Power-of-2-Arms for Bandit Learning With Switching Costs

Ming Shi Purdue University West Lafayette, IN, USA sming@purdue.edu Xiaojun Lin Purdue University West Lafayette, IN, USA linx@ecn.purdue.edu Lei Jiao University of Oregon Eugene, OR, USA jiao@cs.uoregon.edu

ABSTRACT

Motivated by edge computing with artificial intelligence, in this paper we study a bandit-learning problem with switching costs. Existing results in the literature either incur $\Theta(T^{\frac{2}{3}})$ regret with bandit feedback, or rely on free full-feedback in order to reduce the regret to $O(\sqrt{T})$. In contrast, we expand our study to incorporate two new factors. First, full feedback could incur a cost. Second, the player may choose 2 (or more) arms at a time, in which case she is free to use any one of the chosen arms to calculate loss, and switching costs are incurred only when she changes the set of chosen arms. For the setting where the player pulls only one arm at a time, our new regret lower-bound shows that, even when costly full-feedback is added, the $\Theta(T^{\frac{2}{3}})$ regret still cannot be improved. However, the dependence on the number of arms may be improved when the full-feedback cost is small. In contrast, for the setting where the player can choose 2 (or more) arms at a time, we provide a novel online learning algorithm that achieves a lower $O(\sqrt{T})$ regret. Further, our new algorithm does not need any full feedback at all. This sharp difference therefore reveals the surprising power of choosing 2 (or more) arms for this type of bandit-learning problems with switching costs. Both our new algorithm and regret analysis involve several new ideas, which may be of independent interest.

CCS CONCEPTS

• Theory of computation \rightarrow Online learning algorithms; • Networks \rightarrow Network algorithms; Network performance analysis.

KEYWORDS

Bandit learning, switching costs, regret analysis, edge computing with artificial intelligence

ACM Reference Format:

Ming Shi, Xiaojun Lin, and Lei Jiao. 2022. Power-of-2-Arms for Bandit Learning With Switching Costs. In *The Twenty-third International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing (MobiHoc '22), October 17–20, 2022, Seoul, Republic of Korea.* ACM, New York, NY, USA, 10 pages. https://doi.org/10.1145/3492866.3549720

1 INTRODUCTION

In this paper, we are interested in bandit learning with switching costs, which can be used to model many practical decision-making

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

MobiHoc '22, October 17–20, 2022, Seoul, Republic of Korea © 2022 Copyright held by the owner/author(s). ACM ISBN 978-1-4503-9165-8/22/10. https://doi.org/10.1145/3492866.3549720

problems that not only face significant uncertainty, but also incur costs for changing decisions. Consider edge computing with artificial intelligence (Edge AI) [6] as an example, where an edge server close to the end users downloads machine learning (ML) models from the cloud to process incoming inference requests. As the underlying ground-truth model of the data changes in uncertain ways (which is often referred to as concept drift [12]), the best ML model also changes in time. However, because of the limited capability of the edge server, it can often only accommodate a small number of ML models. Thus, the edge server needs to learn which subset of ML models should be used, based on the feedback (e.g., inference losses) observed. Further, downloading an ML model (which is not currently on the edge server) from the cloud incurs communication overhead, which can be modelled by a switching cost β_1 . Hence, the edge server has to carefully select the ML models to reduce the total inference losses and switching costs in the long run, which thus corresponds to a bandit-learning problem with switching costs. Other examples of such problems can be found in transportation networks [2], data-center networks [17], wireless communication [4] and cyber-physical systems [14], etc.

In the online learning literature, it is well-known that the existence of switching costs significantly changes the nature of the regret. Specifically, in an adversarial setting (which we will focus on in this paper), for bandit learning without switching costs, the Exp3 algorithm can attain $O(\sqrt{T})$ regret over a time-horizon T [3]. However, once the switching cost is added, the regret (for the setting where only one arm can be pulled at each time) increases substantially to $\Theta(T^{\frac{2}{3}})$ [1]. A matching lower bound in [9] suggests that such an increased regret is unavoidable. While this result may be somewhat discouraging, it leaves many important open questions, as we explained below. Note that since ML models in Edge AI corresponds to arms in bandit learning, we use the word "model" and "arm" interchangeably in the rest of the paper.

First, in practice, in addition to pulling one arm, there are often other ways to obtain feedback. For example, in Edge AI, the edge server could send the data to the cloud for analysis. In this case, the feedback from all ML models can be obtained, beyond the model already deployed on the edge server. This is somewhat analogous to the full-feedback setting studied in [13]. Reference [13] shows that, if the full feedback can be obtained with zero costs, the regret for bandit learning with switching costs will remain at $O(\sqrt{T})$, which would have been much lower than that of [1] where only bandit feedback is available. However, in practice, feedback from the cloud also incurs non-negative costs due to multiple reasons, e.g., communication costs, latency and privacy issues [6]. Thus, the regret for bandit learning with both switching costs and full-feedback costs remains an open problem.

Second, instead of holding only one ML model at each time, in Edge AI, the edge server can usually accommodate $M \geq 2$ ML

models at each time. In this setting, using any of these *M* models for inference does not incur a switching/downloading cost, and at each time the feedback from all these *M* models (currently on the edge server) can be observed. This setting is thus most similar to a bandit-learning problem with limited advice [16], where $M \ge 2$ arms can be chosen at each time. However, [16] only studied the case without switching costs, where the regret is $O(\sqrt{T})$ regardless of whether one (M = 1) or more $(M \ge 2)$ arms are chosen at each time. Our setting is also related to bandit-learning problems with semi-bandit feedback [7] and side information [2]. The studies for semi-bandit feedback [7] typically do not consider switching costs either. Although the side-information setting [2] has been studied with switching costs, it is somewhat different from ours because switching within the $M \ge 2$ arms also incur switching costs there. Partly due to this difference, the regret [2] remains at $\Theta(T^{\frac{2}{3}})$. In summary, it remains an open problem whether in our setting, choosing $M \ge 2$ arms can improve the regret.

In this paper, we provide new answers to the above-mentioned two open problems. First, we study the case when M=1, i.e., only one arm can be pulled at each time, and there is a switching cost β_1 to change the arm and a full-feedback cost β_2 to obtain feedback from all arms. As we discussed earlier, the latter action corresponds to the edge server sending data to the cloud for analysis. We provide a lower bound of the regret, which grows as $\Theta(T^{\frac{2}{3}})$. In other words, when only one arm can be pulled (M=1), adding costly full-feedback will not fundamentally change how regret depends on T. However, our lower bound does suggest that utilizing costly full-feedback may change the multiplication factor in front of $T^{\frac{2}{3}}$. In some settings, this factor can be reduced from $O(K^{\frac{1}{3}})$ to $O((\ln K)^{\frac{1}{3}})$, where K is the total number of arms. Moreover, we provide an algorithm that achieves a regret that matches the lower bound.

Second, we study the setting when $M \ge 2$, i.e., more than one arm can be chosen at each time and one of them is used to incur the loss, while there are still switching costs and full-feedback costs. Surprisingly, here we provide a new online learning algorithm, called Randomized Online Learning With Working Groups (ROW), that can achieve a regret of $O(\sqrt{T})$ without even using full feedback (see Theorem 4.1), which significantly improves the $\Theta(T^{\frac{2}{3}})$ regret for M = 1. In other words, having the flexibility to accommodate one additional model (i.e., M = 2) almost brings comparable benefit as having free full-feedback [13]. To the best of our knowledge, this sharp transition from M = 1 to $M \ge 2$ has never be reported in the literature for bandit learning with switching costs¹. This may be seen as somewhat analogous to the "power-of-2" routing in load balancing [15] (where sampling two queues can attain comparable reduction to delay as sampling all queues), which is why we refer to it as the "power-of-2-arms". As M increases, the regret of ROW further decreases. Using a trivial lower bound for bandit learning with free full-feedback [5, 13], we conclude that the dependence of the regret of ROW on T must be optimal.

To achieve the improved $O(\sqrt{T})$ regret, ROW employs several new ideas. In order to fully utilize the flexibility of choosing $M \ge 2$ arms and minimize switching costs, the first idea of ROW is to

fix a primary arm over $O(\sqrt{T})$ time-slots (which we refer to as an episode), and switch the secondary arms $\lceil \frac{K-1}{M-1} \rceil$ times during an episode, each time to a new subset of secondary arms that have not yet been chosen in this episode. In this way, ROW only makes a constant number of switches within each episode (and $\Theta(\sqrt{T})$ switches for all the time), but it can obtain not only the feedback of the primary arm for the entire episode, but also the feedback of every other arms for $\frac{1}{K-1}$ fraction of the episode. Intuitively, this way of obtaining feedback incurs much lower costs than using costly full-feedback to obtain the same amount of feedback (for any K and $\beta_2 > 0$ independent of T), which is also the reason that ROW does not use costly full-feedback at all.

