## **Multi-Armed Bandits in Metric Spaces**\*

Robert Kleinberg<sup>†</sup>

Aleksandrs Slivkins<sup>‡</sup>

Eli Upfal<sup>§</sup>

November 2007 Revised: April 2008, September 2008

#### Abstract

In a multi-armed bandit problem, an online algorithm chooses from a set of strategies in a sequence of n trials so as to maximize the total payoff of the chosen strategies. While the performance of bandit algorithms with a small finite strategy set is quite well understood, bandit problems with large strategy sets are still a topic of very active investigation, motivated by practical applications such as online auctions and web advertisement. The goal of such research is to identify broad and natural classes of strategy sets and payoff functions which enable the design of efficient solutions.

In this work we study a very general setting for the multi-armed bandit problem in which the strategies form a metric space, and the payoff function satisfies a Lipschitz condition with respect to the metric. We refer to this problem as the *Lipschitz MAB problem*. We present a solution for the multiarmed problem in this setting. That is, for every metric space (L, X) we define an isometry invariant MaxMinCOV(X) which bounds from below the performance of Lipschitz MAB algorithms for X, and we present an algorithm which comes arbitrarily close to meeting this bound. Furthermore, our technique gives even better results for benign payoff functions.

## **1** Introduction

In a multi-armed bandit problem, an online algorithm must choose from a set of strategies in a sequence of n trials so as to maximize the total payoff of the chosen strategies. These problems are the principal theoretical tool for modeling the exploration/exploitation tradeoffs inherent in sequential decision-making under uncertainty. Studied intensively for the last three decades [7, 8, 13], bandit problems are having an increasingly visible impact on computer science because of their diverse applications including online auctions, adaptive routing, and the theory of learning in games. The performance of a multi-armed bandit algorithm is often evaluated in terms of its *regret*, defined as the gap between the expected payoff of the algorithm and that of an optimal strategy. While the performance of bandit algorithms with a small finite strategy set is quite well understood, bandit problems with exponentially or infinitely large strategy sets are still a topic of very active investigation [1, 3, 4, 5, 6, 9, 10, 11, 12, 14, 15, 16, 19].

Absent any assumptions about the strategies and their payoffs, bandit problems with large strategy sets allow for no non-trivial solutions — any multi-armed bandit algorithm performs as badly, on some inputs,

<sup>\*</sup>The conference version [18] of this paper has appeared in ACM STOC 2008. This is the full version.

<sup>&</sup>lt;sup>†</sup>Computer Science Department, Cornell University, Ithaca, NY 14853. Email: rdk at cs.cornell.edu. Supported in part by NSF awards CCF-0643934 and CCF-0729102.

<sup>&</sup>lt;sup>‡</sup>Microsoft Research, Mountain View, CA 94043. Email: slivkins at microsoft.com. Parts of this work were done while the author was a postdoctoral research associate at Brown University.

<sup>&</sup>lt;sup>§</sup>Computer Science Department, Brown University, Providence, RI 02912. Email: eli at cs.brown.edu. Supported in part by NSF awards CCR-0121154 and DMI-0600384, and ONR Award N000140610607.

as random guessing. But in most applications it is natural to assume a structured class of payoff functions, which often enables the design of efficient learning algorithms [16]. In this paper, we consider a broad and natural class of problems in which the structure is induced by a metric on the space of strategies. While bandit problems have been studied in a few specific metric spaces (such as a one-dimensional interval) [1, 4, 9, 15, 22], the case of general metric spaces has not been treated before, despite being an extremely natural setting for bandit problems. As a motivating example, consider the problem faced by a website choosing from a database of thousands of banner ads to display to users, with the aim of maximizing the click-through rate of the ads displayed by matching ads to users' characterizations and the web content that they are currently watching. Independently experimenting with each advertisement is infeasible, or at least highly inefficient, since the number of ads is too large. Instead, the advertisements are usually organized into a taxonomy based on metadata (such as the category of product being advertised) which allows a similarity measure to be defined. The website can then attempt to optimize its learning algorithm by generalizing from experiments with one ad to make inferences about the performance of similar ads [22, 23]. Abstractly, we have a bandit problem of the following form: there is a strategy set X, with an unknown payoff function  $\mu : X \to [0,1]$  satisfying a set of predefined constraints of the form  $|\mu(u) - \mu(v)| \leq \delta(u,v)$  for some  $u, v \in X$  and  $\delta(u, v) > 0$ . In each period the algorithm chooses a point  $x \in X$  and observes an independent random sample from a payoff distribution whose expectation is  $\mu(x)$ .

A moment's thought reveals that this abstract problem can be regarded as a bandit problem in a metric space. Specifically, if L(u, v) is defined to be the infimum, over all finite sequences  $u = x_0, x_1, \ldots, x_k = v$  in X, of the quantity  $\sum_i \delta(x_i, x_{i+1})$ , then L is a metric<sup>1</sup> and the constraints  $|\mu(u) - \mu(v)| < \delta(u, v)$  may be summarized by stating that  $\mu$  is a Lipschitz function (of Lipschitz constant 1) on the metric space (L, X). We refer to this problem as the *Lipschitz MAB problem* on (L, X), and we refer to the ordered triple  $(L, X, \mu)$  as an *instance* of the Lipschitz MAB problem.<sup>2</sup>

**Prior work.** While our work is the first to treat the Lipschitz MAB problem in general metric spaces, special cases of the problem are implicit in prior work on the continuum-armed bandit problem [1, 4, 9, 15] — which corresponds to the space [0, 1] under the metric  $L_d(x, y) = |x - y|^{1/d}$ ,  $d \ge 1$  — and the experimental work on "bandits for taxonomies" [22], which corresponds to the case in which (L, X) is a tree metric. Before describing our results in greater detail, it is helpful to put them in context by recounting the nearly optimal bounds for the one-dimensional continuum-armed bandit problem, a problem first formulated by R. Agrawal in 1995 [1] and recently solved (up to logarithmic factors) by various authors [4, 9, 15]. In the following theorem and throughout this paper, the *regret* of a multi-armed bandit algorithm  $\mathcal{A}$  running on an instance  $(L, X, \mu)$  is defined to be the function  $R_{\mathcal{A}}(t)$  which measures the difference between its expected payoff at time t and the quantity  $t \sup_{x \in X} \mu(x)$ . The latter quantity is the expected payoff of always playing a strategy  $x \in \operatorname{argmax} \mu(x)$  if such strategy exists.

**Theorem 1.1** ([4, 9, 15]). For any  $d \ge 1$ , consider the Lipschitz MAB problem on  $(L_d, [0, 1])$ . There is an algorithm  $\mathcal{A}$  whose regret on any instance  $\mu$  satisfies  $R_{\mathcal{A}}(t) = \tilde{O}(t^{\gamma})$  for every t, where  $\gamma = \frac{d+1}{d+2}$ . No such algorithm exists for any  $\gamma < \frac{d+1}{d+2}$ .

In fact, if the time horizon t is known in advance, the upper bound in the theorem can be achieved by an extremely naïve algorithm which simply uses an optimal k-armed bandit algorithm (such as the UCB1 algorithm [2]) to choose strategies from the set  $S = \{0, \frac{1}{k}, \frac{2}{k}, \dots, 1\}$ , for a suitable choice of the parameter k. While the regret bound in Theorem 1.1 is essentially optimal for the Lipschitz MAB problem in  $(L_d, [0, 1])$ , it is strikingly odd that it is achieved by such a simple algorithm. In particular, the algorithm

<sup>&</sup>lt;sup>1</sup>More precisely, it is a pseudometric because some pairs of distinct points  $x, y \in X$  may satisfy L(x, y) = 0.

 $<sup>^{2}</sup>$  When the metric space (L, X) is understood from context, we may also refer to  $\mu$  as an instance.

approximates the strategy set by a fixed mesh S and does not refine this mesh as it gains information about the location of the optimal strategy. Moreover, the metric contains seemingly useful proximity information, but the algorithm ignores this information after choosing its initial mesh. Is this really the best algorithm?

A closer examination of the lower bound proof raises further reasons for suspicion: it is based on a contrived, highly singular payoff function  $\mu$  that alternates between being constant on some distance scales and being very steep on other (much smaller) distance scales, to create a multi-scale "needle in haystack" phenomenon which nearly obliterates the usefulness of the proximity information contained in the metric  $L_d$ . Can we expect algorithms to do better when the payoff function is more benign? For the Lipschitz MAB problem on  $(L_1, [0, 1])$ , the question was answered affirmatively in [9, 4] for some classes of instances, with algorithms that are tuned to the specific classes.

**Our results and techniques.** In this paper we consider the Lipschitz MAB problem on arbitrary metric spaces. We are concerned with the following two main questions motivated by the discussion above:

- (i) What is the best possible bound on regret for a given metric space?
- (ii) Can one take advantage of benign payoff functions?

In this paper we give a complete solution to (i), by describing for every metric space X a family of algorithms which come arbitrarily close to achieving the best possible regret bound for X. We also give a satisfactory answer to (ii); our solution is arbitrarily close to optimal in terms of the zooming dimension defined below. In fact, our algorithm for (i) is an extension of the algorithmic technique used to solve (ii).

Our main technical contribution is a new algorithm, the *zooming algorithm*, that combines the upper confidence bound technique used in earlier bandit algorithms such as UCB1 with a novel *adaptive refinement* step that uses past history to zoom in on regions near the apparent maxima of  $\mu$  and to explore a denser mesh of strategies in these regions. This algorithm is a key ingredient in our design of an optimal bandit algorithm for every metric space (L, X). Moreover, we show that the zooming algorithm can perform significantly better on benign problem instances. That is, for every instance  $(L, X, \mu)$  we define a parameter called the *zooming dimension*, and use it to bound the algorithm's performance in a way that is often significantly stronger than the corresponding per-metric bound. Note that the zooming algorithm is *self-tuning*, i.e. it achieves this bound without requiring prior knowledge of the zooming dimension.

To state our theorem on the per-metric optimal solution for (i), we need to sketch a few definitions which arise naturally as one tries to extend the lower bound from [15] to general metric spaces. Let us say that a subset Y in a metric space X has covering dimension d if it can be covered by  $O(\delta^{-d})$  sets of diameter  $\delta$ for all  $\delta > 0$ . A point  $x \in X$  has local covering dimension d if it has an open neighborhood of covering dimension d. The space X has max-min-covering dimension d = MaxMinCOV(X) if it has no subspace whose local covering dimension is uniformly bounded below by a number greater than d.

**Theorem 1.2.** Consider the Lipschitz MAB problem on a compact metric space (L, X). Let d = MaxMinCOV(X). If  $\gamma > \frac{d+1}{d+2}$  then there exists a bandit algorithm  $\mathcal{A}$  such that for every problem instance  $\mathcal{I}$  it satisfies  $R_{\mathcal{A}}(t) = O_{\mathcal{I}}(t^{\gamma})$  for all t. No such algorithm exists if d > 0 and  $\gamma < \frac{d+1}{d+2}$ .

