
Online Learning to Rank in Stochastic Click Models

Masrour Zoghi¹ Tomas Tunys² Mohammad Ghavamzadeh³ Branislav Kveton⁴ Csaba Szepesvari⁵
Zheng Wen⁴

Abstract

Online learning to rank is a core problem in information retrieval and machine learning. Many provably efficient algorithms have been recently proposed for this problem in specific click models. The click model is a model of how the user interacts with a list of documents. Though these results are significant, their impact on practice is limited, because all proposed algorithms are designed for specific click models and lack convergence guarantees in other models. In this work, we propose `BatchRank`, the first online learning to rank algorithm for a broad class of click models. The class encompasses two most fundamental click models, the cascade and position-based models. We derive a gap-dependent upper bound on the T -step regret of `BatchRank` and evaluate it on a range of web search queries. We observe that `BatchRank` outperforms ranked bandits and is more robust than `CascadeKL-UCB`, an existing algorithm for the cascade model.

1. Introduction

Learning to rank (LTR) is a core problem in information retrieval (Liu, 2011) and machine learning; with numerous applications in web search, recommender systems and ad placement. The goal of LTR is to present a list of K documents out of L that maximizes the satisfaction of the user. This problem has been traditionally solved by training supervised learning models on manually annotated relevance judgments. However, strong evidence suggests (Agichtein

et al., 2006; Zoghi et al., 2016) that the feedback of users, that is clicks, can lead to major improvements over supervised LTR methods. In addition, billions of users interact daily with commercial LTR systems, and it is finally feasible to interactively and adaptive maximize the satisfaction of these users from clicks.

These observations motivated numerous papers on online LTR methods, which utilize user feedback to improve the quality of ranked lists. These methods can be divided into two groups: learning the best ranker in a family of rankers (Yue & Joachims, 2009; Hofmann et al., 2013), and learning the best list under some model of user interaction with the list (Radlinski et al., 2008a; Slivkins et al., 2013), such as a *click model* (Chuklin et al., 2015). The click model is a stochastic model of how the user examines and clicks on a list of documents. In this work, we focus on online LTR in click models and address a shortcoming of all past work on this topic.

More precisely, many algorithms have been proposed and analyzed for finding the optimal ranked list in the cascade model (CM) (Kveton et al., 2015a; Combes et al., 2015; Kveton et al., 2015b; Zong et al., 2016; Li et al., 2016), the dependent-click model (DCM) (Katariya et al., 2016), and the position-based model (PBM) (Lagree et al., 2016). The problem is that if the user interacts with ranked lists using a different click model, the theoretical guarantees cease to hold. Then, as we show empirically, these algorithms may converge to suboptimal solutions. This is a grave issue because it is well known that no single click model captures the behavior of an entire population of users (Grotov et al., 2015). Therefore, it is critical to develop efficient learning algorithms for multiple click models, which is the aim of this paper.

We make the following contributions:

- We propose *stochastic click bandits*, a learning framework for maximizing the expected number of clicks in online LTR in a broad class of click models, which includes both the PBM (Richardson et al., 2007) and CM (Craswell et al., 2008).
- We propose the first algorithm, `BatchRank`, that is guaranteed to learn the optimal solution in a diverse class of click models. This is of a great practical significance, as

¹Independent Researcher, Vancouver, BC, Canada (Part of this work was done during an internship at Adobe Research) ²Czech Technical University, Prague, Czech Republic ³DeepMind, Mountain View, CA, USA (This work was done while the author was at Adobe Research) ⁴Adobe Research, San Jose, CA, USA ⁵University of Alberta, Edmonton, AB, Canada. Correspondence to: Branislav Kveton <kveton@adobe.com>, Masrour Zoghi <masrour@zoghi.org>.

it is often difficult or impossible to guess the underlying click model in advance.

- We prove a gap-dependent upper bound on the regret of BatchRank that scales well with all quantities of interest. The key step in our analysis is a KL scaling lemma (Section 5.4), which should be of a broader interest.
- We evaluate BatchRank on both CM and PBM queries. Our results show that BatchRank performs significantly better than RankedExp3 (Radlinski et al., 2008a), an adversarial online LTR algorithm; and is more robust than CascadeKL-UCB (Kveton et al., 2015a), an optimal online LTR algorithm for the CM.

We define $[n] = \{1, \dots, n\}$. For any sets A and B , we denote by A^B the set of all vectors whose entries are indexed by B and take values from A . We use boldface letters to denote important random variables.

2. Background

This section reviews two fundamental *click models* (Chuklin et al., 2015), models of how users click on an ordered list of K documents. The universe of all documents is represented by *ground set* $\mathcal{D} = [L]$ and we refer to the documents in \mathcal{D} as *items*. The user is presented a *ranked list*, an ordered list of K documents out of L . We denote this list by $\mathcal{R} = (d_1, \dots, d_K) \in \Pi_K(\mathcal{D})$, where $\Pi_K(\mathcal{D}) \subset \mathcal{D}^K$ is the set of all K -tuples with distinct elements from \mathcal{D} and d_k is the k -th item in \mathcal{R} . We assume that the click model is parameterized by L *item-dependent attraction probabilities* $\alpha \in [0, 1]^L$, where $\alpha(d)$ is the probability that item d is *attractive*. The items attract the user independently. For simplicity and without loss of generality, we assume that $\alpha(1) \geq \dots \geq \alpha(L)$. The reviewed models differ in how the user examines items, which leads to clicks.

2.1. Position-Based Model

The *position-based model (PBM)* (Richardson et al., 2007) is a model where the probability of clicking on an item depends on both its identity and position. Therefore, in addition to item-dependent attraction probabilities, the PBM is parameterized by K *position-dependent examination probabilities* $\chi \in [0, 1]^K$, where $\chi(k)$ is the examination probability of position k .

The user interacts with a list of items $\mathcal{R} = (d_1, \dots, d_K)$ as follows. The user *examines* position $k \in [K]$ with probability $\chi(k)$ and then *clicks* on item d_k at that position with probability $\alpha(d_k)$. Thus, the expected number of clicks on list \mathcal{R} is

$$r(\mathcal{R}) = \sum_{k=1}^K \chi(k) \alpha(d_k).$$

In practice, it is often observed that $\chi(1) \geq \dots \geq \chi(K)$

(Chuklin et al., 2015), and we adopt this assumption in this work. Under this assumption, the above function is maximized by the list of K most attractive items

$$\mathcal{R}^* = (1, \dots, K), \quad (1)$$

where the k -th most attractive item is placed at position k .

In this paper, we focus on the objective of maximizing the number of clicks. We note that the satisfaction of the user may not increase with the number of clicks, and that other objectives have been proposed in the literature (Radlinski et al., 2008b). The shortcoming of all of these objectives is that none directly measure the satisfaction of the user.

2.2. Cascade Model

In the *cascade model (CM)* (Craswell et al., 2008), the user scans a list of items $\mathcal{R} = (d_1, \dots, d_K)$ from the first item d_1 to the last d_K . If item d_k is *attractive*, the user *clicks* on it and does not examine the remaining items. If item d_k is not attractive, the user *examines* item d_{k+1} . The first item d_1 is examined with probability one.

From the definition of the model, the probability that item d_k is examined is equal to the probability that none of the first $k-1$ items are attractive. Since each item attracts the user independently, this probability is

$$\chi(\mathcal{R}, k) = \prod_{i=1}^{k-1} (1 - \alpha(d_i)). \quad (2)$$

The expected number of clicks on list \mathcal{R} is at most 1, and is equal to the probability of observing any click,

$$r(\mathcal{R}) = \sum_{k=1}^K \chi(\mathcal{R}, k) \alpha(d_k) = 1 - \prod_{k=1}^K (1 - \alpha(d_k)).$$

This function is maximized by the list of K most attractive items \mathcal{R}^* in (1), though any permutation of $[K]$ would be optimal in the CM. Note that the list \mathcal{R}^* is optimal in both the PBM and CM.

3. Online Learning to Rank in Click Models

The PBM and CM (Section 2) are similar in many aspects. First, both models are parameterized by L item-dependent attraction probabilities. The items attract the user independently. Second, the probability of clicking on the item is a product of its attraction probability, which depends on the identity of the item; and the examination probability of its position, which is independent of the identity of the item. Finally, the optimal solution in both models is the list of K most attractive items \mathcal{R}^* in (1), where the k -th most attractive item is placed at position k .

This suggests that the optimal solution in both models can be learned by a single learning algorithm, which does not know the underlying model. We propose this algorithm in Section 4. Before we discuss it, we formalize our learning problem as a *multi-armed bandit* (Auer et al., 2002; Lai & Robbins, 1985).

3.1. Stochastic Click Bandit

We refer to our learning problem as a *stochastic click bandit*. An instance of this problem is a tuple (K, L, P_α, P_χ) , where K is the number of positions, L is the number of items, P_α is a distribution over binary vectors $\{0, 1\}^L$, and P_χ is a distribution over binary matrices $\{0, 1\}^{\Pi_K(\mathcal{D}) \times K}$.

The learning agent interacts with our problem as follows. Let $(\mathbf{A}_t, \mathbf{X}_t)_{t=1}^T$ be T i.i.d. random variables drawn from $P_\alpha \otimes P_\chi$, where $\mathbf{A}_t \in \{0, 1\}^L$ and $\mathbf{A}_t(d)$ is the *attraction indicator* of item d at time t ; and $\mathbf{X}_t \in \{0, 1\}^{\Pi_K(\mathcal{D}) \times K}$ and $\mathbf{X}_t(\mathcal{R}, k)$ is the *examination indicator* of position k in list $\mathcal{R} \in \Pi_K(\mathcal{D})$ at time t . At time t , the agent chooses a list $\mathcal{R}_t = (d_1^t, \dots, d_K^t) \in \Pi_K(\mathcal{D})$, which depends on past observations of the agent, and then observes *clicks*. These clicks are a function of \mathcal{R}_t , \mathbf{A}_t , and \mathbf{X}_t . Let $\mathbf{c}_t \in \{0, 1\}^K$ be the vector of *click indicators* on all positions at time t . Then

$$\mathbf{c}_t(k) = \mathbf{X}_t(\mathcal{R}_t, k) \mathbf{A}_t(d_k^t)$$

for any $k \in [K]$, the item at position k is clicked only if both $\mathbf{X}_t(\mathcal{R}_t, k) = 1$ and $\mathbf{A}_t(d_k^t) = 1$.

