# **Batched Nonparametric Contextual Bandits**

Rong Jiang<sup>1</sup> and Cong Ma<sup>2</sup>

<sup>1</sup>Committee on Computational and Applied Mathematics, University of Chicago <sup>2</sup>Department of Statistics, University of Chicago

February 2024; Revised June 2024

#### **Abstract**

We study nonparametric contextual bandits under batch constraints, where the expected reward for each action is modeled as a smooth function of covariates, and the policy updates are made at the end of each batch of observations. We establish a minimax regret lower bound for this setting and propose a novel batch learning algorithm that achieves the optimal regret (up to logarithmic factors). In essence, our procedure dynamically splits the covariate space into smaller bins, carefully aligning their widths with the batch size. Our theoretical results suggest that for nonparametric contextual bandits, a nearly constant number of policy updates can attain optimal regret in the fully online setting.

## 1 Introduction

Recent years have witnessed substantial progress in the field of sequential decision making under uncertainty. Especially noteworthy are the advancements in personalized decision making, where the decision maker uses side-information to make customized decision for a user. The contextual bandit framework has been widely adopted to model such problems because of its applicability and elegance [35, 53, 6]. In this framework, one interacts with an environment for a number of rounds: at each round, one is given a context, picks an action, and receives a reward. One can update the action-assignment policy based on previous observations and the goal is to maximize the expected cumulative rewards. For example, in online news recommendation, a recommendation algorithm selects an article for each newly arrived user based on the user's contextual information, and observes whether the user clicks the article or not. The goal is to try to maximize the number of clicks received. Apart from news recommendation, contextual bandits have found numerous applications in other fields such as clinical trials, personalized medicine, and online advertising [30, 62, 13].

At the core of designing a contextual bandit algorithm is deciding how to update the policy based on prior observations. A standard metric of performance for bandit algorithms is regret, which is the expected difference between the cumulative rewards obtained by an oracle who knows the optimal action for every context and that obtained by the actual algorithm under consideration. Many existing regret optimal bandit algorithms require a policy update per observation (unit) [4, 1, 39, 34]. At a first glance, such frequent policy updates are needed so that the algorithm can quickly learn the optimal action under each context and reduce regret. However, this kind of algorithm ignores an important concern in the practice of sequential decision making—the batch constraint.

In many real world scenarios, the data often arrive in batches: the statistician can only observe the outcomes of the policy at the end of a batch, and then decides what to do for the next batch. For example, this batch constraint is ubiquitous in clinical trials: statisticians need to divide the participants into batches, determine a treatment allocation policy before the batch starts, and then observe all the outcomes at the end of the batch [49]. Policy updates are made per batch instead of per unit. In fact, it is infeasible to apply unit-wise policy update in this case because observing the effect of a treatment takes time and if one waits for the result before deciding how to treat the next patient, the entire experiment will take too long to complete when the number of participants is huge. The batch constraint also appears in areas such as online marketing, crowdsourcing, and simulations [8, 50, 31, 15]. Clearly, the batch constraint presents additional

challenges to online learning. Indeed, from an information perspective, the statistician's information set is largely restricted since she can only observe all the responses at the end of a batch. The following questions naturally arise:

Given a batch budget M and a total number of T rounds, how should the statistician determine the size of each batch, and how should she update the policy after each batch? Can the statistician design batch learning algorithms that achieve regret performances on par with the fully online setting using as few policy updates as possible?

#### 1.1 Main contributions

In this work, we address the aforementioned questions under a classical framework for personalized decision making—nonparametric contextual bandits [48, 39]. In this framework, the expected reward associated with each treatment (or arm in the language of bandits) is modeled as a nonparametric smooth function of the covariates [59]. In the fully online setup, seminal works [48, 39] establish the minimax optimal regret bounds for the nonparametric contextual bandits. Nevertheless, under the more challenging setting with the batch constraint, the fundamental limits for nonparametric bandits remain unknown. Our paper aims to bridge this gap. More concretely, we make the following contributions:

- First, we establish a minimax regret lower bound for the nonparametric bandits with the batch constraint. Our lower bound holds even when the batch size is adaptively chosen (based on the data observed in prior batches). The proof relies on a simple but useful insight that the worst-case regret over the entire horizon is greater than the worst-case regret over the first i batches for any 1 ≤ i ≤ M. To exploit this insight, for each different batch, we construct different families of hard instances to target it, leading to a maximal regret over this batch.
- In addition, we demonstrate that the aforementioned lower bound is tight by providing a matching upper bound (up to log factors). Specifically, we design a novel algorithm—Batched Successive Elimination with Dynamic Binning (BaSEDB)—for the nonparametric bandits with batch constraints. BaSEDB progressively splits the covariate space into smaller bins whose widths are carefully selected to align well with the corresponding batch size. The delicate interplay between the batch size and the bin width is crucial for obtaining the optimal regret in the batch setting.
- On the other hand, we show the suboptimality of static binning under the batch constraint by proving an algorithm-specific lower bound. Unlike the fully online setting where policies that use a fixed number of bins can attain the optimal regret [39], our lower bound indicates that batched successive elimination with static binning is strictly suboptimal. This highlights the necessity of dynamic binning in some sense under the batch setting, which is uncommon in classical nonparametric estimation.
- Last but not least, we demonstrate the challenge of adapting to the margin parameter in the batch setting. Specifically, we show that when M is small, the price of not knowing the true margin parameter for an algorithm is at least a polynomial increase in terms of the regret.

It is also worth mentioning that an immediate consequence of our results is that M 2 log log T number of batches sufices to achieve the optimal regret in the fully online setting. In other words, we can use a nearly constant number of policy updates in practice to achieve the optimal regret obtained by policies that require one update per round.

#### 1.2 Related work

Nonparametric contextual bandits. [58] introduced the mathematical framework of contextual bandit. The theory of contextual bandits in the fully online setting has been continuously developed in the past few decades. On one hand, [4, 1, 23, 6, 7, 41] obtained learning guarantees for linear contextual bandits in both low and high dimensional settings. On the other hand, [59] introduced the nonparametric approach to model the mean reward function. [48] proved a minimax lower bound on the regret of nonparametric

<sup>&</sup>lt;sup>1</sup>In a certain regime the BSE policy from [39] which uses a fixed number bins could loose by log factors compared to the optimal fully online regret. However, we will show the price of fixed binning is polynomial under the batch setting.

bandit and developed an upper-confidence-bound (UCB) based policy to achieve a near-optimal rate. [39] improved this result and proposed the Adaptively Binned Successive Elimination (ABSE) policy that can also adapt to the unknown margin parameter. Further insights in this nonparametric setting were developed in subsequent works [42, 43, 45, 24, 27, 52, 25, 10, 51, 9]. The smoothness assumption is also adopted in another line of work [37, 36, 33, 11] on the continuum-armed bandit problems. However in contrast to what we study, the reward is assumed to be a Lipschitz function of the action, and the covariates are not taken into considerations.

Batch learning. The batch constraint has received increasing attention in recent years. [40, 21] considered the multi-armed bandit problem under the batch setting and showed that O(log log T) batches are adequate in achieving the rate-optimal regret, compared to the fully online setting. [26, 47] extended batch learning to the (generalized) linear contextual bandits and [46, 56, 17] further studied the setting with high-dimensional covariates. [29, 28] established batch learning guarantees for the Thompson sampling algorithm. [18] considered Lipschitz continuum-armed bandit problem with the batch constraint. Inference for batched bandits was considered in [60]. A concept related to batch learning in literature is called delayed feedback [14, 13, 55, 19]. These works consider the setting where rewards are observed with delay and analyze effects of delay on the regret. [32, 2] studied delayed feedback in nonparametric bandits and the key difference to batch learning is that the batch size is given, whereas in our case, it is a design choice by the statistician. Batch learning's focus is different to that of delayed feedback in the sense that the former gives the decision maker discretion to choose the batch size which makes it possible to approximate the optimal standard online regret with a small number of batches. Finally, the notion switching cost is intimately related to the batch constraint. [12] studied online learning with low switching cost and obtained minimax optimal regret with O(log log T) batches. [5, 61, 20, 57, 44] developed regret guarantees with low switching cost for reinforcement learning. Low switching cost can be interpreted as infrequent policy updates, but it does not require the learner to divide the samples into batches with feedback only becoming available at the end of a batch.

## 2 Problem setup

We begin by introducing the problem setup for nonparametric bandits with the batch constraint.

A two-arm nonparametric bandit with horizon  $T \ge 1$  is specified by a sequence of independent and identically distributed random vectors

$$(X_t, Y_t^{(1)}, Y_t^{(-1)}), \quad \text{for } t = 1, 2, ..., T,$$
 (1)

$$E[Y_t^{(k)} | X_t] = f^{(k)}(X_t).$$

Here  $f^{(k)}$  is the unknown mean reward function for the arm k.

Without the batch constraint, the game of nonparametric bandits plays sequentially. At each step t, the statistician observes the context  $X_t$ , and pulls an action  $A_t \ 2 \ 1, -1 \ according to a rule <math>\pi_t : X \to \{1, -1\}$ . Then she receives the corresponding reward  $Y_t^{(A_t)}$ . In this case, the rule  $\pi_t$  for selecting the action at time t is allowed to depend on all the observations strictly anterior to t.

In an M-batch game, the statistician needs to design an M-batch policy  $(\Gamma,\pi)$ , where  $\Gamma=\{t_0,t_1,...,t_M\}$  is a partition of the entire time horizon T that satisfies  $0=t_0< t_1<...< t_{M-1}< t_M=T$ , and  $\pi=\{\pi_t\}_{t=1}^T$  is a sequence of random functions  $\pi_t:X\to\{1,-1\}$ . The grid  $\Gamma$  can be chosen adaptively, meaning that the statistician can use all information up to  $t_{i-1}$  to determine  $t_i$ . More specifically, prior to the start of the game, she will specify the first batch  $t_1$ , and at the end of  $t_1$ , she will use all observations she have to decide the next batch  $t_2$ , and this process repeats in batches. In contrast to the case without the batch constraint, only the rewards associated with timesteps prior to the current batch are observed and available for making decisions for the current batch. Specifically, let  $\Gamma(t)$  be the batch index for the time t, i.e.,  $\Gamma(t)$  is the unique integer such that  $t_{\Gamma(t)-1} < t \le t_{\Gamma(t)}$ . Then at time t, the available information for  $\pi_t$  is only

 $\{X_l\}_{l=1}^t \mathbb{P}\{Y_l^{(A_l)}\}_{l=1}^{\Gamma(t)-1}$ , which we denote by  $F^t$ . The statistician's policy  $\pi_t$  at time t is allowed to depend on  $F_t$ .

The goal of the statistician is to design an M-batch policy  $(\Gamma, \pi)$  that can compete with an oracle that has perfect knowledge (i.e., the law of  $(X_t, Y_t^{(1)}, Y_t^{(-1)})$ ) of the environment. Formally, we define the cumulative

regret as

$$R_{T}(\pi) := E \Big|_{t=1}^{"X_{T}} f^{\mathbb{Z}}(X_{t}) - f^{(\pi_{t}(X_{t}))}(X_{t})^{\#}, \qquad (2)$$

where  $f^{\mathbb{Z}}(x) = \max_{k \mathbb{Z}\{1,-1\}} f^{(k)}(x)$  is the maximum mean reward one could obtain on the context x. Note here we omit the dependence on  $\Gamma$  for simplicity.

## 2.1 Assumptions

We adopt two standard assumptions in the nonparametric bandits literature [48, 39]. The first assumption is on the smoothness of the mean reward functions.

Assumption 1 (Smoothness). We assume that the reward function for each arm is  $(\beta, L)$ -smooth, that is, there exist  $\beta$  (0,1] and L>0 such that for k (1,-1),

$$|f^{(k)}(x) - f^{(k)}(x')| \le L2x - x'2^{\beta}_{2}$$

holds for all  $x, x' \supseteq X$ .

The second assumption is about the separation between the two reward functions.

Assumption 2 (Margin). We assume that the reward functions satisfy the margin condition with parameter  $\alpha > 0$ , that is there exist  $\delta_0 \ (0,1)$  and  $D_0 > 0$  such that

$$P_X = 0 < f^{(1)}(X) - f^{(-1)}(X) \le \delta \le D_0 \delta^{\alpha}$$

Assumption 2 is related to the margin condition in classification [38, 54, 3] and is introduced to bandits in [22, 48, 39]. The margin parameter affects the complexity of the problem. Intuitively, a small  $\alpha$ , say  $\alpha \approx 0$ , means the two mean functions are entangled with each other in many regions and hence it is challenging to distinguish them; a large  $\alpha$ , on the other hand, means the two reward functions are mostly well-separated.

From now on, we use  $F(\alpha, \beta)$  to denote the class of nonparametric bandit instances (i.e., distributions over (1)) that satisfy Assumptions 1-2.

Remark 1. Throughout the paper, we assume that  $\alpha\beta \leq 1$ . By proposition 2.1 from [48], when  $\alpha\beta > 1$ , one of the arms will dominate the other one for the entire covariate space. The instance is reduced to a multi-armed bandit without covariates which is not the interest of the current paper. Therefore, we focus on the case  $\alpha\beta \leq 1$  hereafter.

## 3 Fundamental limits of batched nonparametric bandits

In this section, we establish minimax lower bounds for the regret achievable by any M-batch policy  $(\Gamma, \pi)$ ; see Theorem 2. To begin with, we state a minimax lower bound, together with its proof, when the grid  $\Gamma$  is prespecified, that is, the statistician divides the horizon [1:T] into M disjoint batches  $[1:t_1]$ ,  $[t_1+1:t_2]$ , ...,  $[t_{M-1}+1,T]$  before the game begins; see Theorem 1. As we will soon see, the proof of the lower bound with fixed grid is not only useful for establishing the lower bound for any general M-batch policy  $(\Gamma,\pi)$ , but also instrumental in our development of an optimal policy to be detailed in Section 4.

Recall that  $F(\alpha, \beta)$  denotes the class of nonparametric bandit instances (i.e., distributions over (1)) that obey Assumptions 1-2. We have the following minimax lower bound for any M-batch policy with a fixed grid, in which we define

$$\gamma := \frac{\beta(1+\alpha)}{2\beta+d} ? (0,1).$$

Theorem 1. Suppose that  $\alpha\beta \le 1$ , and assume that  $P_X$  is the uniform distribution on  $X = [0,1]^d$ . For any Mbatch policy  $(\Gamma, \pi)$  where  $\Gamma$  is prespecified, there exists a nonparametric bandit instance in  $\Gamma(\alpha, \beta)$  such that the regret of  $(\Gamma, \pi)$  on this instance is lower bounded by

$$E[R_T(\pi)] \geq \widetilde{D}T^{\frac{1-\gamma}{1-\gamma M}},$$

where  $\tilde{D} > 0$  is a constant independent of T and M.

See Section 3.1 for the proof of this lower bound.

As a sanity check, one sees that as M increases, the lower bound decreases. This is intuitive, as the policy is more powerful as M increases. As a result, the problem of batched nonparametric bandits becomes easier.

#### Proof of Theorem 1 3.1

Let  $(\Gamma, \pi)$  be the M-batch policy under consideration, with

$$\Gamma = \{t_0 = 0, t_1, t_2, \dots, t_M = T\}.$$

Throughout this proof, we consider Bernoulli reward distributions, that is  $Y_t^{(1)}$ ,  $Y_t^{(-1)}$  are Bernoulli random variables with mean  $f^{(1)}(X_t)$ , and  $f^{(-1)}(X_t)$ , respectively. In addition, we fix  $f^{(-1)}(x) = \frac{1}{2}$ . Let f be the mean reward function of the first arm. To make the dependence on the reward instance clear, we write the cumulative regret up to time n as  $R_n(\pi; f)$ .

Our proof relies on a simple observation: the worst-case regret over [T] is larger than the worst-case regret over the first i batches. Formally, we have

$$\sup_{(f,\frac{1}{2})\mathbb{Z}F(\alpha,\beta)} R_{T}(\pi;f) \geq \max_{1 \leq i \leq M} \sup_{(f,\frac{1}{2})\mathbb{Z}F(\alpha,\beta)} R_{t_{i}}(\pi;f). \tag{3}$$

Though simple, this observation lends us freedom on choosing different families of instances in  $F(\alpha, \beta)$ targeting different batch indices i.