In general MaxMinCOV(X) is bounded above by the covering dimension of X. For metric spaces which are highly homogeneous (in the sense that any two  $\epsilon$ -balls are isometric to one another) the two dimensions are equal, and the upper bound in the theorem can be achieved using a generalization of the naïve algorithm described earlier. The difficulty in Theorem 1.2 lies in dealing with inhomogeneities in the metric space.<sup>3</sup>

<sup>&</sup>lt;sup>3</sup>To appreciate this issue, it is very instructive to consider a concrete example of a metric space (L, X) where MaxMinCOV(X) is strictly less than the covering dimension, and for this specific example design a bandit algorithm whose regret bounds are better than those suggested by the covering dimension. This is further discussed in Section 3.

It is important to treat the problem at this level of generality, because some of the most natural applications of the Lipschitz MAB problem, e.g. the web advertising problem described earlier, are based on highly inhomogeneous metric spaces. (That is, in web taxonomies, it is unreasonable to expect different categories at the same level of a topic hierarchy to have the roughly the same number of descendants.)

The algorithm in Theorem 1.2 combines the zooming algorithm described earlier with a delicate transfinite construction over closed subsets consisting of "fat points" whose local covering dimension exceeds a given threshold d. For the lower bound, we craft a new dimensionality notion, the max-min-covering dimension introduced above, which captures the inhomogeneity of a metric space, and we connect this notion with the transfinite construction that underlies the algorithm.

For "benign" input instances we provide a better performance guarantee for the zooming algorithm. The lower bounds in Theorems 1.1 and 1.2 are based on contrived, highly singular, "needle in haystack" instances in which the set of near-optimal strategies is astronomically larger than the set of precisely optimal strategies. Accordingly, we quantify the tractability of a problem instance in terms of the number of near-optimal strategies. We define the *zooming dimension* of an instance  $(L, X, \mu)$  as the smallest d such that the following covering property holds: for every  $\delta > 0$  we require only  $O(\delta^{-d})$  sets of diameter  $\delta/8$  to cover the set of strategies whose payoff falls short of the maximum by an amount between  $\delta$  and  $2\delta$ .

# **Theorem 1.3.** If d is the zooming dimension of a Lipschitz MAB instance then at any time t the zooming algorithm suffers regret $\tilde{O}(t^{\gamma})$ , $\gamma = \frac{d+1}{d+2}$ . Moreover, this is the best possible exponent $\gamma$ as a function of d.

The zooming dimension can be significantly smaller than the max-min-covering dimension. Let us illustrate this point with two examples (where for simplicity the max-min-covering dimension is equal to the covering dimension). For the first example, consider a metric space consisting of a high-dimensional part and a low-dimensional part. For concreteness, consider a rooted tree T with two top-level branches T' and T'' which are complete infinite k-ary trees, k = 2, 10. Assign edge weights in T that are exponentially decreasing with distance to the root, and let L be the resulting shortest-path metric on the leaf set X.<sup>4</sup> If there is a unique optimal strategy that lies in the low-dimensional part T' then the zooming dimension is bounded above by the covering dimension of T', whereas the "global" covering dimension is that of T''. In the second example, let (L, X) be a homogeneous high-dimensional metric, e.g. the Euclidean metric on the unit k-cube, and the payoff function is  $\mu(x) = 1 - L(x, S)$  for some subset S. Then the zooming dimension is equal to the covering dimension of S, e.g. it is 0 if S is a finite point set.

**Discussion.** In stating the theorems above, we have been imprecise about specifying the model of computation. In particular, we have ignored the thorny issue of how to provide an algorithm with an input containing a metric space which may have an infinite number of points. The simplest way to interpret our theorems is to ignore implementation details and interpret "algorithm" to mean an abstract decision rule, i.e. a (possibly randomized) function mapping a history of past observations  $(x_i, r_i) \in X \times [0, 1]$  to a strategy  $x \in X$  which is played in the current period. All of our theorems are valid under this interpretation, but they can also be made into precise algorithmic results provided that the algorithm is given appropriate oracle access to the metric space. In most cases, our algorithms require only a *covering oracle* which takes a finite collection of open balls and either declares that they cover X or outputs an uncovered point. We refer to this setting as the standard Lipschitz MAB problem. For example, the zooming algorithm uses only a covering oracle for (L, X), and requires only one oracle query per round (with at most t balls in round t). However, the per-metric optimal algorithm in Theorem 1.2 uses more complicated oracles, and we defer the definition of these oracles to Section 3.

While our definitions and results so far have been tailored for the Lipschitz MAB problem on infinite metrics, some of them can be extended to the finite case as well. In particular, for the zooming algorithm

<sup>&</sup>lt;sup>4</sup>Here a *leaf* is defined as an infinite path away from the root.

we obtain sharp results (that are meaningful for both finite and infinite metrics) using a more precise, *non-asymptotic* version of the zooming dimension. Extending the notions in Theorem 1.2 to the finite case is an open question.

**Extensions.** We provide a number of extensions in which we elaborate on our analysis of the zooming algorithm. First, we provide sharper bounds for several examples in which the reward from playing each strategy u is  $\mu(u)$  plus an independent *noise* of a known and "benign" shape. Second, we upgrade the zooming algorithm so that it satisfies the guarantee in Theorem 1.3 and enjoys a better guarantee if the maximal reward is exactly 1. Third, we apply this result to a version where  $\mu(\cdot) = 1 - L(\cdot, S)$  for some *target set* S which is not revealed to the algorithm. Fourth, we relax some assumptions in the analysis of the zooming algorithm, and use this generalization to analyze the version in which  $\mu(\cdot) = 1 - f(L(\cdot, S))$  for some known function f. Finally, we extend our analysis from reward distributions supported on [0, 1] to those with unbounded support and finite absolute third moment.

**Follow-up work.** For metric spaces whose max-min-covering dimension is exactly 0, this paper provides an upper bound  $R(T) = O_{\mathcal{I}}(T^{\gamma})$  for any  $\gamma > \frac{1}{2}$ , but no matching lower bound. Characterizing the optimal regret for such metric spaces remained an open question. Following the publication of the conference version, this question has been settled in [17], revealing the following dichotomy: for every metric space, the optimal regret of a Lipschitz MAB algorithm is either bounded above by any  $f \in \omega(\log t)$ , or bounded below by any  $g \in o(\sqrt{T})$ , depending on whether the completion of the metric space is compact and countable.

#### 1.1 Preliminaries

Given a metric space, B(x, r) denotes an open ball of radius r around point x. Throughout the paper, he constants in the  $O(\cdot)$  notation are absolute unless specified otherwise.

**Definition 1.4.** In the *Lipschitz MAB problem* on (L, X), there is a strategy set X, a metric space (L, X) of diameter  $\leq 1$ , and a payoff function  $\mu : X \to [0, 1]$  such that the following *Lipschitz condition* holds:

$$|\mu(x) - \mu(y)| \le L(x, y) \quad \text{for all } x, y \in X.$$
(1)

Call L is the *similarity function*. The metric space (L, X) is revealed to an algorithm, whereas the payoff function  $\mu$  is not. In each round the algorithm chooses a strategy  $x \in X$  and observes an independent random sample from a payoff distribution  $\mathcal{D}(x)$  with support  $S \subset [0, 1]$  and expectation  $\mu(x)$ .

The *regret* of a bandit algorithm  $\mathcal{A}$  running on a given problem instance is  $R_{\mathcal{A}}(t) = W_{\mathcal{A}}(t) - t\mu^*$ , where  $W_{\mathcal{A}}(t)$  is the expected payoff of  $\mathcal{A}$  at time t and  $\mu^* = \sup_{x \in X} \mu(x)$  is the *maximal expected reward*.

The *C*-zooming dimension of the problem instance  $(L, X, \mu)$  is the smallest *d* such that for every  $r \in (0, 1]$  the set  $X_r = \{x \in X : \frac{r}{2} < \mu^* - \mu(x) \le r\}$  can be covered by  $Cr^{-d}$  sets of diameter at most r/8.

**Definition 1.5.** Fix a metric space on set X. Let N(r) be the smallest number of sets of diameter r required to cover X. The *covering dimension* of X is

$$\operatorname{COV}(X) = \inf \{ d : \exists c \,\forall r > 0 \quad N(r) \le cr^{-d} \}.$$

The *c*-covering dimension of X is defined as the infimum of all d such that  $N(r) \leq cr^{-d}$  for all r > 0.

**Outline of the paper.** In Section 2 we prove Theorem 1.3. In Section 3 we discuss the per-metric optimality and prove Theorem 1.2. Section 4 covers the extensions.

## 2 Adaptive exploration: the zooming algorithm

In this section we introduce the *zooming algorithm* which uses adaptive exploration to take advantage of the "benign" input instances, and prove the main guarantee (Theorem 1.3).

Consider the standard Lipschitz MAB problem on (L, X). The zooming algorithm proceeds in phases i = 1, 2, 3, ... of  $2^i$  rounds each. Let us consider a single phase  $i_{ph}$  of the algorithm. For each strategy  $v \in X$  and time t, let  $n_t(v)$  be the number of times this strategy has been played in this phase before time t, and let  $\mu_t(v)$  be the corresponding average reward. Define  $\mu_t(v) = 0$  if  $n_t(v) = 0$ . Note that at time t both quantities are known to the algorithm. Define the *confidence radius* of v at time t as

$$r_t(v) := \sqrt{8 \, i_{\rm ph} \, / \, (2 + n_t(v))}. \tag{2}$$

Let  $\mu(v)$  be the expected reward of strategy v. Note that  $E[\mu_t(v)] = \mu(v)$ . Using Chernoff Bounds, we can bound  $|\mu_t(v) - \mu(v)|$  in terms of the confidence radius:

**Definition 2.1.** A phase is called *clean* if for each strategy  $v \in X$  that has been played at least once during this phase and each time t we have  $|\mu_t(v) - \mu(v)| \le r_t(v)$ .

**Claim 2.2.** *Phase*  $i_{ph}$  *is clean with probability at least*  $1 - 4^{-i_{ph}}$ .

Throughout the execution of the algorithm, a finite number of strategies are designated *active*. Our algorithm only plays active strategies, among which it chooses a strategy v with the maximal *index* 

$$I_t(v) = \mu_t(v) + 2r_t(v).$$
 (3)

Say that strategy v covers strategy u at time t if  $u \in B(v, r_t(v))$ . Say that a strategy u is covered at time t if at this time it is covered by some active strategy v. Note that the covering oracle (as defined in Section 1) can return a strategy which is not covered if such strategy exists, or else inform the algorithm that all strategies are covered. Now we are ready to state the algorithm:

**Algorithm 2.3** (Zooming Algorithm). Each phase *i* runs for  $2^i$  rounds. In the beginning of the phase no strategies are active. In each round do the following:

- 1. If some strategy is not covered, make it active.
- 2. Play an active strategy with the maximal index (3); break ties arbitrarily.