The goal of the learning agent is to maximize the number of clicks. Therefore, the number of clicks at time t is the *reward* of the agent at time t . We define it as

$$\mathbf{r}_t = \sum_{k=1}^K \mathbf{c}_t(k) = r(\mathcal{R}_t, \mathbf{A}_t, \mathbf{X}_t), \quad (3)$$

where $r : \Pi_K(\mathcal{D}) \times [0, 1]^L \times [0, 1]^{\Pi_K(\mathcal{D}) \times K} \rightarrow [0, K]$ is a *reward function*, which we define for any $\mathcal{R} \in \Pi_K(\mathcal{D})$, $A \in [0, 1]^L$, and $X \in [0, 1]^{\Pi_K(\mathcal{D}) \times K}$ as

$$r(\mathcal{R}, A, X) = \sum_{k=1}^K X(\mathcal{R}, k) A(d_k).$$

We adopt the same independence assumptions as in Section 2. In particular, we assume that items attract the user independently.

Assumption 1. For any $A \in \{0, 1\}^L$,

$$P(\mathbf{A}_t = A) = \prod_{d \in \mathcal{D}} \text{Ber}(A(d); \alpha(d)),$$

where $\text{Ber}(\cdot; \theta)$ denotes the probability mass function of a Bernoulli distribution with mean $\theta \in [0, 1]$, which we define as $\text{Ber}(y; \theta) = \theta^y (1 - \theta)^{1-y}$ for any $y \in \{0, 1\}$.

Moreover, we assume that the attraction of any item is independent of the examination of its position.

Assumption 2. For any list $\mathcal{R} \in \Pi_K(\mathcal{D})$ and position k ,

$$\mathbb{E}[\mathbf{c}_t(k) \mid \mathcal{R}_t = \mathcal{R}] = \chi(\mathcal{R}, k) \alpha(d_k),$$

where $\chi \in [0, 1]^{\Pi_K(\mathcal{D}) \times K}$ and $\chi(\mathcal{R}, k) = \mathbb{E}[\mathbf{X}_t(\mathcal{R}, k)]$ is the examination probability of position k in list \mathcal{R} .

We do not make any independence assumptions among the entries of \mathbf{X}_t , and on other interactions of \mathbf{A}_t and \mathbf{X}_t .

From our independence assumptions and the definition of the reward in (3), the *expected reward* of list \mathcal{R} is

$$\mathbb{E}[r(\mathcal{R}, \mathbf{A}_t, \mathbf{X}_t)] = \sum_{k=1}^K \chi(\mathcal{R}, k) \alpha(d_k) = r(\mathcal{R}, \alpha, \chi).$$

We evaluate the performance of a learning agent by its *expected cumulative regret*

$$R(T) = \mathbb{E} \left[\sum_{t=1}^T R(\mathcal{R}_t, \mathbf{A}_t, \mathbf{X}_t) \right],$$

where $R(\mathcal{R}_t, \mathbf{A}_t, \mathbf{X}_t) = r(\mathcal{R}^*, \mathbf{A}_t, \mathbf{X}_t) - r(\mathcal{R}_t, \mathbf{A}_t, \mathbf{X}_t)$ is the *instantaneous regret* of the agent at time t and

$$\mathcal{R}^* = \arg \max_{\mathcal{R} \in \Pi_K(\mathcal{D})} r(\mathcal{R}, \alpha, \chi)$$

is the *optimal list* of items, the list that maximizes the expected reward. To simplify exposition, we assume that the optimal solution, as a set, is unique.

3.2. Position Bandit

The learning variant of the PBM in Section 2.1 can be formulated in our setting when

$$\forall \mathcal{R}, \mathcal{R}' \in \Pi_K(\mathcal{D}) : \mathbf{X}_t(\mathcal{R}, k) = \mathbf{X}_t(\mathcal{R}', k) \quad (4)$$

at any position $k \in [K]$. Under this assumption, the probability of clicking on item d_k^t at time t is

$$\mathbb{E}[\mathbf{c}_t(k) \mid \mathcal{R}_t] = \chi(k) \alpha(d_k^t),$$

where $\chi(k)$ is defined in Section 2.1. The expected reward of list \mathcal{R}_t at time t is

$$\mathbb{E}[\mathbf{r}_t \mid \mathcal{R}_t] = \sum_{k=1}^K \chi(k) \alpha(d_k^t).$$

3.3. Cascading Bandit

The learning variant of the CM in Section 2.2 can be formulated in our setting when

$$\mathbf{X}_t(\mathcal{R}, k) = \prod_{i=1}^{k-1} (1 - \mathbf{A}_t(d_i)) \quad (5)$$

for any list $\mathcal{R} \in \Pi_K(\mathcal{D})$ and position $k \in [K]$. Under this assumption, the probability of clicking on item d_k^t at time t is

$$\mathbb{E}[\mathbf{c}_t(k) \mid \mathcal{R}_t] = \left[\prod_{i=1}^{k-1} (1 - \alpha(d_i^t)) \right] \alpha(d_k^t).$$

The expected reward of list \mathcal{R}_t at time t is

$$\mathbb{E}[r_t \mid \mathcal{R}_t] = \sum_{k=1}^K \left[\prod_{i=1}^{k-1} (1 - \alpha(d_i^t)) \right] \alpha(d_k^t).$$

3.4. Additional Assumptions

The above assumptions are not sufficient to guarantee that the optimal list \mathcal{R}^* in (1) is learnable. Therefore, we make four additional assumptions, which are quite natural.

Assumption 3 (Order-independent examination). *For any lists $\mathcal{R} \in \Pi_K(\mathcal{D})$ and $\mathcal{R}' \in \Pi_K(\mathcal{D})$, and position $k \in [K]$ such that $d_k = d'_k$ and $\{d_1, \dots, d_{k-1}\} = \{d'_1, \dots, d'_{k-1}\}$, $\mathbf{X}_t(\mathcal{R}, k) = \mathbf{X}_t(\mathcal{R}', k)$.*

The above assumption says that the examination indicator $\mathbf{X}_t(\mathcal{R}, k)$ only depends on the identities of d_1, \dots, d_{k-1} . Both the CM and PBM satisfy this assumption, which can be validated from (4) and (5).

Assumption 4 (Decreasing examination). *For any list $\mathcal{R} \in \Pi_K(\mathcal{D})$ and positions $1 \leq i \leq j \leq K$, $\chi(\mathcal{R}, i) \geq \chi(\mathcal{R}, j)$.*

The above assumption says that a lower position cannot be examined more than a higher position, in any list \mathcal{R} . Both the CM and PBM satisfy this assumption.

Assumption 5 (Correct examination scaling). *For any list $\mathcal{R} \in \Pi_K(\mathcal{D})$ and positions $1 \leq i \leq j \leq K$, let $\alpha(d_i) \leq \alpha(d_j)$ and $\mathcal{R}' \in \Pi_K(\mathcal{D})$ be the same list as \mathcal{R} except that d_i and d_j are exchanged. Then $\chi(\mathcal{R}, j) \geq \chi(\mathcal{R}', j)$.*

The above assumption says that the examination probability of a position cannot increase if the item at that position is swapped for a less-attractive higher-ranked item, in any list \mathcal{R} . Both the CM and PBM satisfy this assumption. In the CM, the inequality follows directly from the definition of examination in (2). In the PBM, $\chi(\mathcal{R}, j) = \chi(\mathcal{R}', j)$.

Assumption 6 (Optimal examination). *For any list $\mathcal{R} \in \Pi_K(\mathcal{D})$ and position $k \in [K]$, $\chi(\mathcal{R}, k) \geq \chi(\mathcal{R}^*, k)$.*

This assumption says that any position k is least examined if the first $k-1$ items are optimal. Both the CM and PBM satisfy this assumption. In the CM, the inequality follows from the definition of examination in (2). In the PBM, we have that $\chi(\mathcal{R}, k) = \chi(\mathcal{R}^*, k)$.

4. Algorithm BatchRank

The design of BatchRank (Algorithm 1) builds on two key ideas. First, we *randomize* the placement of items to avoid

Algorithm 1 BatchRank

```

1: // Initialization
2: for  $b = 1, \dots, 2K$  do
3:   for  $\ell = 0, \dots, T-1$  do
4:     for all  $d \in \mathcal{D}$  do
5:        $\mathbf{c}_{b,\ell}(d) \leftarrow 0$ ,  $\mathbf{n}_{b,\ell}(d) \leftarrow 0$ 
6:    $\mathcal{A} \leftarrow \{1\}$ ,  $b_{\max} \leftarrow 1$ 
7:    $\mathbf{I}_1 \leftarrow (1, K)$ ,  $\mathbf{B}_{1,0} \leftarrow \mathcal{D}$ ,  $\ell_1 \leftarrow 0$ 
8:   for  $t = 1, \dots, T$  do
9:     for all  $b \in \mathcal{A}$  do
10:      DisplayBatch( $b, t$ )
11:    for all  $b \in \mathcal{A}$  do
12:      CollectClicks( $b, t$ )
13:    for all  $b \in \mathcal{A}$  do
14:      UpdateBatch( $b, t$ )

```

biases due to the click model. Second, we *divide and conquer*; recursively divide the batches of items into more and less attractive items. The result is a sorted list of K items, where the k -th most attractive item is placed at position k .

BatchRank explores items in batches, which are indexed by integers $b > 0$. A *batch* b is associated with the initial set of items $\mathbf{B}_{b,0} \subseteq \mathcal{D}$ and a range of *positions* $\mathbf{I}_b \in [K]^2$, where $\mathbf{I}_b(1)$ is the *highest position* in batch b , $\mathbf{I}_b(2)$ is the *lowest position* in batch b , and $\text{len}(b) = \mathbf{I}_b(2) - \mathbf{I}_b(1) + 1$ is *number of positions* in batch b . The batch is explored in *stages*, which we index by integers $\ell > 0$. The remaining items in stage ℓ of batch b are $\mathbf{B}_{b,\ell} \subseteq \mathbf{B}_{b,0}$. The lengths of the stages quadruple. More precisely, any item $d \in \mathbf{B}_{b,\ell}$ in stage ℓ is explored n_ℓ times, where $n_\ell = \lceil 16\tilde{\Delta}_\ell^{-2} \log T \rceil$ and $\tilde{\Delta}_\ell = 2^{-\ell}$. The current stage of batch b is ℓ_b .

Method DisplayBatch (Algorithm 2) explores batches as follows. In stage ℓ of batch b , we randomly choose $\text{len}(b)$ least observed items in $\mathbf{B}_{b,\ell}$ and display them at randomly chosen positions in \mathbf{I}_b . If the number of these items is less than $\text{len}(b)$, we mix them with randomly chosen more observed items, which are not explored. This exploration has two notable properties. First, it is uniform in the sense that no item in $\mathbf{B}_{b,\ell}$ is explored more than once than any other item in $\mathbf{B}_{b,\ell}$. Second, any item in $\mathbf{B}_{b,\ell}$ appears in any list over $\mathbf{B}_{b,\ell}$ with that item with the same probability. This is critical to avoid biases due to click models.

