ , where  $D_{\mathbf{a},k} \stackrel{def}{=} \{t \in [T] : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \mathbf{a}(t) = \mathbf{a}, |\hat{\rho}_{a_k,k}(t) - \rho_{a_k,k}| \geq \frac{\delta}{8}\}$ ,
- $E \stackrel{def}{=} \bigcup_{(\mathbf{a}^+,k) \in \mathcal{N}^+ \times [K]} E_{\mathbf{a}^+,k}$ , where  $E_{\mathbf{a}^+,k} \stackrel{def}{=} \{t \in [T] : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, b_{a_k^+,k}(t) \leq \rho_{a_k^+,k}\}$ ,
- and  $F \stackrel{def}{=} \{t \in [T] : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{\pi}(t) \notin \Pi_{\rho}(\tilde{\mathbf{a}})\}$ .

Let  $t \in B_{t_0}^T$ . By construction, GRAB forces itself to select  $\left\lceil \frac{\tilde{T}_{\tilde{\mathbf{a}}}(t)}{L} \right\rceil$  times the leader  $\tilde{\mathbf{a}}$  between iterations 1 and  $t - 1$ . So,

$$\tilde{T}_{\tilde{\mathbf{a}}}(t) = \left\lceil \frac{\tilde{T}_{\tilde{\mathbf{a}}}(t)}{L} \right\rceil + \sum_{\mathbf{a} \in \mathcal{N} \cup \{\tilde{\mathbf{a}}\}} T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(t)$$

where  $T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(t) = \sum_{s=1}^{t-1} \mathbf{1} \left\{ \tilde{\mathbf{a}}(s) = \tilde{\mathbf{a}}, \mathbf{a}(s) = \mathbf{a}, \tilde{T}_{\tilde{\mathbf{a}}}(s)/L \notin \mathbb{N} \right\}$  is the number of times arm  $\mathbf{a} \in \mathcal{N} \cup \{\tilde{\mathbf{a}}\}$  has been played **normally** (i.e not forced) while  $\tilde{\mathbf{a}}$  was leader, up to time  $t - 1$ . Let prove by contradiction that there is at least one recommendation  $\mathbf{a}$  that has been selected **normally** more than  $\epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t) + 1$  times, namely  $T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(t) \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t) + 1$ .

Assume that for each recommendation  $\mathbf{a}$  in  $\mathcal{N} \cup \{\tilde{\mathbf{a}}\}$ ,  $T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(t) < \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t) + 1$ . Then

$$\begin{aligned} \tilde{T}_{\tilde{\mathbf{a}}}(t) &= \left\lceil \frac{\tilde{T}_{\tilde{\mathbf{a}}}(t)}{L} \right\rceil + \sum_{\mathbf{a} \in \mathcal{N} \cup \{\tilde{\mathbf{a}}\}} T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(t) \\ &< 1 + \frac{\tilde{T}_{\tilde{\mathbf{a}}}(t)}{L} + N(\epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t) + 1). \end{aligned}$$

Therefore  $\tilde{T}_{\tilde{\mathbf{a}}}(t)(1 - \frac{1}{L} - N\epsilon) < 1 + N$ , which contradicts  $t \in B_{t_0}^T$ .

So, there exists a recommendation  $\mathbf{a}$  such that  $T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(t) \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t) + 1$ . Let denote  $s'$  the first iteration such that  $T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(s') \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t) + 1$ . At this iteration,  $T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(s') = T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(s' - 1) + 1$ , meaning that  $\tilde{\mathbf{a}}(s' - 1) = \tilde{\mathbf{a}}, \mathbf{a}(s' - 1) = \mathbf{a}, \tilde{T}_{\tilde{\mathbf{a}}}(s' - 1)/L \notin \mathbb{N}$ , and  $T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(s' - 1) \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t)$ . Therefore, the set  $\{s \in [t] : \tilde{\mathbf{a}}(s) = \tilde{\mathbf{a}}, T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(s) \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t), \tilde{T}_{\tilde{\mathbf{a}}}(s)/L \notin \mathbb{N}\}$  is non-empty. We define  $\psi(t)$  as the minimum on this set

$$\psi(t) \stackrel{def}{=} \min \left\{ s \in [t] : \tilde{\mathbf{a}}(s) = \tilde{\mathbf{a}}, T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(s) \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t), \tilde{T}_{\tilde{\mathbf{a}}}(s)/L \notin \mathbb{N} \right\}.$$