We formulate the main result of this section as follows:

**Theorem 2.4.** Consider the standard Lipschitz MAB problem. Let  $\mathcal{A}$  be Algorithm 2.3. Then  $\forall C > 0$ 

$$R_{\mathcal{A}}(t) \le O(C\log t)^{1/(2+d)} \times t^{1-1/(2+d)} \text{ for all } t,$$
(4)

where d is the C-zooming dimension of the problem instance.

*Remark* The zooming algorithm is *not* parameterized by the C in (4), yet satisfies (4) for all C > 0. For sharper guarantees, C can be tuned to the specific problem instance and specific time t.

Let us prove Theorem 2.4. Note that after step 1 in Algorithm 2.3 all strategies are covered. (Indeed, if some strategy is activated in step 1 then it covers the entire metric.) Let  $\mu^* = \sup_{u \in X} \mu(u)$  be the maximal expected reward; note that we do not assume that the supremum is achieved by some strategy. Let  $\Delta(v) = \mu^* - \mu(v)$ . Let us focus on a given phase  $i_{ph}$  of the algorithm.

**Lemma 2.5.** If phase  $i_{ph}$  is clean then we have  $\Delta(v) \leq 4 r_t(v)$  for any time t and any strategy v. It follows that  $n_t(v) \leq O(i_{ph}) \Delta^{-2}(v)$ .

*Proof.* Suppose strategy v is played at time t. First we claim that  $I_t(v) \ge \mu^*$ . Indeed, fix  $\epsilon > 0$ . By definition of  $\mu^*$  there exists a strategy  $v^*$  such that  $\Delta(v^*) < \epsilon$ . Let  $v_t$  be an active strategy that covers  $v^*$ . By the algorithm specification  $I_t(v) \ge I_t(v_t)$ . Since v is clean at time t, by definition of index we have  $I_t(v_t) \ge \mu(v_t) + r_t(v_t)$ . By the Lipschitz property we have  $\mu(v_t) \ge \mu(v^*) - L(v_t, v^*)$ . Since  $v_t$  covers  $v^*$ , we have  $L(v_t, v^*) \le r_t(v_t)$  Putting all these inequalities together, we have  $I_t(v) \ge \mu(v^*) \ge \mu^* - \epsilon$ . Since this inequality holds for an arbitrary  $\epsilon > 0$ , we in fact have  $I_t(v) \ge \mu^*$ . Claim proved.

Furthermore, note that by the definitions of "clean phase" and "index" we have  $\mu^* \leq I_t(v) \leq \mu(v) + 3r_t(v)$  and therefore  $\Delta(v) \leq 3r_t(v)$ .

Now suppose strategy v is not played at time t. If it has never been played before time t in this phase, then  $r_t(v) > 1$  and thus the lemma is trivial. Else, let s be the last time strategy v has been played before time t. Then by definition of the confidence radius  $r_t(v) = r_{s+1}(v) \ge \sqrt{2/3} r_s(v) \ge \frac{1}{4} \Delta(v)$ .

**Corollary 2.6.** In a clean phase, for any active strategies u, v we have  $L(u, v) > \frac{1}{4} \min(\Delta(u), \Delta(v))$ .

*Proof.* Assume u has been activated before v. Let s be the time when v has been activated. Then by the algorithm specification we have  $L(u, v) > r_s(u)$ . By Lemma 2.5  $r_s(u) \ge \frac{1}{4}\Delta(u)$ .

Let d be the the C-zooming dimension. For a given time t in the current phase, let S(t) be the set of all strategies that are active at time t, and let

$$A(i,t) = \{ v \in S(t) : 2^{i} \le \Delta^{-1}(v) < 2^{i+1} \}.$$

We claim that  $|A(i,t)| \leq C 2^{id}$ . Indeed, set A(i,t) can be covered by  $C 2^{id}$  sets of diameter at most  $2^{-i}/8$ ; by Corollary 2.6 each of these sets contains at most one strategy from A(i,t).

**Claim 2.7.** In a clean phase  $i_{ph}$ , for each time t we have

$$\sum_{v \in S(t)} \Delta(v) \, n_t(v) \le O(C \, i_{\text{ph}})^{1-\gamma} \, t^{\gamma}, \tag{5}$$

where  $\gamma = \frac{d+1}{d+2}$  and d is the C-zooming dimension.

*Proof.* Fix the time horizon t. For a subset  $S \subset X$  of strategies, let  $R_S = \sum_{v \in S} \Delta(v) n_t(v)$ . Let us choose  $\rho \in (0, 1)$  such that

$$\rho t = (\frac{1}{\rho})^{d+1} (C \, i_{\rm ph}) = t^{\gamma} \, (C \, i_{\rm ph})^{1-\gamma}.$$

Define B as the set of all strategies  $v \in S(t)$  such that  $\Delta(v) \leq \rho$ . Recall that by Lemma 2.5 for each  $v \in A(i,t)$  we have  $n_t(v) \leq O(i_{ph}) \Delta^{-2}(v)$ . Then

$$\begin{aligned} R_{A(i,t)} &\leq O(i_{\mathrm{ph}}) \sum_{v \in A(i,t)} \Delta^{-1}(v) \\ &\leq O(2^{i} i_{\mathrm{ph}}) |A(i,t)| \\ &\leq O(C i_{\mathrm{ph}}) 2^{i(d+1)} \\ \sum_{v \in S(t)} \Delta(v) n_{t}(v) &\leq R_{B} + \sum_{i < \log(1/\rho)} R_{A(i,t)} \\ &\leq \rho t + O(C i_{\mathrm{ph}}) (\frac{1}{\rho})^{d+1} \\ &\leq O\left(t^{\gamma} (C i_{\mathrm{ph}})^{1-\gamma}\right). \end{aligned}$$

The left-hand side of (5) is essentially the contribution of the current phase to the overall regret. It remains to sum these contributions over all past phases.

**Proof of Theorem 2.4:** Let  $i_{ph}$  be the current phase, let t be the time spend in this phase, and let T be the total time since the beginning of phase 1. Let  $R_{ph}(i_{ph}, t)$  be the left-hand side of (5). Combining Claim 2.2 and Claim 2.7, we have

$$\begin{split} E[R_{\mathrm{ph}}(i_{\mathrm{ph}},t)] &< O(C\,i_{\mathrm{ph}})^{1-\gamma}\,t^{\gamma},\\ R_{\mathcal{A}}(T) &= E\left[R_{\mathrm{ph}}(i_{\mathrm{ph}},t) + \sum_{i=1}^{i_{\mathrm{ph}}-1} R_{\mathrm{ph}}(i,2^{i})\right]\\ &< O(C\,\log T)^{1-\gamma}\,T^{\gamma}. \end{split}$$

## **3** Attaining the optimal per-metric performance

In this section we ask, "What is the best possible algorithm for the Lipschitz MAB problem on a given metric space?" We consider the *per-metric performance*, which we define as the worst-case performance of a given algorithm over all possible problem instances on a given metric. As everywhere else in this paper, we focus on minimizing the exponent  $\gamma$  such that  $R_A(t) \leq t^{\gamma}$  for all sufficiently large t. Motivated by the shape of the guarantees in Theorem 1.1, let us define the *regret dimension* of an algorithm as follows.

**Definition 3.1.** Consider the Lipschitz MAB problem on a given metric space. For algorithm A and problem instance I let

$$\mathsf{DIM}_{\mathcal{I}}(\mathcal{A}) = \inf_{d \ge 0} \{ \exists t_0 \ \forall t \ge t_0 \ R_{\mathcal{A}}(t) \le t^{1-1/(d+2)} \}$$

The regret dimension of  $\mathcal{A}$  is  $DIM(\mathcal{A}) = \sup_{\mathcal{I}} DIM_{\mathcal{I}}(\mathcal{A})$ , where the supremum is taken over all problem instances  $\mathcal{I}$  on the given metric space.

Then Theorem 1.1 states that for the Lipschitz MAB problem on  $(L_d, [0, 1])$ , the regret dimension of the "naïve algorithm" is at most d. In fact, it is easy to extend the "naïve algorithm" to arbitrary metric spaces. Such algorithm is parameterized by the covering dimension d of the metric space. It divides time into phases of exponentially increasing length, chooses a  $\delta$ -net during each phase,<sup>5</sup> and runs a K-armed bandit algorithm such as UCB1 on the elements of the  $\delta$ -net. The parameter  $\delta$  is tuned optimally given d and the phase length T; the optimal value turns out to be  $\delta = T^{-1/(d+2)}$ . Using the technique from [15] it is easy to prove that the regret dimension of this algorithm is at most d.

**Lemma 3.2.** Consider the Lipschitz MAB problem on a metric space (L, X) of covering dimension d. Let  $\mathcal{A}$  be the naïve algorithm that uses UCB1 in each phase. Then  $DIM(\mathcal{A}) \leq d$ .

*Proof.* Let  $\mathcal{A}$  be the naïve algorithm. For concreteness, assume each phase *i* lasts  $2^i$  rounds. By definition of the covering dimension, it suffices to assume that *d* is a *c*-covering dimension, for some constant c > 0. By definition of the regret dimension, it suffices to prove that  $R_{\mathcal{A}}(t) \leq \tilde{O}(t^{\gamma})$  for all *t*, where  $\gamma = \frac{d+1}{d+2}$ . In order to prove *that*, it suffices to show that for each phase *i* we have  $R_{(\mathcal{A},i)}(2^i) \leq \tilde{O}(2^{i\gamma})$ , where  $R_{(\mathcal{A},i)}(t)$  is the expected regret accumulated in the first *t* rounds of phase *i*.

Let us focus on some phase *i*. In this phase the algorithm chooses a  $\delta$ -net, call it *S*. We claim that  $|S| \leq c \delta^{-d}$ . Indeed, for any  $\delta' < \delta$  the metric space can be covered by  $c \delta^{-d}$  sets of diameter at most  $\delta'$ , each of which can contain only one point from *S*. Claim proved. The algorithm proceeds to run UCB1 on

<sup>&</sup>lt;sup>5</sup>It is easy to see that the cardinality of this  $\delta$ -net is  $K = O(\delta^{-d})$ .

the elements of S. By [2] the expected regret of UCB1 on K arms in t rounds is at most  $O(\sqrt{K t \log t})$ . Since the maximal  $\mu$  on S is at most  $\delta$  off of the maximal  $\mu$  on X, we have

$$R_{(\mathcal{A},i)}(t) \le O(\sqrt{|S|t\log t}) + \delta t \le \tilde{O}(\sqrt{\delta^{-d}t} + \delta t).$$

Plugging in  $t = 2^i$  and  $\delta = t^{-1/(d+2)}$ , we obtain  $R_{(\mathcal{A},i)}(2^i) \leq \tilde{O}(2^{i\gamma})$  as claimed.

Thus we ask: *is it possible to achieve a better regret dimension*, perhaps using a more sophisticated algorithm? We show that this is indeed the case. Moreover, we provide an algorithm such that for any given metric space its regret dimension is arbitrarily close to optimal.

The rest of this section is organized as follows. In Section 3.1 we develop a lower bound on regret dimension. In Section 3.2 we will show that for some metric spaces, there exist algorithms whose regret dimension is smaller than the covering dimension. We develop these ideas further in Section 3.3 and provide an algorithm whose regret dimension is arbitrarily close to optimal.