Method CollectClicks (Algorithm 3) collects feedback. We denote the *number of observations* of item d in stage ℓ of batch b by $\mathbf{n}_{b,\ell}(d)$ and the *number of clicks* on that item by $\mathbf{c}_{b,\ell}(d)$. At the end of the stage, all items $d \in \mathbf{B}_{b,\ell}$ are observed exactly n_ℓ times and we estimate the probability of clicking on item d as

$$\hat{\mathbf{c}}_{b,\ell}(d) = \mathbf{c}_{b,\ell}(d) / n_\ell. \quad (6)$$

Algorithm 2 DisplayBatch

- 1: **Input:** batch index b , time t
 - 2: $\ell \leftarrow \ell_b$, $n_{\min} \leftarrow \min_{d \in \mathcal{B}_{b,\ell}} n_{b,\ell}(d)$
 - 3: Let $d_1, \dots, d_{|\mathcal{B}_{b,\ell}|}$ be a random permutation of items $\mathcal{B}_{b,\ell}$ such that $n_{b,\ell}(d_1) \leq \dots \leq n_{b,\ell}(d_{|\mathcal{B}_{b,\ell}|})$
 - 4: Let $\pi \in \Pi_{\text{len}(b)}([\text{len}(b)])$ be a random permutation of position assignments
 - 5: **for** $k = I_b(1), \dots, I_b(2)$ **do**
 - 6: $d_k^t \leftarrow d_{\pi(k-I_b(1)+1)}$
-

Algorithm 3 CollectClicks

- 1: **Input:** batch index b , time t
 - 2: $\ell \leftarrow \ell_b$, $n_{\min} \leftarrow \min_{d \in \mathcal{B}_{b,\ell}} n_{b,\ell}(d)$
 - 3: **for** $k = I_b(1), \dots, I_b(2)$ **do**
 - 4: **if** $n_{b,\ell}(d_k^t) = n_{\min}$ **then**
 - 5: $c_{b,\ell}(d_k^t) \leftarrow c_{b,\ell}(d_k^t) + c_t(k)$
 - 6: $n_{b,\ell}(d_k^t) \leftarrow n_{b,\ell}(d_k^t) + 1$
-

Method UpdateBatch (Algorithm 4) updates batches and has three main parts. First, we compute KL-UCB upper and lower confidence bounds (Garivier & Cappe, 2011) for all items $d \in \mathcal{B}_{b,\ell}$ (lines 5–6),

$$U_{b,\ell}(d) \leftarrow \arg \max_{q \in [\hat{c}_{b,\ell}(d), 1]} \{n_\ell D_{\text{KL}}(\hat{c}_{b,\ell}(d) \| q) \leq \delta_T\},$$

$$L_{b,\ell}(d) \leftarrow \arg \min_{q \in [0, \hat{c}_{b,\ell}(d)]} \{n_\ell D_{\text{KL}}(\hat{c}_{b,\ell}(d) \| q) \leq \delta_T\},$$

where $D_{\text{KL}}(p \| q)$ denotes the *Kullback-Leibler divergence* between Bernoulli random variables with means p and q , and $\delta_T = \log T + 3 \log \log T$. Then we test whether batch b can be safely divided into s more attractive items and the rest (lines 7–15). If it can, we split the batch into two new batches (lines 21–27). The first batch contains s items and is over positions $I_b(1), \dots, I_b(1) + s - 1$; and the second batch contains the remaining items and is over the remaining positions. The stage indices of new batches are initialized to 0. If the batch can be divided at multiple positions s , we choose the highest s . If the batch is not divided, we eliminate items that cannot be at position $I_b(2)$ or higher with a high probability (lines 17–19).

The set of active batches is denoted by \mathcal{A} , and we explore and update these batches in parallel. The highest index of the latest added batch is b_{\max} . Note that $b_{\max} \leq 2K$, because any batch with at least two items is split at a unique position into two batches. BatchRank is initialized with a single batch over all positions and items (lines 6–7).

Note that by the design of UpdateBatch, the following invariants hold. First, the positions of active batches \mathcal{A} are a partition of $[K]$ at any time t . Second, any batch contains at least as many items as is the number of the positions in

Algorithm 4 UpdateBatch

- 1: **Input:** batch index b , time t
 - 2: // End-of-stage elimination
 - 3: $\ell \leftarrow \ell_b$
 - 4: **if** $\min_{d \in \mathcal{B}_{b,\ell}} n_{b,\ell}(d) = n_\ell$ **then**
 - 5: **for all** $d \in \mathcal{B}_{b,\ell}$ **do**
 - 6: Compute $U_{b,\ell}(d)$ and $L_{b,\ell}(d)$
 - 7: Let $d_1, \dots, d_{|\mathcal{B}_{b,\ell}|}$ be any permutation of items $\mathcal{B}_{b,\ell}$ such that $L_{b,\ell}(d_1) \geq \dots \geq L_{b,\ell}(d_{|\mathcal{B}_{b,\ell}|})$
 - 8: **for** $k = 1, \dots, \text{len}(b)$ **do**
 - 9: $B_k^+ \leftarrow \{d_1, \dots, d_k\}$
 - 10: $B_k^- \leftarrow \mathcal{B}_{b,\ell} \setminus B_k^+$
 - 11: // Find a split at the position with the highest index
 - 12: $s \leftarrow 0$
 - 13: **for** $k = 1, \dots, \text{len}(b) - 1$ **do**
 - 14: **if** $L_{b,\ell}(d_k) > \max_{d \in B_k^-} U_{b,\ell}(d)$ **then**
 - 15: $s \leftarrow k$
 - 16: **if** ($s = 0$) and ($|\mathcal{B}_{b,\ell}| > \text{len}(b)$) **then**
 - 17: // Next elimination stage
 - 18: $\mathcal{B}_{b,\ell+1} \leftarrow \{d \in \mathcal{B}_{b,\ell} : U_{b,\ell}(d) \geq L_{b,\ell}(d_{\text{len}(b)})\}$
 - 19: $\ell_b \leftarrow \ell_b + 1$
 - 20: **else if** $s > 0$ **then**
 - 21: // Split
 - 22: $\mathcal{A} \leftarrow \mathcal{A} \cup \{b_{\max} + 1, b_{\max} + 2\} \setminus \{b\}$
 - 23: $I_{b_{\max}+1} \leftarrow (I_b(1), I_b(1) + s - 1)$
 - 24: $B_{b_{\max}+1,0} \leftarrow B_s^+$, $\ell_{b_{\max}+1} \leftarrow 0$
 - 25: $I_{b_{\max}+2} \leftarrow (I_b(1) + s, I_b(2))$
 - 26: $B_{b_{\max}+2,0} \leftarrow B_s^-$, $\ell_{b_{\max}+2} \leftarrow 0$
 - 27: $b_{\max} \leftarrow b_{\max} + 2$
-

that batch. Finally, when $I_b(2) < K$, the number of items in batch b is equal to the number of its positions.

5. Analysis

In this section, we state our regret bound for BatchRank. Before we do so, we discuss our estimator of clicks in (6). In particular, we show that (6) is the attraction probability of item d scaled by the average examination probability in stage ℓ of batch b . The examination scaling preserves the order of attraction probabilities, and therefore BatchRank can operate on (6) in place of $\alpha(d)$.

5.1. Confidence Radii

Fix batch b , positions I_b , stage ℓ , and items $\mathcal{B}_{b,\ell}$. Then for any item $d \in \mathcal{B}_{b,\ell}$, we can write the estimator in (6) as

$$\hat{c}_{b,\ell}(d) = \frac{1}{n_\ell} \sum_{t \in \mathcal{T}} \sum_{k=I_b(1)}^{I_b(2)} c_t(k) \mathbb{1}\{d_k^t = d\} \quad (7)$$

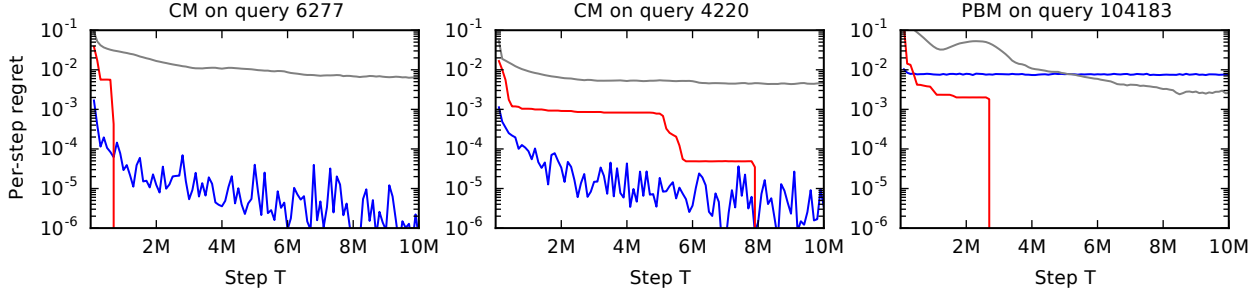


Figure 1. The expected per-step regret of BatchRank (red), CascadeKL-UCB (blue), and RankedExp3 (gray) on three problems. The results are averaged over 10 runs.

for some set of n_ℓ time steps \mathcal{T} and its expected value is

$$\bar{c}_{b,\ell}(d) = \mathbb{E}[\hat{c}_{b,\ell}(d)]. \quad (8)$$

The key step in the design of BatchRank is that we maintain confidence radii around (7). This is sound because the observations in (7),

$$\{\mathbf{X}_t(\mathcal{R}_t, k) \mathbf{A}_t(d)\}_{t \in \{t \in \mathcal{T}: d_k^* = d\}} \quad (9)$$

at any position k , are i.i.d. in time. More precisely, by the design of DisplayBatch, all displayed items from batch b are chosen randomly from $\mathbf{B}_{b,\ell}$; and independently of the realizations of $\mathbf{X}_t(\mathcal{R}_t, k)$ and $\mathbf{A}_t(d)$, which are random as well. The last problem is that the policy for placing items at positions $1, \dots, \mathbf{I}_b(1) - 1$ can change independently of batch b because BatchRank splits batches independently. But this has no effect on $\mathbf{X}_t(\mathcal{R}_t, k)$ because the examination of position k does not depend on the order of higher ranked items (Assumption 3).

