
Almost Optimal Anytime Algorithm for Batched Multi-Armed Bandits

Tianyuan Jin¹ Jing Tang² Pan Xu³ Keke Huang¹ Xiaokui Xiao¹ Quanquan Gu³

Abstract

In batched multi-armed bandit problems, the learner can adaptively pull arms and adjust strategy in batches. In many real applications, not only the regret but also the batch complexity need to be optimized. Existing batched bandit algorithms usually assume that the time horizon T is known in advance. However, many applications involve an unpredictable stopping time. In this paper, we study the anytime batched multi-armed bandit problem. We propose an anytime algorithm that achieves the asymptotically optimal regret for exponential families of reward distributions with $\mathcal{O}(\log \log T \cdot \text{ilog}^\alpha(T))$ ¹ batches, where $\alpha \in \mathcal{O}_T(1)$. Moreover, we prove that for any constant $c > 0$, no algorithm can achieve the asymptotically optimal regret within $c \log \log T$ batches.

1 Introduction

The multi-armed bandit (MAB) problem provides an elementary model for exploration-exploitation tradeoffs and finds many real applications such as medical trials (Thompson, 1933; Perchet et al., 2016), crowdsourcing (Kittur et al., 2008; Zhou et al., 2014), and marketing (Bertsimas & Mersereau, 2007; Vaswani et al., 2017). The problem is typically described as a game between the agent and the environment. The game proceeds in a total of T time steps. At each time step t , the agent pulls an arm A_t from the arm set $[K]$ with the goal of maximizing the accumulated reward over T time steps. Ideally, the agent can observe the immediate feedback of each pull, e.g., reward, and exploit

it to guide the next action. However, this is impractical for many real applications where the number of interactions between the agent and environment is limited. For example, in clinical trials, typically, it takes some time to test the efficacy of a treatment on a patient. It is thus computationally prohibitive to conduct the experiments in fully sequential. Instead, patients are usually grouped into batches, and each batch of patients are tested in a parallel manner. In such a case, the outcomes are unavailable till the end of each batch. As another example, in online advertising, the agent cannot immediately update her strategy upon receiving the feedback, since there may be a large amount of responses in every second.

Perchet et al. (2016) modeled the above problem as the batched multi-armed bandit problem. In such problems, the time horizon T is split into a small number of batches, and the outcomes are only revealed at the end of each batch. Previous batched bandit algorithms (Perchet et al., 2016; Gao et al., 2019; Esfandiari et al., 2019; Jin et al., 2020) all assume that the time T is known in advance. However, many real-world applications involve an unpredictable stopping time. Again, consider the clinical trials example, T may correspond to the number of participated patients in the test for a certain period; and in the online advertisement example, T may correspond to the number of visitors of a website for a certain period. In both cases, designing an anytime bandit algorithm is imperative.

Motivated by the above observations, in this paper, we study the anytime batched multi-armed bandit problem, where the horizon length T is unknown ahead of time. In particular, we have a set $[K] = \{1, 2, \dots, K\}$ of K arms, where each arm i is associated with a reward distribution of some canonical one-dimensional exponential family with mean μ_i . We assume the best arm is unique. Without loss of generality, we assume that arm 1 has the maximum expected reward throughout the paper, i.e., $\mu_1 > \mu_i$ for any $i \in [K] \setminus \{1\}$. The pulls of each arm yield rewards which are independent and identically distributed (i.i.d.) samples from the arm's distribution. Furthermore, the time horizon T is divided into batches represented by a grid $\mathcal{T} = \{t_1, t_2, \dots\}$, which means after j -th batch, the total number of pulls of all arms reaches t_j . In this paper, we study the static grid setting (Perchet et al., 2016; Gao et al., 2019), i.e., t_1, t_2, \dots are predefined numbers. At each time step t , there exists a

¹School of Computing, National University of Singapore, Singapore ²Data Science and Analytics Thrust, The Hong Kong University of Science and Technology ³Department of Computer Science, University of California, Los Angeles, CA 90095, USA. Correspondence to: Xiaokui Xiao <xkxiao@nus.edu.sg>, Quanquan Gu <qgu@cs.ucla.edu>.

¹Notation $\text{ilog}^\alpha(T)$ is the result of iteratively applying the logarithm function on T for α times, e.g., $\text{ilog}^3(T) = \log \log \log T$.

unique j such that $t_{j-1} < t \leq t_j$, and the agent makes a decision on pulling arm A_t based on all the outcomes up to time t_{j-1} . The ultimate goal is to minimize the regret, which is defined as the expected cumulative difference between playing the best arm and playing the arm according to the strategy. The formal definition is given as follows.

$$R_T = T \cdot \mu_1 - \mathbb{E} \left[\sum_{t=1}^T \mu_{A_t} \right].$$

Lai & Robbins (1985) shows that for distributions that are continuously parameterized by their means,

$$\liminf_{T \rightarrow \infty} \frac{R_T}{\log T} \geq \sum_{i \in [K] \setminus \{1\}} \frac{\Delta_i}{\text{kl}(\mu_i, \mu_1)}, \quad (1)$$

where $\Delta_i := \mu_1 - \mu_i$ and $\text{kl}(\mu, \mu')$ is the Kullback-Leibler divergence between two distributions with mean μ and mean μ' . We refer to $\lim_{T \rightarrow \infty} \frac{R_T}{\log T}$ as the asymptotic regret rate, and say that the algorithm is asymptotically optimal if its asymptotic regret rate matches the right hand side of (1).

The well-known algorithms such as KL-UCB (Garivier & Cappé, 2011) and Thompson Sampling (Korda et al., 2013) are shown to be asymptotically optimal in anytime setting. Nevertheless, these algorithms are fully sequential, which require $\mathcal{O}(T)$ batches in batched bandits. A very recent work (Jin et al., 2020) proposes asymptotically optimal algorithms for the 2-armed bandit problem with sub-Gaussian rewards, requiring $\mathcal{O}(1)$ expected batches if T is known. However, in the anytime setting, their algorithm needs at least $\Omega(\log T)$ batches even for 2-armed bandits. Therefore, a natural question is:

How many batches are needed for anytime K -armed bandit algorithms to achieve the asymptotically optimal regret?

On the other hand, Besson & Kaufmann (2018) conjectured that no anytime algorithm can achieve the asymptotically optimal regret with the exponential time grid (i.e., $t_i = a^{b^i}$ for some constants a and b) incurring $\mathcal{O}(\log \log T)$ batches. However, confirming this conjecture theoretically still remains an open problem. This gives rise to another question:

What is the fundamental limit in batch complexity of anytime K -armed bandit algorithms for achieving the asymptotically optimal regret?

In this paper, we answer the above two questions through the lens of both the upper bound and the lower bound of batch complexities for anytime K -armed bandit algorithms that achieve the asymptotically optimal regret.