#### 3.1 Lower bound on regret dimension

Let us develop a lower bound on regret dimension of any algorithm on a given metric space. This bound is equal to the covering dimension for highly homogeneous metric spaces (such as those in which all balls of a given radius are isometric to each other), but in general it can be much smaller.

It is known [3] that a worst-case instance of the K-armed bandit problem consists of K - 1 strategies with identical payoff distributions, and one which is slightly better. We refer to this as a "needle-in-haystack" instance. The known constructions of lower bounds for Lipschitz MAB problems rely on creating a *multiscale* needle-in-haystack instance in which there are K disjoint open sets, and K - 1 of them consist of strategies with identical payoff distributions, but in the remaining open set there are strategies whose payoff is slightly better. Moreover, this special open set contains  $K' \gg K$  disjoint subsets, only one of which contains strategies superior to the others, and so on down through infinitely many levels of recursion. To ensure that this construction can be continued indefinitely, one needs to assume a covering property which ensures that *each* of the open sets arising in the construction has sufficiently many disjoint subsets to continue to the next level of recursion.

**Definition 3.3.** For a metric space (L, X), we say that d is the *min-covering dimension* of X, d = MinCOV(X), if d is the infimum of COV(U) over all non-empty open subsets  $U \subseteq X$ . The *max-min-covering dimension* of X is defined by

$$extsf{MaxMinCOV}(X) = \sup \{ extsf{MinCOV}(Y) \, : \, Y \subseteq X \}.$$

The infimum over open  $U \subseteq X$  in the definition of min-covering dimension ensures that every open set which may arise in the needle-in-haystack construction described above will contain  $\Omega(\delta^{\varepsilon-d})$  disjoint  $\delta$ balls for some sufficiently small  $\delta, \varepsilon$ . Constructing lower bounds for Lipschitz MAB algorithms in a metric space X only requires that X should have *subsets* with large min-covering dimension, which explains the supremum over subsets in the definition of max-min-covering dimension.

We will use the following simple packing lemma.<sup>6</sup>

**Lemma 3.4.** If Y is a metric space of covering dimension d, then for any b < d and  $r_0 > 0$ , there exists  $r \in (0, r_0)$  such that Y contains a collection of at least  $r^{-b}$  disjoint open balls of radius r.

*Proof.* Let  $r < r_0$  be a positive number such that every covering of Y requires more than  $r^{-b}$  balls of radius 2r. Such an r exists, because the covering dimension of Y is strictly greater than b. Now let  $\mathcal{P} =$ 

<sup>&</sup>lt;sup>6</sup>This is a folklore result; we provide the proof for convenience.

 $\{B_1, B_2, \ldots, B_M\}$  be any maximal collection of disjoint *r*-balls. For every  $y \in Y$  there must exist some ball  $B_i$   $(1 \le i \le M)$  whose center is within distance 2r of y, as otherwise B(y, r) would be disjoint from every element of  $\mathcal{P}$  contradicting the maximality of that collection. If we enlarge each ball  $B_i$  to a ball  $B_i^+$  of radius 2r, then every  $y \in Y$  is contained in one of the balls  $\{B_i^+ | 1 \le i \le M\}$ , i.e. they form a covering of Y. Hence  $M \ge r^{-b}$  as desired.

**Theorem 3.5.** If X is a metric space and d is the max-min-covering dimension of X then  $DIM(A) \ge d$  for every bandit algorithm A.

*Proof.* Without loss of generality let us assume that d > 0. Given  $\gamma < \frac{d+1}{d+2}$ , let a < b < c < d be such that  $\gamma < \frac{a+1}{a+2}$ . Let Y be a subset of X such that  $MinCOV(Y) \ge c$ . Using Lemma 3.4 we recursively construct an infinite sequence of sets  $\mathcal{P}_0, \mathcal{P}_1, \ldots$  each consisting of finitely many disjoint open balls in X, centered at points of Y. Let  $\mathcal{P}_0 = \{X\}$  consist of a single ball that contains all of X. If i > 0, for every ball  $B \in \mathcal{P}_{i-1}$ , let r denote the radius of B and choose a number  $r_i(B) \in (0, r/4)$  such that B contains  $n_i(B) = \lceil r_i(B)^{-b} \rceil$  disjoint balls of radius  $r_i(B)$  centered at points of Y. Such a collection of disjoint balls exists, by Lemma 3.4. Let  $\mathcal{P}_i(B)$  denote this collection of disjoint balls and let  $\mathcal{P}_i = \bigcup_{B \in \mathcal{P}_{i-1}} \mathcal{P}_i(B)$ . Now sample a random sequence of balls  $B_1, B_2, \ldots$  by picking  $B_1 \in \mathcal{P}_1$  uniformly at random, and for i > 1 picking  $B_i \in \mathcal{P}_i(B_{i-1})$  uniformly at random.

Given a ball  $B = B(x_*, r_*)$ , let  $f_B(x)$  be a Lipschitz function on X defined by

$$f_B(x) = \begin{cases} \min\{r_* - L(x, x_*), r_*/2\} & \text{if } x \in B\\ 0 & \text{otherwise} \end{cases}.$$
 (6)

Let  $f_i = f_{B_i}$  for  $i \ge 1$ . Define  $f_0$  by setting  $f_0(x) = 1/3$  for all  $x \in X$ . The reader may verify that the sum  $\mu = \sum_{i=0}^{\infty} f_i$  is a Lipschitz function. Define the payoff distribution for  $x \in X$  to be a Bernoulli random variable with expectation  $\mu(x)$ . We have thus specified a randomized construction of an instance  $(L, X, \mu)$ .

We claim that for any algorithm  $\mathcal{A}$  and any constant C,

$$\Pr_{\mu, \mathcal{A}}(\forall t \ R_{\mathcal{A}}(t) < Ct^{\gamma}) = 0.$$
(7)

The proof of this claim is based on a "needle in haystack" lemma (Lemma 3.6 below) which states that for all *i*, conditional on the sequence  $B_1, \ldots, B_{i-1}$ , with probability at least  $1 - O((r_i(B_i))^{(b-a)/2})$ , no more than half of the first  $t_i(B_i) = r_i(B_i)^{-a-2}$  strategies picked by  $\mathcal{A}$  lie inside  $B_i$ . The proof of the lemma is deferred to the end of this section.

Any strategy  $x \notin B_i$  satisfies  $\mu(x) < \mu(x^*) - r_i/2$ , so we may conclude that

$$\Pr_{\mu} \left( R_{\mathcal{A}}(t_i(B_i)) < \frac{1}{4} r_i(B_i)^{-a-1} \,|\, B_1, \dots, B_{i-1} \right) \le O\left( (r_i(B_i))^{(b-a)/2} \right). \tag{8}$$

Denoting  $r_i(B_i)$  and  $t_i(B_i)$  by  $r_i$  and  $t_i$ , respectively, we have  $\frac{1}{4}r_i^{-a-1} = \frac{1}{4}t_i^{(a+1)/(a+2)} > Ct_i^{\gamma}$  for all sufficiently large *i*. As *i* runs through the positive integers, the terms on the right side of (8) are dominated by a geometric progression because  $r_i(B_i) \leq 4^{-i}$ . By the Borel-Cantelli Lemma, almost surely there are only finitely many *i* such that the events on the left side of (8) occur. Thus (7) follows.

*Remark.* To prove Theorem 3.5 it suffices to show that for every given algorithm there exists a "hard" problem instance. In fact we proved a stronger result (7): essentially, we construct a probability distribution over problem instances which is hard, almost surely, for every given algorithm. This seems to be the best possible bound since, obviously, a single problem instance cannot be hard for every algorithm.

In rest of this subsection we prove the "needle in haystack" lemma used in the proof of Theorem 3.5.

**Lemma 3.6.** Consider the randomized construction of an instance  $(L, X, \mu)$  in the proof of Theorem 3.5. Fix a bandit algorithm A. Then for all *i*, conditional on the sequence  $B_1, \ldots, B_{i-1}$ , with probability at least  $1 - O((r_i(B_i))^{(b-a)/2})$ , no more than half of the first  $r_i(B_i)^{-a-2}$  strategies picked by A lie inside  $B_i$ .

Let us introduce some notation needed to prove the lemma. Let us fix an arbitrary Lipschitz MAB algorithm  $\mathcal{A}$ . We will assume that  $\mathcal{A}$  is deterministic; the corresponding result for randomized algorithms follows by conditioning on the algorithm's random bits (so that its behavior, conditional on these bits, is deterministic), invoking the lemma for deterministic algorithms, and then removing the conditioning by averaging over the distribution of random bits. Note that since our construction uses only  $\{0, 1\}$ -valued payoffs, and the algorithm  $\mathcal{A}$  is deterministic, the entire history of play in the first t rounds can be summarized by a binary vector  $\sigma \in \{0, 1\}^t$ , consisting of the payoffs observed by  $\mathcal{A}$  in the first t rounds. Thus a payoff function  $\mu$ determines a probability distribution  $P_{\mu}$  on the set  $\{0, 1\}^t$ , i.e. the distribution on t-step histories realized when using algorithm  $\mathcal{A}$  on instance  $\mu$ .

Let B be any ball in the set  $\mathcal{P}_{i-1}$ , let  $n = n_i(B)$ ,  $r = r_i(B)$ , and  $t = t_i(B) = r_i(B)^{-a-2}$ . Let  $B^1, B^2, \ldots, B^n$  be an enumeration of the balls in  $\mathcal{P}_i(B)$ . Choose an arbitrary sequence of balls  $B_1 \supseteq B_2 \supseteq \ldots \supseteq B_{i-1} = B$  such that  $B_1 \in \mathcal{P}_1$  and for all j > 0  $B_j \in \mathcal{P}(B_{j-1})$ . Similarly, for  $k = 1, 2, \ldots, n$ , choose an arbitrary sequence of balls  $B^k = B_i^k \supseteq B_{i+1}^k \supseteq \ldots$  such that  $B_j^k \in \mathcal{P}(B_{j-1}^k)$  for all  $j \ge i$ . Define functions  $f_j$   $(1 \le j \le i-1)$  and  $f_j^k$   $(j \ge i)$  using the balls  $B_j, B_j^k$ , as in the proof of Theorem 3.5. Specifically, use definition (6) and set  $f_j = f_{B_j}$  and  $f_j^k = f_{B_j^k}$ . Let  $\mu^0 = \sum_{j=0}^{i-1} f_j$  and

$$\mu^k = \mu^0 + \sum_{j=i}^{\infty} f_j^k \quad \text{(for } 1 \le k \le n\text{)}.$$

Note that the instances  $\mu^k$   $(1 \le k \le n)$  are equiprobable under our distribution on input instances  $\mu$ . The instance  $\mu^0$  is not one that could be randomly sampled by our construction, but it is useful as a "reference measure" in the following proof. Note that the functions  $\mu^k$  have the following properties, by construction.