We note  $\mathbf{a}$  the recommendation  $\mathbf{a}(\psi(t))$  at iteration  $\psi(t)$ . We have  $\mathbf{a} \notin \mathcal{N}^+$  since for any recommendation  $\mathbf{a}^+ \in \mathcal{N}^+$ ,  $T_{\mathbf{a}^+}^{\tilde{\mathbf{a}}}(\psi(t)) \leq T_{\mathbf{a}^+}^{\tilde{\mathbf{a}}}(t) \leq T_{\mathbf{a}^+}(t) < \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t)$ . Let  $\mathbf{a}^+$  be one of the best recommendations in  $\mathcal{N}_{\tilde{\pi}(\psi(t))}(\tilde{\mathbf{a}}) \cup \{\tilde{\mathbf{a}}\}$ , meaning  $\mu_{\mathbf{a}^+} = \max_{\mathbf{a}' \in \mathcal{N}_{\tilde{\pi}(\psi(t))}(\tilde{\mathbf{a}}) \cup \{\tilde{\mathbf{a}}\}} \mu_{\mathbf{a}'}$ , and let  $\mathcal{K}$  denote the set of positions for which  $\mathbf{a}$  and  $\mathbf{a}^+$  disagree. As both recommendations are in  $\mathcal{N}_{\tilde{\pi}(\psi(t))}(\tilde{\mathbf{a}}) \cup \{\tilde{\mathbf{a}}\}$ ,  $|\mathcal{K}| \leq 4$ .

Let prove by contradiction that  $\psi(t) \in D \cup E \cup F$ . Assume that  $\psi(t) \notin D \cup E \cup F$ .

Since  $\psi(t) \notin F$ ,  $\tilde{\pi}(\psi(t))$  belongs to  $\Pi_{\rho}(\tilde{\mathbf{a}})$  and hence  $\mathbf{a}^+$  is in  $\mathcal{N}^+$  and  $\sum_k \rho_{a_k^+,k} - \sum_k \rho_{a_k,k} = \mu_{\mathbf{a}^+} - \mu_{\mathbf{a}} \geq \delta$ .

Moreover, since  $\psi(t) \notin D \cup E$ , for each position  $k \in [K]$ ,  $|\hat{\rho}_{a_k,k}(\psi(t)) - \rho_{a_k,k}| < \frac{\delta}{8}$ , and  $b_{a_k^+,k}(\psi(t)) > \rho_{a_k^+,k}$ .

Finally,  $T_{\mathbf{a}}(\psi(t)) \geq T_{\mathbf{a}}^{\tilde{\mathbf{a}}}(\psi(t)) \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t) \geq 1$ , and therefore  $b_{a_k,k}(\psi(t))$  and  $\hat{\rho}_{a_k,k}(\psi(t))$  are properly defined for any position  $k \in [K]$ .

Then, by Pinsker's inequality and the fact that  $\psi(t) \leq t$ ,  $\tilde{T}_{\tilde{\mathbf{a}}}(s)$  is non-decreasing in  $s$ , and  $T_{\mathbf{a}}(\psi(t)) \geq \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t)$ ,

$$\begin{aligned}
 \sum_k b_{a_k, k}(\psi(t)) - \sum_k b_{a_k^+, k}(\psi(t)) &= \sum_{k \in \mathcal{K}} b_{a_k, k}(\psi(t)) - b_{a_k^+, k}(\psi(t)) \\
 &\leq \sum_{k \in \mathcal{K}} \hat{\rho}_{a_k, k}(\psi(t)) + \sqrt{\frac{\log(\tilde{T}_{\tilde{\mathbf{a}}}(\psi(t)) + 1) + 3 \log(\log(\tilde{T}_{\tilde{\mathbf{a}}}(\psi(t)) + 1))}{2T_{\mathbf{a}}(\psi(t))}} - b_{a_k^+, k}(\psi(t)) \\
 &< \sum_{k \in \mathcal{K}} \rho_{a_k, k} + \frac{\delta}{8} + \sqrt{\frac{\log(\tilde{T}_{\tilde{\mathbf{a}}}(t) + 1) + 3 \log(\log(\tilde{T}_{\tilde{\mathbf{a}}}(t) + 1))}{2\epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t)}} - \rho_{a_k^+, k} \\
 &\leq \sum_{k \in \mathcal{K}} \rho_{a_k, k} + \frac{\delta}{8} + \frac{\delta}{8} - \rho_{a_k^+, k} \\
 &\leq \sum_k \rho_{a_k, k} - \sum_k \rho_{a_k^+, k} + |\mathcal{K}| \cdot 2 \frac{\delta}{8} \\
 &\leq -\delta + 8 \frac{\delta}{8} \\
 &= 0,
 \end{aligned}$$