Contributions. Our results can be summarized as follows:

- **(Upper Bound)** We propose an anytime algorithm BABA for batched multi-armed bandits where the reward distributions are from exponential families. We prove

that BABA is asymptotically optimal and only requires $\mathcal{O}(\log \log T \cdot \text{ilog}^\alpha(T))$ batches, where $\alpha \in \mathcal{O}_T(1)$ is a constant and $\text{ilog}^\alpha(T)$ iteratively applies the logarithm function on T for α times.

- **(Lower Bound)** We prove that in the anytime setting, for any positive constant c , no bandit algorithm can achieve the asymptotic optimality within $c \log \log T$ batches. This is the first lower bound for anytime batched bandit algorithms in the literature. Our lower bound is almost tight since it matches the upper bound of our BABA algorithm up to an iterative logarithm factor.
- We empirically evaluate our proposed algorithm and show that it enjoys comparable performance with the fully sequential algorithm KL-UCB in terms of regret while requiring significantly fewer batches.

2 Preliminaries

In this section, we first review the previous work related to ours. We then introduce the definition of exponential families and some useful properties. Finally, we give some notations that will be frequently used.

2.1 Previous Results

MAB. The MAB problem provides an elementary model for a class of sequential optimization problems, which has been extensively studied since the seminal work by Thompson (1933). In the fully adaptive setting, a large body of research has analyzed the regret (Audibert & Bubeck, 2009; Garivier & Cappé, 2011; Korda et al., 2013; Degenne & Perchet, 2016; Agrawal & Goyal, 2017; Kaufmann et al., 2018; Lattimore, 2018). We refer interested readers to the book by Lattimore & Szepesvári (2020) for a comprehensive introduction of bandit algorithms and various applications of MAB.

Batched MAB. Cesa-Bianchi et al. (2013) studied the batched bandit problem under the name of switching cost and show that $\mathcal{O}(\log \log T)$ batches are sufficient for achieving the near-optimal *minimax* regret bound. Perchet et al. (2016) further prove that such batch complexity is sufficient and necessary for achieving the optimal minimax regret for the 2-armed bandit problem, which is later generalized to the K -armed case (Gao et al., 2019). Perchet et al. (2016); Gao et al. (2019); Esfandiari et al. (2019) show that $\mathcal{O}(\log T)$ batches are sufficient for achieving the near-optimal instance-dependent regret bound. However, there exists a multiplicative constant between their regret bounds and the optimal bound, which makes their algorithms sub-optimal in the asymptotic sense. Jin et al. (2020) propose DETC, consisting of two exploration and two exploitation stages, that can achieve the asymptotic optimality with $\mathcal{O}(1)$ expected batches for the 2-armed case with sub-Gaussian rewards. However, the generalization to K -armed bandit and exponential families of reward distributions is non-trivial

and unclear. More importantly, all the aforementioned studies assume that T is known in advance, while we consider anytime batched MAB with unpredictable stopping time T .

Anytime Algorithm. An effective technique to construct an anytime algorithm from a non-anytime algorithm is the *doubling trick* strategy (Auer et al., 1995). At the i -th round/epoch, the doubling trick strategy guesses $T = a^i$ (referred to as geometric doubling trick) or $T = a^{b^i}$ (referred to as exponential doubling trick). For geometric doubling trick, it costs at least $\log T$ batches. For example, the anytime version of DETC requires $\mathcal{O}(\log T)$ exploration and exploitation rounds by guessing $T = 2^i$ at the i -th round, and each round takes $\Omega(1)$ batches. For exponential doubling trick with $\mathcal{O}(\log \log T)$ batches, Besson & Kaufmann (2018) conjecture that it cannot achieve the asymptotically optimal regret, which is confirmed by our lower bound in this paper. Motivated by the deficiencies of the above two doubling tricks, we present a new trick strategy that is asymptotically optimal and takes $\mathcal{O}(\log \log T \cdot \text{ilog}^\alpha(T))$ batches. Compared with the existing tricks, the number of batches incurred by our trick is significantly smaller than geometric doubling trick and is slightly larger than exponential doubling trick.

2.2 Exponential Families

An exponential family is a parametric set of probability distributions $\{\nu_\theta : \theta \in \Theta\}$ dominated by a measure ρ on \mathbb{R} , with density given by

$$\frac{d\nu_\theta}{d\rho}(x) = \exp(x\theta - b(\theta)),$$

where $b(\theta) = \log \int_{\mathbb{R}} e^{x\theta} d\rho(x)$ and $\Theta = \{\theta \in \mathbb{R} : b(\theta) < \infty\}$. Exponential families have the following properties.

$$b'(\theta) = \mathbb{E}[\nu_\theta] \quad \text{and} \quad 0 < b''(\theta) = \text{Var}(\nu_\theta),$$

where $b'(\theta)$ and $b''(\theta)$ are the first derivative and second derivative of $b(\theta)$ with respect to θ , respectively. A direct computation gives the Kullback-Leibler (KL) divergence as

$$\text{KL}(\nu_\theta, \nu_{\theta'}) = b(\theta') - b(\theta) - b'(\theta)(\theta' - \theta).$$

Let $\mu = b'(\theta)$ and $\text{kl}(\mu, \mu') := \text{KL}(\nu_\theta, \nu_{\theta'})$. In this paper, we assume the variance is bounded, i.e.,

$$0 < b''(\theta) \leq V < +\infty.$$

We have the following property on the KL divergence.

Proposition 1. *For all μ and μ' , we have*

$$\text{kl}(\mu, \mu') \geq (\mu - \mu')^2 / (2V). \quad (2)$$

In addition, for $\epsilon > 0$ and $\mu \leq \mu' - \epsilon$, we can obtain that

$$\begin{aligned} \text{kl}(\mu, \mu') &\geq \text{kl}(\mu, \mu' - \epsilon), \\ \text{and } \text{kl}(\mu, \mu') &\leq \text{kl}(\mu - \epsilon, \mu'). \end{aligned} \quad (3)$$

Algorithm 1: Batched Anytime Bandit Alg. (BABA)

Input: a set of K arms and parameters α and I_1

```

1 initialize  $t \leftarrow 0, r \leftarrow 1, c_0 \leftarrow 1$ ;
2 while experiment proceeds do
3   Step I: perform UNIFORMEXPLORATION;
4   Step II: perform INITIALEXPLOITATION;
5   Step III: perform OPTIMISTICEXPLORATION;
6   Step IV: perform CONFIDENTEXPLORATION;
7   Step V: perform CONFIDENTEXPLOITATION;
8    $r \leftarrow r + 1$ ;
9    $I_r \leftarrow f(I_{r-1})$ ;

```

Interested readers are referred to Appendix A for the proof of Proposition 1.