Our proof consists of four steps. In Step 1, we reduce bounding the regret of a policy to lower bounding its inferior sampling rate to be defined. In Step 2, we detail the choice of different families of instances for each different batch index i. Then in Step 3, we apply an Assouad-type of argument to lower bound the average inferior sampling rate of the family of hard instances. Lastly in Step 4, we combine the arguments to complete the proof.

Step 1: Relating regret to inferior sampling rate. Given an M-batch policy, we define its inferior

sampling rate at time n on an instance 
$$(f, \frac{1}{2})$$
 to be 
$$S_n(\pi; f) \coloneqq E \begin{bmatrix} X^n \\ t=1 \end{bmatrix} \{\pi_t(X_t) = \pi^{\mathbb{Z}}(X_t), f(X_t) = \frac{1}{2} \}.$$

In words,  $S_n(\pi;f)$  counts the number of times  $\pi$  selects the strictly suboptimal arm up to time n. Thanks to the following lemma, we can reduce lower bounding the regret to the inferior sampling rate.

Lemma 1 (Lemma 3.1 in [48]). Suppose that  $(f, \frac{1}{2})$   $\mathbb{P}(\alpha, \beta)$ . Then for any  $1 \le n \le T$ , we have

$$S_n(\pi;f) \leq D n^{\frac{1}{1+\alpha}} R_n(\pi;f)^{\frac{\alpha}{1+\alpha}}$$

for some constant D > 0.

As an immediate consequence of the above lemma, we obtain

$$\begin{split} \sup_{(f,\frac{1}{2})\mathbb{B}F(\alpha,\beta)} R_T(\pi;f) &\geq \max_{1\leq i\leq M} \sup_{(f,\frac{1}{2})\mathbb{B}F(\alpha,\beta)} (\frac{1}{D})^{\frac{1+\alpha}{\alpha}} t_i^{-\frac{1}{\alpha}} (S_{t_i}(\pi;f))^{\frac{1+\alpha}{\alpha}} \\ &= (\frac{1}{D})^{\frac{1+\alpha}{\alpha}} \max_{1\leq i\leq M} t_i^{-\frac{1}{\alpha}} \sup_{(f,\frac{1}{2})\mathbb{B}F(\alpha,\beta)} S_{t_i}(\pi;f) \quad . \end{split}$$

From now on, we focus on lower bounding  $\sup_{(f,\frac{1}{\tau})\boxtimes F(\alpha,\beta)} S_{t_i}(\pi;f)$ .

Step 2: Introducing the family of reward instances for  $t_i$ . Our construction of the family of hard instances is adapted from [48]. Define  $z_1 = 1$ , and  $z_i = 2t_{i-1}^{1/(2\beta+d)}$  for i = 2, 3, ..., M. Henceforth, we will fix some i and write  $z_i$  as z. We partition  $[0, 1]^d$  into  $z^d$  bins with equal width. Denote the bins by  $C_j$  for  $j = 1, ..., z^d$ , and let  $q_j$  be the center of  $C_j$ .

Define a set of binary sequences  $\Omega_s = \{\pm 1\}^s$ , with  $s = 2z^{d-\alpha\beta}$ . For each  $\omega 2\Omega_s$  we define a function  $f_\omega : [0,1]^d \to R$ :

$$f_{\omega}(x) = 1 + \partial_{j} \phi_{j}(x), j = 1$$

where  $\phi_j(x) = D_{\phi}z^{-\beta}\phi(2z(x-q_j))\mathbf{1}\{x \ \mathbb{Z} \ C_j\}$  with  $\phi(x) = (1-\mathbb{Z}x\mathbb{Z}_{\infty})^{\beta}\mathbf{1}\{\mathbb{Z}x\mathbb{Z}_{\infty} \le 1\}$ , and  $D_{\phi} = \min(2^{-\beta}L, 1/4)$ . In all, we consider the family of reward instances

$$C_z := f^{(1)}(x) = f_{\omega}(x), f^{(-1)}(x) = \frac{1}{2} | \omega ? \Omega_s .$$
 (4)

With slight abuse of notation, we also use  $C_z$  to denote  $\{f_\omega:\omega\ \mathbb{Z}\ \Omega_s\}$ . It is straightforward to check that  $C_z\ \mathbb{Z}\ F(\alpha,\beta)$ .

Step 3: Lower bounding the inferior sampling rate. Fix some i  $\mathbb{Z}[M]$ , and consider  $z=z_i$ . Since  $C_z \mathbb{Z}[K(\alpha,\beta)]$ , we have

$$\sup_{(f,\frac{1}{2}) \text{@F}(\alpha,\beta)} S_{t_i}(\pi;f) \ge \sup_{f^{\underline{\alpha}}C_z} S_{t_i}(\pi;f).$$

Using the definitions of  $C_z$  and  $S_{t_i}(\pi; f)$ , we have

$$\sup_{f \in C_z} S_{t_i}(\pi; f) = \sup_{\omega \in Q} E_{\pi, f_{\omega}} \left( X_t^{t_i} \right) = \sup_{t = 1} (f_{\omega}(X_t) - \frac{1}{2}), f_{\omega}(X_t) = \frac{1}{2}$$

$$\geq \frac{1}{2^S} \left( X_t^{t_i} \right) = \sum_{t = 1}^{N} (f_{\omega}(X_t) - \frac{1}{2}), f_{\omega}(X_t) = \frac{1}{2}$$

$$= \frac{1}{2^S} \left( X_t^{t_i} \right) = \sum_{t = 1}^{N} (f_{\omega}(X_t) - \frac{1}{2}), f_{\omega}(X_t) = \frac{1}{2}$$

Since  $f_{\omega}(x) = \frac{1}{2}$  for  $x \not \square_{j=1,...s} C_j$ , we further obtain

$$\sup_{f \in C_z} S_{t_i}(\pi; f) \ge \frac{1}{2^s} \sum_{\omega \in Q_{t+1}}^{X} \sum_{j=1}^{X^{t_i}} \sum_{j=1}^{s} \left[ 1\{\pi_t(X_t) = \omega_j, X_t ? C_j\} \right].$$
 (5)

Here we use  $P_{\pi,f_{\omega}}^{t}$  to denote the joint distribution of  $\{X_{l}\}_{l=1}^{t}$   $\mathbb{P}\{Y_{l}^{\pi_{l}(X_{l})}\}_{l=1}^{\Gamma(t)-1}$ , where  $\Gamma(t)$  is the batch index for t, i.e., the unique integer such that  $t_{\Gamma(t)-1} < t \le t_{\Gamma(t)}$ . We use  $E_{\pi,f_{\omega}}^{t}$  to denote the corresponding expectation. Expand the right hand side of (5) to see that

$$\sup_{f \otimes C_{z}} S_{t_{i}}(\pi; f) \geq \frac{1}{2^{s}} \sum_{j=1}^{X^{s}} \sum_{t=1}^{X^{t}} \sum_{\omega_{[-j]} \otimes Q_{-1}} \sum_{h \otimes \{\pm 1\}}^{h \otimes \{\pm 1\}} \underbrace{E_{\pi, f_{q_{-jj}}}^{t}}_{W_{j, t, \omega_{[-j]}}} [1\{\pi_{t}(X_{t}) = h, X_{t} \otimes C_{j}\}], \tag{6}$$

where  $\omega_{[-j]}^h$  is the same as  $\omega$  except for the j-th entry being h. Note that here we use the fact that for f  $_{\omega_{[-j]}^h}^h$ , the optimal arm in the bin  $C_j$  is h. We then relate  $W_{j,t,\omega_{[-j]}}$  to a binary testing error,

$$W_{j,t,\omega_{[-j]}} = \frac{1}{z^{d}} \sum_{h \in \{\pm 1\}}^{X} P_{\pi,f_{q_{-j}^{h}}}^{t} (\pi_{t}(X_{t}) = h \mid X_{t} \boxtimes C_{j})$$

$$\geq \frac{1}{4z^{d}} \exp -KL(P_{\pi,f_{\omega_{[-j]}^{-1}}}^{t}, P_{\pi,f_{\omega_{[-j]}^{1}}}^{t}) , \qquad (7)$$

where the second step invokes Le Cam's method. Under the batch setting, at time t, the available information is only up to  $t_{\Gamma(t)-1}$ . Consequently, we can apply Lemma 5 to obtain

$$\mathsf{KL}(\mathsf{P}_{\pi,\mathsf{f}_{\omega_{[-j]}^{-1}}}^{\mathsf{t}},\mathsf{P}_{\pi,\mathsf{f}_{\omega_{[-j]}^{-1}}}^{\mathsf{t}}) = \mathsf{KL}(\mathsf{P}_{\pi,\mathsf{f}_{\omega_{[-j]}^{-1}}}^{\mathsf{t}_{\Gamma(\mathsf{t})-1}},\mathsf{P}_{\pi,\mathsf{f}_{\omega_{[-j]}^{-1}}}^{\mathsf{t}_{\Gamma(\mathsf{t})-1}}) \leq 2\mathsf{z}^{-(2\beta+\mathsf{d})}\mathsf{t}_{\Gamma(\mathsf{t})-1}. \tag{8}$$

Combining (6), (7), and (8), we arrive at

$$\sup_{f \in C_z} S_{t_i}(\pi; f) \ge \frac{1}{8} \sum_{j=1}^{X^s} \sum_{t=1}^{t_i} \frac{1}{z^d} \exp_{-2z^{-(2\beta+d)}} t_{\Gamma(t)-1}$$

$$\ge \frac{1}{8} \sum_{i=1}^{z^{d-\alpha\beta}} X^i \frac{t_{i-1}t}{t_{i-1}} \exp_{-2z^{-(2\beta+d)}} t_{i-1,j=1}$$

$$\ge \frac{1}{8} \sum_{i=1}^{z^{d-\alpha\beta}} X^i \frac{t_{i-1}t}{t_{i-1}} \exp_{-2z^{-(2\beta+d)}} t_{i-1,j=1}$$

where the second line uses the fact that  $s=\mathbb{Z}z^{d-\alpha\beta}\mathbb{Z}$ , and the last inequality holds since  $t_{i-1}\leq t_{i-1}$  for all  $1\leq i\leq i$ . Now recall that  $z=z_i=\mathbb{Z}(t_{i-1})^{1/(2\beta+d)}\mathbb{Z}$  for  $i\geq 1$ , and z=1 for i=1. We can continue the lower bound to see that

$$\sup_{f \in C_{z_{i}}} S_{t_{i}}(\pi; f) \ge \frac{1}{8} \sum_{j=1}^{z} \sum_{l=1}^{X^{d-\alpha\beta}} X^{i} \frac{t_{l} - t_{l-1}}{z^{d}} \exp -2z^{-(2\beta+d)} t_{i-1}$$

$$\ge c^{\mathbb{Z}} X^{d-\alpha\beta} X^{i} \frac{t_{l} - t_{l}}{dz_{j} = 1}$$

$$\ge c^{\mathbb{Z}} \sum_{l=1}^{z} \frac{t_{l} - t_{l}}{dz_{l} - t_{l}}, \quad i > 1 = 0$$

$$c \cdot \mathbb{Z} \frac{t_{i}}{z^{\alpha\beta}} = \frac{\mathbb{Z}^{2d-\alpha\beta}}{\mathbb{Z}^{2d-\alpha\beta}}, \quad i > 1 = 0$$

$$\mathbb{Z}^{2d-\alpha\beta} = \mathbb{Z}^{2d-\alpha\beta}$$

$$\mathbb{Z}^{2d-\alpha\beta} = \mathbb{Z}^{2d-\alpha\beta}$$

for some  $c^2 > 0$ .

Step 4: Combining bounds together. Combining the previous arguments together leads to the conclusion that

$$\sup_{(f,\frac{1}{2})\boxtimes F(\alpha,\beta)} R_{T}(\pi;f) \geq \max_{1\leq i\leq M} \sup_{f\boxtimes C_{z_{i}}} R_{t_{i}}(\pi;f)$$

$$\geq \left(\frac{1}{D}\right)^{\frac{1+\alpha}{\alpha}} \max_{1\leq i\leq M} t_{i}^{-\frac{1}{\alpha}} \sup_{f\boxtimes C_{z_{i}}} S_{t_{i}}(\pi;f)$$

$$\boxtimes \max_{1\leq i\leq M} t_{1}, t_{2}, \dots, t_{M-1}$$

$$\geq \widetilde{D}T^{\frac{1-\gamma}{1-\gamma M}}.$$
(9)

This finishes the proof.

## 3.2 Lower bound for general M-batch policy

Now we are ready to state the minimax lower bound for any general M-batch policy  $(\Gamma, \pi)$ , i.e., when the grid  $\Gamma$  is allowed to be adaptively chosen.

Theorem 2. Suppose that  $\alpha\beta \leq 1$ , and assume that  $P_X$  is the uniform distribution on  $X = [0,1]^d$ . For any M-batch policy  $(\Gamma,\pi)$ , there exists a nonparametric bandit instance in  $F(\alpha,\beta)$  such that the regret of  $\pi$  on this instance is lower bounded by

$$E[R_{T}(\pi)] \geq \widetilde{D_{1}}(\frac{1}{M})^{\widetilde{D}_{2}}T^{\frac{1-\gamma}{1-\gamma M}},$$

where  $\tilde{D}_1$ ,  $\tilde{D}_2 > 0$  are constants independent of T and M.

See Appendix A for the proof.

Since our focus is on  $M > \log \log T$  (when M?  $\log \log T$ , by Corollary 1 there exists an algorithm whose regret attains the optimal fully online regret), we can see Theorem 1 and Theorem 2 differ at most by poly-log factors in T.

Unlike the fixed grid case where we choose a specific family of hard instances to target the regret in a certain batch, we cannot directly do so when the grid is adaptively selected because the adversary does not know  $\{t_i\}_{i=1}^M$  in advance. Inspired by [21], we overcome this difficulty by using an appropriately defined bad event that happens with sufficient probability to reduce the adaptive case to the fixed grid case. However, the proof under the nonparametric setting is much more challenging because the presence of contexts makes the instances for different batches less indistinguishable with each other. A key ingredient of our proof is to establish tight control of the total variation distance between two mixture distributions. The full proof can be found in Appendix A.

## 3.3 Implications on design of optimal M-batch policy

As we have mentioned, the proof of the lower bound with fixed grid, i.e., Theorem 1 facilitates the design of optimal M-batch policy.

Grid selection. First, the lower bound of the whole horizon is reduced to the worst-case regret over a specific batch; see (3). Consequently, we need to design the grid  $\Gamma = (t_0, t_1, t_2, \ldots, t_{M-1}, t_M)$  such that the total regret is evenly distributed across batches. More concretely, in view of the lower bound (9), one needs to set  $\mathbb{F}_1^{-\gamma} \stackrel{\iota}{\vdash} \mathbb{F}_1^{-\gamma} \stackrel{\iota}{\vdash} \mathbb{F}_$ 

Dynamic binning. In addition, in the proof of the lower bound, for each different batch i, we use different families of hard reward instances, parameterized by the number of bins  $z_i = \mathbb{Z}t_{i-1}^{1/(2\beta+d)}\mathbb{Z}$ . In other words, from the lower bound perspective, the granularity (i.e., the bin width  $1/z_i$ ) at which we investigate the mean reward function depends crucially on the grid points  $\{t_i\}$ : the larger the grid point  $t_i$ , the finer the granularity. This key observation motivates us to consider the batched successive elimination with dynamic binning algorithm to be introduced below.

## 4 Batched successive elimination with dynamic binning

In this section, we present the batched successive elimination with dynamic binning policy (BaSEDB) that nearly attains the minimax lower bound, up to log factors; see Algorithm 1. On a high level, Algorithm 1 gradually partitions the covariate space X into smaller hypercubes (i.e., bins) throughout the batches based on a list of carefully chosen cube widths, and reduces the nonparametric bandit in each cube to a bandit problem without covariates.