5.2. Correct Examination Scaling

Fix batch b , positions \mathbf{I}_b , stage ℓ , and items $\mathbf{B}_{b,\ell}$. Since the examination of a position does not depend on the order of higher ranked items, and does not depend on lower ranked items at all (Assumption 3), we can express the probability of clicking on any item $d \in \mathbf{B}_{b,\ell}$ in (7) as

$$\bar{c}_{b,\ell}(d) = \frac{\alpha(d)}{|\mathcal{S}_d|} \sum_{\mathcal{R} \in \mathcal{S}_d} \sum_{k=\mathbf{I}_b(1)}^{\mathbf{I}_b(2)} \chi(\mathcal{R}, k) \mathbb{1}\{d_k = k\}, \quad (10)$$

where

$$\mathcal{S}_d = \left\{ (e_1, \dots, e_{\mathbf{I}_b(2)}) : d \in \{e_{\mathbf{I}_b(1)}, \dots, e_{\mathbf{I}_b(2)}\}, \right. \\ \left. (e_{\mathbf{I}_b(1)}, \dots, e_{\mathbf{I}_b(2)}) \in \Pi_{\text{len}(b)}(\mathbf{B}_{b,\ell}) \right\}$$

is the set of all lists with permutations of $\mathbf{B}_{b,\ell}$ on positions \mathbf{I}_b that contain item d , for some fixed higher ranked items $e_1, \dots, e_{\mathbf{I}_b(1)-1} \notin \mathbf{B}_{b,\ell}$. Let $d^* \in \mathbf{B}_{b,\ell}$ be any item such that $\alpha(d^*) \geq \alpha(d)$, and $\bar{c}_{b,\ell}(d^*)$ and \mathcal{S}_{d^*} be defined analogously to $\bar{c}_{b,\ell}(d)$ and \mathcal{S}_d above. Then we argue that

$$\bar{c}_{b,\ell}(d^*)/\alpha(d^*) \geq \bar{c}_{b,\ell}(d)/\alpha(d), \quad (11)$$

the examination scaling of a less attractive item d is never higher than that of a more attractive item d^* .

Before we prove (11), note that for any list $\mathcal{R} \in \mathcal{S}_d$, there exists one and only one list in \mathcal{S}_{d^*} that differs from \mathcal{R} only in that items d and d^* are exchanged. Let this list be \mathcal{R}^* . We analyze three cases. First, suppose that list \mathcal{R} does not contain item d^* . Then by Assumption 3, the examination probabilities of d in \mathcal{R} and d^* in \mathcal{R}^* are the same. Second, let item d^* be ranked higher than item d in \mathcal{R} . Then by Assumption 5, the examination probability of d in \mathcal{R} is not higher than that of d^* in \mathcal{R}^* . Third, let item d^* be ranked lower than item d in \mathcal{R} . Then by Assumption 3, the examination probabilities of d in \mathcal{R} and d^* in \mathcal{R}^* are the same, since they do not depend on lower ranked items. Finally, from the definition in (10) and that $|\mathcal{S}_d| = |\mathcal{S}_{d^*}|$, we have that (11) holds.

5.3. Regret Bound

For simplicity of exposition, let $\alpha(1) > \dots > \alpha(L) > 0$. Let $\alpha_{\max} = \alpha(1)$, and $\chi^*(k) = \chi(\mathcal{R}^*, k)$ for all $k \in [K]$. The regret of BatchRank is bounded below.

Theorem 1. *For any stochastic click bandit in Section 3.1 that satisfies Assumptions 1 to 6 and $T \geq 5$, the expected T -step regret of BatchRank is bounded as*

$$R(T) \leq \frac{192K^3L}{(1 - \alpha_{\max})\Delta_{\min}} \log T + 4KL(3e + K),$$

where $\Delta_{\min} = \min_{k \in [K]} \{\alpha(k) - \alpha(k+1)\}$.

Proof. The key idea is to bound the expected T -step regret in any batch (Lemma 7 in Appendix). Since the number of batches is at most $2K$, the regret of BatchRank is at most $2K$ times larger than that of in any batch.

The regret in a batch is bounded as follows. Let all confidence intervals hold, \mathbf{I}_b be the positions of batch b , and the maximum gap in batch b be

$$\Delta_{\max} = \max_{d \in \{\mathbf{I}_b(1), \dots, \mathbf{I}_b(2)-1\}} [\alpha(d) - \alpha(d+1)].$$

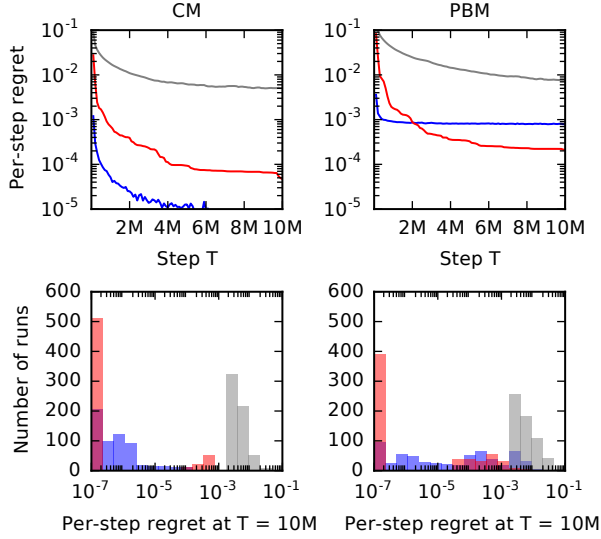


Figure 2. The comparison of BatchRank (red), CascadeKL-UCB (blue), and RankedExp3 (gray) in the CM and PBM. In the top plots, we report the per-step regret as a function of time T , averaged over 60 queries and 10 runs per query. In the bottom plots, we show the distribution of the regret at $T = 10M$.

If the gap of item d in batch b is $O(K\Delta_{\max})$, its regret is dominated by the time that the batch splits, and we bound this time in Lemma 6 in Appendix. Otherwise, the item is likely to be eliminated before the split, and we bound this time in Lemma 5 in Appendix. Now take the maximum of these upper bounds. ■

5.4. Discussion

Our upper bound in Theorem 1 is logarithmic in the number of steps T , linear in the number of items L , and polynomial in the number of positions K . To the best of our knowledge, this is the first gap-dependent upper bound on the regret of a learning algorithm that has sublinear regret in both the CM and PBM. The gap Δ_{\min} characterizes the hardness of sorting $K + 1$ most attractive items, which is sufficient for solving our problem. In practice, the maximum attraction probability α_{\max} is bounded away from 1. Therefore, the dependence on $(1 - \alpha_{\max})^{-1}$ is not critical. For instance, in most queries in Section 6, $\alpha_{\max} \leq 0.9$.

We believe that the cubic dependence on K is not far from being optimal. In particular, consider the problem of learning the most clicked item-position pair in the PBM (Section 2.1), which is easier than our problem. This problem can be solved as a stochastic rank-1 bandit (Katariya et al., 2017b) by Rank1E1im. Now consider the following PBM. The examination probability of the first position is close to one and the examination probabilities of all other positions are close to zero. Then the T -step regret of Rank1E1im is

$O([K^3 + K^2L\Delta_{\min}^{-1}] \log T)$ because $\mu = O(1/K)$, where μ is defined in Katariya et al. (2017b). Note that the gap-dependent term nearly matches our upper bound.

The KL confidence intervals in BatchRank are necessary to achieve sample efficiency. The reason is, as we prove in Lemma 9 in Appendix, that

$$\chi(1 - m)D_{\text{KL}}(\alpha \parallel \alpha^*) \leq D_{\text{KL}}(\chi\alpha \parallel \chi\alpha^*)$$

for any $\chi, \alpha, \alpha^* \in [0, 1]$ and $m = \max\{\alpha, \alpha^*\}$. This implies that any two items with the expected rewards of $\chi\alpha$ and $\chi\alpha^*$ can be distinguished in $O(\chi^{-1}(\alpha^* - \alpha)^{-2})$ observations for any scaling factor χ , when m is bounded away from 1. Suppose that $\alpha < \alpha^*$. Then the expected per-step regret for choosing the suboptimal item is $\chi(\alpha^* - \alpha)$, and the expected cumulative regret is $O((\alpha^* - \alpha)^{-1})$. The key observation is that the regret is independent of χ . This is a major improvement over UCB1 confidence intervals, which only lead to $O(\chi^{-1}(\alpha^* - \alpha)^{-1})$ regret. Because χ can be exponentially small, such a dependence is undesirable.

The elimination of items in UpdateBatch (lines 17–19) is necessary. The regret of BatchRank would be quadratic in L otherwise.

6. Experiments

We experiment with the *Yandex* dataset (Yandex), a dataset of 35 million (M) search sessions, each of which may contain multiple search queries. Each query is associated with displayed documents at positions 1 to 10 and their clicks. We select 60 frequent search queries, and learn their CMs and PBMs using PyClick (Chuklin et al., 2015), which is an open-source library of click models for web search. In each query, our goal is to rerank $L = 10$ most attractive items with the objective of maximizing the expected number of clicks at the first $K = 5$ positions. This resembles a real-world setting, where the learning agent would only be allowed to rerank highly attractive items, and not allowed to explore unattractive items (Zoghi et al., 2016).

BatchRank is compared to two methods, CascadeKL-UCB (Kveton et al., 2015a) and RankedExp3 (Radlinski et al., 2008a). CascadeKL-UCB is an optimal algorithm for learning to rank in the cascade model. RankedExp3 is a variant of ranked bandits (Section 7) where the base bandit algorithm is Exp3 (Auer et al., 1995). This approach is popular in practice and does not make any independence assumptions on the attractions of items.

Many solutions in our queries are near optimal, and therefore the optimal solutions are hard to learn. Therefore, we decided to evaluate the performance of algorithms by their *expected per-step regret*, in up to 10M steps. If a solution is suboptimal and does not improve over time, its expected per-step regret remains constant and is bounded away from

zero, and this can be easily observed even if the gap of the solution is small. We expect this when CascadeKL-UCB is applied to the PBM because CascadeKL-UCB has no guarantees in this model. The reported regret is averaged over periods of 100k steps to reduce randomness.

We report the performance of all compared algorithms on two CMs and one PBM in Figure 1. The plots are chosen to represent general trends in this experiment. In the CM, CascadeKL-UCB performs very well on most queries. This is not surprising since CascadeKL-UCB is designed for the CM. BatchRank often learns the optimal solution quickly (Figure 1a), but sometimes this requires close to $T = 10M$ steps (Figure 1b). In the PBM, CascadeKL-UCB may converge to a suboptimal solution. Then its expected per-step regret remains constant and even RankedExp3 can learn a better solution over time (Figure 1c). We also observe that BatchRank outperforms RankedExp3 in all experiments.