Exponential families include many of the most common distributions, such as Gaussian, Bernoulli, exponential, etc. In particular, for Gaussian distribution with known variance σ^2 by choosing $V = \sigma^2$, $\text{kl}(\mu, \mu') = (\mu - \mu')^2 / (2\sigma^2)$; for Bernoulli distribution by choosing $V = 1/4$, $\text{kl}(\mu, \mu') = \mu \log(\mu/\mu') + (1 - \mu) \log((1 - \mu)/(1 - \mu'))$; and for exponential distribution with known parameter λ by choosing $V = 1/\lambda^2$, $\text{kl}(\mu, \mu') = \log(\mu) - \log(\mu') + \mu'/\mu - 1$.

2.3 Notations

Denote by $\text{ilog}^m(x)$ the result of iteratively applying the logarithm function on x for m times, i.e., $\text{ilog}^m(x) = \max\{\log \text{ilog}^{m-1}(x), 0\}$ for any $x > 0$ and $m \in \mathbb{N}$. We also define $\text{ilog}^0(x) = x$. We define $\text{kl}_+(p, q) := \text{kl}(p, q) \mathbb{1}(p \leq q)$, where $\mathbb{1}(\cdot)$ is the indicator function. We use Δ_i to denote the gap between arm 1 and arm i , i.e., $\Delta_i = \mu_1 - \mu_i$. Let $\hat{\mu}_i(t)$ be the average reward of arm i at time t and $\hat{\mu}_{i,s}$ be the average reward of arm i after its s 's pull. Let $T_i(t)$ be the number of pulls of arm i at time step t , i.e., $T_i(t) = \sum_{\ell=1}^t \mathbb{1}(A_\ell = i)$. Throughout the paper, we adopt the standard asymptotic notations. In particular, we use $f(\cdot) \lesssim g(\cdot)$ to denote $f(\cdot) = \mathcal{O}_T(g(\cdot))$ and $f(\cdot) \gtrsim g(\cdot)$ to denote $f(\cdot) = \Omega_T(g(\cdot))$.

3 The Proposed Algorithm

3.1 Overview

Algorithm 1 presents the framework of our algorithm, referred to as BABA. In particular, our algorithm guesses T in epochs and proceeds in five batches for each epoch.

In the following, we first introduce two core functions that are essential for constructing our time grid \mathcal{T} .

$$\hat{f}(x) = \max\{\lceil x^{1+1/(1+\text{ilog}^\alpha x)} \rceil, 2x\}, \quad (4)$$

$$g(x) = \lceil \log x / \log \log x \rceil, \quad (5)$$

where $\alpha \geq 3$ and $\alpha \in \mathcal{O}_T(1)$ is a constant. That is, at the r -th epoch, we guess $T = f^{(r)}(I_1)$, where $f^{(r)}(\cdot)$ iteratively

Algorithm 2: UNIFORMEXPLORATION

```

1 for  $i = 1, 2, \dots, K$  do
2   while  $T_i(t) \leq g(I_r)$  do
3     pull arm  $i$ ;
4      $t \leftarrow t + 1$ ;
5 while  $t \leq I_{r-1} + K \cdot g(I_r)$  do
6   pull arm  $c_{r-1}$ ;
7    $t \leftarrow t + 1$ ;
    
```

Algorithm 3: INITIALEXPLOITATION

```

1  $a_{1,r} \leftarrow \arg \max_{i \in [K]} \widehat{\mu}_{i,g(I_r)}$ ;
2  $\ell \leftarrow 1$ ;
3 while  $\ell \leq \log^2 I_r$  do
4   pull arm  $a_{1,r}$ ;
5    $t \leftarrow t + 1, \ell \leftarrow \ell + 1$ ;
    
```

applies function f for r times and I_1 is an input parameter satisfying

$$Kg(I_1) + (2K + 1) \log^2 I_1 < I_1. \quad (6)$$

Let $I_r = f^{(r)}(I_1)$. At the r -th epoch, the first four batches pull the arms exactly $Kg(I_r)$, $\log^2 I_r$, $K \cdot \log^2 I_r$ and $K \cdot \log^2 I_r$ times (for ease of presentation, we assume $\log I_r \in \mathbb{N}^+$), respectively, while the fifth batch pulls the arms until a total number of I_r times is pulled. Our time grid \mathcal{T} is given as $\mathcal{T} = \{t_{1,1}, \dots, t_{5,1}, t_{1,2}, \dots, t_{5,2}, \dots\}$, where $t_{j,r}$ denote the checkpoint at j -th step of the r -th epoch for any $j \in [5]$, defined as follows:

$$\begin{aligned}
 t_{1,r} &= I_{r-1} + Kg(I_r), \\
 t_{2,r} &= I_{r-1} + Kg(I_r) + \log^2 I_r, \\
 t_{3,r} &= I_{r-1} + Kg(I_r) + (K + 1) \log^2 I_r, \\
 t_{4,r} &= I_{r-1} + Kg(I_r) + (2K + 1) \log^2 I_r, \\
 \text{and } t_{5,r} &= I_r.
 \end{aligned}$$

It is trivial to see that \mathcal{T} is a static time grid. Note that the trick $\{f^{(r)}(I_1)\}_{r \geq 1}$ grows (i) faster than geometric doubling trick which results in far less number of batches than $\log T$, but (ii) a litter bit slower than exponential doubling trick which ensures the asymptotically optimal regret.

3.2 Detailed Design

Next, we elaborate the details of the five steps in each epoch.

Step I. UNIFORMEXPLORATION (Algorithm 2) shows the first step, which pulls the arms a total of $Kg(I_r)$ times. Specifically, at the r -th epoch, for every arm i that is pulled less than $g(I_r)$ times, we pull arm i till reaching a total of $g(I_r)$ times (Lines 1–4). In addition, we pull the “best arm” c_{r-1} found after the $(r - 1)$ -th epoch until the total number

Algorithm 4: OPTIMISTICEXPLORATION

```

1  $s_1 \leftarrow T_{a_{1,r}}(t)$ ;
2  $\{a_{2,r}, \dots, a_{K,r}\} \leftarrow [K] \setminus \{a_{1,r}\}$ ;
3 for  $i = 2, 3, \dots, K$  do
4   while  $T_{a_{i,r}}(t) \leq \min\{\delta_{i,r}, \log^2 I_r\}$  do
5     pull arm  $a_{i,r}$ ;
6      $t \leftarrow t + 1$ ;
7  $\mathcal{F} \leftarrow \perp$ ;
8 for  $i = 2, 3, \dots, K$  do
9    $s_i \leftarrow \min\{\log^2 I_r, T_{a_{i,r}}(t)\}$ ;
10  if  $\text{kl}_+(\widehat{\mu}_{a_{i,r},s_i}, \widehat{\mu}_{a_{1,r},s_1}) < \frac{\log(I_r \cdot \log^2 I_r)}{s_i}$  then
11     $\mathcal{F} \leftarrow \top$ ;
12    break;
13 while  $t \leq I_{r-1} + K \cdot g(I_r) + \log^2 I_r + K \cdot \log^2 I_r$  do
14   pull arm  $c_{r-1}$ ;
15    $t \leftarrow t + 1$ ;
    
```

of pulls of all arms reaches $I_{r-1} + Kg(I_r)$ so that this batch pulls the arms $Kg(I_r)$ times (Lines 5–7).