which contradicts the fact that  $\mathbf{a}$  is played at iteration  $\psi(t)$ . So  $\psi(t) \in D \cup E \cup F$ .

Overall, for any  $t \in B_{t_0}^T$ ,  $\psi(t) \in D \cup E \cup F$ . So,  $B_{t_0}^T \subseteq \bigcup_{n \in D \cup E \cup F} B_{t_0}^T \cap \{t \in [T] : \psi(t) = n\}$ . Let  $n$  be in  $D \cup E \cup F$ . For any  $t$  in  $B_{t_0}^T \cap \{t \in [T] : \psi(t) = n\}$ ,  $T_{\tilde{\mathbf{a}}(n)}(n) = \lceil \epsilon \tilde{T}_{\tilde{\mathbf{a}}}(t) \rceil$  and  $\tilde{T}_{\tilde{\mathbf{a}}}(t+1) = \tilde{T}_{\tilde{\mathbf{a}}}(t) + 1$ . So  $|B_{t_0}^T \cap \{t \in [T] : \psi(t) = n\}| < 1/\epsilon + 1$ . Overall,

$$\mathbb{E}[|B|] \leq t_0 + \mathbb{E}[|B_{t_0}^T|] \leq t_0 + (1/\epsilon + 1)(\mathbb{E}[|D|] + \mathbb{E}[|E|] + \mathbb{E}[|F|]).$$

It remains to upper-bound  $\mathbb{E}[|D|]$ ,  $\mathbb{E}[|E|]$ , and  $\mathbb{E}[|F|]$  to conclude the proof.

**Bound on  $\mathbb{E}[|D|]$**  The upper-bound on  $\mathbb{E}[|D|]$  is obtained with the same strategy as the last step in the proof of the upper-bound on  $\mathbb{E}[|A_{\mathbf{a}^+}|]$ . Let  $\mathbf{a}$  be a recommendation in  $\mathcal{N} \cup \{\tilde{\mathbf{a}}\} \setminus \mathcal{N}^+$ , and  $k \in [K]$  be a position.  $D_{\mathbf{a}, k} \subseteq \bigcup_{s \in \mathbb{N}} \Lambda_{\mathbf{a}, k}(s)$ , where  $\Lambda_{\mathbf{a}, k}(s) \stackrel{\text{def}}{=} \{t \in D_{\mathbf{a}, k} : T_{\mathbf{a}}(t) = s\}$ .  $|\Lambda_{\mathbf{a}, k}(s)| \leq 1$  as  $T_{\mathbf{a}}(t)$  increases for each  $t \in D_{\mathbf{a}, k}$ . Note that for each  $s \in \mathbb{N}$  and  $n \in \Lambda_{\mathbf{a}, k}(s)$ ,  $T_{a_k, k}(n) \geq T_{\mathbf{a}}(n) = s$ . Then, by Lemma 4

$$\begin{aligned}
 \mathbb{E}[|D_{\mathbf{a}, k}|] &\leq \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\{t \in D_{\mathbf{a}, k}\}\right] \\
 &= \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\left\{t \in D_{\mathbf{a}, k}, |\hat{\rho}_{a_k, k}(t) - \rho_{a_k, k}| \geq \frac{\delta}{8}\right\}\right] \\
 &\leq \frac{64}{\delta^2}
 \end{aligned}$$

Hence,  $\mathbb{E}[|D|] \leq \sum_{(\mathbf{a}, k) \in (\mathcal{N} \cup \{\tilde{\mathbf{a}}\} \setminus \mathcal{N}^+) \times [K]} \mathbb{E}[|D_{\mathbf{a}, k}|] \leq \frac{64(N+1)K}{\delta^2}$ .