- (a)  $1/3 \le \mu^k(x) \le 2/3$  for all  $x \in X$ .
- (b)  $0 \le \mu^k(x) \mu^0(x) \le r$  for all  $x \in X$ .
- (c) If  $x \in X \setminus B^k$ , then  $\mu^k(x) = \mu^0(x)$ .
- (d) If  $x \in X \setminus B^k$ , then there exists some point  $x^k \in B^k$  such that  $\mu^k(x^k) \mu^k(x) \ge r/2$ .

Each of the payoff functions  $\mu^k$   $(0 \le k \le n)$  gives rise to a probability distribution  $P_{\mu^k}$  on  $\{0,1\}^t$ as described in the preceding section. We will use the shorthand notation  $P_k$  instead of  $P_{\mu^k}$ . We will also use  $\mathbf{E}_k$  to denote the expectation of a random variable under distribution  $P_k$ . Finally, we let  $N_k$  denote the random variable defined on  $\{0,1\}^t$  that counts the number of rounds s  $(1 \le s \le t)$  in which algorithm  $\mathcal{A}$ chooses a strategy in  $B^k$  given the history  $\sigma$ .

The following lemma is analogous to Lemma A.1 of [3], and its proof is identical to the proof of that lemma.

**Lemma 3.7.** Let  $f : \{0,1\}^t \to [0,M]$  be any function defined on reward sequences  $\sigma$ . Then for any k,

$$\mathbf{E}_k[f(\sigma)] \le \mathbf{E}_0[f(\sigma)] + \frac{M}{2}\sqrt{-\ln(1-4r^2)\mathbf{E}_0[N_i]}.$$

Applying Lemma 3.7 with  $f = N_k$  and M = t, and averaging over k, we may apply exactly the same reasoning as in the proof of Theorem A.2 of [3] to derive the bound

$$\frac{1}{n}\sum_{k=1}^{n}\mathbf{E}_{k}(N_{k}) \leq \frac{t}{n} + O\left(tr\sqrt{\frac{t}{n}}\right).$$
(9)

Recalling that the actual ball  $B_k$  sampled when randomly constructing  $\mu$  in the proof of Theorem 3.5 is a uniform random sample from  $B^1, B^2, \ldots, B^n$ , we may write  $N_*$  to denote the random variable which counts the number of rounds in which the algorithm plays a strategy in  $B_k$  and the bound (9) implies

$$\mathbf{E}(N_*) = O\left(\frac{t}{n} + tr\sqrt{\frac{t}{n}}\right)$$

Recalling that  $t = r^{-a-2}$  and  $n = r^{-b}$ , we see that the  $O(tr\sqrt{t/n})$  term is the dominant term on the right side, and that it is bounded by  $O(tr^{(b-a)/2})$ . An application of Markov's inequality now yields:

$$\Pr(N_* \ge t/2) = O(r^{(b-a)/2}),$$

completing the proof of Lemma 3.6.

#### **3.2** Beyond the covering dimension

Thus far, we have seen that every metric space X has a bandit algorithm  $\mathcal{A}$  such that  $DIM(\mathcal{A}) = COV(X)$ (the naïve algorithm), and we have seen (via the needle-in-haystack construction, Theorem 3.5) that X can never have a bandit algorithm satisfying  $DIM(\mathcal{A}) < MaxMinCOV(X)$ . When  $COV(X) \neq MaxMinCOV(X)$ , which of these two bounds is correct, or can they both be wrong?

To gain intuition, we will consider two concrete examples. Consider an infinite rooted tree where for each level  $i \in \mathbb{N}$  most nodes have out-degree 2, whereas the remaining nodes (called *fat nodes*) have outdegree x > 2 so that the total number of nodes is  $4^i$ . In our first example, there is exactly one fat node on every level and the fat nodes form a path (called the *fat leaf*). In our second example, there are exactly  $2^i$ fat nodes on every level *i* and the fat nodes form a binary tree (called the *fat subtree*). In both examples, we assign a *weight* of  $2^{-id}$  (for some constant d > 0) to each level-*i* node; this weight encodes the diameter of the set of points contained in the corresponding subtree. An infinite rooted tree induces a metric space (L, X) where X is the set of all infinite paths from the root, and for  $u, v \in X$  we define L(u, v) to be the weight of the least common ancestor of paths u and v. In both examples, the covering dimension is 2d, whereas the max-min-covering dimension is only d because the "fat subset" (i.e. the fat leaf or fat subtree) has covering dimension at most d, and every point outside the fat subset has an open neighborhood of covering dimension d.

In both of the metrics described above, the zooming algorithm (Algorithm 2.3) performs poorly when the optimum  $x^*$  is located inside the fat subset S, because it is too burdensome to keep covering<sup>7</sup> the profusion of strategies located near  $x^*$  as the ball containing  $x^*$  shrinks. An improved algorithm, achieving regret exponent d, modifies the zooming algorithm by imposing quotas on the number of active strategies that lie outside S. At any given time, some strategies outside S may not be covered; however, it is guaranteed that there exists an optimal strategy which eventually becomes covered and remains covered forever afterward. Intuitively, if some optimal strategy lies in S then imposing a quota on active strategies outside S does not hurt. If no optimal strategy lies in S then all of S gets covered eventually and stays covered thereafter, in which case the uncovered part of the strategy set has low covering dimension and (starting after the time when S becomes permanently covered) no quota is ever exceeded.

This use of quotas extends to the following general setting which abstracts the idea of "fat subsets":

**Definition 3.8.** Fix a metric space (L, X). A closed subset  $S \subset X$  is *d*-fat if  $COV(S) \leq d$  and for any open superset U of S we have  $COV(X \setminus U) \leq d$ . More generally, a *d*-fat decomposition of depth k is a decreasing sequence  $X = S_0 \supset \ldots \supset S_k \supset S_{k+1} = \emptyset$  of closed subsets such that  $COV(S_k) \leq d$  and  $COV(S_i \setminus U) \leq d$  whenever  $i \in [k]$  and U is an open superset of  $S_{i+1}$ .

<sup>&</sup>lt;sup>7</sup>Recall that a strategy u is called *covered* at time t if for some active strategy v we have  $L(u, v) \leq r_t(v)$ .

**Example 3.9.** Let (L, X) be the metric space in either of the two "tree with a fat subset" examples. Then the corresponding "fat subset" S is d-fat. For an example of a fat decomposition of depth k = 2, consider the product metric  $(L^*, X \times X)$  defined by

$$L^*((x_1, x_2), (y_1, y_2)) = L(x_1, y_1) + L(x_2, y_2),$$

with a fat decomposition given by  $S_1 = (S \times X) \cup (X \times S)$  and  $S_2 = S \times S$ .

When X is a metric space with a  $d^*$ -fat decomposition  $\mathcal{D}$ , the algorithm described earlier can be modified to achieve regret  $O(t^{\gamma})$  for any  $\gamma > 1 - 1/(d^* + 2)$ , by instituting a separate quota for each subset  $S_i$ . The algorithm requires access to a  $\mathcal{D}$ -covering oracle which for a given *i* and a given finite set of open balls (given by the centers and the radii) either reports that the balls cover  $S_i$ , or returns some strategy in  $S_i$  which is not covered by the balls. No further knowledge of  $\mathcal{D}$  or the metric space is required.

**Theorem 3.10.** Consider the Lipschitz MAB problem on a fixed compact metric space with a  $d^*$ -fat decomposition  $\mathcal{D}$ . Then for any  $d > d^*$  there is an algorithm  $\mathcal{A}_{\mathcal{D}}$  such that  $\text{DIM}(\mathcal{A}_{\mathcal{D}}) \leq d$ .

*Remarks.* (1) We can relax the compactness assumption in Theorem 3.10: instead, we can assume that the *completion* of the metric space is compact and re-define the sets in the d-fat decomposition as subsets of the completion (possibly disjoint with the strategy set). This corresponds to the "fat leaf" which lies outside the strategy set. Such extension requires some minor modifications.

(2) The per-metric guarantee expressed by Theorem 3.10 can be complemented with sharper *per-instance* guarantees. First, for every problem instance  $\mathcal{I}$  the per-instance regret dimension  $\text{DIM}_{\mathcal{I}}(\mathcal{A})$  is upper-bounded by the zooming dimension of  $\mathcal{I}$ . Second, if for some c > 0 the c-covering dimension of X is finite then for some  $\gamma < 1$  and all t we have  $R_{\mathcal{A}}(t) \leq O(ct^{\gamma})$ . However, as this extension is tangential to our main storyline, we focus on analyzing the regret dimension.

**The algorithm.** Our algorithm proceeds in phases i = 1, 2, 3, ... of  $2^i$  rounds each. In a given phase, we run a fresh instance of the following *phase algorithm*  $\mathcal{A}_{ph}(T, d, \mathcal{D})$  parameterized by the phase length  $T = 2^i$ , target dimension  $d > d^*$  and the  $\mathcal{D}$ -covering oracle. The phase algorithm is a version of a single phase of the zooming algorithm (Algorithm 2.3) with very different rules for activating strategies. As in Algorithm 2.3, the confidence radius and the index are defined by (2) and (3), respectively. At the start of each round some strategies are activated, and then an active strategy with the maximal index is played.

Let us specify the activation rules. Let k be the depth of the decomposition  $\mathcal{D}$ , and denote  $\mathcal{D} = \{S_i\}_{i=0}^{k+1}$ . Initially the algorithm constructs  $2^{-j}$ -nets  $\mathcal{N}_j$ ,  $j \in \mathbb{N}$ , using the covering oracle. It finds the largest j such that  $\mathcal{N} = \mathcal{N}_j$  contains at most  $\frac{1}{2} T^{d/(d+2)}$  points, and activates all strategies in  $\mathcal{N}$ . The rest of the active strategies are partitioned into k+1 pools  $P_i \subset S_i$  such that at each time t each pool  $P_i$  satisfies the following quota (that we denote  $Q_i$ ):

$$|\{u \in P_i : r_t(u) \ge \rho\}| \le C_\rho \ \rho^{-d} \tag{10}$$

where  $\rho = T^{-1/(d+2)}$  and  $C_{\rho} = (64k \log \frac{1}{\rho})^{-1}$ . In the beginning of each round the following activation routine is performed. If there exists a set  $S_i$  such that some strategy in  $S_i$  is not covered and *there is room under the corresponding quota*  $Q_i$ , pick one such strategy, activate it, and add it to the corresponding pool  $P_i$ . Since for a given strategy u the confidence radius  $r_t(u)$  is non-increasing in t, the constraint (10) is never violated. Repeat until there are no such sets  $S_i$  left. This completes the description of the algorithm.

**Analysis.** As was the case in Section 2, the analysis of the unbounded-time-horizon algorithm reduces to proving a lemma about the regret of each phase algorithm.