Purpose. Let $a_{1,r}$ be the arm with the largest average reward when every arm is pulled exactly $g(I_r)$ times, which is likely to be the best arm. In fact, we will show that $\mathbb{P}(a_{1,r} = 1) \geq 1 - 1/\log^2 I_r$, which supports us to pull arm $a_{1,r}$ additional $\log^2 I_r$ times while keeping the optimal asymptotic regret.

Step II. INITIALEXPLOITATION (Algorithm 3) shows the second step, which simply pulls $\log^2(I_r)$ times of arm $a_{1,r}$.

Purpose. When $a_{1,r}$ is pulled $\log^2 I_r$ times, the sample average of arm $a_{1,r}$ will concentrate on its true mean. This ensures that when we explore whether other arms have the potential to be the best arm, we do not pull $a_{1,r}$ as its estimated mean is sufficiently accurate.

Step III. OPTIMISTICEXPLORATION (Algorithm 4) shows the third step, which pulls the arms $K \log^2 I_r$ times. Define $\{a_{i,r}\}_{i \geq 2} := [K] \setminus \{a_{1,r}\}$ as the set of other arms except for $a_{1,r}$. Let $\epsilon_r := 1/\log \log I_r$. For $i \geq 2$, define

$$\delta_{i,r} := \frac{\log(I_r \cdot \log^2(I_r))}{\text{kl}(\widehat{\mu}_{a_{i,r},g(I_r)} + \epsilon_r, \widehat{\mu}_{a_{1,r},s_1} - \epsilon_r)}, \quad (7)$$

where s_1 is the number of pulls of $a_{1,r}$ after Step II. Then, we pull $a_{i,r}$ till $T_{a_{i,r}}(t) \geq \min\{\delta_{i,r}, \log^2 I_r\}$ for $i \geq 2$ (Lines 3–6). We further check the following condition

$$\text{kl}_+(\widehat{\mu}_{a_{i,r},s_i}, \widehat{\mu}_{a_{1,r},s_1}) < \frac{\log(I_r \log^2 I_r)}{s_i}, \quad (8)$$

where $s_i = \min\{\log^2 I_r, T_{a_{i,r}}(t)\}$. If for some i , (8) holds, then we set $\mathcal{F} = \top$; otherwise we set $\mathcal{F} = \perp$ (Lines 7–12). Finally, we pull arm c_{r-1} until a total number of $K \cdot \log^2 I_r$ pulls are pulled in this batch (Lines 13–15).

Algorithm 5: CONFIDENTEXPLORATION

```

1 if  $\mathcal{F} = \perp$  then
2    $\ell \leftarrow 1$ ;
3   while  $\ell \leq K \cdot \log^2(I_r)$  do
4     pull arm  $a_{1,r}$ ;
5      $t \leftarrow t + 1$ ;
6    $c_r \leftarrow a_{1,r}$ ;
7 else
8   for  $i = 1, 2, 3, \dots, K$  do
9      $\ell \leftarrow 1$ ;
10    while  $\ell \leq \log^2 I_r$  do
11      pull arm  $i$ ;
12       $t \leftarrow t + 1$ ;
13   $c_r \leftarrow \arg \max_{i \in [K]} \widehat{\mu}_i(t)$ ;

```

Purpose. In this batch, we try to explore whether other arms rather than $a_{1,r}$ have potential to be the best arm. In particular, if $\text{kl}_+(\widehat{\mu}_{a_{i,r},s_i}, \widehat{\mu}_{a_{1,r},s_1})$ is small enough to satisfy (8) for some i , we know that $\widehat{\mu}_{a_{i,r},s_i} > \widehat{\mu}_{a_{1,r},s_1}$ or the difference between $\widehat{\mu}_{a_{i,r},s_i}$ and $\widehat{\mu}_{a_{1,r},s_1}$ is very small, which indicates that $a_{i,r}$ has potential to be the best arm. Then, we will further pull each arm in the future to determine whether $a_{i,r}$ is better than $a_{1,r}$. Otherwise if (8) does not hold for every i , we will confident that $a_{1,r}$ is the best arm. Finally, we pull best arm c_{r-1} found after the $(r-1)$ -th epoch to exhaust the budget of this batch.

Step IV. CONFIDENTEXPLORATION (Algorithm 5) shows the fourth step, which pulls the arms a total of $K \log^2 I_r$ times. In particular, if $\mathcal{F} = \perp$, we directly pull arm $a_{1,r}$ a total of $K \log^2(I_r)$ times (Lines 2–5) and update the new best arm $c_r = a_{1,r}$ (Line 6); otherwise, we pull every arm $\log^2(I_r)$ times (Lines 8–12) and update the new best arm $c_r = \arg \max_{i \in [K]} \widehat{\mu}_i(t)$ (Line 13).

Purpose. Intuitively, if $\mathcal{F} = \perp$, i.e., (8) fails for all $i \geq 2$, we can show that $\mathbb{P}(a_{1,r} = 1) \geq 1 - 1/I_r$, which ensures that the regret of pulling $a_{1,r}$ additional I_r times is bounded in the optimal range. Hence, if $\mathcal{F} = \perp$, all the budget of $K \cdot (I_r)^2$ in this batch is used on arm $a_{1,r}$ and we set $c_r = a_{1,r}$ for future pulls. On the other hand, if $\mathcal{F} = \top$, the arm $a_{1,r}$ may not be the optimal arm. We then pull every arm $\log^2 I_r$ times and update c_r to the arm with largest average reward.

Step V. CONFIDENTEXPLOITATION (Algorithm 6) shows the fifth step, which pulls the “best arm” observed so far until the total number of pulls reaches I_r .

Purpose. After the first four steps, we are confident now that c_r is the best arm. Thus, we keep pulling c_r to optimize regret until the budget of this batch is exhausted.

Algorithm 6: CONFIDENTEXPLOITATION

```

1 while  $t \leq I_r$  do
2   pull arm  $c_r$ ;
3    $t \leftarrow t + 1$ ;

```

4 Main Results

Now, we present our main theoretical results, which consist of an upper bound for the batch complexity of Algorithm 1, and a lower bound for the batch complexity of anytime batched bandit algorithms that attain the asymptotically optimal regret in (1).

4.1 Upper Bound

The batch complexity of Algorithm 1 is given as follows.