We report the average performance of all compared algorithms in both click models in Figure 2. These trends confirm our earlier findings. In the CM, CascadeKL-UCB outperforms BatchRank; while in the PBM, BatchRank outperforms CascadeKL-UCB at $T = 2M$ steps. The regret of CascadeKL-UCB in the PBM flattens and is bounded away from zero. This trend can be explained by the histograms in Figure 2. They show that CascadeKL-UCB converges to suboptimal solutions, whose regret is at least 10^{-3} , in one sixth of its runs. The performance of BatchRank is more robust and we do not observe many runs whose regret is of that magnitude.

We are delighted with the performance of BatchRank. Although it is not designed to be optimal (Section 5.4), it is more robust than CascadeKL-UCB and clearly outperforms RankedExp3. The performance of CascadeKL-UCB is unexpectedly good. Although it does not have guarantees in the PBM, it performs very well on many queries. We plan to investigate this in our future work.

7. Related Work

A popular approach to online learning to rank are *ranked bandits* (Radlinski et al., 2008a; Slivkins et al., 2013). The key idea in ranked bandits is to model each position in the recommended list as an individual bandit problem, which is then solved by a *base bandit algorithm*. This algorithm is typically adversarial (Auer et al., 1995) because the distribution of clicks on lower positions is affected by higher positions. We compare to ranked bandits in Section 6.

Online learning to rank in click models (Craswell et al., 2008; Chuklin et al., 2015) was recently studied in several papers (Kveton et al., 2015a; Combes et al., 2015; Kveton et al., 2015b; Katariya et al., 2016; Zong et al., 2016; Li et al., 2016; Lagree et al., 2016). In all of these papers, the

attraction probabilities of items are estimated from clicks and the click model. The learning agent has no guarantees beyond this model.

The problem of finding the most clicked item-position pair in the PBM, which is arguably easier than our problem of finding K most clicked item-position pairs, can be solved as a stochastic rank-1 bandit (Katariya et al., 2017b;a). We discuss our relation to these works in Section 5.4.

Our problem can be also viewed as an instance of partial monitoring, where the attraction indicators of items are unobserved. General partial-monitoring algorithms (Agrawal et al., 1989; Bartok et al., 2012; Bartok & Szepesvari, 2012; Bartok et al., 2014) are unsuitable for our setting because their computational complexity is polynomial in the number of actions, which is exponential in K .

The click model is a model of how the user interacts with a list of documents (Chuklin et al., 2015), and many such models have been proposed (Becker et al., 2007; Richardson et al., 2007; Craswell et al., 2008; Chapelle & Zhang, 2009; Guo et al., 2009a;b). Two fundamental click models are the CM (Craswell et al., 2008) and PBM (Richardson et al., 2007). These models have been traditionally studied separately. In this work, we show that learning to rank problems in these models can be solved by the same algorithm, under reasonable assumptions.

8. Conclusions

We propose stochastic click bandits, a framework for online learning to rank in a broad class of click models that encompasses two most fundamental click models, the cascade and position-based models. In addition, we propose a computationally and sample efficient algorithm for solving our problems, BatchRank, and derive an upper bound on its T -step regret. Finally, we evaluate BatchRank on web search queries. Our algorithm outperforms ranked bandits (Radlinski et al., 2008a), a popular online learning to rank approach; and is more robust than CascadeKL-UCB (Kveton et al., 2015a), an existing algorithm for online learning to rank in the cascade model.

The goal of this work is not to propose the optimal algorithm for our setting, but to demonstrate that online learning to rank in multiple click models is possible with theoretical guarantees. We strongly believe that the design of BatchRank, as well as its analysis, can be improved. For instance, BatchRank resets its estimators of clicks in each batch, which is wasteful. In addition, based on the discussion in Section 5.4, our analysis may be loose by a factor of K . We hope that the practically relevant setting, which is introduced in this paper, will spawn new enthusiasm in the community and lead to more work in this area.

References

- Agichtein, Eugene, Brill, Eric, and Dumais, Susan. Improving web search ranking by incorporating user behavior information. In *Proceedings of the 29th Annual International ACM SIGIR Conference*, pp. 19–26, 2006.
- Agrawal, Rajeev, Teneketzis, Demosthenis, and Anantharam, Venkatachalam. Asymptotically efficient adaptive allocation schemes for controlled i.i.d. processes: Finite parameter space. *IEEE Transactions on Automatic Control*, 34(3):258–267, 1989.
- Auer, Peter, Cesa-Bianchi, Nicolo, Freund, Yoav, and Schapire, Robert. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *Proceedings of the 36th Annual Symposium on Foundations of Computer Science*, pp. 322–331, 1995.
- Auer, Peter, Cesa-Bianchi, Nicolo, and Fischer, Paul. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47:235–256, 2002.
- Bartok, Gabor and Szepesvari, Csaba. Partial monitoring with side information. In *Proceedings of the 23rd International Conference on Algorithmic Learning Theory*, pp. 305–319, 2012.
- Bartok, Gabor, Zolghadr, Navid, and Szepesvari, Csaba. An adaptive algorithm for finite stochastic partial monitoring. In *Proceedings of the 29th International Conference on Machine Learning*, 2012.
- Bartok, Gabor, Foster, Dean, Pal, David, Rakhlin, Alexander, and Szepesvari, Csaba. Partial monitoring - classification, regret bounds, and algorithms. *Mathematics of Operations Research*, 39(4):967–997, 2014.
- Becker, Hila, Meek, Christopher, and Chickering, David Maxwell. Modeling contextual factors of click rates. In *Proceedings of the 22nd AAAI Conference on Artificial Intelligence*, pp. 1310–1315, 2007.
- Chapelle, Olivier and Zhang, Ya. A dynamic Bayesian network click model for web search ranking. In *Proceedings of the 18th International Conference on World Wide Web*, pp. 1–10, 2009.
- Chuklin, Aleksandr, Markov, Ilya, and de Rijke, Maarten. *Click Models for Web Search*. Morgan & Claypool Publishers, 2015.
- Combes, Richard, Magureanu, Stefan, Proutiere, Alexandre, and Laroche, Cyrille. Learning to rank: Regret lower bounds and efficient algorithms. In *Proceedings of the 2015 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*, 2015.
- Craswell, Nick, Zoeter, Onno, Taylor, Michael, and Ramsey, Bill. An experimental comparison of click position-bias models. In *Proceedings of the 1st ACM International Conference on Web Search and Data Mining*, pp. 87–94, 2008.
- Garivier, Aurelien and Cappe, Olivier. The KL-UCB algorithm for bounded stochastic bandits and beyond. In *Proceeding of the 24th Annual Conference on Learning Theory*, pp. 359–376, 2011.
- Grotov, Artem, Chuklin, Aleksandr, Markov, Ilya, Stout, Luka, Xumara, Finde, and de Rijke, Maarten. A comparative study of click models for web search. In *Proceedings of the 6th International Conference of the CLEF Association*, 2015.
- Guo, Fan, Liu, Chao, Kannan, Anitha, Minka, Tom, Taylor, Michael, Wang, Yi Min, and Faloutsos, Christos. Click chain model in web search. In *Proceedings of the 18th International Conference on World Wide Web*, pp. 11–20, 2009a.
- Guo, Fan, Liu, Chao, and Wang, Yi Min. Efficient multiple-click models in web search. In *Proceedings of the 2nd ACM International Conference on Web Search and Data Mining*, pp. 124–131, 2009b.
- Hoeffding, Wassily. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- Hofmann, Katja, Schuth, Anne, Whiteson, Shimon, and de Rijke, Maarten. Reusing historical interaction data for faster online learning to rank for IR. In *Proceedings of the 6th ACM International Conference on Web Search and Data Mining*, pp. 183–192, 2013.
- Katariya, Sumeet, Kveton, Branislav, Szepesvari, Csaba, and Wen, Zheng. DCM bandits: Learning to rank with multiple clicks. In *Proceedings of the 33rd International Conference on Machine Learning*, pp. 1215–1224, 2016.
- Katariya, Sumeet, Kveton, Branislav, Szepesvari, Csaba, Vernade, Claire, and Wen, Zheng. Bernoulli rank-1 bandits for click feedback. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, 2017a.
- Katariya, Sumeet, Kveton, Branislav, Szepesvari, Csaba, Vernade, Claire, and Wen, Zheng. Stochastic rank-1 bandits. In *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, 2017b.
- Kveton, Branislav, Szepesvari, Csaba, Wen, Zheng, and Ashkan, Azin. Cascading bandits: Learning to rank in the cascade model. In *Proceedings of the 32nd International Conference on Machine Learning*, 2015a.

- Kveton, Branislav, Wen, Zheng, Ashkan, Azin, and Szepesvari, Csaba. Combinatorial cascading bandits. In *Advances in Neural Information Processing Systems 28*, pp. 1450–1458, 2015b.
- Lagree, Paul, Vernade, Claire, and Cappe, Olivier. Multiple-play bandits in the position-based model. In *Advances in Neural Information Processing Systems 29*, pp. 1597–1605, 2016.
- Lai, T. L. and Robbins, Herbert. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1):4–22, 1985.
- Li, Shuai, Wang, Baoxiang, Zhang, Shengyu, and Chen, Wei. Contextual combinatorial cascading bandits. In *Proceedings of the 33rd International Conference on Machine Learning*, pp. 1245–1253, 2016.
- Liu, Tie-Yan. *Learning to Rank for Information Retrieval*. Springer, 2011.
- Radlinski, Filip, Kleinberg, Robert, and Joachims, Thorsten. Learning diverse rankings with multi-armed bandits. In *Proceedings of the 25th International Conference on Machine Learning*, pp. 784–791, 2008a.
- Radlinski, Filip, Kurup, Madhu, and Joachims, Thorsten. How does clickthrough data reflect retrieval quality? In *Proceedings of the 17th ACM Conference on Information and Knowledge Management*, pp. 43–52, 2008b.
- Richardson, Matthew, Dominowska, Ewa, and Ragno, Robert. Predicting clicks: Estimating the click-through rate for new ads. In *Proceedings of the 16th International Conference on World Wide Web*, pp. 521–530, 2007.
- Slivkins, Aleksandrs, Radlinski, Filip, and Gollapudi, Sreenivas. Ranked bandits in metric spaces: Learning diverse rankings over large document collections. *Journal of Machine Learning Research*, 14(1):399–436, 2013.
- Yandex. Yandex personalized web search challenge. <https://www.kaggle.com/c/yandex-personalized-web-search-challenge>, 2013.
- Yue, Yisong and Joachims, Thorsten. Interactively optimizing information retrieval systems as a dueling bandits problem. In *Proceedings of the 26th International Conference on Machine Learning*, pp. 1201–1208, 2009.
- Zoghi, Masrour, Tunys, Tomas, Li, Lihong, Jose, Damien, Chen, Junyan, Chin, Chun Ming, and de Rijke, Maarten. Click-based hot fixes for underperforming torso queries. In *Proceedings of the 39th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 195–204, 2016.
- Zong, Shi, Ni, Hao, Sung, Kenny, Ke, Nan Rosemary, Wen, Zheng, and Kveton, Branislav. Cascading bandits for large-scale recommendation problems. In *Proceedings of the 32nd Conference on Uncertainty in Artificial Intelligence*, 2016.