**Bound on  $\mathbb{E}[|E|]$**  By Theorem 10 of (Garivier & Cappé, 2011),  $\mathbb{E}[|E_{\mathbf{a}^+, k}|] = O(\log(\log(T)))$ , so  $\mathbb{E}[|E|] \leq \sum_{(\mathbf{a}^+, k) \in \mathcal{N}^+ \times [K]} \mathbb{E}[|E_{\mathbf{a}^+, k}|] = O(|\mathcal{N}^+|K \log(\log(T)))$ .

**Bound on  $\mathbb{E}[|F|]$**  By Lemma 3,  $\mathbb{E}[|F|] = \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\{\tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{\pi}(t) \notin \Pi_{\rho}(\tilde{\mathbf{a}})\}\right] = \mathcal{O}(1)$ .

Overall  $\mathbb{E}[\mathbb{1}\{\tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}\}] \leq \frac{|\mathcal{K}^+|}{\epsilon} + \frac{8|\mathcal{K}^+|^3|\mathcal{N}^+|}{\epsilon\delta^2} + t_0 + \left(\frac{1}{\epsilon} + 1\right) \frac{64(N+1)K}{\delta^2} + \mathcal{O}\left(\frac{|\mathcal{N}^+|K}{\epsilon} \log \log T\right) + \mathcal{O}(1) = \mathcal{O}\left(\frac{|\mathcal{N}^+|K}{\epsilon} \log \log T\right)$ , which concludes the proof.  $\square$

## F. Proof of Lemma 3 (Upper-bound on the Number of Iterations of GRAB for which

$$\tilde{\pi}(t) \notin \Pi_{\rho}(\tilde{\mathbf{a}}))$$

*Proof of Theorem 3.* Let  $\tilde{\mathbf{a}}$  be a  $K$ -permutation of  $L$  items. If  $\Pi_{\rho}(\tilde{\mathbf{a}})$  contains all the permutations of  $K$  elements, the set  $\{t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{\pi}(t) \notin \Pi_{\rho}(\tilde{\mathbf{a}})\}$  is empty.

Otherwise, let denote  $\delta$  the smallest non-zero gap between the probability of click at position  $k$  and the probability of click at position  $k' \neq k$ :  $\delta \stackrel{\text{def}}{=} \min \{\rho_{\tilde{\mathbf{a}}_k, k} - \rho_{\tilde{\mathbf{a}}_{k'}, k'} : (k, k') \in [K]^2, \rho_{\tilde{\mathbf{a}}_k, k} - \rho_{\tilde{\mathbf{a}}_{k'}, k'} > 0\}$ . The gap  $\delta$  is the minimum on a finite set, so  $\delta > 0$ .

By definition of  $\tilde{\pi}(t)$ ,  $\hat{\rho}_{\tilde{\mathbf{a}}_{\tilde{\pi}_1(t)}(t), \tilde{\pi}_1(t)}(t) \geq \hat{\rho}_{\tilde{\mathbf{a}}_{\tilde{\pi}_2(t)}(t), \tilde{\pi}_2(t)}(t) \geq \dots \geq \hat{\rho}_{\tilde{\mathbf{a}}_{\tilde{\pi}_K(t)}(t), \tilde{\pi}_K(t)}(t)$ , so,