However, just using the above idea alone is insufficient to produce the $O(\sqrt{T})$ regret. The reason is that the feedback obtained is highly correlated in time. This is because each subset of secondary arms is retained for the whole sub-episode (whose length is also $O(\sqrt{T})$). It is known that such correlation tends to increase the regret (see the counter-examples in Sec. 4.1). ROW utilizes a second crucial idea to overcome this difficulty. Our key observation is that, whenever such a sub-episode with highly-correlated feedback occurs, one of arms in the current working group (either the primary arm or a secondary arm) will likely be consistently better than other arms. Then, ROW will try to switch to the better arm more quickly within the sub-episode, and thus improve the regret. To accomplish this faster switching within a sub-episode, our proposed ROW algorithm will use a larger weight-decay parameter η_2 within each sub-episode, while using a smaller parameter η_1 across episodes. In Sec. 4.2.2, we give a sufficient condition on how much η_2 should be larger than η_1 to strike the right balance. We note that this idea of using two different weight-decay parameters is new and may be of independent interest.

Finally, since in each episode the primary arm will receive much more feedback than the secondary arms, this creates a bias in the overall quality of feedback at the end of each episode. This bias issue is resolved by using instead the loss differences between the primary and secondary arms (please see our Idea 3 in Sec. 4.1). Our proof for the $O(\sqrt{T})$ regret carefully captures the effect of the above ideas by utilizing several new techniques (please see Sec. 4.2 for details).

2 PROBLEM FORMULATION

In this section, we provide the problem formulation for our banditlearning problem with switching costs and full-feedback costs. Moreover, we present a motivating example based on edge computing with artificial intelligence (Edge AI), which has received extensive attention recently [6]. Finally, we introduce the performance metric.

2.1 Bandit Learning With Switching Costs and Full-Feedback Costs

A player interacts with the adversary/environment sequentially in time. Let $\mathcal{K} \triangleq \{1, 2, ..., K\}$ denote the set of all arms and let M be an integer, $1 \leq M < K$. In each time-slot t = 1, ..., T, first the player chooses M arms among all K arms. Let $\hat{\mathbb{k}}(t)$ denote the set of the M arms chosen at time t. The player uses one of the arms in $\hat{\mathbb{k}}(t)$ as the

¹Note that for bandit learning *without* switching costs, choosing $M \ge 2$ arms will improve the regret, but it cannot alter the dependence on T [16].

active arm, which is denoted by $\mathbf{k}(t)$. The loss of this arm, $l_{\mathbf{k}(t)}(t)$, will be used to calculate the loss and regret later. In addition, the losses $l_k(t)$ of all arms $k \in \hat{\mathbf{k}}(t)$ are observed by the player. The loss $l_k(t)$ can be any arbitrary value in [0,1]. In this paper, we study both the cases when M=1 and $1 \leq M < K$. When M=1, $\hat{\mathbf{k}}(t)$ only contains the active arm $\mathbf{k}(t)$ and only the loss of this arm is observed. In this case, we simply say that the player "pulls" the single arm $\mathbf{k}(t)$ at time t. On the other hand, when t0 arms in t1 arms in t2 arms also observed.

Next, for every arm that is newly added to the set $\hat{\mathbb{L}}(t)$, a switching cost $\beta_1 > 0$ will be incurred. Thus, the switching cost at time t is $\beta_1 \sum_{k \in \hat{\mathbb{L}}(t)} \mathbf{1}_{\{k \notin \hat{\mathbb{L}}(t-1)\}}$, where $\mathbf{1}_E$ is an indicator function (i.e., $\mathbf{1}_E = 1$ if the event E is true, and $\mathbf{1}_E = 0$ otherwise). As typically assumed in bandit-learning problems [2, 3, 9, 13, 19], we assume that $\hat{\mathbb{L}}(0) = \Phi$ is empty. In addition to the feedback from the M arms in $\hat{\mathbb{L}}(t)$, at each time t, the player can choose to obtain full feedback of time t for all the arms (including those not in $\hat{\mathbb{L}}(t)$) at a cost β_2 . Let z(t) = 1 if the player chooses to obtain the full feedback at time t, and z(t) = 0 otherwise. Therefore, the total cost is

$$\operatorname{Cost}(1:T) \triangleq \sum_{t=1}^{T} \left\{ l_{\mathbf{k}(t)}(t) + \beta_{1} \sum_{k \in \hat{\mathbb{k}}(t)} \mathbf{1}_{\left\{k \notin \hat{\mathbb{k}}(t-1)\right\}} + \beta_{2} z(t) \right\}. \tag{1}$$

2.2 An Example Motivated by Edge AI

We consider an Edge AI setting where an edge server collaborates with a remote cloud. The edge server runs machine learning (ML) models on an online stream of input data to predict their labels. (For example, in an E-commerce recommendation system, the input data at each time contains the customer data, item data and web shop transactions, etc. The input data will be used by the edge server to return the recommendations, i.e., the predicted labels of what the customer is interested in.) We assume that K ML models are already trained and available in the remote cloud. However, due to the limited capability of the edge server, only M models can be deployed at the edge server at each time. Since it is unknown which ML model works best, the edge server needs to use the feedback (e.g., the actual product picked by the customer) to learn which subset of ML models it should deploy. (In practice, both the underlying distribution of the input data and the mapping from data to labels change in time due to the so-called concept drift [12]. Therefore, the best model(s) also changes in time. As a result, this learning process may be performed again after a concept drift.)

This learning process can be modelled as the bandit-learning problem described above. Each arm corresponds to one of the K ML models. At each time t, the edge server chooses the subset $\hat{\mathbb{k}}(t)$ of M models, which correspond to the M arms chosen in bandit learning. This subset $\hat{\mathbb{k}}(t)$ may be the same as the subset $\hat{\mathbb{k}}(t-1)$ chosen at last time t-1, or it may differ, in which case a switching cost $\beta_1 \sum_{k \in \hat{\mathbb{k}}(t)} \mathbf{1}_{\{k \notin \hat{\mathbb{k}}(t-1)\}}$ for downloading the ML models that are not currently on the edge server will be incurred. Note that this switching cost is assumed to be proportional to the number of ML models (which are not currently on the edge server) downloaded at time t. Then, the input data $\vec{X}(t)$ is revealed. The edge server

will use the models in $\hat{\mathbb{K}}(t)$ to infer the label of $\vec{X}(t)$. Further, it will use the result $\vec{Y}_{k(t)}'(t)$ of one of the models $k(t) \in \hat{\mathbb{K}}(t)$, to return to the end user. This model k(t) then corresponds to the active arm in bandit learning. Next, the true label $\vec{Y}(t)$ of $\vec{X}(t)$ is revealed. The edge server can then calculate the inference loss $l_k(t)$ for each ML model $k \in \hat{\mathbb{K}}(t)$, based on the difference between the inferred label $\vec{Y}_k'(t)$ and the true label $\vec{Y}(t)$, e.g., using the squared loss (i.e., $l_k(t) = \|\vec{Y}(t) - \vec{Y}_k'(t)\|^2$) [10].

At the end of time t, the edge server may also choose to consult the cloud for the quality of all ML models. In that case, it sends the data $\vec{X}(t)$ to the cloud. After the cloud processes this data with all ML models $k \in \mathcal{K}$, the edge server can retrieve the inference-loss $l_k(t)$ of all the ML models. Clearly, it incurs additional computation/communication overhead to obtain such feedback from the cloud, which we model by the full-feedback cost β_2 .

2.3 Performance Metric

We use the regret [1, 3, 9, 13] as the performance metric. For an online learning algorithm π , its total cost $\mathrm{Cost}^\pi(1:T)$ is given by (1), which includes both switching costs and full-feedback costs. For the optimal static solution OPT, it knows the future losses in advance, and hence can choose only one arm/model throughout the time-horizon. The cost of OPT is then given by $\mathrm{Cost}^{\mathrm{OPT}}(1:T) = \min_{k \in \mathcal{K}} \sum_{t=1}^T l_k(t) + \beta_1$, where there is only one switching $\mathrm{cost}\,\beta_1$ at the beginning of the time-horizon, and there is no full-feedback cost. The regret of algorithm π is defined to be the worst-case difference between the expected total cost of algorithm π and the total cost of OPT, i.e.,

$$R^{\pi}(T) \triangleq \sup_{l_{1:K}(1:T)} \left\{ \mathbb{E}_{\pi} \left[\operatorname{Cost}^{\pi}(1:T) \right] - \operatorname{Cost}^{\operatorname{OPT}}(1:T) \right\}, \quad (2)$$

where the expectation is taken over the possible randomness of the algorithm π , and $l_{1:K}(1:T)$ denotes the losses $l_k(t)$ of all arms $k \in [1, K]$ for all time $t \in [1, T]$. Our goal is to design an online learning algorithm with a regret as low as possible.