A. Notation

Symbol	Definition
$\alpha(d)$	Attraction probability of item d
α_{\max}	Highest attraction probability, $\alpha(1)$
A	Binary attraction vector, where $A(d)$ is the attraction indicator of item d
P_α	Distribution over binary attraction vectors
\mathcal{A}	Set of active batches
b_{\max}	Index of the last created batch
$\mathbf{B}_{b,\ell}$	Items in stage ℓ of batch b
$c_t(k)$	Indicator of the click on position k at time t
$c_{b,\ell}(d)$	Number of observed clicks on item d in stage ℓ of batch b
$\hat{c}_{b,\ell}(d)$	Estimated probability of clicking on item d in stage ℓ of batch b
$\bar{c}_{b,\ell}(d)$	Probability of clicking on item d in stage ℓ of batch b , $\mathbb{E}[\hat{c}_{b,\ell}(d)]$
\mathcal{D}	Ground set of items $[L]$ such that $\alpha(1) \geq \dots \geq \alpha(L)$
δ_T	$\log T + 3 \log \log T$
$\tilde{\Delta}_\ell$	$2^{-\ell}$
\mathbf{I}_b	Interval of positions in batch b
K	Number of positions to display items
$\text{len}(b)$	Number of positions to display items in batch b
L	Number of items
$\mathbf{L}_{b,\ell}(d)$	Lower confidence bound of item d , in stage ℓ of batch b
n_ℓ	Number of times that each item is observed in stage ℓ
$\mathbf{n}_{b,\ell}$	Number of observations of item d in stage ℓ of batch b
$\Pi_K(\mathcal{D})$	Set of all K -tuples with distinct elements from \mathcal{D}
$r(\mathcal{R}, A, X)$	Reward of list \mathcal{R} , for attraction and examination indicators A and X
$r(\mathcal{R}, \alpha, \chi)$	Expected reward of list \mathcal{R}
$\mathcal{R} = (d_1, \dots, d_K)$	List of K items, where d_k is the k -th item in \mathcal{R}
$\mathcal{R}^* = (1, \dots, K)$	Optimal list of K items
$R(\mathcal{R}, A, X)$	Regret of list \mathcal{R} , for attraction and examination indicators A and X
$R(T)$	Expected cumulative regret in T steps
T	Horizon of the experiment
$\mathbf{U}_{b,\ell}(d)$	Upper confidence bound of item d , in stage ℓ of batch b
$\chi(\mathcal{R}, k)$	Examination probability of position k in list \mathcal{R}
$\chi^*(k)$	Examination probability of position k in the optimal list \mathcal{R}^*
X	Binary examination matrix, where $X(\mathcal{R}, k)$ is the examination indicator of position k in list \mathcal{R}
P_χ	Distribution over binary examination matrices

B. Proof of Theorem 1

Let $\mathbf{R}_{b,\ell}$ be the stochastic regret associated with stage ℓ of batch b . Then the expected T -step regret of MergeRank can be decomposed as

$$R(T) \leq \mathbb{E} \left[\sum_{b=1}^{2K} \sum_{\ell=0}^{T-1} \mathbf{R}_{b,\ell} \right]$$

because the maximum number of batches is $2K$. Let

$$\bar{\chi}_{b,\ell}(d) = \frac{\bar{c}_{b,\ell}(d)}{\alpha(d)} \quad (12)$$

be the *average examination probability* of item d in stage ℓ of batch b . Let

$$\begin{aligned} \mathcal{E}_{b,\ell} = & \left\{ \begin{array}{l} \text{Event 1: } \forall d \in \mathbf{B}_{b,\ell} : \bar{c}_{b,\ell}(d) \in [\mathbf{L}_{b,\ell}(d), \mathbf{U}_{b,\ell}(d)], \\ \text{Event 2: } \forall \mathbf{I}_b \in [K]^2, d \in \mathbf{B}_{b,\ell}, d^* \in \mathbf{B}_{b,\ell} \cap [K] \text{ s.t. } \Delta = \alpha(d^*) - \alpha(d) > 0 : \\ \quad n_\ell \geq \frac{16K}{\chi^*(\mathbf{I}_b(1))(1 - \alpha_{\max})\Delta^2} \log T \implies \hat{c}_{b,\ell}(d) \leq \bar{\chi}_{b,\ell}(d)[\alpha(d) + \Delta/4], \\ \text{Event 3: } \forall \mathbf{I}_b \in [K]^2, d \in \mathbf{B}_{b,\ell}, d^* \in \mathbf{B}_{b,\ell} \cap [K] \text{ s.t. } \Delta = \alpha(d^*) - \alpha(d) > 0 : \\ \quad n_\ell \geq \frac{16K}{\chi^*(\mathbf{I}_b(1))(1 - \alpha_{\max})\Delta^2} \log T \implies \hat{c}_{b,\ell}(d^*) \geq \bar{\chi}_{b,\ell}(d^*)[\alpha(d^*) - \Delta/4] \end{array} \right\} \end{aligned}$$

be the “good event” in stage ℓ of batch b , where $\bar{c}_{b,\ell}(d)$ is the probability of clicking on item d in stage ℓ of batch b , which is defined in (8); $\hat{c}_{b,\ell}(d)$ is its estimate, which is defined in (7); and both χ^* and α_{\max} are defined in Section 5.3. Let $\bar{\mathcal{E}}_{b,\ell}$ be the complement of event $\mathcal{E}_{b,\ell}$. Let \mathcal{E} be the “good event” that all events $\mathcal{E}_{b,\ell}$ happen; and $\bar{\mathcal{E}}$ be its complement, the “bad event” that at least one event $\mathcal{E}_{b,\ell}$ does not happen. Then the expected T -step regret can be bounded from above as

$$R(T) \leq \mathbb{E} \left[\sum_{b=1}^{2K} \sum_{\ell=0}^{T-1} \mathbf{R}_{b,\ell} \mathbb{1}\{\mathcal{E}\} \right] + TP(\bar{\mathcal{E}}) \leq \sum_{b=1}^{2K} \mathbb{E} \left[\sum_{\ell=0}^{T-1} \mathbf{R}_{b,\ell} \mathbb{1}\{\mathcal{E}\} \right] + 4KL(3e + K),$$

where the second inequality is from Lemma 2. Now we apply Lemma 7 to each batch b and get that

$$\sum_{b=1}^{2K} \mathbb{E} \left[\sum_{\ell=0}^{T-1} \mathbf{R}_{b,\ell} \mathbb{1}\{\mathcal{E}\} \right] \leq \frac{192K^3L}{(1 - \alpha_{\max})\Delta_{\min}} \log T.$$

This concludes our proof.

C. Upper Bound on the Probability of Bad Event $\bar{\mathcal{E}}$

Lemma 2. Let $\bar{\mathcal{E}}$ be defined as in the proof of Theorem 1 and $T \geq 5$. Then

$$P(\bar{\mathcal{E}}) \leq \frac{4KL(3e + K)}{T}.$$

Proof. By the union bound,

$$P(\bar{\mathcal{E}}) \leq \sum_{b=1}^{2K} \sum_{\ell=0}^{T-1} P(\bar{\mathcal{E}}_{b,\ell}).$$

Now we bound the probability of each event in $\bar{\mathcal{E}}_{b,\ell}$ and then sum them up.

Event 1

The probability that event 1 in $\mathcal{E}_{b,\ell}$ does not happen is bounded as follows. Fix I_b and $B_{b,\ell}$. For any $d \in B_{b,\ell}$,

$$\begin{aligned} P(\bar{c}_{b,\ell}(d) \notin [L_{b,\ell}(d), U_{b,\ell}(d)]) &\leq P(\bar{c}_{b,\ell}(d) < L_{b,\ell}(d)) + P(\bar{c}_{b,\ell}(d) > U_{b,\ell}(d)) \\ &\leq \frac{2e \lceil \log(T \log^3 T) \log n_\ell \rceil}{T \log^3 T} \\ &\leq \frac{2e \lceil \log^2 T + \log(\log^3 T) \log T \rceil}{T \log^3 T} \\ &\leq \frac{2e \lceil 2 \log^2 T \rceil}{T \log^3 T} \\ &\leq \frac{6e}{T \log T}, \end{aligned}$$

where the second inequality is by Theorem 10 of [Garivier & Cappé \(2011\)](#), the third inequality is from $T \geq n_\ell$, the fourth inequality is from $\log(\log^3 T) \leq \log T$ for $T \geq 5$, and the last inequality is from $\lceil 2 \log^2 T \rceil \leq 3 \log^2 T$ for $T \geq 3$. By the union bound,

$$P(\exists d \in B_{b,\ell} \text{ s.t. } \bar{c}_{b,\ell}(d) \notin [L_{b,\ell}(d), U_{b,\ell}(d)]) \leq \frac{6eL}{T \log T}$$

for any $B_{b,\ell}$. Finally, since the above inequality holds for any $B_{b,\ell}$, the probability that event 1 in $\mathcal{E}_{b,\ell}$ does not happen is bounded as above.