A tree-based interpretation. The process is best illustrated with the notion of a tree T of depth M; see Figure 1. Each layer of of the tree T is a set of bins that form a regular partition of X using hypercubes with equal widths. And the common width of the bins  $B_i$  in layer i is dictated by a list  $\{g_i\}^{M-1}$  of  $\underbrace{sp}_i$  lit factors. More precisely, we let

$$w_{i} := (\int_{|z|}^{i} \gamma^{1} g_{1})^{-1}$$
 (10)

be the width of the cubes in the i-th layer  $B_i$  for  $i \ge 1$ , and  $w_0 = 1$ . In other words,  $B_i$  contains all the cubes

#### Algorithm 1 Batched successive elimination with dynamic binning (BaSEDB)

```
Input: Batch size M, grid \Gamma = \{t_i\}_{i=0}^M, split factors \{g_i\}_{i=0}^{M-1}.
L \leftarrow B_1
for C 2 L do
     I_{C} = I
for i = 1, ..., M - 1 do
      for t = t_{i-1} + 1, ..., t_i do
           C \leftarrow L(X_{t})
           Pull an arm from I_{C} in a round-robin way.
          if t = t_i then
               Update L and \{I_C\}_{C \square L} by Algorithm 2 (L, \{I_C\}_{C \square L}, i, g_i).
for t = t_{M-1} + 1, ..., T do
      C \leftarrow L(X_t)
      Pull any arm from I_C.
```

$$C_{i,v} = \{x \ ? \ X : (v_j - 1)w_i \le x_j < v_j w_i, 1 \le j \le d\},$$

corresponding active arms I c for each C 2L; see Figure 1 for an example. Specifically, prior to the game (i.e., prior to the first batch), L is set to be  $B_1$ , all bins in layer 1, and  $I_C = \{1, -1\}$  for all  $C \supseteq L$ . Within this batch, the statistician tries the arms in Ic equally likely for all bins in L. Then at the end of the batch, given the revealed rewards in this batch, we update the active arms I C for each C I L via successive elimination. If no arm were eliminated from I<sub>C</sub>, this suggests that the current bin is not fine enough for the statistician to tell the difference between the two arms. As a result, she splits the bin C 2 L into its children child(C) in T. All the child nodes will be included in L, while the parent C stops being active (i.e., C is removed from L). The whole process repeats in a batch fashion. <sup>2</sup>

Grid  $\Gamma$  and split factors  $\{g_i\}_{i=0}^{M-1}$ . As one can see, the split factor  $g_i$  controls how many children a node at layer i can have and its appropriate choice is crucial for obtaining small regret. Intuitively, gi should be selected in a way such that a node  $C_{i+1}$  with width  $w_i$  can fully leverage the number of samples allocated to it during the (i + 1)-th batch. With these goals in mind, we design the grid  $\Gamma = \{t_i\}$  and split factors  $\{g_i\}$  as follows. Recall that  $v = {\beta(1+\alpha) \over \alpha}$ . We set

$$b = \Theta T^{\frac{1-\gamma}{2}}$$

The split factors are chosen according to

$$g_0 = \mathbb{P}b^{\frac{1}{2\beta+d}}\mathbb{P}$$
, and  $g_i = \mathbb{P}g_{i-1}^{y}\mathbb{P}$ ,  $i = 1, ..., M-2$ . (11)

In addition, the grid is chosen such that

$$t_i - t_{i-1} = 2 I_i w_i^{-(2\beta+d)} \log(T w_i^d) 2, 1 \le i \le M - 1,$$
 (12)

where  $l_i > 0$  is a constant to be specified later. It is easy to check that with these choices, we have

$$t_1 \ \ \ T^{\frac{1}{1-\gamma}M}$$
, and  $t_i = \ \ \ \ b(t_{i-1})^{\gamma} \ \ \ \ \$  for  $i = 2,...,M$ .

<sup>&</sup>lt;sup>2</sup> For the final batch M, the split factor  $g_{M-1} = 1$  by default because there is no need to further partition the nodes for estimation.

#### Algorithm 2 Tree growing subroutine

In particular, we set b properly to make  $t_M = T$ . Indeed, these choices taken together meet the expectation laid out in Section 3.3: we need to choose the grid and the split factors appropriately so that (1) the total regret spreads out across different batches, and (2) the granularity becomes finer as we move further to later batches.

When to eliminate arms? Now we zoom in on the elimination process described in Algorithm 2. The basic idea follows from successive elimination in the bandit literature [16, 39, 21]: the statistician eliminates an arm from  $I_C$  if she expects the arm to be suboptimal in the bin C given the rewards collected in C. Specifically, for any node  $C \ T$ , define

$$U(\tau,T,C) := 4 \frac{\frac{\log(2T|C|^d)}{\tau}}{\tau},$$

where |C| denotes the width of the bin. Let  $m_{C,i} := \bigcap_{t=t_{i-1}+1}^{t_i} \mathbf{1}\{X_t \ \mathbb{Z} \ C\}$  be the number of times we observe contexts from C in batch i. We then define for  $k \ \mathbb{Z} \{1,-1\}$  that

$$\bar{Y}_{C,i}^{(k)} := \frac{P_{\substack{t_i \\ t = t_{i-1}+1}} Y_t \cdot 1\{X_t ? C, A_t = k\}}{P_{\substack{t_i \\ t = t_{i-1}+1}} 1\{X_t ? C, A_t = k\}},$$

which is the empirical mean reward of arm k in node C during the i-th batch. It is easy to check that  $\bar{Y}_{C,i}^{(k)}$  has expectation  $f_{C}^{(k)}$  given by

$$f_{C}^{(k)} := E[f^{(k)}(X) | X \boxtimes C] = \frac{1}{P_{X}(C)} Z_{C} f^{(k)}(x)dP_{X}(x).$$

Similarly, we define the average optimal reward in bin C to be

$$\bar{f_C} := \frac{1}{P_X(C)} \sum_{C}^{Z} f^{\mathbb{R}}(x) dP_X(x).$$

The elimination threshold  $U(m_{C,i},T,C)$  is chosen such that an arm k with  $f^{\boxtimes}$   $-\overline{t}_{C}^{(k)}$   $\overline{\mathbb{Q}}$   $|C|^{\beta}$  is eliminated with high probability at the end of batch i. Therefore, when  $|I_{C}| > 1$ , the remaining arms are statistically

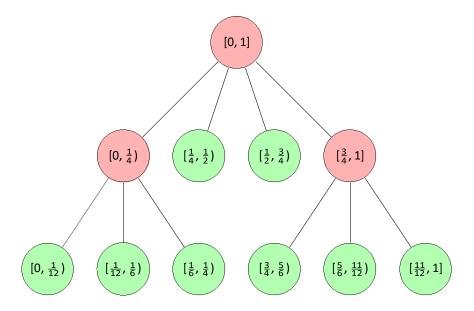


Figure 1: An example of the tree growing process for d = 1, M = 3, G = {4,3,1}. The root node is at depth 0. For the first batch, the 4 nodes located at depth 1 of the tree were used. Both  $[\frac{1}{2}, \frac{1}{2}]$  and  $[\frac{1}{2}, \frac{3}{2}]$  only had one active arm remaining so they were not further split and remained in the set of active nodes (green). Meanwhile,  $|I_{[0,\frac{1}{2}]}| = |I_{[\frac{3}{2},1]}| = 2$  so aeach of them was split into 3 smaller nodes, and both nodes were marked as inactive (red). For the second batch, all the green nodes were actively used but arm elimination was performed at the end of batch 2 only for nodes located at depth 2 (the green nodes at depth 1 already have 1 active arm remaining so there is no need to eliminate again).

indistinguishable from each other, so C is split into smaller nodes to estimate those arms more accurately using samples from future batches. On the other hand, when  $|I_C| = 1$ , the remaining arm is optimal in C with high probability—a consequence of the smoothness condition, and it will be exploited in the later batches.

Connections and differences with ABSE in [39]. In appearance, BaSEDB (Algorithm 1) looks quite similar to the Adaptively Binned Successive Elimination (ABSE) proposed in [39]. However, we would like to emphasize several fundamental differences. First, the motivations for the algorithms are completely different. [39] designs ABSE to adapt to the unknown margin condition  $\alpha$ , while our focus is to tackle the batch constraint. In fact, without the batch constraints, if  $\alpha$  is known, adaptive binning is not needed to achieve the optimal regret [39]. This is certainly not the case in the batched setting. Fixing the number of bins used across different batches is suboptimal because one can construct instances that cause the regret incurred during a certain batch to explode. We will expand on this phenomenon in Section 4.3. Secondly, the algorithm in [39] partitions a bin into a fixed number  $2^d$  of smaller ones once the original bin is unable to distinguish the remaining arms. In this way, the algorithm can adapt to the difference in the local dificulty of the problem. In comparison, one of our main contributions is to carefully design the list of varying split factors that allows the new cubes to maximally utilize the number of samples allocated to it during the next batch.

#### 4.1 Regret guarantees

Now we are ready to present the regret performance of BaSEDB (Algorithm 1).

Theorem 3. Suppose that  $\alpha\beta \leq 1$ . Fix any constant  $D_1 > 0$  and suppose that  $M \leq D_1 \log T$ . Equipped with the grid and split factors list that satisfy (12) and (11), the policy  $\hat{\pi}$  given by Algorithm 1 obeys

$$E[R_T(\hat{\pi})] \leq \tilde{C}(\log T)^2 \cdot T^{\frac{1-\gamma}{1-\gamma M}},$$

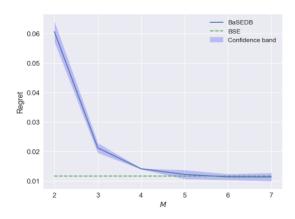


Figure 2: Regret vs. batch budget M.

where  $\tilde{C} > 0$  is a constant independent of T and M.

See Appendix B for the proof.

While Theorem 3 requires  $M \ @ \log T$ , we see from the corollary below that it is in fact sufficient to show the optimality of Algorithm 1.

Corollary 1. As long as  $M \ge D_2 \log \log(T)$ , where  $D_2$  depends on  $\gamma = \frac{\beta(1+\alpha)}{2\beta+d}$ , Algorithm 1 achieves

$$E[R_T(\hat{\pi})] \leq \tilde{C}(\log T)^2 \cdot T^{1-\gamma}$$

where  $\tilde{C} > 0$  is a constant independent of T and M.

Theorem 3, together with Corollary 1 and Theorem 2 establish the fundamental limits of batch learning for the nonparametric bandits with covariates, as well as the optimality of BaSEDB, up to logarithmic factors. To see this, when M  $2 \log \log(T)$ , the upper bound in Theorem 3 matches the lower bounds in Theorem 1 and Theorem 2, apart from log factors. On the other end, when M  $2 \log \log(T)$ , Algorithm 1, while splitting the horizon into M batches, achieves the optimal regret (up to log factors) for the setting without the batch constraint [39]. It is evident that Algorithm 1 is optimal in this case.

### 4.2 Numerical experiments

In this section, we provide some experiments on the empirical performance of Algorithm 1. We set T = 50000, d =  $\beta$  = 1,  $\alpha$  = 0.2. We let P<sub>X</sub> be the uniform distribution on [0,1]. Denote q<sub>j</sub> = (j - 1/2)/4 and C<sub>j</sub> = [q<sub>j</sub> - 1/8, q<sub>j</sub> + 1/8] for 1 ≤ j ≤ 4. For the mean reward functions, we choose f<sup>(1)</sup>, f<sup>(-1)</sup> : [0,1]  $\rightarrow$  R such that

$$f^{(1)}(x) = \frac{1}{2} + \chi^4 \omega_j \phi_j(x), \qquad f^{(-1)}(x) = \frac{1}{2},$$

where  $\omega_j$ 's are sampled i.i.d. from Rad( $_2$ \frac{1}{2},  $\varphi_j(x) = _4 \varphi(8(x - q_j))1\{x \ 2 \ C_j\}$  and  $\varphi(x) = (1 - |x|)1\{|x| \le 1\}$ . We let  $Y^{(k)}$   $2 \ Bernoulli(f^{(k)}(x))$ . To illustrate the performance of Algorithm 1, we compare it with the Binned Successive Elimination (BSE) policy from [39], which is shown to be minimax optimal in the fully online case. Figure 2 shows the regret of Algorithm 1 under different batch budegts. One can see that it is suficient to have M = 5 batches to achieve the fully online eficiency.

## 4.3 Failure of static binning

We have seen the power of dynamic binning in solving batched nonparametric bandits by establishing its rateoptimality in minimizing regret. Now we turn to a complimentary but intriguing question: is it necessary to use dynamic binning to achieve optimal regret under the batch constraint? To formally address this question, we investigate the performance of successive elimination with static binning, i.e., Algorithm 1 with  $g_0 = g$ , and  $g_1 = g_2 = \cdots g_{M-2} = 1$ . Although static binning works when M is large (e.g., a single choice of g attains the optimal regret [48, 39] in the fully online setting), we show that it must fail when M is small.

To bring the failure mode of static binning into focus, we consider the simplest scenario when M=3, and  $\alpha=\beta=d=1$ . Note that the successive elimination with static binning algorithm is parameterized by the grid choice  $\Gamma=\{t_0=0,t_1,t_2,t_3=T\}$  and the fixed number g of bins. The following theorem formalizes the failure of static binning in achieving optimal regret when M=3.

Theorem 4. Consider M = 3, and  $\alpha$  =  $\beta$  = d = 1. For any choice of  $1 \le t_1 < t_2 \le T - 1$ , and any choice of g, there exists a nonparametric bandit instance in F(1,1) such that the resulting successive elimination with static binning algorithm  $\hat{\pi}_{static}$  satisfies

$$E[R_T(\hat{\pi}_{static})] \geq \tilde{C_1}T^{\frac{9}{19}+\kappa},$$

for some  $\kappa$ ,  $\tilde{C_1} > 0$  that are independent of T. Here T  $\frac{9}{19}$  is the optimal regret achieved by BaSEDB—an successive elimination algorithm with dynamic binning.

While the formal proof is deferred to Appendix C, we would like to immediately point out the intuition underlying the failure of static binning.

Necessary choice of grid  $\Gamma$ . It is evident from the proof of the minimax lower bound (Theorem 1) that one needs to set  $t_1 \ \ T^{9/19}$ , and  $t_2 \ \ T^{15/19}$ . Otherwise, the inequality (9) guarantees the worst-case regret of  $\hat{\pi}_{static}$  exceeds the optimal one  $T_{19}$ . Consequently, we can focus on the algorithm with  $t_1 \ \ T^{9/19}$ ,  $t_2 \ \ T^{15/19}$ , and only consider the design choice g.

Why fixed g fails. As a baseline for comparison, recall that in the optimal algorithm with dynamic binning, we set  $g_0 \ \mathbb{Z} \ \mathsf{T}^{3/19}$ , and  $g_0g_1 \ \mathbb{Z} \ \mathsf{T}^{5/19}$  so that the worst case regret in three batches are all on the order of  $\mathsf{T}_{19}$ . In view of this, we split the choice of g into three cases.

- Suppose that g  $2 T^{3/19}$ . In this case, we can construct an instance such that the reward difference only appears on an interval with length 1/z 2 1/g; see Figure 3. In other words, the static binning is finer than that in the reward instance. As a result, the number of pulls in the smaller bin (used by the algorithm) in the first batch is not sufficient to tell the two arms apart, that is with constant probability, arm elimination will not happen after the first batch. This necessarily yields the blowup of the regret in the second batch.
- Suppose that g  $\[mathbb{T}\]$  In this case, we can construct an instance such that the reward difference only appears on an interval with length 1/z  $\[mathbb{T}\]$  1/g; see Figure 4. In other words, the static binning is coarser than that in the reward instance. Since the aggregated reward difference on the larger bin is so small, the number of pulls in the larger bin (used by the algorithm) in the first batch is still not sufficient to result in successful arm elimination. Again, the regret on the second batch blows up.
- Suppose that g  $\[ \]$  T  $\[ \]$  Since this choices matches  $\[ \]$  used in the optimal dynamic binning algorithm, there is no reward instance that can blow up the regret in the first two batches. Nevertheless, since  $\[ \]$  g  $\[ \]$  g  $\[ \]$  g  $\[ \]$  one can construct the instance similar to the previous case (i.e., Figure 4) such that the regret on the third batch blows up.

## 5 Adaptivity to margin parameter $\alpha$

In this section, we provide some discussions on the possibility of adapting to the margin parameter  $\alpha$  if it is unknown. Recall in Section 4, the grid choice of Algorithm 1 requires knowledge of  $\alpha$ . One may ask if such knowledge is essential in obtaining small regret. Unfortunately, the following theorem demonstrates that the price of not knowing  $\alpha$  is at least a polynomial increase in regret.