**Lemma 3.11.** Fix a problem instance in the setting of Theorem 3.10. Let  $\mathcal{A}_{ph}(T) = \mathcal{A}_{ph}(T, d, \mathcal{D})$ . Then

$$(\exists t_{\min} < \infty) \ (\forall T \ge t_{\min}) \quad R_{\mathcal{A}_{ph}(T)}(T) \le T^{1-1/(d+2)}.$$

$$(11)$$

Note that the lemma bounds the regret of  $\mathcal{A}_{ph}(T)$  for time  $T \ge t_{\min}$  only. Proving Theorem 3.10 is now straightforward:

**Proof of Theorem 3.10:** Let  $\mathcal{A}_{ph}(T)$  be the phase algorithm from Lemma 3.11. Recall that in each phase *i* in the overall algorithm  $\mathcal{A}$  we simply run a fresh instance of algorithm  $\mathcal{A}_{ph}(2^i)$  for  $2^i$  steps.

Let  $t_0$  be the  $t_{\min}$  from (11) rounded up to the nearest end-of-phase time. Let  $i_0$  be the phase starting at time  $t_0 + 1$ . Note that  $R_{\mathcal{A}}(t_0) \leq t_0$ . Let  $R_i$  be the regret accumulated by  $\mathcal{A}$  during phase i. Let  $\gamma = \frac{d+1}{d+2}$ . Then for any time  $t \geq t_0^{1/\gamma}$  in phase i we have  $R_{\mathcal{A}}(t) \leq t_0 + \sum_{j=i_0}^{i} R_j \leq t_0 + \sum_{j=i_0}^{i} (2^j)^{\gamma} \leq O(t^{\gamma})$ .  $\Box$ 

In the remainder of this section we prove Lemma 3.11. Let us fix a problem instance of the Lipschitz MAB problem on a compact metric space (L, X) with a depth- $k d^*$ -fat decomposition  $\mathcal{D} = \{S_i\}_{i=0}^{k+1}$ . Fix  $d > d^*$  and let  $\mathcal{A}_{ph}(T) = \mathcal{A}_{ph}(T, d, \mathcal{D})$  be the phase algorithm. Let  $\mu$  be the expected reward function and let  $\mu^* = \sup_{u \in X} \mu(u)$  be the optimal reward. Let  $\Delta(u) = \mu^* - \mu(u)$ .

By definition of the Lipschitz MAB problem,  $\mu$  is a continuous function on the metric space (L, X). Therefore the supremum  $\mu^*$  is achieved by some strategy (call such strategies *optimal*). Say that a run of algorithm  $\mathcal{A}_{ph}(T)$  is *well-covered* if at every time  $t \leq T$  some optimal strategy is covered.

Say that a run of algorithm  $\mathcal{A}_{ph}(T)$  is *clean* if the property in Claim 2.2 holds for all times  $t \leq T$ . Note that a given run is clean with probability at least  $1 - T^{-2}$ . The following lemma adapts the technique from Lemma 2.5 to the present setting:

**Claim 3.12.** Consider a clean run of algorithm  $\mathcal{A}_{ph}(T)$ .

- (a) If strategies u, v are active at time  $t \leq T$  then  $\Delta(v) \Delta(u) \leq 4r_t(v)$ .
- (b) if the run is well-covered and strategy v is active at time  $t \leq T$  then  $\Delta(v) \leq 4r_t(v)$ .

The quotas (10) are chosen so that the regret computation in Claim 2.7 works out for a clean and wellcovered run of algorithm  $\mathcal{A}_{ph}(T)$ .

**Claim 3.13.**  $R_{\mathcal{A}}(T) \leq T^{1-1/(d+2)}$  for any clean well-covered run of algorithm  $\mathcal{A} = \mathcal{A}_{ph}(T)$ .

Sketch. Let  $A_t(\delta)$  be the set of all strategies  $u \in X$  such that u is active at time  $t \leq T$  and  $\delta \leq r_t(u) < 2\delta$ . Note that for any such strategy we have  $n_t(u) \leq O(\log T) \delta^{-2}$  and  $\Delta(u) \leq 4r_t(u) < 8\delta$ . Write

$$R^*(T) := \sum_{u \in X} \Delta(u) \, n_T(u) \le \rho T + \sum_{i=0}^{\lceil \log 1/\rho \rceil} \sum_{u \in A_T(2^{-i})} \Delta(u) \, n_T(u),$$

where  $\rho = T^{-1/(d+2)}$  and apply the quotas (10).

Let  $S_{\ell}$  be the smallest set in  $\mathcal{D}$  which contains some optimal strategy. Then there is an optimal strategy contained in  $S_{\ell} \setminus S_{\ell+1}$ ; let  $u^*$  be one such strategy. The following claim essentially shows that the irrelevant high-dimensional subset  $S_{\ell+1}$  is eventually pruned away.

**Claim 3.14.** There exists an open set U containing  $S_{\ell+1}$  such that  $u^* \notin U$  and U is always covered throughout the first T steps of any clean run of algorithm  $\mathcal{A}_{ph}(T)$ , provided that T is sufficiently large.

*Proof.*  $S_{\ell+1}$  is a compact set since it is a closed subset of a compact metric space. Since function  $\mu$  is continuous, it assumes a maximum value on  $S_{\ell+1}$ . By construction, this maximum value is strictly less than  $\mu^*$ . So there exists  $\epsilon > 0$  such that  $\Delta(w) > 8\epsilon$  for any  $w \in S_{\ell+1}$ . Define  $U = B(S_{\ell+1}, \epsilon/2)$ . Note that  $u^* \notin U$  since  $8\epsilon < \Delta(w) \le L(u^*, w)$  for any  $w \in S_{\ell+1}$ .

Recall that in the beginning of algorithm  $\mathcal{A}(T)$  all strategies in some  $2^{-j}$ -net  $\mathcal{N}$  are activated. Suppose T is large enough so that  $2^{-j} \leq \epsilon$ .

Consider a clean run of algorithm  $\mathcal{A}_{ph}(T)$ . We claim that U is covered at any given time  $t \leq T$ . Indeed, fix  $u \in U$ . By definition of U there exists a strategy  $w \in S_{\ell+1}$  such that  $L(u, w) < \epsilon/2$ . By definition of  $\mathcal{N}$  there exist  $v, v^* \in \mathcal{N}$  such that  $L(v, w) \leq \epsilon$  and  $L(u^*, v^*) \leq \epsilon$ . Note that:

- (a)  $\Delta(v^*) = \mu(u^*) \mu(v^*) \le L(u^*, v^*) \le \epsilon.$
- (b) Since  $L(v, w) \leq \epsilon$  and  $\Delta(w) > 8\epsilon$ , we have  $\Delta(v) > 7\epsilon$ .
- (c) By Claim 3.12 we have  $\Delta(v) \Delta(v^*) \leq 4r_t(v^*)$ .

Combining (a-c), it follows that  $r_t(v) \ge \frac{3}{2} \epsilon \ge L(u, v)$ , so v covers u. Claim proved.

**Proof of Lemma 3.11:** By Claim 3.13 it suffices to show that if T is sufficiently large then any clean run of algorithm  $\mathcal{A}_{ph}(T)$  is well-covered. (Runs that are not clean contribute only O(1/T) to the expected regret of  $\mathcal{A}_{ph}(T)$ , because the probability that a run is not clean is at most  $T^{-2}$  and the regret of such a run is at most T.) Specifically, we will show that  $u^*$  is covered at any time  $t \leq T$  during a clean run of  $\mathcal{A}_{ph}(T)$ . It suffices to show that at any time  $t \leq T$  there is room under the corresponding quota  $Q_{\ell}$  in (10).

Let U be the open set from Claim 3.14. Since U is an open neighborhood of  $S_{\ell+1}$ , by definition of the fat decomposition it follows that  $COV(S_{\ell} \setminus U) \leq d^*$ . Define  $\rho$  and  $C_{\rho}$  as in (10) and fix  $d' \in (d^*, d)$ . Then for any sufficiently large T it is the case that (i)  $S_{\ell} \setminus U$  can be covered with  $(\frac{1}{\rho})^{d'}$  sets of diameter  $< \rho$  and moreover (ii) that  $(\frac{1}{\rho})^{d'} \leq \frac{1}{2} C_{\rho} \rho^{-d}$ . Fix time  $t \leq T$  and let  $A_t$  be the set of all strategies u such that u is in the pool  $P_{\ell}$  at time t and

Fix time  $t \leq \dot{T}$  and let  $A_t$  be the set of all strategies u such that u is in the pool  $P_{\ell}$  at time t and  $r_t(u) \geq \rho$ . Note that  $A_t \subset S_{\ell} \setminus U$  since U is always covered, and by the specification of  $\mathcal{A}_{ph}$  only active uncovered strategies in  $S_{\ell}$  are added to pool  $P_{\ell}$ . Moreover,  $A_t$  is  $\rho$ -separated. (Indeed, let  $u, v \in A_t$  and assume u has been activated before v. Then  $L(u, v) > r_s(u) \geq r_t(u) \geq \rho$ , where s is the time when v was activated.) It follows that  $|A_t| \leq \frac{1}{2} C_{\rho} \rho^{-d}$ , so there is room under the corresponding quota  $Q_{\ell}$  in (10).

#### **3.3** The per-metric optimal algorithm

The algorithm in Theorem 3.10 requires a fat decomposition of finite depth, which in general might not exist. To extend the ideas of the preceding section to arbitrary metric spaces, we must generalize Definition 3.8 to *transfinitely infinite* depth.

**Definition 3.15.** Fix a metric space (L, X). Let  $\beta$  denote an arbitrary ordinal. A *transfinite d-fat decomposition* of depth  $\beta$  is a transfinite sequence  $\{S_{\lambda}\}_{0 \le \lambda \le \beta}$  of closed subsets of X such that:

- (a)  $S_0 = X$ ,  $S_\beta = \emptyset$ , and  $S_\nu \supseteq S_\lambda$  whenever  $\nu < \lambda$ .
- (b) if  $V \subset X$  is closed, then the set {ordinals  $\nu \leq \beta$ : V intersects  $S_{\nu}$ } has a maximum element.
- (c) for any ordinal  $\lambda \leq \beta$  and any open set  $U \subset X$  containing  $S_{\lambda+1}$  we have  $COV(S_{\lambda} \setminus U) \leq d$ .

Note that for a finite depth  $\beta$  the above definition is equivalent to Definition 3.8. In Theorem 3.17 below, we will show how to modify the "quota algorithms" from the previous section to achieve regret dimension d in any metric with a transfinite  $d^*$ -fat decomposition for  $d^* < d$ . This gives an optimal algorithm for every metric space X because of the following surprising relation between the max-min-covering dimension and transfinite fat decompositions.

**Proposition 3.16.** For every compact metric space (L, X), the max-min-covering dimension of X is equal to the infimum of all d such that X has a transfinite d-fat decomposition.

*Proof.* If  $\emptyset \neq Y \subseteq X$  and MinCOV(Y) > d then, by transfinite induction,  $Y \subseteq S_{\lambda}$  for all  $\lambda$  in any transfinite d-fat decomposition, contradicting the fact that  $S_{\beta} = \emptyset$ . Thus, the existence of a transfinite d-fat decomposition of X implies  $d \geq MaxMinCOV(X)$ . To complete the proof we will construct, given any

d > MaxMinCOV(X), a transfinite d-fat decomposition of depth  $\beta$ , where  $\beta$  is any ordinal whose cardinality exceeds that of X. For a metric space Y, define the set of d-thin points TP(Y, d) to be the union of all open sets  $U \subseteq Y$  satisfying COV(U) < d. Its complement, the set of d-fat points, is denoted by FP(Y, d). Note that it is a closed subset of Y.