Event 2

The probability that event 2 in $\mathcal{E}_{b,\ell}$ does not happen is bounded as follows. Fix I_b and $B_{b,\ell}$, and let $k = I_b(1)$. If the event does not happen for items d and d^* , then it must be true that

$$n_\ell \geq \frac{16K}{\chi^*(k)(1 - \alpha_{\max})\Delta^2} \log T, \quad \hat{c}_{b,\ell}(d) > \bar{\chi}_{b,\ell}(d)[\alpha(d) + \Delta/4].$$

From the definition of the average examination probability in (12) and a variant of Hoeffding's inequality in Lemma 8, we have that

$$P(\hat{c}_{b,\ell}(d) > \bar{\chi}_{b,\ell}(d)[\alpha(d) + \Delta/4]) \leq \exp[-n_\ell D_{\text{KL}}(\bar{\chi}_{b,\ell}(d)[\alpha(d) + \Delta/4] \parallel \bar{c}_{b,\ell}(d))].$$

From Lemma 9, $\bar{\chi}_{b,\ell}(d) \geq \chi^*(k)/K$ (Lemma 3), and Pinsker's inequality, we have that

$$\begin{aligned} \exp[-n_\ell D_{\text{KL}}(\bar{\chi}_{b,\ell}(d)[\alpha(d) + \Delta/4] \parallel \bar{c}_{b,\ell}(d))] &\leq \exp[-n_\ell \bar{\chi}_{b,\ell}(d)(1 - \alpha_{\max}) D_{\text{KL}}(\alpha(d) + \Delta/4 \parallel \alpha(d))] \\ &\leq \exp\left[-n_\ell \frac{\chi^*(k)(1 - \alpha_{\max})\Delta^2}{8K}\right]. \end{aligned}$$

From our assumption on n_ℓ , we conclude that

$$\exp\left[-n_\ell \frac{\chi^*(k)(1 - \alpha_{\max})\Delta^2}{8K}\right] \leq \exp[-2 \log T] = \frac{1}{T^2}.$$

Finally, we chain all above inequalities and get that event 2 in $\mathcal{E}_{b,\ell}$ does not happen for any fixed \mathbf{I}_b , $\mathbf{B}_{b,\ell}$, d , and d^* with probability of at most T^{-2} . Since the maximum numbers of items d and d^* are L and K , respectively, the event does not happen for any fixed \mathbf{I}_b and $\mathbf{B}_{b,\ell}$ with probability of at most $KL T^{-2}$. In turn, the probability that event 2 in $\mathcal{E}_{b,\ell}$ does not happen is bounded by $KL T^{-2}$.

Event 3

This bound is analogous to that of event 2.

Total probability

The maximum number of stages in any batch in BatchRank is $\log T$ and the maximum number of batches is $2K$. Hence, by the union bound,

$$P(\bar{\mathcal{E}}) \leq \left(\frac{6eL}{T \log T} + \frac{KL}{T^2} + \frac{KL}{T^2} \right) (2K \log T) \leq \frac{4KL(3e + K)}{T}.$$

This concludes our proof. ■

D. Upper Bound on the Regret in Individual Batches

Lemma 3. For any batch b , positions \mathbf{I}_b , stage ℓ , set $\mathbf{B}_{b,\ell}$, and item $d \in \mathbf{B}_{b,\ell}$,

$$\frac{\chi^*(k)}{K} \leq \bar{\chi}_{b,\ell}(d),$$

where $k = \mathbf{I}_b(1)$ is the highest position in batch b .

Proof. The proof follows from two observations. First, by Assumption 6, $\chi^*(k)$ is the lowest examination probability of position k . Second, by the design of `DisplayBatch`, item d is placed at position k with probability of at least $1/K$. ■

Lemma 4. Let event \mathcal{E} happen and $T \geq 5$. For any batch b , positions \mathbf{I}_b , set $\mathbf{B}_{b,0}$, item $d \in \mathbf{B}_{b,0}$, and item $d^* \in \mathbf{B}_{b,0} \cap [K]$ such that $\Delta = \alpha(d^*) - \alpha(d) > 0$, let $k = \mathbf{I}_b(1)$ be the highest position in batch b and ℓ be the first stage where

$$\tilde{\Delta}_\ell < \sqrt{\frac{\chi^*(k)(1 - \alpha_{\max})}{K}} \Delta.$$

Then $\mathbf{U}_{b,\ell}(d) < \mathbf{L}_{b,\ell}(d^*)$.

Proof. From the definition of n_ℓ in `BatchRank` and our assumption on $\tilde{\Delta}_\ell$,

$$n_\ell \geq \frac{16}{\tilde{\Delta}_\ell^2} \log T > \frac{16K}{\chi^*(k)(1 - \alpha_{\max})\Delta^2} \log T. \quad (13)$$

Let $\mu = \bar{\chi}_{b,\ell}(d)$ and suppose that $\mathbf{U}_{b,\ell}(d) \geq \mu[\alpha(d) + \Delta/2]$ holds. Then from this assumption, the definition of $\mathbf{U}_{b,\ell}(d)$, and event 2 in $\mathcal{E}_{b,\ell}$,

$$\begin{aligned} D_{\text{KL}}(\hat{\mathbf{c}}_{b,\ell}(d) \parallel \mathbf{U}_{b,\ell}(d)) &\geq D_{\text{KL}}(\hat{\mathbf{c}}_{b,\ell}(d) \parallel \mu[\alpha(d) + \Delta/2]) \mathbb{1}\{\hat{\mathbf{c}}_{b,\ell}(d) \leq \mu[\alpha(d) + \Delta/2]\} \\ &\geq D_{\text{KL}}(\mu[\alpha(d) + \Delta/4] \parallel \mu[\alpha(d) + \Delta/2]). \end{aligned}$$

From Lemma 9, $\mu \geq \chi^*(k)/K$ (Lemma 3), and Pinsker's inequality, we have that

$$\begin{aligned} D_{\text{KL}}(\mu[\alpha(d) + \Delta/4] \parallel \mu[\alpha(d) + \Delta/2]) &\geq \mu(1 - \alpha_{\max})D_{\text{KL}}(\alpha(d) + \Delta/4 \parallel \alpha(d) + \Delta/2) \\ &\geq \frac{\chi^*(k)(1 - \alpha_{\max})\Delta^2}{8K}. \end{aligned}$$

From the definition of $\mathbf{U}_{b,\ell}(d)$, $T \geq 5$, and above inequalities,

$$n_\ell = \frac{\log T + 3 \log \log T}{D_{\text{KL}}(\hat{\mathbf{c}}_{b,\ell}(d) \parallel \mathbf{U}_{b,\ell}(d))} \leq \frac{2 \log T}{D_{\text{KL}}(\hat{\mathbf{c}}_{b,\ell}(d) \parallel \mathbf{U}_{b,\ell}(d))} \leq \frac{16K \log T}{\chi^*(k)(1 - \alpha_{\max})\Delta^2}.$$

This contradicts to (13), and therefore it must be true that $\mathbf{U}_{b,\ell}(d) < \mu[\alpha(d) + \Delta/2]$ holds.

On the other hand, let $\mu^* = \bar{\chi}_{b,\ell}(d^*)$ and suppose that $\mathbf{L}_{b,\ell}(d^*) \leq \mu^*[\alpha(d^*) - \Delta/2]$ holds. Then from this assumption, the definition of $\mathbf{L}_{b,\ell}(d^*)$, and event 3 in $\mathcal{E}_{b,\ell}$,

$$\begin{aligned} D_{\text{KL}}(\hat{\mathbf{c}}_{b,\ell}(d^*) \parallel \mathbf{L}_{b,\ell}(d^*)) &\geq D_{\text{KL}}(\hat{\mathbf{c}}_{b,\ell}(d^*) \parallel \mu^*[\alpha(d^*) - \Delta/2]) \mathbb{1}\{\hat{\mathbf{c}}_{b,\ell}(d^*) \geq \mu^*[\alpha(d^*) - \Delta/2]\} \\ &\geq D_{\text{KL}}(\mu^*[\alpha(d^*) - \Delta/4] \parallel \mu^*[\alpha(d^*) - \Delta/2]). \end{aligned}$$

From Lemma 9, $\mu^* \geq \chi^*(k)/K$ (Lemma 3), and Pinsker's inequality, we have that

$$\begin{aligned} D_{\text{KL}}(\mu^*[\alpha(d^*) - \Delta/4] \parallel \mu^*[\alpha(d^*) - \Delta/2]) &\geq \mu^*(1 - \alpha_{\max})D_{\text{KL}}(\alpha(d^*) - \Delta/4 \parallel \alpha(d^*) - \Delta/2) \\ &\geq \frac{\chi^*(k)(1 - \alpha_{\max})\Delta^2}{8K}. \end{aligned}$$

From the definition of $\mathbf{L}_{b,\ell}(d^*)$, $T \geq 5$, and above inequalities,

$$n_\ell = \frac{\log T + 3 \log \log T}{D_{\text{KL}}(\hat{\mathbf{c}}_{b,\ell}(d) \parallel \mathbf{L}_{b,\ell}(d^*))} \leq \frac{2 \log T}{D_{\text{KL}}(\hat{\mathbf{c}}_{b,\ell}(d^*) \parallel \mathbf{L}_{b,\ell}(d^*))} \leq \frac{16K \log T}{\chi^*(k)(1 - \alpha_{\max})\Delta^2}.$$

This contradicts to (13), and therefore it must be true that $\mathbf{L}_{b,\ell}(d^*) > \mu^*[\alpha(d^*) - \Delta/2]$ holds.

Finally, based on inequality (11),

$$\mu^* = \frac{\bar{\mathbf{c}}_{b,\ell}(d^*)}{\alpha(d^*)} \geq \frac{\bar{\mathbf{c}}_{b,\ell}(d)}{\alpha(d)} = \mu,$$

and item d is guaranteed to be eliminated by the end of stage ℓ because

$$\begin{aligned} \mathbf{U}_{b,\ell}(d) &< \mu[\alpha(d) + \Delta/2] \\ &\leq \mu\alpha(d) + \frac{\mu^*\alpha(d^*) - \mu\alpha(d)}{2} \\ &= \mu^*\alpha(d^*) - \frac{\mu^*\alpha(d^*) - \mu\alpha(d)}{2} \\ &\leq \mu^*[\alpha(d^*) - \Delta/2] \\ &< \mathbf{L}_{b,\ell}(d^*). \end{aligned}$$

This concludes our proof. ■

Lemma 5. *Let event \mathcal{E} happen and $T \geq 5$. For any batch b , positions \mathbf{I}_b where $\mathbf{I}_b(2) = K$, set $\mathbf{B}_{b,0}$, and item $d \in \mathbf{B}_{b,0}$ such that $d > K$, let $k = \mathbf{I}_b(1)$ be the highest position in batch b and ℓ be the first stage where*

$$\tilde{\Delta}_\ell < \sqrt{\frac{\chi^*(k)(1 - \alpha_{\max})}{K}} \Delta$$

for $\Delta = \alpha(K) - \alpha(d)$. Then item d is eliminated by the end of stage ℓ .