3 THE CASE OF M = 1

In this section, we focus on the case when M=1, i.e., the player (e.g., edge-server) can pull only one arm (e.g., model) at each time. We are interested in studying whether adding full feedback with a cost β_2 can alter the regret of bandit learning with switching costs. Recall that in this case, the active arm $\mathbf{k}(t)$ is the only arm in $\hat{\mathbf{k}}(t)$. As we mentioned in the introduction, when full feedback is free, it has been shown in [13] that using full feedback will improve the regret from $\Theta(T^{\frac{2}{3}})$ to $O(\sqrt{T})$. However, since in our model the full feedback incurs a cost, it is no longer clear whether the regret can still be improved.

3.1 A Lower Bound on the Regret (M=1)

Our first main result shows that adding costly full-feedback will not change the dependence of the regret on T, but may change the multiplication factor as a function of K.

Theorem 3.1. Consider bandit learning with switching costs and full-feedback costs introduced in Sec. 2.1. When M = 1, the regret of

any online algorithm π must be lower-bounded as follows,

$$R^{\pi}(T) \ge \underline{R}^{\pi}(T) \triangleq \max \left\{ C_1 \beta_a^{\frac{1}{3}} \left(\log_2 K \right)^{\frac{1}{3}} T^{\frac{2}{3}}, C_2 \beta_b^{\frac{1}{3}} T^{\frac{2}{3}} \right\}, \quad (3)$$

where

$$\begin{split} \beta_a &= \min \left\{ \frac{3}{2}\beta_1, \beta_2 \right\}, \beta_b = \min \left\{ \frac{3}{4}K\beta_1, \beta_2 \right\}, \\ C_1 &= \Theta \left(\frac{1}{\log_2 T - \log_2 \log_2 K} \right), \text{ and } C_2 = \Theta \left(\frac{1}{\log_2 T} \right). \end{split}$$

We can see from Theorem 3.1 that, even when the costly full-feedback is added, as long as M=1, $\Theta(T^{\frac{2}{3}})$ is still the optimal regret for bandit learning with switching costs. This is in sharp contrast to the case of free full-feedback [13], where the regret can be improved to $O(\sqrt{T})$. While this result may be somewhat discouraging, the costly full-feedback does play some role in the multiplication factor in front of $T^{\frac{2}{3}}$, which depends on the relative magnitude of β_1 and β_2 . Intuitively, when the full-feedback cost β_2 is large, the online learning algorithm would rather switch to obtain feedback than using costly full-feedback. On the other hand, when β_2 is small, the online learning algorithm should avoid switching and obtain feedback from costly full-feedback. Thus, we expect that costly full-feedback will be more useful in the latter case than in the former case. The conclusion of Theorem 3.1 shows this difference precisely. Specifically, we can make the following observations.

(i) When $\beta_2 \geq \frac{3}{4}K\beta_1$, the lower bound $\underline{R}^{\pi}(T)$ in (3) is equal to

$$\max \left\{ C_1 \left(\frac{3}{2} \beta_1 \right)^{\frac{1}{3}} \left(\log_2 K \right)^{\frac{1}{3}} T^{\frac{2}{3}}, C_2 \left(\frac{3}{4} \beta_1 \right)^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}} \right\}. \tag{4}$$

As K increases, the second term in (4) quickly dominates. This means that, when the full-feedback cost β_2 is high, the regret of any online learning algorithm π will at least increase as $\beta_1^{\frac{1}{3}}K^{\frac{1}{3}}T^{\frac{2}{3}}$. Note that this expression is the same as the regret (for bandit learning with switching costs) when there is no full feedback at all [9]. This observation is consistent with our intuition that, when β_2 is large, the online algorithm cannot benefit from costly full-feedback.

(ii) When $\beta_2 < \frac{3}{4}K\beta_1$, the lower bound $\underline{R}^{\pi}(T)$ in (3) is equal to

$$\max \left\{ C_1 \beta_a^{\frac{1}{3}} \left(\log_2 K \right)^{\frac{1}{3}} T^{\frac{2}{3}}, C_2 \beta_2^{\frac{1}{3}} T^{\frac{2}{3}} \right\}. \tag{5}$$

As K increases, the first term in (5) quickly dominates. This means that, when the full-feedback cost β_2 is not high, the regret of any online algorithm π will at least increase as $\beta_a^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}}$. If in addition $\beta_2 \leq \frac{3}{2}\beta_1$, we have $\beta_a^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}} = \beta_2^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}}$, which is smaller than $\beta_1^{\frac{1}{3}} K^{\frac{1}{3}} T^{\frac{2}{3}}$. Compared with the earlier case (with large β_2), our regret expression here has the same dependence on T, but now increases more slowly as a function of the total number K of arms. This observation is also consistent with our intuition that, when β_2 is small, the online algorithm can benefit from costly full-feedback more.

Finally, we note that the division of the two cases depends on the value of $K\beta_1$ and β_2 . The intuition is that, with K switches, an online algorithm may also attain the feedback from all K arms. Thus, it makes sense to compare $K\beta_1$ with β_2 to determine which type of feedback is more effective.

Algorithm 1 The Multivariate Hidden Markov (MHM) adversary

Parameters: Choose ϵ and σ according to (9).

Initialization: Choose k^* uniformly from \mathcal{K} .

for t = 1 : T **do**

Step 1: Generate the value of G(t) according to (7).

Step 2: Generate the losses of each arm $k \in \mathcal{K}$ as follows,

$$l_k(t) = G(t) + \frac{1}{2} - \epsilon \cdot \mathbf{1}_{\{k=k^*\}} + \gamma_k(t), \tag{6}$$

where $\gamma_k(t) \sim \mathcal{N}(0, \sigma^2)$ are *i.i.d.* Gaussian random variables with zero-mean and σ^2 -variance.

end for

3.2 Lower Bound Analysis

To prove Theorem 3.1, we design two important adversaries. For both adversaries, we make use of Yao's principle [20] that the worst-case expected regret $R^{\pi}(T)$ of a randomized online algorithm π is lower-bounded by the expected regret of the best deterministic online algorithm against a randomized adversary. Thus, in the following we focus on designing randomized adversaries, and studying the regret of deterministic online algorithms. Recall that $\mathcal{K} = \{1, ..., K\}$.

3.2.1 Multivariate Hidden Markov (MHM) Adversary. In this section, we provide the first randomized adversary, called Multivariate Hidden Markov (MHM) adversary, which generalizes the idea in [9]. Please see Algorithm 1.

Specifically, Step 1 in Algorithm 1 is the same as that used by the adversary introduced in [9]. That is, for each time t, define the parent time of t as $\rho(t) \triangleq t - 2^{\delta(t)}$, where $\delta(t) \triangleq \max\{\delta \mid t \equiv 0 \pmod{2^{\delta}}\}$. The main reason that the parent time $\rho(t)$ is $2^{\delta(t)}$ time-slot ahead of time t is to guarantee that with high probability, the generated losses $l_k(t)$ are in [0,1]. Please see our technical report [18] for the concrete proof of this guarantee. Then, Step 1 of MHM generates a Gaussian process G(t) in the following way,

$$G(t) = G(\rho(t)) + \xi(t), \text{ for all time } t \in [1, T], \tag{7}$$

where G(0)=0, and $\xi(t)\sim\mathcal{N}(0,\sigma^2)$ are *i.i.d.* Gaussian random variables with zero-mean and σ^2 -variance. As in [9], this process G(t) creates a common uncertainty across all arms. In Step 2, the first three terms² in (6) are also the same as that used in [9]. However, (6) differs from the adversary of [9] in the fourth term. This additional term adds a Gaussian noise $\gamma_k(t)$ to the loss $l_k(t)$ of each arm at each time. This additional noise is critical because our online algorithm π can use costly full-feedback, which is not considered in [9]. Intuitively, without this noise $\gamma_k(t)$, by using one round of costly full-feedback, the online algorithm can know the losses of all arms in the same time-slot. Then, the online algorithm will immediately know which arm is the optimal one (i.e., the arm with a loss that is ϵ lower). In contrast, the additional noise in (6) eliminates the possibility for such a trivial solution.