$$\begin{aligned} \{t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{\pi}(t) \notin \Pi_{\rho}(\tilde{\mathbf{a}})\} &= \bigcup_{\tilde{\pi} \in \mathcal{P}_K^K} \bigcup_{k \in [K-1]} \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{\pi}(t) = \tilde{\pi}, \rho_{\tilde{\mathbf{a}}_{\tilde{\pi}_k}, \tilde{\pi}_k} < \rho_{\tilde{\mathbf{a}}_{\tilde{\pi}_{k+1}}, \tilde{\pi}_{k+1}} \right\} \\ &\subseteq \bigcup_{\tilde{\pi} \in \mathcal{P}_K^K} \bigcup_{k \in [K-1]} \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{\pi}(t) = \tilde{\pi}, \text{ or } |\hat{\rho}_{\tilde{\mathbf{a}}_{\tilde{\pi}_k}, \tilde{\pi}_k}(t) - \rho_{\tilde{\mathbf{a}}_{\tilde{\pi}_k}, \tilde{\pi}_k}| > \frac{\delta}{2} \right. \\ &\quad \left. \text{ or } |\hat{\rho}_{\tilde{\mathbf{a}}_{\tilde{\pi}_{k+1}}, \tilde{\pi}_{k+1}}(t) - \rho_{\tilde{\mathbf{a}}_{\tilde{\pi}_{k+1}}, \tilde{\pi}_{k+1}}| > \frac{\delta}{2} \right\} \\ &= \bigcup_{\tilde{\pi} \in \mathcal{P}_K^K} \bigcup_{k \in [K]} \Lambda_{\tilde{\pi}, k}, \end{aligned}$$

with  $\Lambda_{\tilde{\pi}, k} \stackrel{\text{def}}{=} \left\{ t : \tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{\pi}(t) = \tilde{\pi}, |\hat{\rho}_{\tilde{\mathbf{a}}_{\tilde{\pi}_k}, \tilde{\pi}_k}(t) - \rho_{\tilde{\mathbf{a}}_{\tilde{\pi}_k}, \tilde{\pi}_k}| > \frac{\delta}{2} \right\}$ , for any ranking of positions  $\tilde{\pi} \in \mathcal{P}_K^K$  and any rank  $k \in [K]$ .

Let  $\tilde{\pi} \in \mathcal{P}_K^K$  be a ranking of positions, and  $k \in [K]$  be a rank.  $\Lambda_{\tilde{\pi}, k} \subseteq \bigcup_{s \in \mathbb{N}} \Lambda_{\tilde{\pi}, k}(s)$ , with  $\Lambda_{\tilde{\pi}, k}(s) \stackrel{\text{def}}{=} \{t \in \Lambda_{\tilde{\pi}, k} : \tilde{T}_{\tilde{\mathbf{a}}}(t) = s\}$ .  $|\Lambda_{\tilde{\pi}, k}(s)| \leq 1$  as  $\tilde{T}_{\tilde{\mathbf{a}}}(t)$  increases for each  $t \in \Lambda_{\tilde{\pi}, k}$ . Note that for each  $s \in \mathbb{N}$  and  $n \in \Lambda_{\tilde{\pi}, k}(s)$ ,  $T_{\tilde{\mathbf{a}}_{\tilde{\pi}_k}, \tilde{\pi}_k}(n) \geq T_{\tilde{\mathbf{a}}}(n) \geq \tilde{T}_{\tilde{\mathbf{a}}}(n)/L = s/L$ . Then, by Lemma 4

$$\begin{aligned} \mathbb{E}[|\Lambda_{\tilde{\pi}, k}|] &= \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\{t \in \Lambda_{\tilde{\pi}, k}\}\right] \\ &= \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\left\{t \in \Lambda_{\tilde{\pi}, k}, |\hat{\rho}_{\tilde{\mathbf{a}}_{\tilde{\pi}_k}, \tilde{\pi}_k}(t) - \rho_{\tilde{\mathbf{a}}_{\tilde{\pi}_k}, \tilde{\pi}_k}| > \frac{\delta}{2}\right\}\right] \\ &\leq \frac{4L}{\delta^2} \end{aligned}$$

Hence,

$$\begin{aligned} \mathbb{E}\left[\sum_{t=1}^T \mathbb{1}\{\tilde{\mathbf{a}}(t) = \tilde{\mathbf{a}}, \tilde{\pi}(t) \notin \Pi_{\rho}(\tilde{\mathbf{a}})\}\right] &\leq \sum_{\tilde{\pi} \in \mathcal{P}_K^K} \sum_{k \in [K]} \mathbb{E}[\Lambda_{\tilde{\pi}, k}] \\ &\leq \frac{4LKK!}{\delta^2} \\ &= \mathcal{O}(LKK!), \end{aligned}$$

which concludes the proof. □

## G. KL-CombUCB and its Application to PBM Setting

In this section we first define the generic combinatorial semi-bandit algorithm KL-CombUCB and we compare two upper-bounds on its regret. Then, we present the application of KL-CombUCB to PBM setting and discuss its relation to GRAB.