Figure 3: Instance with g > z. Each bin B produced by  $\hat{\pi}_{static}$  has width 1/g.

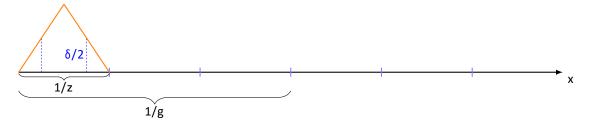


Figure 4: Instance with g < z. Each bin B produced by  $\hat{\pi}_{static}$  has width 1/g.

Theorem 5. Consider M = 2 (or 3) and  $\beta$  = d = 1. For any algorithm that does not know the true margin parameter  $\alpha^{\mathbb{Z}}$ , there exists a choice of  $\alpha^{\mathbb{Z}}$  such that

$$\sup_{+\kappa_{1}} E[R_{T}(\pi)] \ge D_{3}T^{\frac{1-\frac{\alpha^{2}+1}{\alpha^{2}+1}}} \cdot M$$

for some  $\tilde{D_3} > 0$ ,  $\kappa_1$  that are independent of T.

See Appendix D for the proof.

Theorem 5 says for any algorithm that does not have knowledge of  $\alpha^{\mathbb{Z}}$ , its regret is at least a polynomial factor larger than the optimal regret attained by Algorithm 1. This result shows batch learning for nonparametric bandits is much harder than the fully online case to some extent, where adaptivity to  $\alpha^{\mathbb{Z}}$  could be achieved for free [39].

The intuition behind the proof of Theorem 5 is that since the algorithm does not know  $\alpha^{\mathbb{Z}}$ , it has little hope to pick the first batch size  $t_1$  optimally. If  $t_1$  is too large, then the adversary can choose a big  $\alpha^{\mathbb{Z}}$ , which corresponds to the family of reward functions with larger gaps, so that the algorithm's regret during the first batch explodes. If  $t_1$  is too small, then the adversary can choose a little  $\alpha^{\mathbb{Z}}$ , which corresponds to the family of reward functions with smaller gaps, so that the algorithm's knowledge gathered during the first batch is not enough to distinguish the arms and its regret will explode in later batches.

### 6 Conclusions

In this paper, we characterize the fundamental limits of batch learning in nonparametric contextual bandits. In particular, our optimal batch learning algorithm (i.e., Algorithm 1) is able to match the optimal regret in the fully online setting with only O(log log T) policy updates. Our work open a few interesting avenues to explore in the future.

Extensions to multiple arms. With slight modification, our algorithm works for nonparametric contextual bandits with more than two arms. However, it remains unclear what the fundamental limits of batch learning are in this multi-armed case (i.e., when K is large).

Improving the log factor. Comparing the upper and lower bounds, it is evident that Algorithm 1 is near-optimal up to log factors. It is certainly interesting to improve this log factor, either by strengthening the lower bound, or making the upper bound more efficient.

Adapting to margin parameters. While we have shown that adaptivity to the margin parameter is not possible when the batch constraint is stringent, i.e., when M = 2 (or 3), it leaves open the question of designing optimal adaptive algorithm when M is large, say M  $\square$  log log T. In fact, when M = T, i.e., in the fully online setting, [39] provides an adaptively binned successive elimination algorithm that is capable of adapting to the margin parameter optimally.

## A cknowledgements

CM is partially supported by the National Science Foundation via grant DMS-2311127.

## References

- [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. Advances in neural information processing systems, 24, 2011.
- [2] Sakshi Arya and Yuhong Yang. Randomized allocation with nonparametric estimation for contextual multi-armed bandits with delayed rewards. Statistics & Probability Letters, 164:108818, 2020.
- [3] Jean-Yves Audibert and Alexandre B Tsybakov. Fast learning rates for plug-in classifiers. The Annals of Statistics, 35(2):608–633, 2007.
- [4] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. Journal of Machine Learning Research, 3(Nov):397–422, 2002.
- [5] Yu Bai, Tengyang Xie, Nan Jiang, and Yu-Xiang Wang. Provably efficient q-learning with low switching cost. Advances in Neural Information Processing Systems, 32, 2019.
- [6] Hamsa Bastani and Mohsen Bayati. Online decision making with high-dimensional covariates. Operations Research, 68(1):276–294, 2020.
- [7] Hamsa Bastani, Mohsen Bayati, and Khashayar Khosravi. Mostly exploration-free algorithms for contextual bandits. Management Science, 67(3):1329–1349, 2021.
- [8] Dimitris Bertsimas and Adam J Mersereau. A learning approach for interactive marketing to a customer segment. Operations Research, 55(6):1120–1135, 2007.
- [9] Moise Blanchard, Steve Hanneke, and Patrick Jaillet. Non-stationary contextual bandits and universal learning. arXiv preprint arXiv:2302.07186, 2023.
- [10] Changxiao Cai, T Tony Cai, and Hongzhe Li. Transfer learning for contextual multi-armed bandits. arXiv preprint arXiv:2211.12612, 2022.
- [11] T Tony Cai and Hongming Pu. Stochastic continuum-armed bandits with additive models: Minimax regrets and adaptive algorithm. The Annals of Statistics, 50(4):2179–2204, 2022.
- [12] Nicolo Cesa-Bianchi, Ofer Dekel, and Ohad Shamir. Online learning with switching costs and other adaptive adversaries. Advances in Neural Information Processing Systems, 26, 2013.
- [13] Olivier Chapelle. Modeling delayed feedback in display advertising. In Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 1097–1105, 2014.
- [14] Olivier Chapelle and Lihong Li. An empirical evaluation of thompson sampling. Advances in neural information processing systems, 24, 2011.

- [15] Stephen E Chick and Noah Gans. Economic analysis of simulation selection problems. Management Science, 55(3):421–437, 2009.
- [16] Eyal Even-Dar, Shie Mannor, Yishay Mansour, and Sridhar Mahadevan. Action elimination and stopping conditions for the multi-armed bandit and reinforcement learning problems. Journal of machine learning research, 7(6):1079–1105, 2006.
- [17] Jianqing Fan, Zhaoran Wang, Zhuoran Yang, and Chenlu Ye. Provably eficient high-dimensional bandit learning with batched feedbacks. arXiv preprint arXiv:2311.13180, 2023.
- [18] Yasong Feng, Zengfeng Huang, and Tianyu Wang. Lipschitz bandits with batched feedback. Advances in Neural Information Processing Systems, 35:19836–19848, 2022.
- [19] Manegueu Anne Gael, Claire Vernade, Alexandra Carpentier, and Michal Valko. Stochastic bandits with arm-dependent delays. In International Conference on Machine Learning, pages 3348–3356. PMLR, 2020.
- [20] Minbo Gao, Tianle Xie, Simon S Du, and Lin F Yang. A provably eficient algorithm for linear markov decision process with low switching cost. arXiv preprint arXiv:2101.00494, 2021.
- [21] Zijun Gao, Yanjun Han, Zhimei Ren, and Zhengqing Zhou. Batched multi-armed bandits problem. Advances in Neural Information Processing Systems, 32, 2019.
- [22] Alexander Goldenshluger and Assaf Zeevi. Woodroofe's one-armed bandit problem revisited. The Annals of Applied Probability, 19(4):1603–1633, 2009.
- [23] Alexander Goldenshluger and Assaf Zeevi. A linear response bandit problem. Stochastic Systems, 3(1):230–261, 2013.
- [24] Melody Guan and Heinrich Jiang. Nonparametric stochastic contextual bandits. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 32, 2018.
- [25] Yonatan Gur, Ahmadreza Momeni, and Stefan Wager. Smoothness-adaptive contextual bandits. Operations Research, 70(6):3198–3216, 2022.
- [26] Yanjun Han, Zhengqing Zhou, Zhengyuan Zhou, Jose Blanchet, Peter W Glynn, and Yinyu Ye. Sequential batch learning in finite-action linear contextual bandits. arXiv preprint arXiv:2004.06321, 2020.
- [27] Yichun Hu, Nathan Kallus, and Xiaojie Mao. Smooth contextual bandits: Bridging the parametric and nondifferentiable regret regimes. Operations Research, 70(6):3261–3281, 2022.
- [28] Cem Kalkanli and Ayfer Ozgur. Batched thompson sampling. Advances in Neural Information Processing Systems, 34:29984–29994, 2021.
- [29] Amin Karbasi, Vahab Mirrokni, and Mohammad Shadravan. Parallelizing thompson sampling. Advances in Neural Information Processing Systems, 34:10535–10548, 2021.
- [30] Edward S Kim, Roy S Herbst, Ignacio I Wistuba, J Jack Lee, George R Blumenschein Jr, Anne Tsao, David J Stewart, Marshall E Hicks, Jeremy Erasmus Jr, Sanjay Gupta, et al. The battle trial: personalizing therapy for lung cancer. Cancer discovery, 1(1):44–53, 2011.
- [31] Aniket Kittur, Ed H Chi, and Bongwon Suh. Crowdsourcing user studies with mechanical turk. In Proceedings of the SIGCHI conference on human factors in computing systems, pages 453–456, 2008.
- [32] Anders Bredahl Kock and Martin Thyrsgaard. Optimal sequential treatment allocation. arXiv preprint arXiv:1705.09952, 2017.
- [33] Akshay Krishnamurthy, John Langford, Aleksandrs Slivkins, and Chicheng Zhang. Contextual bandits with continuous actions: Smoothing, zooming, and adapting. The Journal of Machine Learning Research, 21(1):5402–5446, 2020.

- [34] Tor Lattimore and Csaba Szepesvári. Bandit algorithms. Cambridge University Press, 2020.
- [35] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In Proceedings of the 19th international conference on World wide web, pages 661–670, 2010.
- [36] Andrea Locatelli and Alexandra Carpentier. Adaptivity to smoothness in x-armed bandits. In Conference on Learning Theory, pages 1463–1492. PMLR, 2018.
- [37] Tyler Lu, Dávid Pál, and Martin Pál. Showing relevant ads via context multi-armed bandits. In Proceedings of AISTATS, 2009.
- [38] Enno Mammen and Alexandre B Tsybakov. Smooth discrimination analysis. The Annals of Statistics, 27(6):1808–1829, 1999.
- [39] Vianney Perchet and Philippe Rigollet. The multi-armed bandit problem with covariates. Ann. Statist., 41(2):693–721, 2013.
- [40] Vianney Perchet, Philippe Rigollet, Sylvain Chassang, and Erik Snowberg. Batched bandit problems. Ann. Statist., 44(2):660–681, 2016.
- [41] Wei Qian, Ching-Kang Ing, and Ji Liu. Adaptive algorithm for multi-armed bandit problem with high-dimensional covariates. Journal of the American Statistical Association, pages 1–13, 2023.
- [42] Wei Qian and Yuhong Yang. Kernel estimation and model combination in a bandit problem with covariates. Journal of Machine Learning Research, 17(149), 2016.
- [43] Wei Qian and Yuhong Yang. Randomized allocation with arm elimination in a bandit problem with covariates. Electronic Journal of Statistics, 10(1):242–270, 2016.
- [44] Dan Qiao, Ming Yin, Ming Min, and Yu-Xiang Wang. Sample-eficient reinforcement learning with loglog (t) switching cost. In International Conference on Machine Learning, pages 18031–18061. PMLR, 2022.
- [45] Henry Reeve, Joe Mellor, and Gavin Brown. The k-nearest neighbour ucb algorithm for multi-armed bandits with covariates. In Algorithmic Learning Theory, pages 725–752. PMLR, 2018.
- [46] Zhimei Ren and Zhengyuan Zhou. Dynamic batch learning in high-dimensional sparse linear contextual bandits. Management Science, 2023.
- [47] Zhimei Ren, Zhengyuan Zhou, and Jayant R Kalagnanam. Batched learning in generalized linear contextual bandits with general decision sets. IEEE Control Systems Letters, 6:37–42, 2020.
- [48] Philippe Rigollet and Assaf Zeevi. Nonparametric bandits with covariates. arXiv preprint arXiv:1003.1630, 2010.
- [49] Herbert E. Robbins. Some aspects of the sequential design of experiments. Bulletin of the American Mathematical Society, 58:527–535, 1952.
- [50] Eric M Schwartz, Eric T Bradlow, and Peter S Fader. Customer acquisition via display advertising using multi-armed bandit experiments. Marketing Science, 36(4):500–522, 2017.
- [51] Joe Suk and Samory Kpotufe. Tracking most significant shifts in nonparametric contextual bandits. arXiv preprint arXiv:2307.05341, 2023.
- [52] Joseph Suk and Samory Kpotufe. Self-tuning bandits over unknown covariate-shifts. In Algorithmic Learning Theory, pages 1114–1156. PMLR, 2021.
- [53] Ambuj Tewari and Susan A Murphy. From ads to interventions: Contextual bandits in mobile health. In Mobile Health, pages 495–517. Springer, 2017.
- [54] Alexander B Tsybakov. Optimal aggregation of classifiers in statistical learning. The Annals of Statistics, 32(1):135–166, 2004.

- [55] Claire Vernade, Olivier Cappé, and Vianney Perchet. Stochastic bandit models for delayed conversions. arXiv preprint arXiv:1706.09186, 2017.
- [56] Chi-Hua Wang and Guang Cheng. Online batch decision-making with high-dimensional covariates. In International Conference on Artificial Intelligence and Statistics, pages 3848–3857. PMLR, 2020.
- [57] Tianhao Wang, Dongruo Zhou, and Quanquan Gu. Provably efficient reinforcement learning with linear function approximation under adaptivity constraints. Advances in Neural Information Processing Systems, 34:13524–13536, 2021.
- [58] Michael Woodroofe. A one-armed bandit problem with a concomitant variable. Journal of the American Statistical Association, 74(368):799–806, 1979.
- [59] Yuhong Yang and Dan Zhu. Randomized allocation with nonparametric estimation for a multi-armed bandit problem with covariates. Ann. Statist., 30(1):100–121, 2002.
- [60] Kelly Zhang, Lucas Janson, and Susan Murphy. Inference for batched bandits. Advances in neural information processing systems, 33:9818–9829, 2020.
- [61] Zihan Zhang, Yuan Zhou, and Xiangyang Ji. Almost optimal model-free reinforcement learning via reference-advantage decomposition. Advances in Neural Information Processing Systems, 33:15198–15207, 2020.
- [62] Zhijin Zhou, Yingfei Wang, Hamed Mamani, and David G Coffey. How do tumor cytogenetics in-form cancer treatments? dynamic risk stratification and precision medicine using multi-armed bandits. Dynamic Risk Stratification and Precision Medicine Using Multi-armed Bandits (June 17, 2019).

### A Proof of Theorem 2

It is worth emphasizing that Theorem 2 aims to establish the hardness of batched nonparametric bandits even when the grid  $\Gamma$  is allowed to be adaptively chosen. Nevertheless, the proof of Theorem 1 in the fixed-grid case is still useful.

Define b  $\mathbb{Z}$  T  $(1-\gamma)/(1-\gamma^M)$ . For each  $1 \le i \le M$ , we set  $T_i = \mathbb{Z}b^{(1-\gamma^i)/(1-\gamma)}\mathbb{Z}$ ,  $z_i = \mathbb{Z}(36T_{i-1}M^2)^{1/(2\beta+d)}\mathbb{Z}$ , and  $s_i := \mathbb{Z}z^{d-\alpha\beta}\mathbb{Z}$ . We reuse the family of hard instances  $C_{z_i}$  as defined in (4), and define the mixture distribution

$$Q_{i}(\cdot) = \frac{1}{s_{i} 2^{s_{i}-1}} \sum_{j=1}^{X^{s_{i}}} \frac{X}{\omega_{[-j]} \mathbb{Z}} P_{\pi,i,f_{\omega_{[-j]}^{-1}}}(\cdot) = \frac{X}{\omega \mathbb{Z}} q_{i}(\omega) P_{\pi,i,f_{\omega}}(\cdot),$$
(13)

where  $q_i:\Omega_{s_i}\to [0,1]$  is selected so that the above equality holds,  $P_{\pi,i,f_\omega}$  is the distribution of the observations when  $f_\omega$   $\mathbb{Z}$   $C_{z_i}$ . It is easy to see  $C_{\omega\mathbb{Z}_{Q_i}}$   $Q_i(\omega)=1$ .