As we explain below, this additional noise $\gamma_k(t)$ causes new difficulties in the proof of the lower bound. We follow the approach in [9] to derive the regret lower-bound of any deterministic online algorithm π against the MHM adversary. Specifically, let $\mathcal{P}_{k^*}(\cdot)$

²The first three terms in (6) guarantees that the expected values of the losses are $\frac{1}{2}$ and $\frac{1}{2} - \epsilon$ for the sub-optimal arms $k \neq k^*$ and the optimal arm k^* , respectively.

denote the probability measure under the setting where one optimal arm k^* incurs ϵ lower cost than other arms, as in (6). Let $\mathcal{P}_0(\cdot)$ denotes the probability measure when $\epsilon=0$, i.e., the arm k^* is statistically the same as other arms. In addition, let $l^{ob}(\cdot)$ denote the observed losses of the online learning algorithm. Then, the analysis in [9] focuses on estimating the Kullback-Leibler (KL) divergence $D_{\text{KL}}(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T))||\mathcal{P}_0(l^{\text{ob}}(1:T))$, which then leads to the lower bound on the regret. However, for our MHM adversary, the additional noise $\gamma_k(t)$ incurs a new difficulty. Recall that $\rho(t)$ is the parent (time) of t, and thus t is the child (time) of $\rho(t)$. Let $\bar{\rho}(t)$ denote the set of the predecessors of time t, i.e. its parent, parent's parent, etc. Similarly, let $\rho(t)$ denote the set of the descendants of time t. Note that without $\gamma_k(t)$, the observed loss $l^{\text{ob}}(t)$ would have been a Gaussian process G(t) plus a fixed constant $\frac{1}{2}$ or $\frac{1}{2} - \epsilon$. Thus, $l^{ob}(t)$ would have satisfied a form of Markov property [11, p. 235], i.e., conditioned on current observed losses, the conditional probability distribution of future losses at a descendant time in $\rho(t)$ is independent of past losses at any predecessor time in $\bar{\rho}(t)$. Then, the proof could use the chain rule of KL divergence [8, p. 23]. In contrast, with the additional noise $\gamma_k(t)$, the observed loss $l^{\text{ob}}(t)$ does not satisfy the Markov property any more. This is because, conditioned on the observed losses at time t, past observed losses still provide information for the statistics of the future losses. For example, by taking the average of the losses observed at all predecessors in $\bar{\rho}(t)$, we can average out $\gamma_k(t)$ across time, and thus estimate the mean value of the loss at a descendant time in $\rho(t)$ with a higher accuracy. Therefore, we cannot use the chain rule directly, and must find a new way to bound the KL divergence.

To overcome this new difficulty, we develop a result on the KL divergence of hidden Markov models [8, p. 69]. Specifically, notice that the hidden loss $l^{\text{hi}}(t) \triangleq l_{\mathbf{k}(t)}(t) - \gamma_{\mathbf{k}(t)}(t)$, i.e., the loss in (6) but with $\gamma_{\mathbf{k}(t)}(t)$ removed, satisfies the Markov property. Then, using the chain rule of probability, we can show that

$$\begin{split} D_{\text{KL}} \left(& \mathcal{P}_{k^*}(l^{\text{ob}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T)) \right) \\ & \leq D_{\text{KL}} \left(\mathcal{P}_{k^*}(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T)) \| \mathcal{P}_0(l^{\text{ob}}(1:T) | l^{\text{hi}}(1:T)) \right) \\ & + D_{\text{KL}} \left(\mathcal{P}_{k^*}(l^{\text{hi}}(1:T)) \| \mathcal{P}_0(l^{\text{hi}}(1:T)) \right). \end{split} \tag{8}$$

The first term on the right-hand-side of (8) can be easily calculated at each time, since conditioned on the hidden loss $l^{\rm hi}(t)$, the observed loss $l^{\rm ob}(t)$ is only due to *i.i.d.* Gaussian variables $\gamma_k(t)$. The second term on the right-hand-side of (8) can be calculated by using the chain rule of the KL divergence, since the hidden loss $l^{\rm hi}(t)$ satisfies the Markov property. We can then obtain Lemma 3.2 below for the regret lower-bound against the MHM adversary.

Lemma 3.2. Consider bandit learning with switching costs and full-feedback costs introduced in Sec. 2.1. When M = 1, by choosing

$$\epsilon = \sqrt[3]{\frac{4}{9 \ln 2}} \cdot \frac{1}{9 \log_2 T} \cdot \beta_b^{\frac{1}{3}} T^{-\frac{1}{3}} \text{ and } \sigma = \frac{1}{9 \log_2 T},$$
 (9)

the regret of any online learning algorithm π against the MHM adversary is lower-bounded as follows: for $T \ge \max\{\beta_b, 6K\}$,

$$R^{\pi}(T) \ge C_2 \beta_b^{\frac{1}{3}} T^{\frac{2}{3}},$$
 (10)

where
$$\beta_b = \min\left\{\frac{3}{4}K\beta_1, \beta_2\right\}$$
 and $C_2 = \Theta\left(\frac{1}{\log_2 T}\right)$.

Please see our technical report [18] for the complete proof of Lemma 3.2. From Lemma 3.2, we can see that the regret lower-bound produced by MHM corresponds to the second term in (3). Note that it correctly captures the dependence of the regret on T, but the dependence on K still needs to be refined.

To further refine the dependence on K, we provide a second adversary, called Dividing Set (DS) adversary. Please see our technical report [18] for details. Finally, we design an online learning algorithm, called Randomized Online Learning With Costly Full-Feedback (ROCF), that attains the following regret for large T,

$$R^{\text{ROCF}}(T) \le \begin{cases} 4\beta_1^{\frac{1}{3}} (K \ln K)^{\frac{1}{3}} T^{\frac{2}{3}}, & \text{if } \beta_2 \ge \frac{3K}{4} \beta_1, \\ \frac{7}{2} \beta_2^{\frac{1}{3}} (\ln K)^{\frac{1}{3}} T^{\frac{2}{3}} + O(1), & \text{if } \beta_2 < \frac{3K}{4} \beta_1, \end{cases}$$
(11)

which matches the lower bound in Theorem 3.1. ROCF essentially uses episodic versions of either Exp3 [3] (when β_2 is large) or the shrinking dartboard algorithm [13] (when β_2 is small). Due to page limits, we refer the readers to our technical report [18].

4 THE POWER-OF-2-ARMS ($M \ge 2$)

In this section, we proceed to the case when $M \geq 2$. In contrast to the previous section where we show that adding costly full-feedback does not change the $\Theta(T^{\frac{2}{3}})$ regret, here we provide a new algorithm that utilizes the flexibility of having 2 (or more) arms and successfully improves the regret to $O(\sqrt{T})$.

4.1 Randomized Online Learning With Working Groups (ROW)

We call our new algorithm Randomized Online Learning With Working Groups (ROW). Please see Algorithm 2. We start with describing the high-level skeleton of ROW. Recall that $\mathcal{K} = \{1, ..., K\}$.

Idea 1: Note that in order to obtain the $O(\sqrt{T})$ regret, we can switch or use costly full-feedback at most $O(\sqrt{T})$ number of times. The first idea of ROW is thus to design an effective way to rotate a working group (of M arms) through all K arms, so that plenty of feedback can be obtained for all the arms, while incurring $O(\sqrt{T})$ switching costs and zero full-feedback costs. Specifically, ROW divides the entire time-horizon into $U = \begin{bmatrix} T \\ \tau_1 \end{bmatrix}$ episodes, each with $\tau_1 = \Theta(\sqrt{T})$ time-slots. At the beginning of the first time-slot $t_u = (u-1)\tau_1 + 1$ of the u-th (u=1,...,U) episode, each arm $k \in \mathcal{K}$ is associated with a weight $w_k^{\text{ROW}}[u]$, which is initialized to be $w_k^{\text{ROW}}[1] = 1$ (we will describe how to update $w_k^{\text{ROW}}[u]$ from $w_k^{\text{ROW}}[u-1]$ shortly). Then, from all arms $k \in \mathcal{K}$, ROW chooses a primary arm $k_0^{\text{ROW}}[u]$ with probability (i.e., Step 1 in Algorithm 2)

$$p_k^{\text{ROW}}[u] = \frac{w_k^{\text{ROW}}[u]}{\sum_{k=1}^K w_k^{\text{ROW}}[u]}.$$
 (12)

This primary arm $k_0^{\rm ROW}[u]$ will be fixed for the entire episode u. In addition, ROW divides each episode into $V = \left\lceil \frac{K-1}{M-1} \right\rceil$ sub-episodes, each with $\tau_2 = \frac{\tau_1}{V}$ time-slots. In the rest of this paper, we refer to the v-th sub-episode in the u-th episode as sub-episode (u,v). At the beginning of the first time-slot $t_{u,v} = (u-1)\tau_1 + (v-1)\tau_2 + 1$

Algorithm 2 Randomized Online Learning With Working Groups (ROW)

for $t = t_{u,v} : t_{u,v} + \tau_2 - 1$ **do**

and (14), respectively.

Step 4: Use an arm $k \in \hat{k}^{\text{ROW}}[u, v]$ as the active arm according to the updated probability $\hat{p}_k^{\text{ROW}}(t)$.