Theorem 2 (Batch Complexity). *For any input $\alpha \in \mathcal{O}_T(1)$ and I_1 satisfying (6), the number of batches for Algorithm 1 is $O(\log \log T \cdot \text{ilog}^\alpha(T))$.*

Furthermore, the following theorem shows that Algorithm 1 is asymptotically optimal for an unknown horizon T .

Theorem 3 (Regret). *For any input $\alpha \in \mathcal{O}_T(1)$ and I_1 satisfying (6), Algorithm 1 achieves the asymptotically optimal regret, i.e., $\lim_{T \rightarrow \infty} \frac{R_T}{\log T} = \sum_{i=2}^K \frac{\Delta_i}{\text{kl}(\mu_i, \mu_1)}$.*

Comparison with Previous Work. Compared with the fully sequential algorithms, such as KL-UCB (Garivier & Cappé, 2011) and Thompson Sampling (Korda et al., 2013) that achieve the asymptotically optimal regret with $O(T)$ batches, our algorithm only needs $O(\log \log T \cdot \text{ilog}^\alpha(T))$ batches while maintaining the asymptotically optimal regret. Compared with Anytime-DETC (Algorithm 5 in Jin et al. (2020)) that incurs at least $\Omega(\log T)$ batches, our algorithm not only significantly improves the batch complexity but also expands the applicability since Anytime-DETC only applies to 2-armed bandit with sub-Gaussian rewards, whereas our algorithm generalizes to K -armed bandit with exponential families of reward distributions.

4.2 Lower Bound

Besson & Kaufmann (2018) conjectures that *the geometric doubling trick can never bring the right constant in asymptotic regret bound* in (1). We present a lower bound that theoretically confirms this conjecture.

Theorem 4. *For any static grid and constant $c > 0$, no anytime algorithm can achieve the asymptotically optimal regret in (1) within $c \log \log T$ batches.*

It is worth noting that this is the first lower bound for anytime batched bandit problem in the literature. We observe that the lower bound in Theorem 4 matches the upper bound in Theorem 2 within an $\text{ilog}^\alpha(T)$ factor, which shows that our batch complexity is almost optimal for achieving the

asymptotic optimality in the anytime setting.

We note that [Perchet et al. \(2016\)](#); [Gao et al. \(2019\)](#) also proved certain lower bounds for the batched bandit problem. However, their lower bounds are significantly different from ours, since we focus on anytime algorithm tailored for minimizing the asymptotic regret with an unknown T , whereas they consider minimax regret or problem-dependent regret with a finite known T .

5 Theoretical Analysis

Now we present the proof of the upper bound results. The proof the lower bound can be found in [Appendix C](#).

5.1 Analysis of Batch Complexity

Proof of Theorem 2. According to the definition of f , it is easy to verify that

$$\begin{aligned} I_{r+n} &\geq (I_{r+n-1})^{1+1/(1+\text{ilog}^\alpha I_{r+n-1})} \\ &\geq (I_{r+n-1})^{1+1/(1+\text{ilog}^\alpha I_r)} \geq (I_r)^{(1+1/(1+\text{ilog}^\alpha I_r))^n}. \end{aligned}$$

For $n = \lceil 1 + \text{ilog}^\alpha I_r \rceil$, we have $I_{r+n} \geq (I_r)^2$. Therefore,

$$f^{\lceil 1 + \text{ilog}^\alpha T \rceil}(I_r) \geq f^{\lceil 1 + \text{ilog}^\alpha I_r \rceil}(I_r) \geq (I_r)^2. \quad (9)$$

This implies that it needs at most $\lceil 1 + \text{ilog}^\alpha T \rceil$ epochs for increasing I_r from $(I_1)^{2^\ell}$ to $(I_1)^{2^{\ell+1}}$. Moreover, when $I_1 \geq 2$, $\ell^* = \log_2 \log_2 T$ suffices to ensure $(I_1)^{2^{\ell^*}} \geq T$, which indicates that the algorithm runs at most $\ell^* \cdot \lceil 1 + \text{ilog}^\alpha T \rceil$ epochs. Besides, the number of epochs is proportional to the number of batches. Therefore, the total number of batches is $\mathcal{O}(\log T \log T \cdot \text{ilog}^\alpha T)$. \square

5.2 Analysis of Regret

Proof of Theorem 3. Let T_i be the total number of pulls of arm i in [Algorithm 1](#). Then, the regret can be rewritten as

$$R_T = \sum_{i \geq 2} \mathbb{E}[T_i \Delta_i].$$

Therefore, it suffices to prove the elementary result such that for each $i \geq 2$

$$\lim_{T \rightarrow \infty} \frac{\mathbb{E}[T_i]}{\log T} = \frac{1}{\text{kl}(\mu_i, \mu_1)}. \quad (10)$$

For ease of presentation, for the r -th epoch, we fix a suboptimal arm i and define some notations as follows.

- $Y_1(r)$: the number of pulls of arm i , excluding pulling c_{r-1} , in [Step I](#) (Lines 1–4 in [Algorithm 2](#)).
- $Y_2(r)$: the number of pulls of arm i in [Step II](#) ([Algorithm 3](#)).
- $Y_3(r)$: the number of pulls of arm i in [Step III](#) (Line 1–12 in [Algorithm 4](#));

- $Y_4(r)$: the number of pulls of arm i in [Step IV](#) ([Algorithm 5](#)) and [V](#) ([Algorithm 6](#)).
- $Z(r)$: the total number of pulls of arm c_{r-1} when $c_{r-1} = i$ in [Step I](#) (Line 5–7 in [Algorithm 2](#)) and [Step III](#) (Line 13–15 in [Algorithm 4](#)).

In addition, since for $I_r \leq \sqrt{\log T}$, the algorithm plays suboptimal arms at most $\sqrt{\log T}$ times. We use $\sqrt{\log T}$ to bound the number of pulls of arm i when $I_r \leq \sqrt{\log T}$. Let $r' := \min\{r: I_r > \sqrt{\log T}\}$ and r^o be the total number of epochs. Then, we have

$$\begin{aligned} \mathbb{E}[T_i] &= \sum_{r=1}^{r^o} \left(\sum_{j=1}^4 \mathbb{E}[Y_j(r)] + \mathbb{E}[Z(r)] \right) \\ &\leq \sqrt{\log T} + \sum_{r=r'}^{r^o} \left(\sum_{j=1}^4 \mathbb{E}[Y_j(r)] + \mathbb{E}[Z(r)] \right) \\ &= \sqrt{\log T} + \sum_{j=1}^4 \mathbb{E}[Y_j] + \mathbb{E}[Z], \end{aligned}$$

where $Y_j := \sum_{r=r'}^{r^o} Y_j(r)$ and $Z := \sum_{r=r'}^{r^o} Z(r)$. According to [Lemmas 2–6](#), which shall be given later, we have

$$\frac{\mathbb{E}[T_i]}{\log T} = \frac{\sqrt{\log T}}{\log T} + \frac{1 + 1/(1 + \text{ilog}^\alpha T)}{\text{kl}(\mu_i - \epsilon_{r'}, \mu_1 + \epsilon_{r'})} + o_T(1).$$

Note that $\frac{\sqrt{\log T}}{\log T} \rightarrow 0$, $1/(1 + \text{ilog}^\alpha T) \rightarrow 0$, $\epsilon_{r'} \rightarrow 0$ and $o_T(1) \rightarrow 0$ when $T \rightarrow \infty$. Therefore,

$$\lim_{T \rightarrow \infty} \frac{\mathbb{E}[T_i]}{\log T} = \frac{1}{\text{kl}(\mu_i, \mu_1)}.$$

This completes the proof. \square

In the following, we bound $\mathbb{E}[Y_i]$ for $i \in [4]$ and $\mathbb{E}[Z]$, which requires the following concentration bounds for exponential families.