For an ordinal  $\lambda \leq \beta$ , we define a set  $S_{\lambda}$  using transfinite induction as follows:

- 1.  $S_0 = X$  and  $S_{\lambda+1} = FP(S_{\lambda}, d)$  for each ordinal  $\lambda$ .
- 2. If  $\lambda$  is a limit ordinal then  $S_{\lambda} = \bigcap_{\nu < \lambda} S_{\nu}$ .

Note that each  $S_{\lambda}$  is closed, by transfinite induction. It remains to show that  $\mathcal{D} = \{S_{\lambda}\}_{\lambda \in \mathcal{O}}$  satisfies the properties (a-c) in Definition 3.15. It follows immediately from the construction that  $S_0 = X$  and  $S_{\nu} \supseteq S_{\lambda}$  when  $\nu < \lambda$ . To prove that  $S_{\beta} = \emptyset$ , observe first that the sets  $S_{\lambda} \setminus S_{\lambda+1}$  (for  $0 \le \lambda < \beta$ ) are disjoint subsets of X, and the number of such sets is greater than the cardinality of X, so at least one of them is empty. This means that  $S_{\lambda} = S_{\lambda+1}$  for some  $\lambda < \beta$ . If  $S_{\lambda} = \emptyset$  then  $S_{\beta} = \emptyset$  as desired. Otherwise, the relation  $FP(S_{\lambda}, d) = S_{\lambda}$  implies that  $MinCOV(S_{\lambda}) \ge d$  contradicting the assumption that MaxMinCOV(X) < d. This completes the proof of property (a). To prove property (b), suppose  $\{\nu_i \mid i \in \mathcal{I}\}$  is a set of ordinals such that  $S_{\nu_i}$  intersects V for every i. Let  $\nu = \sup\{\nu_i\}$ . Then  $S_{\nu} \cap V = \bigcap_{i \in \mathcal{I}} (S_{\nu_i} \cap V)$ , and the latter set is nonempty because X is compact and the closed sets  $\{S_{\nu_i} \cap V \mid i \in \mathcal{I}\}$  have the finite intersection property. Finally, to prove property (c), note that if U is an open neighborhood of  $S_{\lambda+1}$  then the set  $T = S_{\lambda} \setminus U$  is closed (hence compact) and is contained in  $TP(S_{\lambda}, d)$ . Consequently T can be covered by open sets V satisfying COV(V) < d. By compactness of T, this covering has a finite subcover  $V_1, \ldots, V_m$ , and consequently  $COV(T) = max_{1 \le i \le m} COV(V_i) < d$ .

**Theorem 3.17.** Consider the Lipschitz MAB problem on a compact metric space (L, X). For any d > MaxMinCOV(X) there exists an algorithm  $A_d$  such that  $DIM(A_d) \leq d$ .

Note that Theorem 1.2 follows immediately by combining Theorem 3.17 with Theorem 3.5.

We next describe an algorithm  $\mathcal{A}_d$  satisfying Theorem 3.17. The algorithm requires two oracles: a depth oracle  $\mathtt{Depth}(\cdot)$  and a  $\mathcal{D}$ -covering oracle  $\mathcal{D}$ - $\mathtt{Cov}(\cdot)$ . For any finite set of open balls  $B_0, B_1, \ldots, B_n$  (given via the centers and the radii) whose union is denoted by B,  $\mathtt{Depth}(B_0, B_1, \ldots, B_n)$  returns the maximum ordinal  $\lambda$  such that  $S_{\lambda}$  intersects the closure  $\overline{B}$ ; such an ordinal exists by Definition 3.15(b).<sup>8</sup> Given a finite set of open balls  $B_0, B_1, \ldots, B_n$  with union B as above, and an ordinal  $\lambda$ ,  $\mathcal{D}$ - $\mathtt{Cov}(\lambda, B_0, B_1, \ldots, B_n)$  either reports that B covers  $S_{\lambda}$ , or it returns a strategy  $x \in S_{\lambda} \setminus B$ .

**The algorithm.** Our algorithm proceeds in phases i = 1, 2, 3, ... of  $2^i$  rounds each. In any given phase i, there is a "target ordinal"  $\lambda(i)$  (defined at the end of the preceding phase), and we run an algorithm during the phase which: (i) activates some nodes initially; (ii) plays a version of the zooming algorithm which only activates strategies in  $S_{\lambda(i)}$ ; (iii) concludes the phase by computing  $\lambda(i + 1)$ . The details are as follows. In a given phase we run a fresh instance of a phase algorithm  $\mathcal{A}_{ph}(T, d, \lambda)$  where  $T = 2^i$  and  $\lambda = \lambda(i)$  is a *target ordinal* for phase *i*, defined below when we give the full description of  $\mathcal{A}_{ph}(T, d, \lambda)$ . The goal of  $\mathcal{A}_{ph}(T, d, \lambda)$  is to satisfy the per-phase bound

$$R_{\mathcal{A}_{\mathrm{ph}}(T,d,\lambda)}(T) = \widetilde{O}(T^{\gamma}) \tag{12}$$

for all  $T > T_0$ , where  $\gamma = 1 - 1/(d+2)$  and  $T_0$  is a number which may depend on the instance  $\mu$ . Then, to derive the bound  $R_{\mathcal{A}_d}(t) = \widetilde{O}(t^{\gamma})$  for all t we simply sum per-phase bounds over all phases ending before time 2t.

<sup>&</sup>lt;sup>8</sup>To avoid the question of how arbitrary ordinals are represented on the oracle's output tape, we can instead say that the oracle outputs a point  $u \in S_{\lambda}$  instead of outputting  $\lambda$ . In this case, the definition of  $\mathcal{D}$ -Cov should be modified so that its first argument is a point of  $S_{\lambda}$  rather than  $\lambda$  itself.

Initially  $\mathcal{A}_{ph}(T, d, \lambda)$  uses the covering oracle to construct  $2^{-j}$ -nets  $\mathcal{N}_j$ , j = 0, 1, 2, ..., until it finds the largest j such that  $\mathcal{N} = \mathcal{N}_j$  contains at most  $\frac{1}{2} T^{d/(d+2)} \log(T)$  points. It activates all strategies in  $\mathcal{N}$  and sets

$$\varepsilon(i) = \max\{2^{-j}, 32 \, T^{-1/(d+2)} \log(T)\}$$

After this initialization step, for every active strategy v we define the confidence radius

$$r_t(v) := \max\left\{T^{-1/(d+2)}, \sqrt{\frac{8\log T}{2 + n_t(v)}}\right\},$$

where  $n_t(v)$  is the number of times v has been played by the phase algorithm  $\mathcal{A}_{ph}(T, d, \lambda)$  before time t. Let  $B_0, B_1, \ldots, B_n$  be an enumeration of the open balls belonging to the collection

 $\{B(v, r_t(v)) \mid v \text{ active at time } t\}.$ 

If  $n < \frac{1}{2} T^{d/(d+2)} \log(T)$  then we perform the oracle call  $\mathcal{D}$ -Cov $(\lambda, B_0, \ldots, B_n)$ , and if it reports that a point  $x \in S_{\lambda}$  is uncovered, we activate x and set  $n_t(x) = 0$ . The index of an active strategy v is defined as  $\mu_t(v) + 4r_t(v)$  — note the slight difference from the index defined in Algorithm 2.3 — and we always play the active strategy with maximum index. To complete the description of the algorithm, it remains to explain how the ordinals  $\lambda(i)$  are defined. The definition is recursive, beginning with  $\lambda(1) = 0$ . At the end of phase i  $(i \ge 1)$ , we let  $B_0, B_1, \ldots, B_m$  be an enumeration of the open balls in the set  $\{B(v, \varepsilon(i)) | v \text{ active}, r_T(v) < \varepsilon(i)/2\}$ . Finally, we set  $\lambda(i + 1) = \text{Depth}(B_0, B_1, \ldots, B_m)$ .

**Proof of Theorem 3.17:** Since we have modified the definition of index, we must prove a variant of Claim 3.12 which asserts the following:

In a clean run of 
$$\mathcal{A}_{ph}$$
, if  $u, v$  are active at time  $t$  then  $\Delta(v) - \Delta(u) \le 5r_t(v)$ . (13)

To prove it, let s be the latest round in  $\{1, 2, ..., t\}$  when v was played. We have  $r_t(v) = r_s(v)$ , and  $\Delta(v) - \Delta(u) = \mu(u) - \mu(v)$ , so it remains to prove that

$$\mu(u) - \mu(v) \le 5r_s(v). \tag{14}$$

From the fact that v was played instead of u at time s, together with the fact that both strategies are clean,

$$\mu_s(u) + 4r_s(u) \le \mu_s(v) + 4r_s(v)$$
(15)

$$\mu(u) - \mu_s(u) \le r_s(u) \tag{16}$$

$$\mu_s(v) - \mu(v) \le r_s(v). \tag{17}$$

We obtain (14) by adding (15)-(17), noting that  $r_s(u) > 0$ . This completes the proof of (13).

Let  $\lambda$  be the maximum ordinal such that  $S_{\lambda}$  contains an optimal strategy  $u^*$ ; such an ordinal exists by Definition 3.15(b). We will prove that for sufficiently large i, if the i-th phase is clean, then  $\lambda(i) = \lambda$ . The set  $S_{\lambda+1}$  is compact, and the function  $\mu$  is continuous, so it assumes a maximum value on  $S_{\lambda+1}$  which is, by construction, strictly less than  $\mu^*$ . Choose  $\varepsilon > 0$  such that  $\Delta(w) > 5\varepsilon$  for all  $w \in S_{\lambda+1}$ , and choose  $T_0 = 2^{i_0}$  such that  $\varepsilon(i_0) \leq \varepsilon$ . We shall prove that for all  $T = 2^i \geq T_0$  and all ordinals  $\nu$ , a clean run of  $\mathcal{A}_{\rm ph}(T, d, \nu)$  results in setting  $\lambda(i + 1) = \lambda$ . First, let  $v^* \in \mathcal{N}$  be such that  $L(u^*, v^*) \leq \varepsilon(i)$ . If vis active and  $r_T(v) < \varepsilon(i)/2$  then (13) implies that  $\Delta(v) - \Delta(v^*) \leq \frac{5}{2}\varepsilon(i)$  hence  $\Delta(v) \leq \frac{7}{2}\varepsilon(i)$ . As  $\Delta(w) > 5\varepsilon \geq 5\varepsilon(i)$  for all  $w \in S_{\lambda+1}$ , it follows that the closure of  $B(v, \varepsilon(i))$  does not intersect  $S_{\lambda+1}$ . This guarantees that  $\text{Depth}(B_0, B_1, \ldots, B_m)$  returns an ordinal less than or equal to  $\lambda$ . Next we must prove that this ordinal is greater than or equal to  $\lambda$ . Note that the total number of strategies activated by  $\mathcal{A}_{ph}(T, d, \nu)$  is bounded above by  $T^{d/(d+2)}\log(T)$ . Let  $A_T$  denote the set of strategies active at time T and let