Step 5: Update the weights $\hat{w}_k^{\text{ROW}}(t)$ and probabilities $\hat{p}_k^{\text{ROW}}(t)$ of all arms $k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]$ according to (15) and (14), respectively.

end for

end for

Step 6: At the end of the last time-slot of the *u*-th episode, update the weights $w_k^{\text{ROW}}[u+1]$ and probabilities $p_k^{\text{ROW}}[u+1]$ of all arms $k \in \mathcal{K}$ according to (17) and (12), respectively.

end for

of sub-episode (u,v), ROW uniformly chooses M-1 secondary arms from the arms that have not yet been chosen in the u-th episode³ (i.e., Step 2 in Algorithm 2). We let $\hat{\mathbb{k}}_{M-1}^{ROW}[u,v]$ denote the set of the M-1 secondary arms chosen in sub-episode (u,v). Let $\hat{\mathbb{k}}^{ROW}[u,v] = \{k_0^{ROW}[u]\} \cup \hat{\mathbb{k}}_{M-1}^{ROW}[u,v]$ denote the working group formed by the primary arm and secondary arms. The working group $\hat{\mathbb{k}}^{ROW}[u,v]$ will be fixed for the whole sub-episode (u,v).

Notice that by using this idea, ROW only switches at the boundaries of sub-episodes and never uses full feedback. Therefore, by tuning τ_2 to be $\Theta(\sqrt{T})$, the total switching cost is guaranteed to be $\Theta(\sqrt{T})$, and the total full-feedback cost is 0. More importantly, with this idea, we not only have the feedback for the primary arm for the entire episode, but also have the feedback for each secondary arm for $\frac{1}{V}$ fraction of each episode. Intuitively, this way of obtaining feedback incurs much lower costs than using costly full-feedback. For example, if we want to obtain the same amount of feedback

by using costly full-feedback alone, we would have to incur a full-feedback cost equal to $\Theta(\sqrt{T})$ in every episode! This is also the reason that ROW does not use full feedback at all.

We now describe the rest of the details of ROW. After choosing the working group $\hat{\mathbb{R}}^{ROW}[u,v]$ as we discussed above, within each sub-episode (u,v) we solve a bandit-learning problem with the set of arms restricted to the chosen working group. Note that this restricted version of the bandit-learning problem has no switching cost (since any arm $k \in \hat{\mathbb{R}}^{ROW}[u,v]$ can be used as the active arm without incurring switching costs), and also has full feedback (from all the arms $k \in \hat{\mathbb{R}}^{ROW}[u,v]$). Specifically, in the first time-slot $t_{u,v}$ of sub-episode (u,v), ROW initializes the weights of all the arms $k \in \mathcal{K}$ as follows (i.e., Step 3 in Algorithm 2),

$$\hat{w}_k^{\text{ROW}}(t_{u,v}) = w_k^{\text{ROW}}[u], \tag{13}$$

i.e., to be the values of the weights at the beginning of the entire episode u. Then, for each time $t = t_{u,v},...,t_{u,v} + \tau_2 - 1$, each arm $k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]$ is used as the active arm $k^{\text{ROW}}(t)$ with probability (i.e., Step 4 and Step 5 in Algorithm 2)

$$\hat{p}_k^{\text{ROW}}(t) = \frac{\hat{w}_k^{\text{ROW}}(t)}{\sum_{k \in \hat{\mathbb{R}}^{\text{ROW}}[u,v]} \hat{w}_k^{\text{ROW}}(t)}.$$
 (14)

After the losses $l_k(t)$ of all the arms $k \in \hat{k}^{ROW}[u, v]$ are obtained for time t, ROW updates their weights with a tunable parameter η_2 as follows (i.e., Step 5 in Algorithm 2),

$$\hat{w}_k^{\text{ROW}}(t+1) = \hat{w}_k^{\text{ROW}}(t) \cdot e^{-\eta_2 l_k(t)}, \tag{15}$$

and then proceeds to the next time-slot t+1. Note that the weights $\hat{w}_k^{\text{ROW}}(t)$ are reset by (13) in the first time-slot $t=t_{u,v}$ of each sub-episode (u,v).

Finally, in the last time-slot of episode u, ROW subtracts the loss of the primary arm from the corresponding loss of arm k in the subepisodes that k was observed. Then, the resulting value is divided by the conditional probability that k is chosen as a secondary arm (conditioned on k not being the primary arm), i.e., $\frac{M-1}{K-1}$. Precisely, we let $v_u(k) \triangleq \left\{v \mid v=1,...,V,k \in \hat{\mathbb{k}}^{\text{ROW}}[u,v]\right\}$ denote the subepisodes (u,v) when the arm k was chosen in the working group. Let $L_k[u,v_u(k)] \triangleq \sum_{v \in v_u(k)} \sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-2} l_k(t)$ denote the sum of the

losses of arm k in sub-episodes (u,v) (except the last time-slot $t=t_{u,v}+\tau_2-1$) for all $v\in v_u(k)$. Then, ROW computes the loss difference of each arm $k\in\mathcal{K}$ as follows,

$$\tilde{L}_{k}^{\text{ROW}}[u] = \frac{L_{k}[u, v_{u}(k)] - L_{k_{0}^{\text{ROW}}[u]}[u, v_{u}(k)]}{\frac{M-1}{K-1}}.$$
 (16)

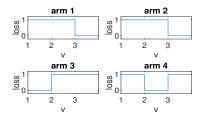
Note that for the primary arm $k_0^{\text{ROW}}[u]$, the loss difference is $\tilde{L}_{k_0^{\text{ROW}}[u]}^{\text{ROW}}[u] = 0$, which is also consistent with (16). Then, ROW updates the weights for all the arms $k \in \mathcal{K}$ with a tunable parameter η_1 as follows (i.e., Step 6 in Algorithm 2),

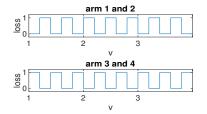
$$w_k^{\text{ROW}}[u+1] = w_k^{\text{ROW}}[u] \cdot e^{-\eta_1 \tilde{L}_k[u]}, \tag{17}$$

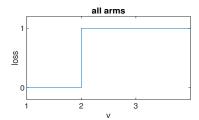
which becomes the initial weights for the next episode u + 1.

Readers familiar with bandit-learning algorithms may have already noticed two other crucial differences in ROW. First, a different weight-decay parameter η_2 is used to update weights in (15) within

 $^{^3}$ When K-1 is not divisible by M-1, the number of the remaining unchosen arms in the last (i.e., V-th) sub-episode may be less than M-1. In this case, after choosing all those unchosen arms, ROW uniformly chooses the secondary arms from the arms that have not yet been chosen for the V-th sub-episode.







(a) Trace in counter-example 1 (i.i.d. across arms k and sub-episodes [u, v]).

(b) Trace in counter-example 2 (repeats every 2 time-slots).

(c) Trace in counter-example 3 (repeats every episode *u*).

Figure 1: One realization of the counter-example traces in one episode.

the episode, compared with the parameter η_1 that is used in (17) across episodes. Second, when updating the weights across episodes in (17), we use the difference between the loss of an arm and that of the primary arm, instead of using the absolute loss of the arm directly. In the following, we explain why these two differences (i.e., our idea 2 and idea 3) are crucial for achieving the $O(\sqrt{T})$ regret.

Idea 2: Use different weight-decay parameters η_2 and η_1 . Recall that in every episode, ROW can obtain at least $\frac{1}{V}$ fraction of feedback from every arm. We would have hoped that this amount of feedback is sufficient for attaining a low $O(\sqrt{T})$ regret. Indeed, consider an alternate bandit-learning problem where the feedback of each arm is obtained independently with probability $\frac{1}{V}$ in every time-slot. It is not difficult to show that Exp3 [3] using this amount of feedback will attain the $O(\sqrt{T})$ regret.

However, compared with the above alternate problem, the difficulty we are facing here is that in ROW the feedback becomes highly correlated in time. Indeed, the secondary arms are fixed during the whole sub-episode. Thus, we either have all feedback of an arm, or have none for the whole sub-episode. Below, we construct two counter-examples to illustrate the difficulties in dealing with such correlation. For ease of exposition, we use $l(t_1:t_2) \triangleq [l(t), \text{ for all } t=t_1,t_1+1,...,t_2]$ to collect l(t) from $t=t_1$ to $t=t_2$.