We pause here to state a useful claim regarding the family  $\{Q_i\}_{i=1}^M$  of mixture distributions, namely, they are all close to each other under the total variation distance.

Lemma 2. For any 
$$1 \le i \le M$$
, one has  $TV(Q_M^{T_{i-1}}, Q_i^{T_{i-1}}) \le \frac{1}{2} \frac{q}{T_{i-1}z_i^{-(2\beta+d)}}$ .

Define the event

$$A_i = \{t_{i-1} < T_{i-1} < T_i \le t_i\}.$$

Intuitively,  $A_i$  models the event when the algorithm's selected grid points  $t_{i-1}$  and  $t_i$  are suboptimal. When  $A_i$  happens, the goal is to design a problem instance such that using observations up  $t_{i-1}$  cannot distinguish the optimal arm and therefore the policy must incur a large regret between  $t_{i-1}$  and  $t_i$ .

The following lemma ensures that at least one of the bad events  $A_i$ 's happens with suficiently large probability under the mixture distribution.

The next lemma indeed shows that when  $A_i$  happens, the regret must be large.

Lemma 4. If  $Q_i(A_i) \ge 1/(2M)$ , then

$$\sup_{f \, \boxtimes C_{Z_i}} \, R_{\,T_{\,_i}}(\pi;f) \, ? \ T_{\,_i} z_i^{\,-\beta\,(1+\alpha)} \, M^{\,-\frac{1+\alpha}{\alpha}}.$$

Now we are ready to establish the desired claim in the theorem. It is straightforward to see that

$$\sup_{(f,\frac{1}{2})\boxtimes F(\alpha,\beta)} R_T(\pi;f) \geq \sup_{\widetilde{\mathbb{R}}C_{z_i \mathbb{Z}}} R_{T_i^{\mathbb{Z}}}(\pi;f) \ ? \ T_{i \mathbb{Z}} Z_{i \mathbb{Z}}^{-\beta\,(1+\alpha)} M^{-\frac{1+\alpha}{\alpha}} = D_1^{\sim}(\frac{1}{M})^{D_2^{\sim}} T^{\frac{1+\gamma}{1-\gamma-M}},$$

where the second inequality uses Lemma 4, and the last one arises from the definitions of  $T_i$  and  $z_i$ .

#### A.1 Proof of Lemma 2

It sufices to bound their KL divergence. By the standard decomposition of KL divergence and Bernoulli reward structure,

where  $f_{i,\omega}$  denotes an instance from  $C_{z_i}$ . To control  $\Delta_t$ , we further decompose it as

$$\Delta_{t} = \begin{array}{c} X^{c_{i}} & X \\ \Delta_{t} = & q_{M}(\omega)f_{M,\omega}(X_{t}) - X \\ \downarrow_{j=1} & \omega \otimes Q_{j} \end{array} q_{i}(\omega)f_{i,\omega}(X_{t}) \mathbf{1}\{X_{t} \otimes C_{i,j}\},$$

where  $C_{i,j}$  is the  $j^{th}$  bin corresponding to the instance family  $C_z$ . Here, the difference between the two sums can be restricted to  $\mathbb{F}_{j=1} C_{i,j}$  because  $\mathbb{F}_{j=1} C_{i,j}$ . Indeed, the area of effective bins for an instance in  $C_{z_i}$  is  $s_i z_i^{d} = z^{-\alpha\beta}$ , which decreases as i increases. Notice for  $X_t \mathbb{F} C_{i,j}$ ,

$$\left| \begin{array}{c} X \\ q_{i}(\omega)f_{i,\omega}(X_{t}) - \frac{1}{2} \right| \leq \frac{2^{s_{i}-1}}{s_{i} \cdot 2^{s_{i}-1}} \cdot \frac{z_{i}^{-\beta}}{4} = \frac{z_{i}^{-\beta}}{4s_{i}},$$

because all the  $\omega_{[-j]}^{-1}$ 's have a negative sign in the j<sup>th</sup> bin and there are  $2^{s_i-1}$  of them, while for each k=j, the positive and negative spikes within the j<sup>th</sup> bin cancel out each other when summing over  $\omega_{[-k]}^{-1}$ 's due to symmetry. Therefore,

$$\left| \Delta_t \right| \, \leq \, \, \frac{z_i^{-\beta}}{4s_i} \, \frac{\chi^{\! S_i}}{j \, = \, 1} \, \, \mathbf{1} \, \{ \, X_t \, \, \mathbb{P} \, \, C_{i,j} \, \} \, = \, \, \frac{1}{4} z_i^{-\beta - d + \alpha \, \beta} \, \, \frac{\chi^{\! S_i}}{j \, = \, 1} \, \, \mathbf{1} \, \{ \, X_t \, \, \mathbb{P} \, \, C_{i,j} \, \}.$$

Plugging the above back to (14) we obtain

$$\begin{split} \text{KL}\big(Q_{M}^{T_{i-1}},Q_{i}^{T_{i-1}}\big) &\leq \frac{1}{2} \frac{\bar{X}^{-1}}{t=1} \, E_{Q_{M}} \, \big[z_{i}^{-2(\beta+d-\alpha\beta)} \, \overset{\hat{X}^{i}}{X^{i}} \, 1 \big\{ X_{t} \, \mathbb{P} \, C_{i,j} \big\} 1 \big\{ \pi_{t}(X_{t}) = 1 \big\} \big] \\ &= \frac{1}{2} z_{i}^{-2(\beta+d-\alpha\beta)} \, \frac{\bar{X}^{-1}}{t=1} \, \overset{\hat{X}^{i}}{j=1} \, P_{Q_{M}} \, \big( X_{t} \, \mathbb{P} \, C_{i,j}, \pi_{t}(X_{t}) = 1 \big) \big] \\ &\leq \frac{1}{2} z_{i}^{-2(\beta+d-\alpha\beta)} \, \frac{\bar{X}^{-1}}{t=1} \, \overset{\hat{X}^{i}}{j=1} \, z_{i}^{-d} = \frac{1}{2} T_{i-1} z_{i}^{-(2\beta+d+(d-\alpha\beta))}, \end{split}$$

where the last inequality is because  $P_{Q_M}\left(X_t \ \ \ C_{i,j}, \pi_t(X_t) = 1\right) = z_i^{-d} P_{Q_M}(\pi_t(X_t) = 1 \ | \ X_t \ \ C_{i,j}) \leq z_i^{-d}$ . Since  $\alpha\beta \leq 1$ ,

$$\mathsf{KL}\big(\,Q_{\,\mathsf{M}}^{\,\mathsf{T}_{\,i-1}}\,,\,Q_{\!i}^{\,\mathsf{T}_{\,i-1}}\,\big)\,\leq\,\,\frac{1}{2}\,\mathsf{T}_{\,i-1}\,\mathsf{z}_{\,i}^{\,-(\,2\,\beta\,+\,d\,+\,(\,d\,-\,\alpha\,\beta\,)\,)}\,\leq\,\,\frac{1}{2}\,\mathsf{T}_{\,i-1}\,\mathsf{z}_{\,i}^{\,-(\,2\,\beta\,+\,d\,)}\,.$$

By Pinsker's inequality, we can conclude

## A.2 Proof of Lemma 3

For any  $1 \le i \le M$ , we have

$$|Q_{M}(A_{i}) - Q_{i}(A_{i})| \stackrel{(i)}{=} |Q_{M}^{T_{i-1}}(A_{i}) - Q_{i}^{T_{i-1}}(A_{i})| \stackrel{(ii)}{\leq} TV(Q_{M}^{T_{i-1}}, Q_{i}^{T_{i-1}}) \stackrel{(iii)}{\leq} \frac{1}{2M}, \tag{15}$$

where step (i) is because  $A_i$  can be determined by observations up to  $T_{i-1}$ , step (ii) uses the definition of TV, and step (iii) applies Lemma 2 and the definition of  $z_i$ . Consequently,

$$\begin{split} \chi^{M} & Q_{i}\left(A_{i}\right) = Q_{M}\left(A_{M}\right) + \sum_{i=1}^{N\chi-1} Q_{i}\left(A_{i}\right) \\ & = Q_{M}\left(A_{M}\right) + \sum_{i=1}^{N\chi-1} \left(Q_{i}\left(A_{i}\right) - Q_{M}\left(A_{i}\right) + Q_{M}\left(A_{i}\right)\right) \\ & \stackrel{(iv)}{\geq} Q_{M}\left(A_{M}\right) + \sum_{i=1}^{N\chi-1} \left(Q_{M}\left(A_{i}\right) - \frac{1}{2M}\right) \geq \sum_{i=1}^{\chi^{M}} Q_{M}\left(A_{i}\right) - \frac{1}{2} \stackrel{(v)}{=} \frac{1}{2}, \end{split}$$

where step (iv) uses inequality (15), and step (v) uses the fact that  $P_{i=1}^{M} Q_{M}(A_{i}) = 1$ .

### A.3 Proof of Lemma 4

We try to lower-bound the number of mistakes we make up to Ti. By inequality (5),

where the second inequality invokes Le Cam's method, and the last inequality holds since  ${}^R\min\{dP,dQ\}=1-TV(P,Q), \text{ and } TV(P_{\pi,f_{\omega_{[-j]}^{-1}}}^t,P_{\pi,f_{\omega_{[-j]}^{-1}}}^t,P_{\pi,f_{\omega_{[-j]}^{-1}}}^{T_i}) \text{ for } t\leq T_i. \text{ We continue the lower-bound to see that}$ 

$$\sup_{f \boxtimes C_{z}} S_{T_{i}}(\pi, f, \frac{1}{2}) \geq \frac{1}{2^{s}} \sum_{j=1}^{X^{s}} \sum_{\omega_{[-j]} \boxtimes \Omega_{-1}} \frac{T_{i}}{z^{d}} \sum_{min\{dP_{\pi, f_{\omega_{-j}}^{-1}}^{T_{i}}, dP_{\pi, f_{\omega_{[-j]}}^{1}}^{T_{i}}\}$$

$$\geq \frac{1}{2^{s}} \sum_{j=1}^{X^{s}} \sum_{\omega_{[-j]} \boxtimes \Omega_{-1}} \frac{T_{i}}{z^{d}} \sum_{A_{i}} \min\{dP_{\pi, f_{\omega_{-j}}^{-1}}^{T_{i}}, dP_{\pi, f_{\omega_{[-j]}}^{1}}^{T_{i}}\}$$

$$= \frac{1}{2^{s}} \sum_{j=1}^{X^{s}} \sum_{\omega_{[-j]} \boxtimes \Omega_{-1}} \frac{T_{i}}{z^{d}} \sum_{A_{i}} \min\{dP_{\pi, f_{\omega_{-j}}^{-1}}^{T_{i-1}}, dP_{\pi, f_{\omega_{-j}}^{1}}^{T_{i-1}}\},$$

where the last step uses the fact that under  $A_i$ , the available observations for  $\pi$  at  $T_i$  are the same as those at  $T_{i-1}$ . Using properties of TV, we reach

$$\sup_{f \ni C_z} S_{T_i}(\pi, f, \frac{1}{2}) \ge \frac{1}{2^s} \frac{X^s}{j=1} \frac{X}{\omega_{[-j]} \ni \Omega_{-1}} \frac{T_i}{z^d} \sum_{A_i}^{Z_i} \min\{dP_{\pi, f_{\omega_{-j}}^{-1}}^{T_{i-1}}, dP_{\pi, f_{\omega_{-j}}^{-1}}^{T_{i-1}}\}$$

$$\ge \frac{1}{2^s} \cdot \frac{T_i}{z^d} \sum_{j=1}^{X^s} \frac{X}{\omega_{[-j]} \ni \Omega_{-1}} P_{\pi, f_{\omega_{-j}}^{-1}}(A_i) - \frac{3}{2} TV(P_{\pi, f_{\omega_{-j}}^{-1}}^{T_{i-1}}, P_{\pi, f_{\omega_{-j}}^{-1}}^{T_{i-1}})$$

$$\ge \frac{1}{2^s} \cdot \frac{T_i}{z^d} \sum_{j=1}^{X^s} \frac{X}{\omega_{[-j]} \ni \Omega_{-1}} P_{\pi, f_{\omega_{-j}}^{-1}}(A_i) - \frac{3}{2} \frac{s}{2} \frac{1}{2} KL(P_{\pi, f_{\omega_{-1}}}^{T_{i-1}}, P_{\pi, f_{\omega_{1}}}^{T_{i-1}})$$

where the second inequality applies Lemma 6, and the third inequality is due to Pinsker's inequality. Take  $z = z_i$  and use Lemma 5, we have

$$\begin{split} \sup_{f \not \in C_{z_{i}}} S_{T_{i}}(\pi, f, \frac{1}{2}) &\geq \frac{1}{2^{s_{i}}} \cdot \frac{T_{i}}{z_{i}^{d}} \sum_{j=1}^{X_{S_{i}}} \frac{X}{Q_{i-1}} P_{\pi, i, f} \\ &= \frac{1}{3} \cdot \frac{T_{i}}{z_{i}^{d}} S_{i} 2^{s_{i}-1} Q_{i}(A_{i}) - \frac{3}{2} S_{i} 2^{s_{i}-1} r \frac{1}{36M^{2}} \\ &\geq \frac{1}{8} T_{i} z_{i}^{-\alpha \beta} \frac{1}{M}, \end{split}$$

where the last inequality uses the assumption  $Q_i(A_i) \ge 1/(2M)$ . Therefore, we arrive at

$$\sup_{f \boxtimes C_{Z_i}} R_{T_i}(\pi;f) ? \ T_i^{-\frac{1}{\alpha}} \sup_{f \boxtimes C_{Z_i}} S_{t_i}(\pi;f) \qquad ? \ T_i z_i^{-\beta(1+\alpha)} M^{-\frac{1+\alpha}{\alpha}}.$$

Lemma 5. Fix z > 0 and suppose  $f_{\omega} \ \mathbb{C}_z$ . For any  $n \ \mathbb{C}[T]$  and any policy  $\pi$ , one has

$$\text{KL}(P^{\,n}_{\pi,f_{\omega^{-1}_{[-j]}}},P^{n}_{\pi,f_{\omega^{1}_{[-j]}}}) \leq \, 2nz^{-(2\beta+d)}.$$

Proof. We can compute

$$\begin{split} \text{KL}(P_{\pi^{\prime}f_{\omega_{i}^{-1}}^{-1}}^{n},P_{\pi,f_{\omega^{-1}}}^{n}) &\overset{(i)}{\leq} 8E_{\pi,f_{\omega^{-1}}} \left[ \overset{X}{} (f_{\omega_{i}^{-1}}^{-1}(X_{t}) - f_{\omega_{i}^{-1}}^{-1}(X_{t}))^{2} \mathbf{1} \{\pi_{t}(X_{t}) = 1\}] \\ &\overset{(ii)}{\leq} 32D_{\phi}^{2} z^{-2\beta} E_{\pi,f_{\omega_{i}^{-1}}} \left[ \overset{X}{} 1 \{\pi_{t}(X_{t}) = 1, X_{t} \ \mathbb{Z} \ C_{j}\}] \right] \\ &\overset{(iii)}{\equiv} 32D_{\phi}^{2} z^{-(2\beta+d)} \overset{X^{n}}{P}_{\pi^{t},f_{\omega_{i}^{-1}}}^{T} (\pi_{t}(X_{t}) = 1 \ | \ X_{t} \ \mathbb{Z} \ C_{j}) \\ &\overset{(iv)}{\leq} 32D_{\phi}^{2} z^{-(2\beta+d)} n \leq 2nz^{-(2\beta+d)}. \end{split}$$

Here, step (i) uses the standard decomposition of K L divergence and Bernoulli reward structure; step (ii) is due to the definition of  $f_{\omega}$ ; step (iii) uses  $P(X_t \ \ \ C_j) = 1/z^d$ , and step (iv) arises from  $P_{\pi,f_{\frac{-1}{\omega-1}}}^t(\pi_t(X_t) = 1/z^d)$ 

$$1 \mid X_t \supseteq C_j \le 1$$
 for any  $1 \le t \le n$ .