Lemma 1 (Maximal Inequality ([Ménard & Garivier, 2017](#))). *Let N and M be two real numbers in $\mathbb{R}^+ \times \overline{\mathbb{R}^+}$, let $\gamma > 0$, and $\hat{\mu}_n$ be the empirical mean of n random variables i.i.d. according to the distribution $\nu_{b^{-1}(\mu)}$. Then*

$$\mathbb{P}(\exists N \leq n \leq M, \text{kl}_+(\hat{\mu}_n, \mu) \geq \gamma) \leq e^{-N\gamma}. \quad (11)$$

As a consequence, for every $x \leq \mu$,

$$\mathbb{P}(\exists N \leq n \leq M, \hat{\mu}_n \leq x) \leq e^{-N(x-\mu)^2/(2V)}. \quad (12)$$

Meanwhile, for every $x \geq \mu$,

$$\mathbb{P}(\exists N \leq n \leq M, \hat{\mu}_n \geq x) \leq e^{-N(x-\mu)^2/(2V)}. \quad (13)$$

Lemma 2. $\mathbb{E}[Y_1]/\log T = o_T(1)$.

Proof. Be definition, we have

$$I_{r^o-1} < T \leq I_{r^o} = f(I_{r^o-1}) \leq f(T) \leq T^{1+1/(1+i\log^\alpha T)}.$$

Note that the total number of pulls of arm i contributed by Y_1 is at most $g(I_{r^o})$. Then, we have

$$\begin{aligned} Y_1 &\leq g(I_{r^o}) \leq \frac{2 \log(I_{r^o})}{\log \log(I_{r^o})} \\ &\leq \frac{2(1+1/(1+i\log^\alpha T)) \cdot \log T}{\log \log T}. \end{aligned}$$

Since $\alpha = \mathcal{O}_T(1)$, we obtain

$$\frac{\mathbb{E}[Y_1]}{\log T} \leq \frac{2(1+1/(1+i\log^\alpha T))}{\log \log T} = o_T(1). \quad \square$$

Lemma 3. $\mathbb{E}[Y_2]/\log T = o_T(1)$.

Proof. Intuitively, when I_r is sufficiently large, $a_{1,r}$ is arm 1 with high probability, i.e., $\mathbb{P}(a_{1,r} = 1) \geq 1 - \frac{1}{\log^2 I_r}$. Hence, the expected number of pulls of arm i in Step II is small. In the following, we first bound $\mathbb{P}(a_{1,r} = i)$.

Let $\Delta_{\min} := \min_{i \geq 2} \Delta_i$. After pulling $g(I_r)$ times of arm 1, by (12), we have

$$\mathbb{P}(\hat{\mu}_{1,g(I_r)} \leq \mu_1 - \Delta_{\min}/2) \leq e^{-g(I_r)\Delta_{\min}^2/(8V)}. \quad (14)$$

Meanwhile,

$$g(I_r) \geq \frac{\log I_r}{\log \log I_r} \gtrsim \frac{8V \cdot (2 \log \log I_r + 2 \log K)}{\Delta_{\min}^2}.$$

Combining with (14) gives

$$\mathbb{P}(\hat{\mu}_{1,g(I_r)} \leq \mu_1 - \Delta_{\min}/2) \lesssim \frac{1}{2K \log^2 I_r}.$$

Similarly, for arm i , we can get that

$$\mathbb{P}(\hat{\mu}_{i,g(I_r)} \geq \mu_i + \Delta_{\min}/2) \lesssim \frac{1}{2K \log^2 I_r}.$$

As a result,

$$\mathbb{P}(\hat{\mu}_{1,g(I_r)} \leq \hat{\mu}_{i,g(I_r)}) \lesssim \frac{1}{K \log^2 I_r}. \quad (15)$$

Furthermore, we can obtain that

$$\begin{aligned} \mathbb{E}[Y_2] &= \sum_{r=r'}^{r^o} \mathbb{E}[Y_2(r)] \leq \sum_{r=r'}^{r^o} (\log^2 I_r \cdot \mathbb{P}(a_{1,r} = i)) \\ &\leq \sum_{r=r'}^{r^o} (\log^2 I_r \cdot \mathbb{P}(\hat{\mu}_{1,g(I_r)} \leq \hat{\mu}_{i,g(I_r)})) \lesssim r^o \end{aligned}$$

Note that $r^o = \mathcal{O}(\log \log T \cdot i \log^\alpha T)$ by Theorem 2. This implies that

$$\frac{\mathbb{E}[Y_2]}{\log T} = \mathcal{O}_T\left(\frac{\log \log T \cdot i \log^\alpha T}{\log T}\right) = o_T(1). \quad \square$$

²By choosing sufficiently large T , all \lesssim hold simultaneously.

Lemma 4. $\frac{\mathbb{E}[Y_3]}{\log T} = \frac{1+1/(1+i\log^\alpha T)}{\text{kl}(\mu_i - 2\epsilon_{r'}, \mu_1 + 2\epsilon_{r'})} + o_T(1)$.