Counter-example 1: Consider K = 4 arms and M = 2. For each arm k, in each sub-episode (u, v), $l_k(t_{u,v}: t_{u,v} + \tau_2 - 1) = 0$ with probability $\frac{1}{2}$, and $l_k(t_{u,v}:t_{u,v}+\tau_2-1)=1$ with probability $\frac{1}{2}$. The losses are independent across arms k and across sub-episodes [u, v]. Please see Fig. 1a for this loss trace in one episode. Using this counter-example, we show why existing bandit-learning method, Exp3 [1], could lead to a poor regret. Let us consider the optimal static loss. First, the expected total loss of each arm is trivially $\mathbb{E}[L] = \frac{T}{2}$. Second, let us estimate the variance of the total loss of each arm. Since the loss is a constant within a sub-episode, the higher correlation in time leads to a higher variance in the total loss of each arm. Specifically, for each arm, the variance of its total loss in a sub-episode⁴ is $\Theta(\tau_2^2)$. Thus the variance of its total loss across T time-slots is $\mathrm{Var}(L)=\frac{T}{\tau_2}\cdot\Theta(\tau_2^2)=\Theta(T^{\frac{3}{2}})$. Thus, one of the K arms may incur a total loss that is smaller than the average by $\Theta(\sqrt{\text{Var}(L)})$. As a result, the total loss of the optimal static decision OPT is $\mathbb{E}[L] - \Theta(\sqrt{\text{Var}(L)}) = \frac{T}{2} - \Theta(T^{\frac{3}{4}})$. (This estimate can also be obtained by applying the random walk analysis [19, p. 111].) Next, we consider the total loss of the episodic version of Exp3 [1]. Such

version of Exp3 picks an arm k_0 at the beginning of an episode, and use it as the active arm for the entire episode. Since the loss in each episode is independent, the total loss of such Exp3 will be the average loss of each arm in this counter-example, i.e., $\frac{T}{2}$. Therefore, the regret would be $\Theta(T^{\frac{3}{4}})$.

Counter-example 1 clearly illustrates why the higher correlation in time leads to a higher regret for the episodic version of Exp3. To overcome this difficulty, we make an important observation. In this setting with highly correlated losses, we observe that one arm (with losses 0) will be consistently better than the other arms (with losses 1) in each sub-episode. We may then beat the average loss by switching to the better arm within a sub-episode. Indeed, with M=2, the chance that one of the two arms incurs zero loss is $\frac{3}{4}$. Thus, if we can switch to the better arm (with losses 0) quickly within a sub-episode, we may attain a total loss approximately equals to $\frac{T}{4}$, which would have beaten the optimal static decision OPT. This counter-example thus suggests why it is important to use Exp3 [3] inside each sub-episode (in addition to across episodes).

However, it is still highly non-trivial to choose the parameter η of Exp3 within each sub-episode. One possible thought is that, we can think of each sub-episode as a bandit-learning problem with $\tau_2 = \Theta(\sqrt{T})$ time-slots. Then, if we view the better arm within the sub-episode as the static optimal arm, we would have to use $\eta = \Theta(T^{-\frac{1}{4}})$ in order to attain the minimal regret against the better arm. However, this choice of η would have been too large, as can be seen in the counter-example below.

Counter-example 2: Consider K=4 arms and M=2. For arms $k=1,2, l_k(t)=0$ for all odd time-slots t, and $l_k(t)=1$ for all even time-slots t. For arms $k=3,4, l_k(t)=1$ for all odd time-slots t, and $l_k(t)=0$ for all even time-slots t. Please see Fig. 1b for this loss trace in one episode. Using this counter-example, we can see why using Exp3 [3] with a parameter $\eta=\Theta(T^{-\frac{1}{4}})$ could lead to a poor regret. Let us consider the optimal static loss. Since the total loss of every arm is $\frac{T}{2}$, the optimal static loss is $\frac{T}{2}$. Next, we consider the total loss of Exp3. Notice that the probability of each arm is initialized to be the same, i.e., $\frac{1}{K}$, at time t=1. Then, at each time, suppose that all arms have been observed almost the same number of times. Thus, the probabilities of all arms would be about the same. However, whenever an arm with loss $l_{k_1}(t)=0$ and an arm with loss $l_{k_2}(t)=1$ are observed simultaneously, at the next time t+1 Exp3 will use the arm k_1 as the active arm with a

 $^{^4}$ In contrast, if the losses were *i.i.d.* in time, the variance should have been $\Theta(au_2)$.

probability higher by approximately $\Theta(\eta)$. According to counter-example 2, $l_{k_1}(t+1)=1$. Thus, Exp3 will suffer an additional loss $\Theta(\eta)$ approximately at each time. Hence, the total loss of Exp3 will be $\frac{T}{2}+\Theta(\eta T)=\frac{T}{2}+\Theta(T^{\frac{3}{4}})$. Therefore, the regret would be $\Theta(T^{\frac{3}{4}})$. Counter-example 2 clearly indicates that, in order to attain the

Counter-example 2 clearly indicates that, in order to attain the $O(\sqrt{T})$ regret, the parameter η_2 should be no larger than $O(T^{-\frac{1}{2}})$. However, since a sub-episode is of length much smaller than T, we conjecture that η_2 still needs to be larger than η_1 (the latter is used across episodes), so that ROW converges fast to the better arm inside the chosen working group. Lemma 4.3 in Sec. 4.2.2 will provide the exact condition on how η_2 and η_1 should be tuned to obtain the $O(\sqrt{T})$ regret.

Idea 3: Use the loss difference from the primary arm to update weights across episodes. We next describe why it is also crucial to use the loss difference in (16) instead of the absolute loss of each arm. Recall that at the end of each episode, we receive τ_1 feedback from the primary arm, but only $\tau_2 = \frac{\tau_1}{V}$ feedback from each secondary arm. Intuitively, this bias will also increase the variance of the total losses accumulated in the past, which again leads to a higher regret. The following counter-example illustrates this difficulty.

Counter-example 3: Consider K = 4 arms and M = 2. In the first sub-episode of each episode, the loss of each arm at each time is 0. For all subsequent sub-episodes of each episode, the loss of each arm at each time is 1. Please see Fig. 1c for this loss trace in one episode. In the literature, the standard way to deal with this bias in the amount of feedback is to divide the observed loss by the probability that the arm is observed [1, 3, 5]. For each arm, this probability is $p_k[u] + (1-p_k[u]) \frac{M-1}{K-1}$, where $p_k[u]$ is the probability that arm k is chosen as the primary arm, and $(1 - p_k[u]) \frac{M-1}{K-1}$ is the probability that arm k is chosen as the secondary arm in a subepisode. With this mechanism, the estimated losses will be $\tilde{L}_k[u] =$ $\frac{1}{p_k[u]+(1-p_k[u])\frac{M-1}{K-1}} \text{ when } k \text{ is the primary arm, } \tilde{L}_k[u]=0 \text{ when } k$ is a secondary arm that is chosen in the first (v=1) sub-episode, and $\tilde{L}_k[u] = \frac{\tau_2}{p_k[u] + (1-p_k[u])\frac{M-1}{K-1}}$ when k is a secondary arms that is chosen in the subsequent (v=2,3) sub-episodes. Suppose that $p_k[u] = \frac{1}{K}$ is the same across all arms. Then, the denominator is actually the same across all arms, but the numerator will still lead to a significant variance. Indeed. since the primary arm is chosen randomly with probability $p_k[u] = \frac{1}{K}$, it is not hard to verify that the total estimated loss of each arm over an episode will have a variance of $\Theta(\tau_2^2)$. In contrast, if full feedback was available, all arms would have a total loss equal to $2\tau_2$ in an episode, and the variance would have been zero. It is easy to show that, with this additional $\Theta(\tau_2^2)$ gap in the variance, the regret of Exp3 [1] is still $O(T^{\frac{2}{3}})$, which is much larger than $O(\sqrt{T})$.

Counter-example 3 thus suggests that, instead of dividing the loss by the probability of observing an arm, we need some new ways to deal with the above bias issue. Precisely, in (16), ROW updates the estimated loss by the difference of the loss of each secondary arm and that of the primary arm. In addition, the loss difference of the primary arm is simply 0. Returning to counter-example 3, the new estimated loss will then be $\tilde{L}_k[u] = 0$ for all the arms $k \in \mathcal{K}$. Thus, the additional variance $\Theta(\tau_2^2)$ of the estimated losses has been eliminated, which is also crucial for attaining the $O(\sqrt{T})$ regret.

4.2 Regret Analysis

In Theorem 4.1 below, we show the upper bound of the regret attained by ROW. For ease of exposition, we focus on the case when K-1 is divisible by M-1. (It is not difficult to extend to the case when K-1 is not divisible by M-1. Please see our technical report [18] for details.)

Theorem 4.1. Consider bandit learning with switching costs and full-feedback costs introduced in Sec. 2.1. When $M \geq 2$, the regret of ROW can be upper-bounded as follows, for $T \geq \frac{448(K-1)^2 \ln K}{\frac{5}{2} + 4\beta_1}$,

$$R^{ROW}(T) \le 8b_1 \frac{K-1}{M-1} \sqrt{\ln K} \sqrt{T} + b_2,$$
 (18)

where
$$b_1 = \sqrt{\frac{5}{2} + 2b_3\beta_1}$$
, $b_2 = b_3\beta_1 + 1$ and $b_3 = \min\{M, K - M\}$.

In Sec. 3 when M=1, the optimal regret is $\Theta(T^{\frac{2}{3}})$ for bandit learning with switching costs and full-feedback costs. In sharp contrast, now with $M\geq 2$, ROW achieves a significantly lower regret equals to $O(\sqrt{T})$. Moreover, ROW never uses full feedback. Further, as M increases, the regret of ROW can be further reduced. To the best of our knowledge, this is the first result in the literature to utilize the flexibility of choosing $M\geq 2$ arms to improved the regret to $O(\sqrt{T})$ for bandit learning with switching costs. Furthermore, using a trivial lower bound for bandit learning with free full-feedback [5, 13], we can conclude that the $O(\sqrt{T})$ regret cannot be further improved.