Lemma 6. For any i 2 [M], one has

$$Z \\ \min \{ dP_{\pi,f_{\omega_{-j}}^{-1}}^{T_{i-1-1}}, dP_{\pi,f_{\omega_{-j}^{1}}}^{T_{i-1}} \} \ge P_{\pi,f_{\omega_{-j}^{-1}}}^{[-]} (A_i) - \frac{3}{2} TV(P_{\pi,f_{\omega_{-j}^{-1}}}^{T_{i-1-1}}, P_{\pi,f_{\omega_{-j}^{1}}}^{T_{i-1}}).$$

Proof. We can compute

$$\begin{split} Z & \min\{dP_{\pi,f_{\omega_{[-j]}}^{T_{i-1}}}^{T_{i-1}},dP_{\pi,f_{\omega_{[-j]}}^{T_{i-1}}}^{T_{i-1}}\} = & \frac{dP_{\pi,f_{\omega_{[-j]}}^{T_{i-1}}}^{T_{i-1}} + dP_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}} - |dP_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}} - dP_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}}| \\ & \geq & \frac{1}{2} \big(P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}} \big(A_{i}\big) + P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}} \big(A_{i}\big) - TV \big(P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}},P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}}\big) \\ & = & \frac{1}{2} \big(2P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}} \big(A_{i}\big) + P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}} \big(A_{i}\big) - P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}} \big(A_{i}\big) - TV \big(P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}},P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}}\big) \\ & \geq & P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}} \big(A_{i}\big) - \frac{1}{2} TV \big(P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}},P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}}\big) - TV \big(P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}},P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}}\big) \\ & = & P_{\pi,f_{\omega_{[-j]}}}^{-1} \big(A_{i}\big) - \frac{3}{2} TV \big(P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}},P_{\pi,f_{\omega_{[-j]}}}^{T_{i-1}}\big), \end{split}$$

where step (i) is due to  $|P_{\pi,f_{\omega_{[-j]}^{-1}}}^{T_{i-1}}(A_i) - P_{\pi,f_{\omega_{[-j]}^{-1}}}^{T_{i-1}}(A_i)| \le TV(P_{\pi,f_{\omega_{[-j]}^{-1}}}^{T_{i-1}},P_{\pi,f_{\omega_{[-j]}^{1}}}^{T_{i-1}})$ , and step (ii) uses the fact

that 
$$P_{\pi,f_{\omega}}^{-1}$$
 and  $P_{\pi^{'}f_{\omega^{-1}}}^{T_{f_{\omega^{-1}}}}$  are equivalent on  $A_i$ .

#### B Proof of Theorem 3

Our proof of Theorem 3 is inspired by the framework developed in [39]. Our setting presents additional technical dificulty due to the batch constraint.

We begin with introducing some useful notations. Recall the tree growing process described in section 4, where we have defined a tree T of depth M. The root (depth 0) of the tree is the whole space X. In depth 1, X has  $g^d$  children, each of which is a bin of width  $1/g_0$ . For each bin in depth 1, it has  $g^d$  children, each of which is a bin of width  $1/(g_0g_1)$ . These children form the depth 2 nodes of the tree T. We form the tree recursively until depth M.

For a bin C ② T, we define its parent by  $p(C) = \{C^{'} ②$  T : C ② child $(C^{'})\}$ . Moreover, we let  $p^{1}(C) = p(C)$  and define  $p^{k}(C) = p(p^{k-1}(C))$  for  $k \ge 2$  recursively. In all, we denote by  $P(C) = \{C^{'} ②$  T :  $C^{'} = p^{k}(C)$  for some  $k \ge 1\}$  all the ancestors of the bin C.

We also define  $L_t$  to be the set of active bins at time t, with the dummy case  $L_0 = \{X\}$ . Clearly, for  $1 \le t \le t_1$ , one has  $L_1 = B_1$ , where  $B_1$  are all the bins in the first layer.

#### B.1 Two clean events

$$m_{C,i} := \sum_{\substack{t=t_{i-1}+1}}^{X^{t_i}} 1\{X_t ? C\}.$$

Clearly, it has expectation

$$m_{C,i}^{2} = E[m_{C,i}] = (t_{i} - t_{i-1}) P_{X}(X ? C).$$

The first clean event claims that  $m_{C,i}$  concentrates well around its expectation  $m_{C,i}^{\mathbb{Z}}$  uniformly over all C  $\mathbb{Z}$  T. We denote this event by E.

Lemma 7. Suppose that  $M \le D_1 \log(T)$  for some constant  $D_1 > 0$ . With probability at least 1 - 1/T, for all  $1 \le i \le M$ , and  $C \supseteq L_{i_{i-1}+1}$ , we have

$$\frac{1}{2} m_{C,i}^{\mathbb{Z}} \leq m_{C,i} \leq \frac{3}{2} m_{C,i}^{\mathbb{Z}}.$$

See Section B.5.1 for the proof.

Since  $M \le D_1 \log(T)$  by assumption, we can apply Lemma 7 to obtain

$$E[R_T(\hat{\pi})1(E^c)] \le TP(E^c) = 1.$$

Therefore, in the remaining proof, we condition on E and focus on bounding  $E[R_T(\hat{\pi})1(E)]$ .

The second clean event is on the elimination process. Since we use successive elimination in each bin, it is natural to expect that the optimal arm in each bin is not eliminated during the process. To mathematically specify this event, we need a few notations.

For each bin C  $2L_i$ , let  $L_i'$  be the set of remaining arms at the end of batch i, i.e., after Algorithm 2 is invoked. Define

$$\bar{I}_{C} = k \ 2 \{1, -1\} : \sup_{x \in C} f^{2}(x) - f^{(k)}(x) \le c_{1} |C|^{\beta}$$
,

$$\underline{I}_{C} = k \ [ \{1, -1\} : \sup_{x \in C} f^{(k)}(x) - f^{(k)}(x) \le c_0 |C|^{\beta} ,$$

where  $c_0 = 2Ld^{\beta/2} + 1$  and  $c_1 = 8c_0$ . Clearly, we have

Define a good event  $A_C = \{ \underline{I}_C \ \ \underline{C} \ \ \underline{I}_C' \ \ \underline{C} \ \ \underline{I}_C' \}$ , which is the event that the remaining arms in C have gaps of correct order. In addition, define  $G_C = \bigcap_{C' \ \underline{C} \ P(C)} A_{C'}$ . Recall  $B_i$  is the set of bins C with  $|C| = (\bigcap_{i \in C_1} \underline{g}_i)^{-1} = w_i$  for  $i \ge 1$ .

Lemma 8. For any  $1 \le i \le M - 1$  and  $C \supseteq B_i$ , we have

$$P(E \cap G_C \cap A_C^c) \leq \frac{4m_{C,i}^{\mathbb{Z}}}{T C|d}.$$

In words, Lemma 8 guarantees that  $A_C$  happens with high probability if E holds and  $A_C$  holds for all the ancestors C' of C. See Section B.5.2 for the proof.

## B.2 Regret decomposition

In this section, we decompose the regret into three terms. First, for a bin C, we define

$$r_T^{live}(C) := \int_{t=1}^{X^T} f^{2}(X_t) - f^{(\pi_t(X_t))}(X_t) \mathbf{1}(X_t 2 C) \mathbf{1}(C 2 L_t).$$

In addition, define  $J_t := \mathbb{Z}_{s \le t} L_s$  to be the set of bins that have been live up until time t. Correspondingly we define

$$r_T^{born}(C) \coloneqq \underset{t=1}{\overset{X_T}{=}} f^{\textcircled{2}}(X_t) - f^{(\pi_t(X_t))}(X_t) \quad \mathbf{1}(X_t \ \textcircled{2} \ C) \mathbf{1}(C \ \textcircled{2} \ J_t).$$

$$\begin{split} r_{T}^{born}(C) &= r_{T}^{live}(C) + \underset{C' \boxtimes child(C)}{X} \\ &= r_{T}^{born}(C) \mathbf{1}(A_{C}^{c}) + r_{T}^{live}(C) \mathbf{1}(A_{C}) + \underset{C' \boxtimes child(C)}{X} r_{T}^{born}(C') \mathbf{1}(A_{C}). \end{split}$$

Applying this relation recursively leads to the following regret decomposition:

$$\begin{split} R_{T}\left(\pi\right) &= r_{T}^{born}(X) \\ &= r_{T}^{live}\left\{\frac{X}{2}\right\} + \sum_{C' \supseteq child(X)} r_{T}^{born}(C') \\ &= \frac{X}{1 \le i < M} \frac{?}{?} X r_{T}^{born}(C) \mathbf{1}(G_{C} \cap A_{C}^{c}) + X r_{T}^{live}(C) \mathbf{1}(G_{C} \cap A_{C}) \frac{?}{?} \\ &+ \sum_{T} r_{T}^{live}(C) \mathbf{1}(G_{C}), \end{split}$$

where the second equality arises from the fact that  $r_T^{live}(X) = 0$ . Indeed,  $X \not \square L_t$  for any  $1 \le t \le T$ .

### B.3 Controlling three terms

In what follows, we control V<sub>i</sub>, U<sub>i</sub> and the last batch separately.

#### B.3.1 Controlling V<sub>i</sub>

Fix some  $1 \le i \le M-1$ , and some bin C  $2B_i$ . On the event  $G_C$  we have  $I_{p(C)}$   $2I_{p(C)}$ , that is, for any k  $2I_{p(C)}$ ,

$$\sup_{x \ni p(C)} f^{\mathbb{Q}}(x) - f^{(k)}(x) \leq c_1 |p(C)|^{\beta}.$$

This implies that for any  $x \ 2 \ C$ , and  $k \ 2 \ I_{p(C)}$ ,

$$f^{\mathbb{B}}(x) - f^{(k)}(x) \quad \mathbf{1}\{G_C\} \le c_1 |p(C)|^{\beta} \mathbf{1}(0 < f^{(1)}(x) - f^{(-1)}(x) \le c_1 |p(C)|^{\beta}). \tag{16}$$

As a result, we obtain

$$E[r_{T}^{live}(C)1(G_{C} \cap A_{C})] = E \int_{t=1}^{T} f^{2}(X_{t}) - f^{(\pi_{t}(X_{t}))}(X_{t}) 1(X_{t} 2C)1(C 2L_{t})1(G_{C} \cap A_{C})$$

$$\begin{array}{l} \text{\#} \\ \leq E \\ & \sum_{t=1}^{|T|} c_1 |p(C)|^{\beta} \mathbf{1}(0 < f^{(1)}(X_t) - f^{(-1)}(X_t) \leq c_1 |p(C)|^{\beta}) \mathbf{1}(X_t \ \mathbb{Z} \ C, C \ \mathbb{Z} \ L_t) \mathbf{1}(G_C \cap A_C) \\ \geq & \sum_{t=1}^{|G|} \mathbb{Z} \\ \text{\#} \\ \text{(S)} \ c_1 |p(C)|^{\beta} E \ \mathbb{Z} \\ & \sum_{t=t_{i-1}+1}^{|T|} \mathbb{Z} \\ \leq & c_1 |p(C)|^{\beta} \\ & \sum_{t=t_{i-1}+1}^{|T|} P(0 < f^{(1)}(X_t) - f^{(-1)}(X_t) \leq c_1 |p(C)|^{\beta}, X_t \ \mathbb{Z} \ C) \\ = & c_1 |p(C)|^{\beta} (t_i - t_{i-1}) P(0 < f^{(1)}(X) - f^{(-1)}(X) \leq c_1 |p(C)|^{\beta}, X \ \mathbb{Z} \ C). \end{array}$$

Taking the sum over all bins in  $B_i$  and using the fact that  $|p(C)| = w_{i-1}$ , we obtain

$$\begin{array}{l} X \\ \subset \mathbb{B}_{i} \end{array} E [r_{T}^{live}(C) \mathbf{1}(G_{C} \cap A_{C})] \leq & X \\ \subset \mathbb{B}_{i} \end{array} c_{1} w_{i-1}^{\beta}(t_{i} - t_{i-1}) P(0 < f^{(1)}(X) - f^{(-1)}(X) \leq c_{1} |p(C)|^{\beta}, X \ \mathbb{C}C) \\ & = c_{1} w_{i-1}^{\beta}(t_{i} - t_{i-1}) X \\ \subset \mathbb{B}_{i} \end{array} P(0 < f^{(1)}(X) - f^{(-1)}(X) \leq c_{1} w_{i-1}^{\beta}, X \ \mathbb{C}C). \end{aligned}$$

Note that

where the last inequality follows from the margin condition. Combining relations (18) and (17), we reach

X
$$E[r_{T}^{live}(C)1(G_{C} \cap A_{C})] \leq (t_{i} - t_{i-1}) \cdot [c_{1}w_{i-1}^{\beta}]^{1+\alpha} \cdot D_{0}.$$
C@B<sub>i</sub>

#### B.3.2 Controlling Ui

Fix some  $1 \le i \le M - 1$ , and some bin C  $\mathbb{Z}$   $B_i$ . Again, using the definition of  $G_C$ , we obtain

$$\begin{split} E[r_T^{born}(C)\mathbf{1}(G_C \cap A_C^c)] &= E & f^{\textcircled{2}}(X_t) - f^{(\pi_t(X_t))}(X_t) \quad \mathbf{1}(X_t \textcircled{2} C)\mathbf{1}(C \textcircled{2} J_t)\mathbf{1}(G_C \cap A_C^c) \\ & \overset{t=1}{X^T} & \\ &\leq E & c_1|p(C)|^{\beta}\mathbf{1}(0 < f^{(1)}(X_t) - f^{(-1)}(X_t) \leq c_1|p(C)|^{\beta})\mathbf{1}(X_t \textcircled{2} C, C \textcircled{2} J_t)\mathbf{1}(G_C \cap A^c)_C \\ &\leq c_1|p(C)|^{\beta}TP(0 < f^{(1)}(X) - f^{(-1)}(X) \leq c_1|p(C)|^{\beta}, X \textcircled{2} C)P(G_C \cap A^c)_C \end{split}$$

Apply Lemma 8 to see that

$$E[r^{born}(C)\mathbf{1}(G_{C}\cap A^{c})] \leq c_{1}|p(C)|^{\beta}TP(0 < f^{(1)}(X) - f^{(-1)}(X) \leq c_{1}|p(C)|^{\beta}, X \ \boxdot C)_{T}|\underbrace{c_{C}}_{d}^{\frac{1}{2}}$$

$$= c_1 w_{i-1}^{\beta} P(0 < f^{(1)}(X) - f^{(-1)}(X) \le c_1 w_{i-1}^{\beta}, X ? C)^{4(t_i - t_{i-1})} P_X(X ? C) \frac{|C|^d}{|C|^d}$$

$$\le 4\bar{c}c_1 w_{i-1}^{\beta} P(0 < f^{(1)}(X) - f^{(-1)}(X) \le c_1 w_{i-1}^{\beta}, X ? C)(t_i - t_{i-1}),$$

where we use the fact that  $P_X(X \ \mathbb{Z} \ C) \le \bar{c} |C|^d$  in the second inequality. Summing over all bins in  $B_i$ , we obtain

$$\begin{split} & \underset{C \boxtimes B_{i}}{X} & E[r_{T}^{born}(C)\mathbf{1}(G_{C} \cap A_{C}^{c})] \leq \, 4\bar{c}c_{1}w_{i-1}^{\beta}(t_{i}-t_{i-1}) & P(0 < \, f^{(1)}(X) - \, f^{(-1)}(X) \leq \, c_{1}w_{i-1}^{\,\,\beta}, X \boxtimes C) \\ & \overset{C \boxtimes B}{h}_{i} \, i \\ & \leq \, 4\bar{c}c_{1}w_{i-1}^{\beta}(t_{i}-t_{i-1})D_{0} \cdot \, c_{1}w_{i-1}^{\beta} \\ & = \, 4D_{0}\bar{c}(t_{i}-t_{i-1})[c_{1}w_{i-1}^{\beta}]^{1+\alpha}, \end{split}$$

where the second inequality reuses the bound in (18).