Proof. Let $B^+ = \{k, \dots, K\}$. Now note that $\alpha(d^*) - \alpha(d) \geq \Delta$ for any $d^* \in B^+$. By Lemma 4, $\mathbf{L}_{b,\ell}(d^*) > \mathbf{U}_{b,\ell}(d)$ for any $d^* \in B^+$; and therefore item d is eliminated by the end of stage ℓ . ■

Lemma 6. *Let \mathcal{E} happen and $T \geq 5$. For any batch b , positions \mathbf{I}_b , and set $\mathbf{B}_{b,0}$, let $k = \mathbf{I}_b(1)$ be the highest position in batch b and ℓ be the first stage where*

$$\tilde{\Delta}_\ell < \sqrt{\frac{\chi^*(k)(1 - \alpha_{\max})}{K}} \Delta_{\max}$$

for $\Delta_{\max} = \alpha(s) - \alpha(s+1)$ and $s = \arg \max_{d \in \{\mathbf{I}_b(1), \dots, \mathbf{I}_b(2)-1\}} [\alpha(d) - \alpha(d+1)]$. Then batch b is split by the end of stage ℓ .

Proof. Let $B^+ = \{k, \dots, s\}$ and $B^- = \mathbf{B}_{b,0} \setminus B^+$. Now note that $\alpha(d^*) - \alpha(d) \geq \Delta_{\max}$ for any $(d^*, d) \in B^+ \times B^-$. By Lemma 4, $\mathbf{L}_{b,\ell}(d^*) > \mathbf{U}_{b,\ell}(d)$ for any $(d^*, d) \in B^+ \times B^-$; and therefore batch b is split by the end of stage ℓ . ■

Lemma 7. *Let event \mathcal{E} happen and $T \geq 5$. Then the expected T -step regret in any batch b is bounded as*

$$\mathbb{E} \left[\sum_{\ell=0}^{T-1} \mathbf{R}_{b,\ell} \right] \leq \frac{96K^2L}{(1 - \alpha_{\max})\Delta_{\max}} \log T.$$

Proof. Let $k = \mathbf{I}_b(1)$ be the highest position in batch b . Choose any item $d \in \mathbf{B}_{b,0}$ and let $\Delta = \alpha(k) - \alpha(d)$.

First, we show that the expected per-step regret of any item d is bounded by $\chi^*(k)\Delta$ when event \mathcal{E} happens. Since event \mathcal{E} happens, all eliminations and splits up to any stage ℓ of batch b are correct. Therefore, items $1, \dots, k-1$ are at positions $1, \dots, k-1$; and position k is examined with probability $\chi^*(k)$. Note that this is the highest examination probability in batch b (Assumption 4). Our upper bound follows from the fact that the reward is linear in individual items (Section 3.1).

We analyze two cases. First, suppose that $\Delta \leq 2K\Delta_{\max}$ for Δ_{\max} in Lemma 6. Then by Lemma 6, batch b splits when the number of steps in a stage is at most

$$\frac{16K}{\chi^*(k)(1 - \alpha_{\max})\Delta_{\max}^2} \log T.$$

By the design of `DisplayBatch`, any item in stage ℓ of batch b is displayed at most $2n_\ell$ times. Therefore, the maximum regret due to item d in the last stage before the split is

$$\frac{32K\chi^*(k)\Delta}{\chi^*(k)(1-\alpha_{\max})\Delta_{\max}^2} \log T \leq \frac{64K^2\Delta_{\max}}{(1-\alpha_{\max})\Delta_{\max}^2} \log T = \frac{64K^2}{(1-\alpha_{\max})\Delta_{\max}} \log T.$$

Now suppose that $\Delta > 2K\Delta_{\max}$. This implies that item d is easy to distinguish from item K . In particular,

$$\alpha(K) - \alpha(d) = \Delta - (\alpha(k) - \alpha(K)) \geq \Delta - K\Delta_{\max} \geq \frac{\Delta}{2},$$

where the equality is from the identity

$$\Delta = \alpha(k) - \alpha(d) = \alpha(k) - \alpha(K) + \alpha(K) - \alpha(d);$$

the first inequality is from $\alpha(k) - \alpha(K) \leq K\Delta_{\max}$, which follows from the definition of Δ_{\max} and $k \in [K]$; and the last inequality is from our assumption that $K\Delta_{\max} < \Delta/2$. Now we apply the derived inequality and, by Lemma 5 and from the design of `DisplayBatch`, the maximum regret due to item d in the stage where that item is eliminated is

$$\frac{32K\chi^*(k)\Delta}{\chi^*(k)(1-\alpha_{\max})(\alpha(K) - \alpha(d))^2} \log T \leq \frac{128K}{(1-\alpha_{\max})\Delta} \log T \leq \frac{64}{(1-\alpha_{\max})\Delta_{\max}} \log T.$$

The last inequality is from our assumption that $\Delta > 2K\Delta_{\max}$.

Because the lengths of the stages quadruple and `BatchRank` resets all click estimators at the beginning of each stage, the maximum expected regret due to any item d in batch b is at most 1.5 times higher than that in the last stage, and hence

$$\mathbb{E} \left[\sum_{\ell=0}^{T-1} \mathbf{R}_{b,\ell} \right] \leq \frac{96K^2 |\mathbf{B}_{b,0}|}{(1-\alpha_{\max})\Delta_{\max}} \log T.$$

This concludes our proof. ■

E. Technical Lemmas

Lemma 8. Let $(X_1)_{i=1}^n$ be n i.i.d. Bernoulli random variables, $\bar{\mu} = \sum_{i=1}^n X_i$, and $\mu = \mathbb{E}[\bar{\mu}]$. Then

$$P(\bar{\mu} \geq \mu + \varepsilon) \leq \exp[-nD_{\text{KL}}(\mu + \varepsilon \parallel \mu)]$$

for any $\varepsilon \in [0, 1 - \mu]$, and

$$P(\bar{\mu} \leq \mu - \varepsilon) \leq \exp[-nD_{\text{KL}}(\mu - \varepsilon \parallel \mu)]$$

for any $\varepsilon \in [0, \mu]$.

Proof. We only prove the first claim. The other claim follows from symmetry.

From inequality (2.1) of [Hoeffding \(1963\)](#), we have that

$$P(\bar{\mu} \geq \mu + \varepsilon) \leq \left[\left(\frac{\mu}{\mu + \varepsilon} \right)^{\mu + \varepsilon} \left(\frac{1 - \mu}{1 - (\mu + \varepsilon)} \right)^{1 - (\mu + \varepsilon)} \right]^n$$

for any $\varepsilon \in [0, 1 - \mu]$. Now note that

$$\begin{aligned} \left[\left(\frac{\mu}{\mu + \varepsilon} \right)^{\mu + \varepsilon} \left(\frac{1 - \mu}{1 - (\mu + \varepsilon)} \right)^{1 - (\mu + \varepsilon)} \right]^n &= \exp \left[n \left[(\mu + \varepsilon) \log \frac{\mu}{\mu + \varepsilon} + (1 - (\mu + \varepsilon)) \log \frac{1 - \mu}{1 - (\mu + \varepsilon)} \right] \right] \\ &= \exp \left[-n \left[(\mu + \varepsilon) \log \frac{\mu + \varepsilon}{\mu} + (1 - (\mu + \varepsilon)) \log \frac{1 - (\mu + \varepsilon)}{1 - \mu} \right] \right] \\ &= \exp[-nD_{\text{KL}}(\mu + \varepsilon \parallel \mu)]. \end{aligned}$$

This concludes the proof. ■

Lemma 9. For any $c, p, q \in [0, 1]$,

$$c(1 - \max\{p, q\})D_{\text{KL}}(p \parallel q) \leq D_{\text{KL}}(cp \parallel cq) \leq cD_{\text{KL}}(p \parallel q). \quad (14)$$

Proof. The proof is based on differentiation. The first two derivatives of $D_{\text{KL}}(cp \parallel cq)$ with respect to q are

$$\frac{\partial}{\partial q} D_{\text{KL}}(cp \parallel cq) = \frac{c(q - p)}{q(1 - cq)}, \quad \frac{\partial^2}{\partial q^2} D_{\text{KL}}(cp \parallel cq) = \frac{c^2(q - p)^2 + cp(1 - cp)}{q^2(1 - cq)^2};$$

and the first two derivatives of $cD_{\text{KL}}(p \parallel q)$ with respect to q are

$$\frac{\partial}{\partial q} [cD_{\text{KL}}(p \parallel q)] = \frac{c(q - p)}{q(1 - q)}, \quad \frac{\partial^2}{\partial q^2} [cD_{\text{KL}}(p \parallel q)] = \frac{c(q - p)^2 + cp(1 - p)}{q^2(1 - q)^2}.$$

The second derivatives show that $D_{\text{KL}}(cp \parallel cq)$ and $cD_{\text{KL}}(p \parallel q)$ are convex in q for any p . Their minima are at $q = p$.

Now we fix p and c , and prove (14) for any q . The upper bound is derived as follows. Since

$$D_{\text{KL}}(cp \parallel cx) = cD_{\text{KL}}(p \parallel x) = 0$$

when $x = p$, the upper bound holds when $cD_{\text{KL}}(p \parallel x)$ increases faster than $D_{\text{KL}}(cp \parallel cx)$ for any $p < x \leq q$, and when $cD_{\text{KL}}(p \parallel x)$ decreases faster than $D_{\text{KL}}(cp \parallel cx)$ for any $q \leq x < p$. This follows from the definitions of $\frac{\partial}{\partial x} D_{\text{KL}}(cp \parallel cx)$ and $\frac{\partial}{\partial x} [cD_{\text{KL}}(p \parallel x)]$. In particular, both derivatives have the same sign and $|\frac{\partial}{\partial x} D_{\text{KL}}(cp \parallel cx)| \leq |\frac{\partial}{\partial x} [cD_{\text{KL}}(p \parallel x)]|$ for any feasible $x \in [\min\{p, q\}, \max\{p, q\}]$.

The lower bound is derived as follows. The ratio of $\frac{\partial}{\partial x} [cD_{\text{KL}}(p \parallel x)]$ and $\frac{\partial}{\partial x} D_{\text{KL}}(cp \parallel cx)$ is bounded from above as

$$\frac{\frac{\partial}{\partial x} [cD_{\text{KL}}(p \parallel x)]}{\frac{\partial}{\partial x} D_{\text{KL}}(cp \parallel cx)} = \frac{1 - cx}{1 - x} \leq \frac{1}{1 - x} \leq \frac{1}{1 - \max\{p, q\}}$$

for any $x \in [\min\{p, q\}, \max\{p, q\}]$. Therefore, we get a lower bound on $D_{\text{KL}}(cp \parallel cx)$ when we multiply $cD_{\text{KL}}(p \parallel x)$ by $1 - \max\{p, q\}$. ■