The rest of this section is devoted to the proof of Theorem 4.1. Due to the three new ideas in ROW, new analytical techniques are needed to capture the evolution of the weights, which are also of independent interest. In order to relate the loss of ROW to that of the optimal static loss, our analysis below is carried out in three steps, first for inside each sub-episode, second for the end of each episode, and third for across all episodes. In the following, we let $\mathcal{H}[u-1]$ denotes the σ -algebra generated by the observation of ROW from time t=1 to $t=(u-1)\tau_1$. Let $L_k[u,v]\triangleq\sum_{t=t,u}^{t_{u,v}+\tau_2-2}l_k(t)$.

4.2.1 Inside each sub-episode. We start by relating the expected loss of ROW inside each sub-episode to a log-sum-exp function $g_2[u,v]$ (see Lemma 4.2). Recall that in (13), the weights in the first time-slots of all sub-episodes are initialized to be the weights at the beginning of the episode. Thus, given a same working group, the probabilities $\hat{p}_k^{\rm ROW}(t_{u,v})$ are also the same at the beginning of all sub-episode v. We let $\hat{p}_k^{\rm ROW}[u]$ denote this common probability.

LEMMA 4.2. For each sub-episode (u, v), given the history $\mathcal{H}[u-1]$ and the chosen working group $\hat{k}[u, v]$, we have

$$\sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-1} \sum_{k \in \hat{\mathbb{k}}^{ROW}[u,v]} \hat{p}_k^{ROW}(t) l_k(t) \le g_2[u,v] + \frac{1}{2} \eta_2 \tau_2 + 1, \quad (19)$$

where

$$g_2[u,v] \triangleq -\frac{1}{\eta_2} \ln \left(\sum_{k \in \hat{\mathbb{k}}^{ROW}[u,v]} \hat{p}_k^{ROW}[u] e^{-\eta_2 L_k[u,v]} \right).$$
 (20)

On the left-hand-side of (19), the probability $\hat{p}_k^{\text{ROW}}(t)$ is the probability of using arm k as the active arm. Thus, the left-hand-side of (19) represents the conditional (conditioned on the working

group $\hat{k}^{\text{ROW}}[u,v]$ and history $\mathcal{H}[u-1]$) expected loss of ROW in sub-episode (u,v). Hence, (19) upper-bounds the conditional expected loss of ROW by a log-sum-exp function $g_2[u,v]$ and the term $\frac{1}{2}\eta_2\tau_2+1$. We make two important comments. First, the value of $g_2[u,v]$ is approximated dominated by the arm with the smallest loss $L_k[u,v]$ (as long as the corresponding probability $\hat{p}_k^{\text{ROW}}[u]$ is non-zero). (19) thus confirms that ROW is trying to switch to the "better" arm in the working group. Second, the gap $\frac{1}{2}\eta_2\tau_2$ is much smaller than the gap $\frac{1}{2}\eta\tau_2^2$ incurred by the episodic version of Exp3 [1]. Note that the above-mentioned two conclusions precisely capture our ideas 1 and 2, which together allow ROW to converge quickly to the better arm in the working group. Please see our technical report [18] for the complete proof of Lemma 4.2.

4.2.2 Relating the loss upper-bound at the end of a sub-episode to the weights across episodes. Lemma 4.2 provides an upper bound on the loss of ROW at the end of each sub-episode (u, v). Note that this upper bound depends on η_2 . On the other hand, at the end of each episode u, we calculate the weights according to (17). Notice that not only is $\tilde{L}_k^{\rm ROW}[u]$ in (17) different from $L_k[u, v]$ in (20), the parameter η_2 is also different from η_1 . Thus, we need a way to convert the loss upper-bound in Lemma 4.2 for each sub-episode to a form that depends on the weights calculated by (17). This is accomplished by Lemma 4.3 below. Further, this lemma gives a sufficient condition on how to tune the parameters η_2 and η_1 .

Specifically, notice that the loss difference $\tilde{L}_k^{\text{ROW}}[u]$ calculated in (16) is a difference from the loss of the primary arm $k_0^{\text{ROW}}[u]$. We let $g_2[u]$ denote the sum of $g_2[u,v]$ for all sub-episodes v, minus a term that corresponds to the loss of the primary arm, i.e.,

$$g_{2}[u] \triangleq \sum_{v=1}^{V} g_{2}[u, v] - \sum_{v=1}^{V} L_{k_{0}^{\text{ROW}}[u]}[u, v]$$

$$= -\frac{1}{\eta_{2}} \sum_{v=1}^{V} \ln \left(\sum_{k \in \hat{\mathbb{k}}^{\text{ROW}}[u, v]} \hat{p}_{k}^{\text{ROW}}[u] e^{-\eta_{2} \mathcal{L}_{k}^{\text{ROW}}[u, v]} \right), \quad (21)$$

where $\mathcal{L}_k^{\text{ROW}}[u, v] = L_k[u, v] - L_{k_\alpha^{\text{ROW}}[u]}[u, v].$

Lemma 4.3. If the parameters η_2 , τ_2 , η_1 and τ_1 satisfy that

$$\eta_2 \ge 16 \left(\frac{K-1}{M-1}\right)^2 \cdot \eta_1, \ \eta_2 \tau_2 \le \ln 2 \ and \ \eta_1 \tau_1 \le \ln 2,$$
(22)

we have

$$\mathbb{E}_{\hat{\mathbb{k}}^{ROW}[u,1:V]} \left[g_2[u] \middle| \mathcal{H}[u-1] \right]$$

$$\leq \mathbb{E}_{\hat{\mathbb{k}}^{ROW}[u,1:V]} \left[g_1[u] \middle| \mathcal{H}[u-1] \right],$$
(23)

where the expectation is taken with respect to the randomness in the working groups, and

$$g_1[u] \triangleq -\frac{1}{\eta_1} \ln \left(\sum_{k=1}^K p_k^{ROW}[u] e^{-\eta_1 \tilde{L}_k^{ROW}[u]} \right).$$
 (24)

The log-sum-exp function $g_2[u]$ on the left-hand-side of (23) is related to $g_2[u,v]$ through (21), which is then related to the loss of ROW in each sub-episode through (19). The log-sum-exp function $g_1[u]$ on the right-hand-side of (23) is related to the weights calculated at the end of the episode. Thus, Lemma 4.3 relates the loss

upper-bound at the end of each sub-episode to the weights across episodes, and (22) confirms our conjecture that η_2 should be larger than η_1 .

The proof of Lemma 4.3 first relates the function $g_1[u]$ and $g_2[u]$ to the variances of the working-group feedback and the loss differences, respectively, and then bounds these variances. Please see our technical report [18] for the complete proof of Lemma 4.3.

Up to now, by combining (19), (21) and (23) for all sub-episode v and episode u, we can relate the total loss of ROW to $g_1[u]$ as follows.

$$\sum_{u=1}^{U} \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{\mathbb{E}}^{\text{ROW}}[u,1:V]} \left[\sum_{v=1}^{V} \sum_{t=t_{u,v}}^{t_{u,v}+\tau_2-1} \sum_{k \in \hat{\mathbb{E}}^{\text{ROW}}[u,v]} \hat{p}_k^{\text{ROW}}(t) \right] \right. \\ \left. \cdot l_k(t) - \sum_{v=1}^{V} L_{k_0^{\text{ROW}}[u]}[u,v] \left| \mathcal{H}[u-1] \right| \right\}$$

$$\leq \sum_{u=1}^{U} \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{\mathbb{E}}^{\text{ROW}}[u,1:V]} \left[g_1[u] \middle| \mathcal{H}[u-1] \right] \right\} + \frac{1}{2} \eta_2 T + VU.$$

$$(25)$$

In the next section, we show how to relate the first term on the right-hand-side of (25) to the optimal static loss.

4.2.3 Relating the upper-bound of the total loss of ROW to the optimal static loss. Lemma 4.4 below relates the sum of $g_1[u]$ on the right-hand-side of (25) to the optimal static loss of OPT.

LEMMA 4.4. We have the following inequality,

$$\sum_{u=1}^{U} \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E}_{\hat{\mathbb{k}}^{ROW}[u,1:V]} \left[g_1[u] \middle| \mathcal{H}[u-1] \right] \right\} - Cost^{OPT}(1:T) \\
\leq \frac{\ln K}{\eta_1} - \sum_{u=1}^{U} \mathbb{E}_{\mathcal{H}[u-1]} \left\{ \mathbb{E} \left[\sum_{v=1}^{V} L_{k_0^{ROW}[u]}[u,v] \middle| \mathcal{H}[u-1] \right] \right\}. \quad (26)$$

In (26), the first term on the left-hand-side is the first term on the right-hand-side of (25). The first term on the right-hand-side of (26) can be obtained by following the Exp3 analysis [3]. The second term on the right-hand-side of (26) is because the loss of the primary arm is subtracted in $g_2[u]$ (see (21)). This term also appears on the left-hand-side of (25), which will eventually be cancelled. Please see our technical report [18] for the complete proof of Lemma 4.4.