**Algorithm 2** KL-ComUCB1 (generic version)

---

**Input:** set of elements  $E$ , set of arms  $\mathcal{A}$

$t \leftarrow 1$

**while**  $\{e \in E : T_e(t) = 0\} \neq \emptyset$  **do**

$\tilde{E} \leftarrow \{e \in E : T_e(t) = 0\}$

$\tilde{\mathcal{A}} \leftarrow \{\mathbf{a} \in \mathcal{A} : \mathbf{a} \cap \tilde{E} \neq \emptyset\}$

recommend  $\mathbf{a}(t) = \operatorname{argmax}_{\mathbf{a} \in \tilde{\mathcal{A}}} \sum_{e \in \mathbf{a}} b_e(t)$

observe the weights  $[w_e(t) : e \in \mathbf{a}]$

$t \leftarrow t + 1$

**end while**

$t_0 \leftarrow t$

**for**  $t = t_0, t_0 + 1, \dots$  **do**

recommend  $\mathbf{a}(t) = \operatorname{argmax}_{\mathbf{a} \in \mathcal{A}} \sum_{e \in \mathbf{a}} b_e(t)$

observe the weights  $[w_e(t) : e \in \mathbf{a}]$

**end for**

---

**G.1. KL-CombUCB for Generic Setting**

CombUCB1 (Kveton et al., 2015b) is a bandit algorithm handling the following combinatorial setting. Let  $E$  be a set of elements and  $\mathcal{A} \subseteq \{0, 1\}^E$  be a set of arms, where each arm  $\mathbf{a}$  is a subset of  $E$ . Following the terminology used in (Kveton et al., 2015b),  $E$  is the *ground set* and  $\mathcal{A}$  the *feasible set*. At each iteration, the bandit algorithm chooses a subset of elements  $\mathbf{a} \in \mathcal{A}$  and receives the reward  $\sum_{e \in \mathbf{a}} w_e$ , where  $\mathbf{w}$  is an independent draw of a distribution  $\nu$  on  $[0, 1]^E$ . Given these assumptions, CombUCB1 chooses an arm  $\mathbf{a}(t)$  at each iteration, aiming at minimizing the total regret defined as usual.

We denote  $\rho_e \stackrel{\text{def}}{=} \mathbb{E}_{\mathbf{w} \sim \nu} [w_e]$  the expected reward associated to element  $e$ ,  $\mu_{\mathbf{a}} \stackrel{\text{def}}{=} \mathbb{E}_{\mathbf{w} \sim \nu} [\sum_{e \in \mathbf{a}} w_e] = \sum_{e \in \mathbf{a}} \rho_e$  the expected reward when choosing the arm  $\mathbf{a} \in \mathcal{A}$ , and  $\mu^* \stackrel{\text{def}}{=} \max_{\mathbf{a} \in \mathcal{A}} \mu_{\mathbf{a}}$  the best expected reward. We also denote  $\Delta_{\mathbf{a}} \stackrel{\text{def}}{=} \mu^* - \mu_{\mathbf{a}}$  the gap between the best expected reward and the reward of an arm  $\mathbf{a}$ , and  $\Delta_{\min} \stackrel{\text{def}}{=} \min_{\mathbf{a} \in \mathcal{A} : \Delta_{\mathbf{a}} > 0} \Delta_{\mathbf{a}}$  the smallest gap of a suboptimal arm. Finally,  $K \stackrel{\text{def}}{=} \max_{\mathbf{a} \in \mathcal{A}} |\mathbf{a}|$  denotes the maximum size of an arm (meaning the maximum number of chosen elements),  $K_{\mathbf{a}} \stackrel{\text{def}}{=} \min_{\mathbf{a}^* \in \mathcal{A} : \mu_{\mathbf{a}^*} = \mu^*} |\mathbf{a} \setminus \mathbf{a}^*|$  is the smallest number of elements to remove from  $\mathbf{a}$  to get an optimal arm, and  $K_{\max} \stackrel{\text{def}}{=} \max_{\mathbf{a} \in \mathcal{A} : \mu_{\mathbf{a}} \neq \mu^*} K_{\mathbf{a}}$  is its larger value.