#### B.3.3 Last Batch

$$E[r_{\tau}^{live}(C)1(G_C)] \le c_1|p(C)|^{\beta}(T-t_{M-1})P(0 < f^{(1)}(X) - f^{(-1)}(X) \le c_1|p(C)|^{\beta}, X \ @ C).$$

Consequently, summing over C B BM yields

$$\begin{array}{l} X \\ \in [r_T^{live}(C)1(G_C)] \leq & X \\ c_{\mathbb{R}B_M} & c_1|p(C)|^{\beta}(T-t_{M-1})P(0 < f^{(1)}(X) - f^{(-1)}(X) \leq c_1|p(C)|^{\beta}, X \ \hline{?}\ C) \\ & & b \\ & \leq c_1 w_{M-1}^{\beta}(T-t_{M-1})D_0 \cdot c_1 w_{M-1}^{\beta} \\ & = D_0(T-t_{M-1})[c_1 w_{M-1}^{\beta}]^{1+\alpha}. \end{array}$$

### B.4 Putting things together

In sum, the total regret is bounded by

$$E\left[R_{T}\left(\pi\right)\right] \leq c \quad t_{1} + \sum_{i=2}^{M} t_{i-1}^{1} \cdot w_{i-1}^{\beta+\alpha\beta} + (T - t_{M-1})w_{M-1}^{\beta+\alpha\beta} \quad ,$$

where c is a constant that depends on  $(\alpha, \beta, D, L)$ . Recall that  $w_i = (Q_{i-1} g_i)^{-1}$ , and the choices for the batch size and the split factors (12)-(11). We then obtain

$$\begin{split} t_1 & ? T^{\frac{1-\gamma}{1-\gamma - M}} log \, T, \\ (t_i - t_{i-1}) & \cdot w_{i-1}^{\beta + \alpha \beta} & ? T^{\frac{1-\gamma}{1-\gamma - M}} log \, T, \qquad \text{for } 2 \leq i \leq M-1, \\ (T - t_{M-1}) w_{M-1}^{\beta + \alpha \beta} & \le T w_{M-1}^{\beta + \alpha \beta} & ? T^{\frac{1-\gamma}{1-\gamma - M}} log \, T. \end{split}$$

The proof is finished by combining the above three bounds.

## B.5 Proofs for the clean e<sub>v</sub>ents

We are left with proving that the two clean events happen with high probability.

#### B.5.1 Proof of Lemma 7

Fix the batch index i, and a node C in layer-i of the tree T. By relation (12), we have

$$\begin{split} m_{C,i}^{2} &= (t_{i} - t_{i-1}) P_{X}(X ? C) \\ &? |C|^{-(2\beta+d)} \log(T |C|^{d}) P_{X}(X ? C) \\ ? |C|^{-2\beta} &\geq g_{0}^{2\beta} ? (T^{\frac{1+\gamma}{1-\gamma}M} \cdot \frac{2^{2\beta}}{2^{\beta+d}}), \end{split}$$

where the last step uses the fact that  $P_X(X \ \ C) \ge \underline{c}|C|^d$ . Therefore,  $m_{C,i}^2 \ge \frac{3}{4}\log(2T^2)$  for all i and C, as long as T is suficiently large. This allows us to invoke Chernoff's bound to obtain that with probability at most  $1/T^2$ 

$$X \xrightarrow[t=t_{i-1}+1]{} 1\{X_t ? C\} - m_{C,i}^{2} \ge q \frac{1}{3 \log(2T^2)m_{C,i}^{2}}$$

Denote  $E^c = \{ 21 \le i \le M, C \ge L_{t_{i-1}+1} \text{ such that } \mid P_{t=t_{i-1}+1}^{p} 1\{X_t \ge C\} - m_{C,i}^2 \mid \ge q \frac{q}{3 \log(2T^2) m_{C,i}^2} \}$ . Applying union bound to reach

$$P(E^{c}) \leq X \frac{1}{T^{2}} \leq \frac{1}{T^{2}} (i) \frac{1}{T^{2}} (g_{I})^{d} \leq \frac{1}{T^{2}} \cdot M \cdot (g_{I})^{d},$$

where step (i) sums over all possible nodes of T across batches, and step (ii) is due to  $(Q_{i-1} g_i)^d \le (Q_{M-1} g_i)^d$  for any  $1 \le i \le M$ . Since  $g_{M-1} = 1$ , we further obtain

$$P(E^{c}) \leq \frac{1}{T^{2}} \cdot M \cdot (\int_{1-\Omega}^{M_{Y}-2} g_{I})^{d} \stackrel{(iii)}{\leq} \frac{1}{T^{2}} \cdot M \cdot t_{M-1}^{\frac{d}{2\beta+d}} \stackrel{(iv)}{\leq} D_{1} \frac{1}{T^{2}} \cdot \log T \cdot T^{\frac{d}{2\beta+d}} \leq \frac{1}{T},$$

where step (iii) invokes relation (12), and step (iv) uses the assumption  $M \le D_1 \log T$ . This completes the proof.

## B.5.2 Proof of Lemma 8

To simplify notation, for any event F, we define  $P^{G_C}(F) = P(E \cap G_C \cap F)$ .

Let  $D_C^1$  be the event that an arm  $k \ \underline{\square} \ \underline{I}_C$  is eliminated at the end of batch i, and  $D_C^2$  be the event that an arm  $k \ \underline{\square} \ \underline{I}_C$  is not eliminated at the end of batch i. Consequently, we have

$$\mathsf{P}^{\mathsf{G}_{\mathsf{C}}}(\mathsf{A}_{\mathsf{C}}^{\mathsf{c}}) = \mathsf{P}^{\mathsf{G}_{\mathsf{C}}}(\mathsf{D}_{\mathsf{C}}^{\mathsf{1}}) + \mathsf{P}^{\mathsf{G}_{\mathsf{C}}}((\mathsf{D}_{\mathsf{C}}^{\mathsf{1}})^{\mathsf{c}} \cap \mathsf{D}_{\mathsf{C}}^{\mathsf{2}}).$$

Recall  $U(\tau,T,C)=4^{q}\frac{q}{\frac{\log(2T|C|^{d})}{\tau}}$  . By relation (12), we can write

$$m_{C,i}^{2} = (t_{i} - t_{i-1}) P_{X}(X \ 2 \ C)$$
  
=  $|i| C |^{-(2\beta+d)} |og(T | C |^{d}) P_{X}(X \ 2 \ C)$ .

where  $I_i > 0$  is a constant chosen such that  $U(2m_{C,i}^{\mathbb{Z}}, T, C) = 2c_0|C|^{\beta}$ . Under E, we have  $U(m_{C,i}, T, C) \le 4c_0|C|^{\beta}$  because  $m_{C,i} \ge \frac{1}{2}m_{C,i}^{\mathbb{Z}}$ .

1. Upper bounding  $P^{G_C}(D_C^1)$ : when  $D_C^1$  occurs, an arm  $k \ \underline{\ } \ \underline{I}_C$  is eliminated by some  $k' \ \underline{\ } \ I_{p(C)}$  at the end of batch i. This means  $\bar{Y}_{C,i}^{(k')} - \bar{Y}_{C,i}^{(k)} > U(m_{C,i},T,C)$ . Meanwhile,

$$\bar{f_{C}^{(k')}} - \bar{f_{C}^{(k)}} \leq \bar{f_{C}^{(k)}} - \bar{f_{C}^{(k)}} \leq \bar{f_{C}^{(k)}} \leq c_{0} |C|^{\beta} \leq \frac{1}{2} U(2m_{C,i}^{2}, T, C),$$

where step (i) uses the definition of  $\underline{L}_C$ . Consequently,  $|\bar{Y}_{C,i}^{(k')} - \bar{f}_C^{(k')}| \le U(m_{C,i},T,C)/4$  and  $|\bar{Y}_{C,i}^{(k)} - \bar{f}_C^{(k)}| \le U(m_{C,i},T,C)/4$  cannot hold simultaneously. Otherwise, this would contradict with  $\bar{Y}_{C,i}^{(k')} - \bar{Y}_{C,i}^{(k)} > U(m_{C,i},T,C)$  because  $m_{C,i} \le 2m_{C,i}^{\mathbb{Z}}$  under E. Therefore,

$$P^{G_{C}} \big( D^{1}_{C} \big) \leq \ P \quad \mathbb{E} k \, \mathbb{E} \, I_{p(C)}^{\;\; \prime}, \, m_{C,i} \leq \ 2 m_{C,i}^{\,\mathbb{B}} : \, |\bar{Y}_{C,i}^{\,(k)} - \, f_{C}^{(k)}| \, \geq \, \frac{1}{4} U \left( m_{C,i}, T, C \right) \ .$$

2. Upper bounding  $P^{G_c}((D_C^1)^c \cap D_C^2)$ : when  $(D_C^1)^c \cap D_C^2$  happens, no arm in  $\underline{L}_C$  is eliminated while some  $k \not \square I_C$  remains in the active arm set. By definition, there exists  $x^{(k)}$  such that  $f^{\not \square}(x^{(k)}) - f^{(k)}(x^{(k)}) > 8c_0 |C|^{\beta}$ . Let  $\eta(k)$  be any arm that satisfies  $f^{\not \square}(x^{(k)}) = f^{(\eta(k))}(x^{(k)})$ , and one can easily verify  $\eta(k) \not \square \underline{L}_C$ . Since k is not eliminated, we have  $\bar{Y}_{C,i}^{(\eta(k))} - \bar{Y}_{C,i}^{(k)} \le U(m_{C,i}, T, C)$ . On the other hand,

$$f_{C}^{(\eta(k))} \stackrel{(iii)}{\geq} f^{(\eta(k))}(x^{(k)}) - c_{0}|C|^{\beta} 
\geq f^{(k)}(x^{(k)}) + 8c_{0}|C|^{\beta} - c_{0}|C|^{\beta} 
= f^{(k)}(x^{(k)}) + 7c_{0}|C|^{\beta} 
\stackrel{(iv)}{\geq} f_{C}^{(k)} + 6c_{0}|C|^{\beta} \geq f_{C}^{(k)} + \frac{3}{2}U(m_{C,i}, T, C),$$
(19)

where steps (iii) and (iv) use Lemma 10. Inequality (19) together with the fact that  $\bar{Y}_{C,i}^{(\eta(k))} - \bar{Y}_{C,i}^{(k)} \leq U(m_{C,i},T,C)$  imply  $|\bar{Y}_{C,i}^{(k_0)} - \bar{f}_{C}^{(k_0)}| \geq U(m_{C,i},T,C)/4$  for either  $k_0 = k$  or  $k_0 = \eta(k)$ . Consequently,

$$P^{G_{C}}((D^{1})^{c} \cap D^{2}) \leq P \mathbb{Z}k \mathbb{Z}I^{'}_{(C)}, m_{C,i} \leq 2m^{\mathbb{Z}}_{i} : |Y_{C,1}(\overline{k})|^{1/2} + \sum_{C} U(m_{C,i}, T, C).$$

Combining the two parts we obtain

$$\begin{split} P^{G_{C}}(A_{C}^{c}) &= P^{G_{C}}(D_{C}^{1}) + P^{G_{C}}((D_{C}^{1})^{c} \cap D_{C}^{2}) \\ &\leq 2 \cdot P \quad \mathbb{Z}k \, \mathbb{Z} \, I_{p(C)}, \, m_{C,i} \leq 2 m_{\S_{p}i} : |Y_{-}^{(k)} - f_{-}^{(k)}| \geq \underbrace{1}_{C,i} U(m_{C,i}, T, C) \\ &\leq \frac{4 m_{C,i}^{\mathbb{Z}}}{T \, |C|^{d}}, \end{split}$$

where the last inequality applies Lemma 9.

## B.6 Auxiliary lemmas

Lemma 9. For any  $1 \le i \le M - 1$  and  $C \supseteq B_i$ , one has

$$P \ \ \mathbb{E} k \ \mathbb{E} \ I_{p(C)}^{'}, m_{C,i} \leq 2m_{C,i}^{\mathbb{S}} : |\bar{Y}_{C,i}^{(k)} - f_{C}^{(k)}| \geq \frac{1}{4} U(m_{C,i}, T, C) \leq \frac{2m_{C,i}^{\mathbb{S}}}{T |C|^{d}}.$$

Proof. Recall in Algorithm 1 we pull each arm in a round-robin fashion within a bin during batch i. Fix  $\tau > 0$ . Let  $Y^{(k)} = \frac{\tau}{\tau} Y^{(k)}_j / \tau$  where  $Y^{(k)}_j$ 's are i.i.d. random variables with  $Y^{(k)}_j \supseteq [0,1]$  and  $E[Y^{(k)}_j] = f^{(k)} - By$  Hoeffding's inequality, with probability  $1/(T \mid C \mid^d)$ , we have

$$|\tilde{Y}_{\tau}^{(k)} - f_{C}^{(k)}| \ge \frac{r}{\frac{\log(2T|C|^d)}{2\tau}}.$$

Applying union bound to get

which completes the proof.

Lemma 10. Fix  $k \ 2 \{1, -1\}$  and  $C \ 2 T$ , for any  $x \ 2 C$ , one has

$$|\bar{f}_{C}^{(k)} - f^{(k)}(x)| \le c_{0}|C|^{\beta},$$

where  $c_0 = 2Ld^{\beta/2} + 1$ .

Proof. For notation simplicity, we write f for  $f^{(k)}$  in the following proof. By definition,

$$|\bar{f}_C - f(x)| = \left| \frac{1}{P(C)} Z^C (f(y) - f(x)) dP(y) \right|$$

$$\leq \frac{1}{P(C)} Z^C |f(y) - f(x)| dP(y)$$

$$\leq \frac{1}{P(C)} Z^C L^2 x - y \mathbb{Z}_2^{\beta} dP(y),$$

where the first inequality uses the triangle inequality, and the second inequality is due to the smoothness condition. Since  $x \ 2 \ C$ , we further have

$$\begin{split} |\bar{f}_C - f(x)| &\leq \frac{1}{P(C)} Z^C \\ &\leq \frac{1}{P(C)} Z^C \\ &\leq \frac{1}{C} Ld^{\beta/2} |C|^{\beta} dP(y) \\ &\leq c_0 |C|^{\beta}. \end{split}$$

This completes the proof.

## C Proof of Theorem 4

$$f(x) = \frac{1}{2} + \phi_1(x).$$

The problem instance of interest is  $v = (f^{(1)}(x) = f(x), f^{(-1)}(x) = \frac{1}{2})$ . It is easy to verify v P(1, 1). Throughout the proof, we condition on the event E specified by Lemma 7, which says the number of samples allocated to a bin concentrates well around its expectation. We will show even under this good event, there exists a choice of z that makes successive elimination fail to remove the suboptimal arms at the end of a batch with constant probability.

### C.1 A helper lemma

We begin with presenting a helper lemma that will be used extensively in the later part of the proof. The claim is intuitive: if the sample size is small, it is not suficient to tell apart two Bernoulli distributions with similar means. Then, in our context, arm elimination will not occur.

Lemma 11. Assume  $m_{B,i} \le 2m_{B,i}^{\mathbb{Z}}$ . For any B  $\mathbb{Z}[0,1]$  and i  $\mathbb{Z}\{1,2\}$ . If  $f^{(\frac{1}{B})} - f^{(\frac{-1}{B})} \le \delta \le 1/\frac{p}{m_{B,i}^{\mathbb{Z}}}$  for some  $\delta > 0$ , then

$$P \ \ \bar{Y}_{B,\,i}^{\,(1)} - \ \bar{Y}_{B,\,i}^{\,(-1)} > \ U\left(m_{B,\,i},\,T\,,\,B\right) \ \le \ \frac{t_i}{T}.$$

Proof. Fix  $0 < \tau \le m_{B,i}^{\mathbb{Z}}$ . Let  $\bar{Y}_{\tau}^{(k)} = \prod_{l=1}^{P} Y_{l}^{(k)}/\tau$  where  $Y_{l}^{(k)}$  and  $Y_{l}^{(k)} = \bar{Y}_{l}^{(k)}$  for  $\mathbb{Z}[1, -1]$ . Recall  $U(\tau, T, B) = 4$   $\frac{\log(2T|B|)^{3}}{\tau}$ . Then,

$$PY_{\tau}^{-1()} - Y_{\tau}^{(-1)} > U(2\tau, T, B) \le P_{\tau}^{(i)} Y_{\tau}^{(1)} - Y_{\tau}^{(-1)} > \delta + r_{\tau}^{r} \frac{!}{2\tau}$$

 $<sup>\</sup>frac{3 \text{ We remark the constant 4 is not essential for the proof to work. For any c > 0 , clog(2T|B|) = log((2T|B|)^c)$  so the final success probability is still tiny as long as T is sufficiently large.

$$\stackrel{\text{(iii)}}{\leq} P \quad \bar{Y}_{\tau}^{(1)} - \bar{Y}_{\tau}^{(-1)} > f_{B}^{(1)} - f_{B}^{(-1)} + \frac{r \log(2T/g)}{2\tau}$$

$$\stackrel{\text{(iiii)}}{\leq} \frac{g}{T},$$

where step (i) is because  $\delta \leq 1/\frac{D}{m_{B,i}^{\frac{m}{B}}} \leq 1/\sqrt{\tau}$ , step (ii) is due to  $f_B^{(1)} - f_B^{(-1)} \leq \delta$ , and step (iii) uses Hoeffding's inequality. Applying union bound to get

$$P \, \mathbb{P} 0 < \tau \leq \, m^{\mathbb{P}_{\gamma_{\dot{B}}^{\dot{i}}}} : Y_{\tau}^{\, (\!\! 1)} - \, Y_{\tau}^{\, (\!\! -1)} > \, U \, (2\tau, T, B) \leq \, \frac{m_{B, i} \, \mathcal{B}}{T} \frac{\leq}{T} \, T \, . \, \frac{t_{i}}{T} \, . \, \frac{t_$$

This finishes the proof.