Proof. Define events

$$\begin{aligned} \mathcal{E}_0(r) &:= \left\{ |\hat{\mu}_{a_{1,r},s_1} - \mu_{a_{1,r}}| < \epsilon_r \right\}, \\ \mathcal{E}_1(r) &:= \{a_{1,r} = 1\}, \\ \mathcal{E}_2(r) &:= \left\{ \forall k \in [K] \setminus \{a_{1,r}\}: |\hat{\mu}_{k,g(I_r)} - \mu_k| < \epsilon_r \right\}, \\ \mathcal{E}(r) &:= \mathcal{E}_0(r) \cap \mathcal{E}_1(r) \cap \mathcal{E}_2(r). \end{aligned}$$

Based on $\mathcal{E}(r)$, we category the epochs into two sets

$$\begin{aligned} S_1 &:= \{r \in [r', r^o]: \mathbb{1}(\mathcal{E}(r)) = 1\}, \\ \text{and } S_2 &:= \{r \in [r', r^o]: \mathbb{1}(\mathcal{E}^c(r)) = 1\}. \end{aligned}$$

Let $a_{i',r} = i$. Thus, for all epochs in S_1 , arm i is pulled at most $\max_{r \in S_1} \delta_{i',r}$ times in Y_3 , i.e.,

$$\sum_{r \in S_1} Y_3(r) \leq \max_{r \in S_1} \delta_{i',r}.$$

Meanwhile, when T is sufficiently large, by definition, for every $r \in [r', r^o]$, we have

$$\epsilon_r = \frac{1}{\log \log I_r} \leq \frac{1}{\log \log \sqrt{\log T}} \leq \frac{\Delta_{\min}}{4}. \quad (16)$$

Thus, for any $r \in S_1$, we have

$$\hat{\mu}_{i,g(I_r)} + \epsilon_r < \mu_i + 2\epsilon_r \leq \mu_1 - 2\epsilon_r < \hat{\mu}_{1,s_1} - \epsilon_r.$$

Then, we can get that

$$\delta_{i',r} \leq \frac{\log(I_r \cdot \log^2 I_r)}{\text{kl}(\mu_i + 2\epsilon_r, \mu_1 - 2\epsilon_r)} \leq \frac{\log(I_{r^o} \cdot \log^2 I_{r^o})}{\text{kl}(\mu_i + 2\epsilon_{r'}, \mu_1 - 2\epsilon_{r'})}.$$

As a result, we have

$$\begin{aligned} \sum_{r \in S_1} Y_3(r) &\leq \max_{r \in S_1} \delta_{i',r} \leq \frac{\log(I_{r^o} \cdot \log^2 I_{r^o})}{\text{kl}(\mu_i + 2\epsilon_{r'}, \mu_1 - 2\epsilon_{r'})} \\ &\leq \frac{(1+1/(1+i\log^\alpha T)) \cdot \log T + 2 \log \log f(T)}{\text{kl}(\mu_i + 2\epsilon_{r'}, \mu_1 - 2\epsilon_{r'})}, \quad (17) \end{aligned}$$

where the last inequality is due to the fact that $I_{r^o} \leq f(T)$.

On the other hand,

$$\sum_{r \in S_2} Y_3(r) \leq \sum_{r \in S_2} \log^2 I_r,$$

since arm i is pulled at most $\log^2 I_r$ time in $Y_3(r)$. Hence,

$$\mathbb{E}\left[\sum_{r \in S_2} Y_3(r)\right] \leq \sum_{r=r'}^{r^o} \mathbb{P}(\mathcal{E}^c(r)) \log^2 I_r.$$

According to the definition of $\mathcal{E}^c(r)$, we have

$$\mathbb{P}(\mathcal{E}^c(r)) \leq \mathbb{P}(\mathcal{E}_0^c(r)) + \mathbb{P}(\mathcal{E}_1^c(r)) + \mathbb{P}(\mathcal{E}_2^c(r)).$$

Note that after first two steps, $a_{1,r}$ has been pulled at least $\log^2 I_r$ times. Then, by Lemma 1, we have

$$\begin{aligned} \mathbb{P}(\mathcal{E}_0^c(r)) &\leq \sum_{k=1}^K \mathbb{P}(\exists s \geq \log^2 I_r, |\hat{\mu}_{k,s} - \mu_k| \geq \epsilon_r) \\ &\leq 2K \cdot \exp\left(-\frac{\log^2 I_r}{2V(\log \log I_r)^2}\right) \lesssim \frac{1}{I_r}, \end{aligned} \quad (18)$$

where the first inequality follows from the union bound over all possible event $a_{1,r} = k$. Meanwhile, by (15), we have

$$\mathbb{P}(\mathcal{E}_1^c(r)) \leq \sum_{i=2}^K \mathbb{P}(\hat{\mu}_{1,g(I_r)} \leq \hat{\mu}_{i,g(I_r)}) \lesssim \frac{1}{\log^2 I_r}.$$

Finally, observing that after first two steps, arm i is pulled at least $g(I_r)$ times, by Lemma 1, we have

$$\begin{aligned} \mathbb{P}(\mathcal{E}_2^c(r)) &\leq \sum_{k=1}^K \mathbb{P}(\exists s \geq g(I_r), |\hat{\mu}_{k,s} - \mu_k| \geq \epsilon_r) \\ &\leq 2K \cdot \exp\left(-\frac{g(I_r)}{2V(\log \log I_r)^2}\right) \lesssim \frac{1}{\log^2 I_r}. \end{aligned}$$

As a result, we have

$$\mathbb{P}(\mathcal{E}^c(r)) \leq \sum_{s=1}^3 \mathbb{P}(\mathcal{E}_s^c(r)) \lesssim \frac{1}{\log^2 I_r}. \quad (19)$$

Therefore,

$$\mathbb{E}\left[\sum_{r \in \mathcal{S}_2} Y_3(r)\right] \leq \sum_{r=r'}^{r^o} \sum_{n=1}^3 \mathbb{P}(\mathcal{E}_n^c(r)) \log^2(I_r) \lesssim r^o. \quad (20)$$

Combining (17) and (20), we obtain

$$\frac{\mathbb{E}[Y_3]}{\log T} = \frac{1 + 1/(1 + i \log^\alpha T)}{\text{kl}(\mu_i - 2\epsilon_{r'}, \mu_1 + 2\epsilon_{r'})} + o_T(1). \quad \square$$

Lemma 5. $\mathbb{E}[Y_4]/\log T = o_T(1)$.

Lemma 6. $\mathbb{E}[Z]/\log T = o_T(1)$.

Due to the space limit, we refer readers to Appendix B for proofs of Lemma 5 and Lemma 6.

6 Experiment

In this section, we compare our algorithm BABA with KL-UCB (Garivier & Cappé, 2011) under two reward distributions, i.e., *Gaussian distribution* and *Bernoulli distribution*. For each distribution, we test BABA and KL-UCB with 2 arms and 5 arms respectively. Specifically, for 2-arm setting,

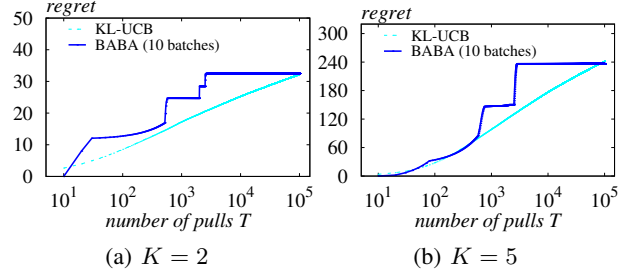


Figure 1. Regrets over Gaussian distributions. The experiments are averaged over 2000 repetitions