In our paper, we use the Kullback-Leibler variation of CombUCB1 which chooses the arm based on the index  $b_e(t)$  (defined hereafter) instead of the usual confidence upper-bound derived from the Hoeffding's inequality. The corresponding algorithm (KL-CombUCB) also assumes that the weight-vector  $\mathbf{w}(t)$  is in  $\{0, 1\}^E$ . KL-CombUCB is depicted by Algorithm 2 which uses the following notations. At each iteration  $t$ , we denote

$$\hat{\rho}_e(t) \stackrel{\text{def}}{=} \frac{1}{T_e(t)} \sum_{s=1}^{t-1} \mathbb{1}\{e \in \mathbf{a}(s)\} w_e(s)$$

the average number of clicks obtained by the element  $e$ , where

$$T_e(t) \stackrel{\text{def}}{=} \sum_{s=1}^{t-1} \mathbb{1}\{e \in \mathbf{a}(s)\}$$

is the number of times element  $e$  has been selected;  $\hat{\rho}_e(t) \stackrel{\text{def}}{=} 0$  when  $T_e(t) = 0$ . The statistics  $\hat{\rho}_e(t)$  are paired with their respective *indices*

$$b_e(t) \stackrel{\text{def}}{=} f(\hat{\rho}_e(t), T_e(t), t),$$

where  $f(\hat{\rho}, s, t)$  stands for

$$\sup\{p \in [\hat{\rho}, 1] : s \times \text{kl}(\hat{\rho}, p) \leq \log(t) + 3 \log(\log(t))\},$$

---

**Algorithm 3** KL-ComUCB1 (applied to PBM)
 

---

**Input:** number of items  $L$ , number of positions  $K$   
**for**  $t = 1, 2, \dots, L$  **do**  
     recommend  $\mathbf{a}(t) = (((t-1)\%L) + 1, (t\%L) + 1, \dots, ((t+K-2)\%L) + 1)$   
     observe the clicks-vector  $\mathbf{c}(t)$   
**end for**  
**for**  $t = L+1, L+2, \dots$  **do**  
     recommend  $\mathbf{a}(t) = \underset{\mathbf{a} \in \mathcal{P}_K^L}{\operatorname{argmax}} \sum_{k=1}^K b_{a_k, k}(t)$   
     observe the clicks-vector  $\mathbf{c}(t)$   
**end for**

---

with

$$\operatorname{kl}(p, q) \stackrel{\text{def}}{=} p \log \left( \frac{p}{q} \right) + (1-p) \log \left( \frac{1-p}{1-q} \right)$$

the *Kullback-Leibler divergence* from a Bernoulli distribution of mean  $p$  to a Bernoulli distribution of mean  $q$ ;  $f(\hat{\rho}, s, t) \stackrel{\text{def}}{=} 1$  when  $\hat{\rho} = 1$ ,  $s = 0$ , or  $t = 0$ .

Kveton et al. prove that the regret of CombUCB1 is upper-bounded by  $\mathcal{O}(|E|K/\Delta_{\min} \log T)$ , and a similar proof would lead to the same upper-bound for KL-CombUCB. In our paper we prove in Theorem 2 a completely different regret upper-bound for KL-CombUCB:  $\mathcal{O}(|\mathcal{A}|K_{\max}^2/\Delta_{\min} \log T)$ . For most combinatorial bandit settings, this new bound is useless since  $|\mathcal{A}| \gg |E|$ , and  $K_{\max} \approx K$ . However, the analysis of GRAB involves an application of KL-CombUCB to a setting where the new bound is smaller than the standard one as  $|\mathcal{A}| = |E| - 1$  and  $K_{\max} = 2$ .

## G.2. KL-CombUCB Applied to PBM Setting

In the experiments (Section 6), we apply KL-CombUCB to PBM bandit setting by choosing the *ground set*  $E = [L] \times [K]$ , the *feasible set*  $\Theta = \{(a_k, k) : k \in [K]\} : \mathbf{a} \in \mathcal{P}_K^L$ , and the *expected weights*  $\rho_{(i,k)} = \theta_i \kappa_k$  for any “element”  $(i, k) \in E$ . Note that the observed weights of the generic setting correspond to the clicks-vector in the PBM setting.