4.2.4 The final regret. Since ROW only switches at the boundaries of the sub-episodes, the total switching cost of ROW can be upper-bounded by min $\{M, K - M\} \cdot \beta_1 \left\lceil \frac{T}{\tau_2} \right\rceil$. Next, since ROW never asks for full feedback, the total full-feedback cost of ROW is 0. Then, together with (25), (26), we can see that the regret of ROW is upper-bounded as follows,

$$R^{\text{ROW}}(T) \le \frac{\ln K}{\eta_1} + \frac{1}{2}\eta_2 T + \min\{M, K-M\} \cdot \beta_1 \left[\frac{T}{\tau_2}\right] + \left[\frac{T}{\tau_2}\right]. \tag{27}$$
Then, by choosing $(c_1 = \sqrt{\frac{\ln K}{\frac{5}{2} + \min\{M, K-M\} \cdot 2\beta_1}} \text{ and } c_2 = \frac{4(K-1)}{M-1})$

$$\begin{cases}
\eta_2 = \frac{c_1 c_2}{\sqrt{T}}, & \tau_2 = \left\lfloor \frac{\ln 2}{c_1 c_2} \sqrt{T} \right\rfloor, \\
\eta_1 = \frac{c_1}{c_2 \sqrt{T}}, & \tau_1 = \left\lceil \frac{K-1}{M-1} \right\rceil \left\lfloor \frac{\ln 2}{c_1 c_2} \sqrt{T} \right\rfloor,
\end{cases} (28)$$

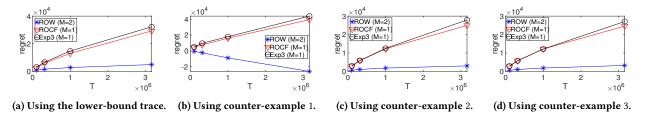


Figure 2: Compare the regrets of ROW, ROCF and the episodic version of Exp3.

we get the final regret of ROW in Theorem 4.1. Please see our technical report [18] for the complete proof of Theorem 4.1.

5 NUMERICAL RESULTS

In this section, we present numerical results comparing the regret of our algorithm ROW introduced in Sec. 4.1 (for $M \ge 2$) with that of the episodic version of Exp3 proposed in [1]. (We also show the regret of our ROCF algorithm from Sec. 3.2 for M=1. Please see our technical report [18] for more numerical results for ROW and ROCF.) According to [1], the theoretical regret of the episodic version of Exp3 is $\Theta(K^{\frac{1}{3}}T^{\frac{2}{3}})$.

In Fig. 2, we use both the lower-bound trace that we designed in Sec. 3.2 and the three counter-example traces that we designed in Sec. 4.1. We consider K = 4 arms, and M = 2 for ROW. (When Mincreases, the gap between the regret of ROW and that of Exp3 will further increase.) We let the switching cost and full-feedback cost be $\beta_1 = \beta_2 = 1$. We compare how the regret increases with the time length T. From Fig. 2, we can see that for all 4 traces, the regret of ROW (with M = 2) is much smaller than that of Exp3 (and ROCF). For example, when using counter-example 3 and $T = \sqrt{10} \times 10^6$, the regret of Exp3 is around 2.61×10⁴. In contrast, the regret of ROW is only about 3.22×10^3 , confirming the power of using 2 arms. For M = 1, the regret of ROCF is also smaller than that of Exp3. This is because the choice of β_1 and β_2 here satisfies $\beta_2 \leq \frac{3}{4}K\beta_1$. As we show in (3) and (11), this is the range where costly full-feedback is helpful for reducing the regret when M = 1. When β_2 increases to values larger than $\frac{3}{4}K\beta_1$, the gap between the regret of ROCF and that of Exp3 will diminish. See [18] for additional numerical results.

6 CONCLUSION

In this paper, we investigate bandit-learning problems with switching costs and full-feedback costs. Although we show that adding costly full-feedback will not alter the $\Theta(T^{\frac{2}{3}})$ regret (for M=1), we provide a novel online learning algorithm ROW that utilizes the flexibility of choosing $M\geq 2$ arms at each time to improve the regret to $O(\sqrt{T})$. Our result reveals that having 2 (or more) arms is surprisingly as powerful as having free full-feedback, for obtaining a low regret in bandit-learning problems with switching costs. Our algorithm ROW and regret analysis involve several new ideas, e.g., using different weight-decay parameters inside and across episodes. Our numerical results confirm that the regret of our algorithm ROW is much smaller than that of the episodic version of Exp3.

There are several interesting directions of future work. First, notice that we study the static regret. It would be interesting to extend our study to the dynamic regret, where the optimal arm

changes in time. Second, ROW assumes the knowledge of the time length T. It would be useful to extend ROW to the setting where T is not known in advance.

ACKNOWLEDGMENTS

This work has been partially supported by NSF grants CNS-2113893 and CNS-2047719.

REFERENCES

- Raman Arora, Ofer Dekel, and Ambuj Tewari. 2012. Online bandit learning against an adaptive adversary: from regret to policy regret. In Proceedings of 29th International Conference on Machine Learning. 1747–1754.
- [2] Raman Arora, Teodor Vanislavov Marinov, and Mehryar Mohri. 2019. Bandits with feedback graphs and switching costs. In Advances in Neural Information Processing Systems. 10397–10407.
- [3] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. 2002. The nonstochastic multiarmed bandit problem. SIAM journal on computing 32, 1 (2002), 48–77.
- [4] Dirk Bergemann and Juuso Välimäki. 2006. Bandit problems. Cowles Foundation discussion paper (2006).
- [5] Avrim Blum and Yishay Monsour. 2007. Learning, regret minimization, and equilibria. Algorithmic Game Theory (2007).
- [6] Jiasi Chen and Xukan Ran. 2019. Deep learning with edge computing: A review. Proc. IEEE 107, 8 (2019), 1655–1674.
- [7] Richard Combes, Mohammad Sadegh Talebi Mazraeh Shahi, and Alexandre Proutiere. 2015. Combinatorial bandits revisited. In Advances in Neural Information Processing Systems. 2116–2124.
- [8] Thomas M Cover. 1999. Elements of information theory. John Wiley & Sons.
- [9] Ofer Dekel, Jian Ding, Tomer Koren, and Yuval Peres. 2014. Bandits with switching costs: T^{2/3} regret. In Proceedings of 46th annual ACM symposium on Theory of computing. 459–467.
- [10] Pedro Domingos. 2000. A unified bias-variance decomposition. In Proceedings of 17th International Conference on Machine Learning. 231–238.
- [11] Rick Durrett. 2019. Probability: theory and examples. Cambridge university press.
- [12] João Gama, Indré Žliobaité, Albert Bifet, Mykola Pechenizkiy, and Abdelhamid Bouchachia. 2014. A survey on concept drift adaptation. ACM computing surveys 46. 4 (2014). 1–37
- [13] Sascha Geulen, Berthold Vöcking, and Melanie Winkler. 2010. Regret minimization for online buffering problems using the weighted majority algorithm. In Conference on Learning Theory. Citeseer, 132–143.
- [14] Sudipto Guha and Kamesh Munagala. 2009. Multi-armed bandits with metric switching costs. In *International Colloquium on Automata, Languages, and Programming*. Springer, 496–507.
- [15] Michael Mitzenmacher. 2001. The power of two choices in randomized load balancing. IEEE Transactions on Parallel and Distributed Systems 12, 10 (2001), 1094–1104.
- [16] Yevgeny Seldin, Peter Bartlett, Koby Crammer, and Yasin Abbasi-Yadkori. 2014. Prediction with limited advice and multiarmed bandits with paid observations. In Proceedings of 31st International Conference on Machine Learning. 280–287.
- [17] Shai Shalev-Shwartz. 2011. Online learning and online convex optimization. Foundations and trends in Machine Learning 4, 2 (2011), 107–194.
- [18] Ming Shi, Xiaojun Lin, and Lei Jiao. 2022. Available at https://engineering.purdue. edu/%7elinx/papers.html. Power-of-2-arms for bandit learning with switching costs. Technical Report. Purdue University.
- [19] Aleksandrs Slivkins. 2019. Introduction to multi-armed bandits. Foundations and Trends® in Machine Learning 12, 1-2 (2019), 1-286.
- [20] Andrew Chi-Chin Yao. 1977. Probabilistic computations: Toward a unified measure of complexity. In 18th Annual Symposium on Foundations of Computer Science. IEEE Computer Society, 222–227.