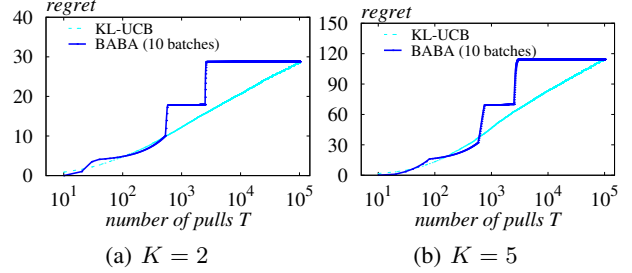


Figure 2. Regrets over Bernoulli distributions. The experiments are averaged over 2000 repetitions

we set $\mu \in \{1, 0\}$ and $\sigma = 1$ for Gaussian distribution; we set $p \in \{0.5, 0.25\}$ for Bernoulli distribution. For 5-arm setting, we set $\mu \in \{1, 0.5, 0.5, 0.5, 0.5\}$ and $\sigma = 1$ for Gaussian distribution; we set $p \in \{0.5, 0.25, 0.25, 0.25, 0.25\}$ for Bernoulli distribution. For our BABA algorithm, we set $\alpha = 3$ and $I_1 = 2000$. All the experiments are averaged over 2000 repetitions.

For Gaussian rewards, Figure 1(a) and Figure 1(b) report the regret when $K = 2$ and $K = 5$, respectively. As we can see, when T approaches 10^5 , BABA achieves the same regret as KL-UCB while requiring 10 batches opposed to 10^5 . The regret of BABA increases rapidly at some time steps. For example, in Figure 1(a), the regret increases from 17.0 to 24.2 from time steps 520 to 560. The reason is that in BABA, the suboptimal arms are mostly pulled during the exploration stages. In addition, as shown in Figure 1(a) and 1(b), when $T = 2000$, the regret of BABA is larger than KL-UCB. The reason is that for small T , BABA may not reach the optimal performance since asymptotic optimality holds only for sufficiently large T .

For Bernoulli rewards, Figure 2(a) and Figure 2(b) report the regret when $K = 2$ and $K = 5$, respectively. Again, the BABA achieves the comparable regret with that of KL-UCB while requiring 10 batches opposed to 10^5 .

7 Conclusion

We study the anytime bathed multi-armed bandit problem. We propose an algorithm BABA that achieves the asymptotically optimal regret using only $\mathcal{O}(\log \log T \cdot i \log^\alpha T)$

batches. We also show a lower bound on the batch complexity of anytime bandit algorithms, which theoretically confirms the conjecture in Besson & Kaufmann (2018) that no algorithm using static time grid can achieve the asymptotic optimality within $c \log \log T$ batches for any constant c . Moreover, we conduct experiments to show that our algorithm achieves the comparable regret with that of KL-UCB while using significantly fewer batches.

Acknowledgement

We thank the anonymous reviewers for their helpful comments. X. Xiao is supported by the Ministry of Education, Singapore, under Tier-2 Grant R-252-000-A70-112. T. Jin is supported by the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No: AISG-PhD/2021-01-004[T]). P. Xu and Q. Gu are partially supported by the National Science Foundation IIS-1904183. The views and conclusions contained in this paper are those of the authors and should not be interpreted as representing any funding agencies.

References

- Agrawal, S. and Goyal, N. Near-optimal regret bounds for thompson sampling. *Journal of the ACM*, 64(5):1–24, 2017.
- Audibert, J.-Y. and Bubeck, S. Minimax policies for adversarial and stochastic bandits. In *Proc. COLT*, 2009.
- Auer, P., Cesa-Bianchi, N., Freund, Y., and Schapire, R. E. Gambling in a rigged casino: The adversarial multi-armed bandit problem. In *Proc. IEEE FOCS*, pp. 322–331, 1995.
- Bertsimas, D. and Mersereau, A. J. A learning approach for interactive marketing to a customer segment. *Operations Research*, 55(6):1120–1135, 2007.
- Besson, L. and Kaufmann, E. What doubling tricks can and can’t do for multi-armed bandits. *arXiv preprint arXiv:1803.06971*, 2018.
- Cesa-Bianchi, N., Dekel, O., and Shamir, O. Online learning with switching costs and other adaptive adversaries. In *Proc. NeurIPS*, 2013.
- Degenne, R. and Perchet, V. Anytime optimal algorithms in stochastic multi-armed bandits. In *Proc. ICML*, pp. 1587–1595, 2016.
- Esfandiari, H., Karbasi, A., Mehrabian, A., and Mirrokni, V. Batched multi-armed bandits with optimal regret. *arXiv preprint arXiv:1910.04959*, 2019.
- Gao, Z., Han, Y., Ren, Z., and Zhou, Z. Batched multi-armed bandits problem. In *Proc. NeurIPS*, pp. 501–511, 2019.
- Garivier, A. and Cappé, O. The kl-ucb algorithm for bounded stochastic bandits and beyond. In *Proc. COLT*, pp. 359–376, 2011.
- Harremoës, P. Bounds on tail probabilities in exponential families. *arXiv preprint arXiv:1601.05179*, 2016.
- Jin, T., Xu, P., Xiao, X., and Gu, Q. Double explore-then-commit: Asymptotic optimality and beyond. *arXiv preprint arXiv:2002.09174*, 2020.
- Kaufmann, E. et al. On bayesian index policies for sequential resource allocation. *The Annals of Statistics*, 46(2): 842–865, 2018.
- Kittur, A., Chi, E. H., and Suh, B. Crowdsourcing user studies with mechanical turk. In *Proc. CHI*, pp. 453–456, 2008.
- Korda, N., Kaufmann, E., and Munos, R. Thompson sampling for one-dimensional exponential family bandits. In *Proc. NeurIPS*, 2013.
- Lai, T. L. and Robbins, H. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1): 4–22, 1985.
- Lattimore, T. Refining the confidence level for optimistic bandit strategies. *Journal of Machine Learning Research*, 19(1):765–796, 2018.
- Lattimore, T. and Szepesvári, C. *Bandit algorithms*. Cambridge University Press, 2020.
- Ménard, P. and Garivier, A. A minimax and asymptotically optimal algorithm for stochastic bandits. In *Proc. ALT*, pp. 223–237, 2017.
- Perchet, V., Rigollet, P., Chassang, S., and Snowberg, E. Batched bandit problems. *The Annals of Statistics*, 44(2): 660–681, 2016.
- Thompson, W. R. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- Vaswani, S., Kveton, B., Wen, Z., Ghavamzadeh, M., Lakshmanan, L. V., and Schmidt, M. Model-independent online learning for influence maximization. In *Proc. ICML*, pp. 3530–3539, 2017.
- Zhou, Y., Chen, X., and Li, J. Optimal pac multiple arm identification with applications to crowdsourcing. In *Proc. ICML*, pp. 217–225, 2014.