The corresponding algorithm, depicted by Algorithm 3, recommends at each iteration  $t$  the best permutation given the indices  $b_{i,k}(t)$  defined for GRAB. This optimization problem is a *linear sum assignment problem* which is solvable in  $\mathcal{O}(K^2(L + \log K))$  time (Ramshaw & Tarjan, 2012). Note the close relationship with GRAB:

- both algorithms solve a linear sum assignment problem, they only differ from the metric to optimize:  $\sum_{k=1}^K \hat{\rho}_{a_k, k}(t)$  for GRAB vs.  $\sum_{k=1}^K b_{a_k, k}(t)$  for KL-CombUCB;
- both algorithms recommend the best permutation  $\mathbf{a}$  regarding  $\sum_{k=1}^K b_{a_k, k}(t)$ , they only differ from the considered set of permutations:  $\{\tilde{\mathbf{a}}(t)\} \cup \mathcal{N}_{\tilde{\pi}(t)}(\tilde{\mathbf{a}}(t))$  for GRAB vs.  $\mathcal{P}_K^L$  for KL-CombUCB.

By considering a larger set of permutations, KL-ComUCB1 suffers a  $\mathcal{O}(LK^2/\Delta_{\min} \log T)$  regret (by applying (Kveton et al., 2015b) bound), which is higher than the upper-bound on the regret of GRAB by a factor  $K^2$ .

## H. S-GRAB: OSUB on a Static Graph

The algorithm S-GRAB, depicted in Algorithm 4, is similar to GRAB except that it explores a static graph  $G = (E, V)$  defined by

$$\begin{aligned}
 V &\stackrel{\text{def}}{=} \mathcal{P}_K^L, \\
 E &\stackrel{\text{def}}{=} \{(\mathbf{a}, \mathbf{a} \circ (k, k')) : k, k' \in [K]^2, k > k'\} \cup \{(\mathbf{a}, \mathbf{a}[k := i]) : k \in [K], i \in [L] \setminus \mathbf{a}([K])\}.
 \end{aligned}$$

This graph is chosen to ensure that with PBM setting any sub-optimal recommendation has a strictly better recommendation

**Algorithm 4** S-GRAB: Static Graph for unimodal RAnking Bandit

---

**Input:** number of items  $L$ , number of positions  $K$

$$\gamma \leftarrow K(2L - K - 1)/2$$

**for**  $t = 1, 2, \dots$  **do**

$$\tilde{\mathbf{a}}(t) \leftarrow \operatorname{argmax}_{\mathbf{a} \in \mathcal{P}_K^L} \sum_{k=1}^K \hat{\rho}_{a_k, k}(t)$$

$$\text{recommend } \mathbf{a}(t) = \begin{cases} \tilde{\mathbf{a}}(t) & , \text{ if } \frac{\tilde{T}_{\tilde{\mathbf{a}}(t)}(t)}{\gamma+1} \in \mathbb{N}, \\ \operatorname{argmax}_{\mathbf{a} \in \{\tilde{\mathbf{a}}(t)\} \cup \mathcal{N}_G(\tilde{\mathbf{a}}(t))} \sum_{k=1}^K b_{a_k, k}(t) & , \text{ otherwise} \end{cases}$$

where  $\mathcal{N}_G(\mathbf{a}) = \{\mathbf{a} \circ (k, k') : k, k' \in [K]^2, k > k'\} \cup \{\mathbf{a}[k := i] : k \in [K], i \in [L] \setminus \mathbf{a}([K])\}$   
 observe the clicks vector  $\mathbf{c}(t)$

**end for**

---

in its neighborhood given  $G$ . This graph is fixed and does not require the knowledge of a mapping  $\mathcal{P}$ , but its degree is also about  $K$  times larger than the degree of the graphs handled by GRAB.

As for GRAB, any recommendation in the neighborhood of the leader given  $G$  differs with the leader at, at most two positions. Therefore a proof similar to the one of Theorem 1 ensures that S-GRAB's regret is upper-bounded by  $\mathcal{O}(LK/\Delta_{\min} \log T)$ . This regret upper-bound is higher than GRAB's one by a factor  $K$  due to the larger size of the considered neighborhoods. However, this regret remains smaller than KL-CombUCB's one by a factor  $K$  thanks to the bounded number of differences between the leader and the arm played.