## C.2 Three failure cases for g

Fix some small constant  $\epsilon > 0$  to be specified later. From now on, we use  $\hat{\pi}$  to denote  $\hat{\pi}_{static}$  for simplicity. We split the proof into three cases: (1)  $g \geq T^{3/19+\epsilon}$ ; (2)  $g \leq T^{3/19-\epsilon}$ ; (3) and  $g \supseteq (T^{3/19-\epsilon}, T^{3/19+\epsilon})$ .

Case 1:  $g \ge T^{3/19+\epsilon}$ . Set  $z = T^{3/19-\epsilon/2}$ . Assume without loss of generality that  $g = H \cdot z$  for some  $H \ge 4$ ; see Figure 3 for an illustration of the instance. Suppose  $C_1 = \mathbb{Z}^H$   $B_{1}$  where  $B_1$ 's are the bins produced by  $\hat{\pi}$  that lie in  $C_1$ . It is clear that

$$E[R_{T}(\hat{\pi})] \stackrel{(i)}{\geq} E \qquad f^{\mathbb{Z}}(X_{t}) - f^{\hat{\pi}_{t}(X_{t})}(X_{t}) \qquad \#$$

$$= \frac{\sum_{t=t_{1}+1}^{t} \qquad \#}{E[X_{t}) - f^{\hat{\pi}_{t}(X_{t})}(X_{t}) \quad 1\{X_{t} \ \mathbb{Z} \ C_{1}\} \ t = t_{1}+1} \qquad i$$

$$= \frac{\sum_{t=t_{1}+1}^{t} \frac{t_{2}}{1 = H/4}}{E[X_{t}) - f^{\hat{\pi}_{t}(X_{t})}(X_{t}) \quad 1\{X_{t} \ \mathbb{Z} \ B_{1}\} , \qquad (20)$$

where step (i) is because the total regret is greater than the regret incurred during the second batch, step (ii) uses the fact that under the instance v, the mean rewards of the two arms differ only in  $C_1$ , and step (iii) arises since  $C_1 = \mathbb{Z}^H$   $B_{I_{t=1}}$ Now we turn to lower bounding E  $f^{\mathbb{Z}}(X_t) - f^{\hat{\pi}}t^{(X_t)}(X_t)$   $1\{X_t \mathbb{Z} B_I\}$  for each  $H/4 \le I \le 3H/4$ .

Consider any such  $B_1$ . We drop the subscripts and write B instead for simplicity. By the design of v, we have  $f_B^{(1)} - f_B^{(-1)} \le D_\phi z^{-1} = \delta$ , which obeys  $D_\phi z^{-1} \le 1/p \, \overline{m_{B,1}}$ —a consequence of the choice of z. Additionally, we have  $m_{B,1} \le 2m_{B,1}^{\mathbb{Z}}$  under E. Therefore, we can invoke Lemma 11 to obtain

$$PY_{-B,1}^{(1)} - Y_{-B,1}^{(-1)} > U(m_{B,1}, T, B) \le T \stackrel{\xi_1}{=} 2 \cdot \frac{1}{-B}$$

In words, with probability exceeding 1/2, no elimination will happen for the bin B. As a result, we obtain

$$\begin{split} E\left[R_{T}\left(\hat{\pi}\right)\right] & \geq & E & f^{\mathbb{B}}(X_{t}) - f^{\hat{\pi}_{t}(X_{t})}(X_{t}) & \mathbf{1}\{X_{t} \ \mathbb{B}_{I}\} \\ & \stackrel{t=t_{12}+1}{=}_{I} + 1 \\ & \mathbb{B}_{I} + \frac{t_{2}}{g} \cdot z^{-1} \, \mathbb{B}_{I} \frac{\dot{x}}{z} = & T^{19} \\ & \frac{\dot{x}}{z} = & T^{19} \\ \end{split}$$

where we have used the choice of z. So Theorem 4 holds with  $\kappa = \epsilon$ .

Case 2:  $g \le T^{3/19-\epsilon}$ . Set  $z = T^{3/19-\epsilon/8}$ . We have g < z and there exists H > 1 such that  $z = H \cdot g$ ; see Figure 4 for an illustration of the instance. Let B be the bin produced by  $\hat{\pi}$  such that  $C_1 \supseteq B$ . By the design of v, we have

$$\bar{f_B^{(1)}} - \bar{f_B^{(-1)}} \leq \frac{1}{H} (1/2 + D_{\phi} z^{-1}) + (1 - \frac{1}{H}) \frac{1}{2} - \frac{1}{2} = \frac{D_{\phi} z^{-1}}{H}.$$

Let  $\delta = \frac{D_{\phi}z^{-1}}{H}$ , we have  $\delta \leq 1/\frac{p}{m_{B,1}^{2}}$  due to our choice of z. Additionally, we have  $m_{B,1} \leq 2m_{B,1}^{2}$  under E. Therefore, we can invoke Lemma 11 to obtain

Thus, with probability exceeding 1/2, the suboptimal arm is not eliminated in B. Similar to the previous case, we obtain

$$\begin{split} E\left[R_{T}\left(\hat{\pi}\right)\right] &\geq E & \overset{*}{\underset{t_{2}}{\overset{t_{2}}{\overset{*}{\longrightarrow}}}} f^{2}\left(X_{t}\right) - f^{\hat{\pi}_{t}\left(X_{t}\right)}\left(X_{t}\right) \\ & \overset{*}{\overset{*}{\overset{*}{\longrightarrow}}} \frac{X^{t_{2}}}{\overset{*}{\longrightarrow}} & \# \\ &= E & f^{2}\left(X_{t}\right) - f^{\hat{\pi}_{t}\left(X_{t}\right)}\left(X_{t}\right) & \mathbf{1}\left\{X_{t} \ \mathbf{?} \ C_{1}\right\} \ t = t_{1} + 1 \\ &? & \frac{1}{2^{2}} = T^{\frac{9}{19} + \frac{6}{4}}. \end{split}$$

So Theorem 4 holds with  $\kappa = \epsilon/4$ .

1 such that  $z = H \cdot g$ ; see Figure 4 for an illustration of the instance. Let B be the bin produced by  $\hat{\pi}$  such that  $C_1 ext{ } extstyle{!} ext$ 

$$\bar{f}_{B}^{(1)} - \bar{f}_{B}^{(-1)} \leq \frac{1}{H}(1/2 + D_{\phi}z^{-1}) + (1 - \frac{1}{H})\frac{1}{2} - \frac{1}{2} = \frac{D_{\phi}z^{-1}}{H}.$$

Let  $\delta = \frac{D_{\phi} z^{-1}}{H}$ , we have  $\delta \leq 1/\frac{p_{B,1}}{B_{B,1}}$  due to our choice of z. Additionally, we have  $m_{B,1} \leq 2m_{B,1}^{\mathbb{Z}}$  under E. Therefore, we can invoke Lemma 11 to obtain

$$PY_{-}^{(1)} - Y_{-}^{(-1)} > U(m_{B,1}, T, B) \le T \stackrel{\xi_1}{=} 4 \cdot \frac{1}{2}$$

 $PY_{-}^{(1)}-Y_{-}^{(-1)}>U\left(m_{B,1},T,B\right)\leq T\xrightarrow{\frac{E_{1}}{B}}4\cdot\frac{1}{-}$  This means with probability at least 3/4, arm elimination does not occur in B after the first batch. Moreover, since  $\delta\leq 1/P$   $\overline{m_{B,2}^{\mathbb{B}}}$  by the choice of z, and  $m_{B,2}\leq 2m_{B,2}^{\mathbb{B}}$  under E, we can apply Lemma 11 again to get

$$PY_{-}^{(1)} - Y_{-}^{(-1)} > U(m_{B,2}, T, B) \le T \le 4 \cdot 1$$

In all, with probability at least 1/2, arm elimination does not occur in B after the second batch. Similar to before, we reach the conclusion that

$$E[R_{T}(\hat{\pi})] \geq E \int_{\|t^{-1}x^{T}\|}^{T} f^{2}(X_{t}) - f^{\hat{\pi}_{t}(X_{t})}(X_{t})$$

$$= E \left(f^{2}(X_{t}) - f^{\hat{\pi}_{t}(X_{t})}(X_{t})\right) 1\{X_{t} ? C_{1}\} t = t_{2} + 1$$

$$? \frac{T}{z^{2}} = T^{2^{\frac{1}{2}}}.$$

We see that Theorem 4 holds with  $\kappa = 1/38$ .

### D Proof of Theorem 5

When  $\beta$  = d = 1, we denote  $\gamma(\alpha^{\tiny{\square}})$  =  $(\alpha^{\tiny{\square}} + 1)/3$ . Fix some small constant  $\epsilon$  > 0 to be specified later. We deal with M = 2 and M = 3 separately. In both cases, we use the fact that the algorithm needs to provide its first batch size  $t_1$  prior to the game, and design  $\alpha^{\tiny{\square}}$  such that any choice of  $t_1$  would fail.

#### D.1 When M = 2

Under M = 2, the theoretical optimal rate is  $T^{(1-\gamma(\alpha^{\mathbb{Z}}))/(1-\gamma(\alpha^{\mathbb{Z}})^2)} = T^{3/(\alpha^{\mathbb{Z}}+4)}$ .

Case of  $t_1 = T^{3/5+\epsilon}$ . Take  $\alpha^{\mathbb{Z}} = 1$ . Since fixed grid and adaptive grid are the same when M = 2, by relation (9), we have

$$\sup_{f \supseteq F\left(\alpha^{\mathbb{S}}, \beta\right)} \ E\left[R_{T}\left(\pi\right)\right] \supseteq \max_{\stackrel{?}{?}} \ \frac{T}{t_{1}, \frac{T}{\frac{\alpha^{\mathbb{S}+1}}{3}}} \nearrow \geq t_{1} = T^{\frac{3}{5} + \kappa_{1}},$$

where  $\kappa_1 = \epsilon$ .

Case of  $t_1 = T^{3/4-\epsilon}$ . Take  $\alpha^{2} = o(1)$ . By relation (9), we have

$$\sup_{f \in F\left(\alpha^{\mathbb{Z}},\beta\right)} E\left[R_{T}\left(\pi\right)\right] ? \max_{\left[\frac{1}{2},\frac{1}{3}\right]} t_{1}, \quad \frac{T}{t_{1}^{\frac{\alpha^{\mathbb{Z}}+1}{3}}} \ge \frac{T}{t_{1}^{\frac{\alpha^{\mathbb{Z}}+1}{3}}} = T^{\frac{1-\left(\frac{3}{4}-\epsilon\right)}{3} = T^{\frac{3}{4}+\kappa_{1}},$$

where  $\kappa_1 = (\alpha^2 + 1)\epsilon/3 - \alpha^2/4 > 0$ .

#### D.2 When M = 3

Under M = 3, the theoretical optimal rate is  $T^{(1-\gamma(\alpha^{\boxtimes}))/(1-\gamma(\alpha^{\boxtimes})^3)} = T^{9/((\alpha^{\boxtimes})^2+5\alpha^{\boxtimes}+13)}$ .

Case of  $t_1 = T^{9/19+\epsilon}$ . Take  $\alpha^{\mathbb{Z}} = 1$ . During the first batch, the learner can do no better than pull an arm uniformly at random. We can use the instance  $(f_1(x) = 1, f_2(x) = 0)$  so that

$$\sup_{f \, \mathbb{B} \, F \, (\alpha^{\underline{\alpha}}, \beta)} \, E \left[ \, \mathsf{R}_{\, \mathsf{T}} \, \left( \pi \right) \right] \, \mathbb{P} \, \mathsf{t}_{1} = \, \mathsf{T}^{\, \frac{9}{19} + \, \mathsf{K}_{\, \mathsf{1}}} \, ,$$

where  $\kappa_1 = \epsilon$ .

Case of  $t_1=T^{9/13-\varepsilon}$ . Take  $\alpha^\mathbb{Z}=o(1)$ . Denote  $T_2=\int_{\omega\mathbb{Z},Q}^{1/(\gamma(\alpha^\mathbb{Z})+1)}t^{\gamma(\alpha^\mathbb{Z})/(\gamma(\alpha^\mathbb{Z})+1)}$ . Define the events  $E_2=\{T_2<\ t_2\}$  and  $E_3=\{t_2\le T_2\}$ . Recall  $Q_i(\cdot)=\bigcup_{\omega\mathbb{Z},Q}^{1/(2\beta+d)}\mathbb{Z}$  and  $Z_3=\mathbb{Z}(36T_22^2)^{1/(2\beta+d)}\mathbb{Z}$ . Since  $E_2$  can be determined by observations up to  $t_1$ , we have

$$|Q_{2}(E_{2}) - Q_{3}(E_{2})| = |Q_{2}^{t_{1}}(E_{2}) - Q_{3}^{t_{1}}(E_{2})| \le TV(Q_{2}^{t_{1}}, Q_{3}^{t_{1}}) \le \frac{1}{2} q \frac{q}{t_{1}z_{2}^{-(2\beta+d)}} \le \frac{1}{4},$$
 (21)

where step (i) applies Lemma 2, and step (ii) is due to the definition of z<sub>2</sub>. Consequently,

$$\begin{split} Q_2(E_2) + \ Q_3(E_3) &= \ Q_2(E_2) - \ Q_3(E_2) + \ Q_3(E_2) + \ Q_3(E_3) \\ &\geq \ -\frac{1}{4} + \ Q_3(E_2) + \ Q_3(E_3) = \ \frac{3}{4}, \end{split}$$

where the second step uses inequality (21) and the last equality is because  $E_2$  and  $E_3$  form a whole partition of the probability space. Then we would have at least one of  $Q_2(E_2) \ge 1/4$  or  $Q_3(E_3) \ge 1/4$ . If  $Q_2(E_2) \ge 1/4$ , by Lemma 4 we obtain

$$\sup_{f \ni F(\alpha^{\mathbb{Z}},\beta)} \ E[R_T(\pi)] \ge \sup_{f \ni C_{z_2}} \ R_{T_2}(\pi;f) \ ? \ T_2 Z_2^{-\beta(1+\alpha^{\mathbb{Z}})} \boxdot T_{13}^{\frac{9}{2}+\kappa_1},$$

for some  $\kappa_1 > 0$ . If  $Q_3(E_3) \ge 1/4$ , we similarly have

$$\sup_{f \, \mathbb{B} \, F \, (\alpha^{\mathbb{B}}, \beta)} \, \, E \, [\, R_{\,T} \, (\pi) \,] \geq \sup_{f \, \mathbb{B} \, C_{\,z_{_{_{3}}}}} \, R_{\,T_{_{_{3}}}}(\pi; f) \, \, ? \, \, T \, z_{_{_{3}}}^{-\beta \, (1+\alpha^{\mathbb{B}})} \, \, \mathbb{P} \, T_{\,13}^{\,\underline{9}+\,K_{\,1}} \, ,$$

for some  $\kappa_1 > 0$ .