$$v^0 = \arg \max_{v \in A_T} n_T(v).$$

By the pigeonhole principle,  $n_T(v^0) \ge T^{2/(d+2)}/\log(T)$  and hence  $r_T(v^0) < 3T^{-1/(d+2)}\log(T)$ . If t denotes the last time at which  $v^0$  was played, then we have

$$I_t(v^0) = \mu_t(v^0) + 4r_t(v^0) \le \mu^* + 5r_t(v^0) \le \mu^* + 15T^{-1/(d+2)}\log(T) < \mu^* + \varepsilon(i)/2,$$

provided that the phase is clean and that  $T \ge T_0$ . Since  $v^0$  had maximum index at time t, we deduce that  $I_t(v^*) < \mu^* + \varepsilon(i)/2$  as well. As  $L(u^*, v^*) \le \varepsilon(i)$  we have  $\mu_t(v^*) \ge \mu^* - \varepsilon(i) - r_t(v^*)$  provided the phase is clean. To finish the proof we observe that

$$\mu^* + \varepsilon(i)/2 > I_t(v^*) \ge \mu^* - \varepsilon(i) + 3r_t(v^*)$$

which implies  $r_t(v^*) < \varepsilon(i)/2$ . Since the confidence radius does not increase over time, we have  $r_T(v^*) < \varepsilon(i)/2$  so  $B(v^*, \varepsilon(i))$  is one of the balls  $B_0, B_1, \ldots, B_m$ . Since  $u^*$  is contained in the closure of this ball, we may conclude that  $\text{Depth}(B_0, B_1, \ldots, B_m)$  returns the ordinal  $\lambda$  as desired.

Let  $U = B(S_{\lambda+1}, \varepsilon(i)/2)$ . As in Claim 3.14 it holds that in any clean phase, U is covered throughout the phase by balls centered at points of  $\mathcal{N}$ . Hence for any pair of consecutive clean phases, in the second phase of the pair our algorithm only calls the covering oracle  $\mathcal{D}$ -Cov with the proper ordinal  $\lambda$  (i.e. the maximum  $\lambda$  such that  $S_{\lambda}$  contains an optimal strategy) and with a set of balls  $B_0, B_1, \ldots, B_n$  that covers U. Also, note that an active strategy v during a run of  $\mathcal{A}_{ph}(T, d, \lambda)$  never has a confidence radius  $r_t(v)$  less than  $\delta = T^{-1/(d+2)}$ , so the strategies activated by the covering oracle form a  $\delta$ -net in the space  $S_{\lambda} \setminus U$ . By Definition 3.15(c), a  $\delta$ -net in  $S_{\lambda} \setminus U$  contains fewer than  $O(\delta^{-d})$  points. Hence for sufficiently large T the "quota" of  $\frac{1}{2} T^{d/(d+2)}$  active strategies is never reached, which implies that every point of  $S_{\lambda}$  — including  $u^*$  — is covered throughout the phase. The upper bound on the regret of  $\mathcal{A}_{ph}(T, d, \lambda)$  concludes as in the proof of Theorem 2.4.

## **4** Zooming algorithm: extensions and examples

We extend the analysis in Section 2 in several directions, and follow up with examples.

- In Section 4.1 we note that our analysis works under a more abstract notion of the confidence radius: essentially, it can be any function of the history of playing a given strategy such that Claim 2.2 holds. This observation leads to sharper results if the reward from playing each strategy u is  $\mu(u)$  plus an independent *noise* of a known and "benign" shape; we provide several concrete examples.
- In Section 4.2 we provide an improved version of the confidence radius such that the zooming algorithm satisfies the guarantee in Theorem 2.4 and achieves a better regret exponent  $\frac{d}{d+1}$  if the maximal reward is exactly 1. The analysis builds on a novel Chernoff-style bound which, to the best of our knowledge, has not appeared in the literature.
- In Section 4.3 we consider the an example which show-cases both the notion of the zooming dimension and the improved algorithm from Section 4.2. It is the *target MAB* problem, a version of the Lipschitz MAB problem in which the expected reward of a given strategy is a equal to its distance to some (unknown) *target set S*. We show that the zooming algorithm performs much better in this setting; in particular, if the metric is doubling and S is finite, it achieves *poly-logarithmic* regret.

- In Section 4.4 we relax some of the assumptions in the Lipschitz MAB problem: we do not require the similarity function L to satisfy the triangle inequality, and we need the Lipschitz condition (1) to hold only if one of the two strategies is optimal. We use this extension to analyze a generalization of the target MAB problem in which  $\mu(u) = f(L(u, S))$  for some known function f.
- Finally, in Section 4.5 we extend the analysis in Section 2 from reward distributions with bounded support<sup>9</sup> to arbitrary reward distributions with a finite absolute third moment. Our analysis relies on the extension of Azuma inequality known as the *non-uniform Berry-Esseen theorem* [21].

Let us recap some conventions we'll be using throughout this section. The zooming algorithm proceeds in phases i = 1, 2, 3, ... of  $2^i$  rounds each. Within a given phase, for each strategy  $v \in X$  and time t,  $n_t(v)$ is the number of times v has been played before time t, and  $\mu_t(v)$  is the corresponding average reward. Also, we denote  $\Delta(v) = \mu^* - \mu(v)$ , where  $\mu^* = \sup_{v \in X} \mu(v)$  is the maximal reward.

#### 4.1 Abstract confidence radius and noisy rewards

In Section 4.1 the confidence radius of a given strategy was defined by (2). Here we generalize this definition to any function of the history of playing this strategy that satisfies certain properties.

**Definition 4.1.** Consider a single phase  $i_{ph}$  of the algorithm. For each strategy v and any time t within this phase, let  $\hat{r}_t(v)$  and  $\hat{\mu}_t(v)$  be non-negative functions of  $i_{ph}$ , t, and the history of playing v up to round t. Call  $\hat{r}_t(v)$  a *confidence radius* with respect to  $\hat{\mu}_t(v)$  if

(i)  $|\hat{\mu}_t(v) - \mu(t)| \leq \hat{r}_t(v)$  with probability at least  $1 - 8^{-i_{\text{ph}}}$ .

(ii) 
$$\frac{3}{4}\hat{r}_t(v) \le \hat{r}_{t+1}(v) \le \hat{r}_t(v)$$
.

The confidence radius is  $(\beta, C)$ -good if  $n_t(v) \leq (C i_{\text{ph}}) \Delta^{-\beta}(v)$  whenever  $\Delta(v) \leq 4\hat{r}_t(v)$ .

*Remark.* Property (i) says that Claim 2.2 holds for the appropriately redefined *clean phase.* Property (ii) is a "smoothness" condition:  $\hat{r}_t(v)$  does not increase with time, and does not decrease too fast. It is needed for the last line of the proof of Lemma 2.5.

Given such confidence radius, we can carry out the proof of Theorem 2.4 with very minor modifications.

**Theorem 4.2.** Consider an instance of the standard Lipschitz MAB problem for which there exists a  $(\beta, c_0)$ -good confidence radius,  $\beta \ge 0$ . Let  $\mathcal{A}$  be an instance of Algorithm 2.3 defined with respect to this confidence radius. Suppose the problem instance has c-zooming dimension d. Then:

(a) If 
$$d + \beta > 1$$
 then  $R_{\mathcal{A}}(t) \le a(t) t^{1-1/(d+\beta)}$  for all t, where  $a(t) = O(c c_0 \log^2 t)^{1/(d+\beta)}$ .

(b) If  $d + \beta \leq 1$  then  $R_{\mathcal{A}}(t) \leq O(c c_0 \log^2 t)$ .

*Remark.* A new feature of this theorem (as compared to Theorem 2.4) is the *poly-logarithmic* bound on regret in part (b). For better intuition on this, note that the exponent in part (a) becomes negative if  $d+\beta < 1$ . Since the regret bound should not be *decreasing* in t, one would expect this term to vanish from the "correct" bound. Indeed, it is easy to check that the computation in the proof of Claim 2.7 results in part (b).

A natural application of Theorem 4.2 if a setting in which the reward from playing each strategy u is  $\mu(u)$  plus an independent *noise* of known shape.

<sup>&</sup>lt;sup>9</sup>In Section 4.1 we also consider *stochastically bounded* distributions such as Gaussians.

**Definition 4.3.** The *Noisy Lipschitz MAB problem* is a standard Lipschitz MAB problem such that every time any strategy u is played, the reward is  $\mu(u)$  plus an independent random sample from some fixed distribution  $\mathcal{P}$  (called the *noise distribution*) which is revealed to the algorithm.

We present several examples in which we take advantage of a "benign" shape of  $\mathcal{P}$ . Interestingly, in these examples the payoff distributions are not restricted to have bounded support.<sup>10</sup> Technically the results are simple corollaries of Theorem 4.2.

We start with perhaps the most natural example when the noise distribution is normal.

**Corollary 4.4.** Consider the Noisy Lipschitz MAB problem with normal noise distribution  $\mathcal{P} = \mathcal{N}(0, \sigma^2)$ . Then there exists an algorithm  $\mathcal{A}$  which enjous guarantee (4) with the right-hand side multiplied by  $\sigma$ .

*Proof.* Define the confidence radius as (2) with the right-hand side multiplied by  $\sigma$ . It is easy to see that this is a  $(2, O(\sigma))$ -good confidence radius. The result follows from Theorem 4.2(a).

*Remark.* In fact, Corollary 4.4 can be extended to noise distributions of a somewhat more general form: let us say that a random variable X is *stochastically*  $(\rho, \sigma)$ -*bounded* if its moment-generating function satisfies

$$E[e^{r(X-E[x])}] \le e^{r^2 \sigma^2/2} \text{ for all } r \in [-\rho, \rho].$$
 (18)

Note that a normal distribution  $\mathcal{N}(0, \sigma^2)$  is  $(\infty, \sigma)$ -bounded, and any distribution with support  $[-\sigma, \sigma]$  is  $(1, \sigma)$ -bounded. The meaning of (18) is that it is precisely the condition needed to establish an Azuma-type inequality: if S is the sum of n independent stochastically  $(\rho, \sigma)$ -bounded random variables with zero mean, then with high probability  $S \leq \tilde{O}(\sigma_i \sqrt{n})$ :

$$\Pr\left[S > \lambda \sigma \sqrt{n}\right] \le \exp(-\lambda^2/2) \quad \text{for any } \lambda \le \frac{1}{2} \rho \, \sigma \sqrt{n}. \tag{19}$$

The derivation and the theorem statement needs to be modified slightly to account for the parameter  $\rho$ ; we omit the details from this version.

Second, we consider the *noiseless* case when all probability mass in  $\mathcal{P}$  is concentrated at 0. Our result holds more generally, when  $\mathcal{P}$  has at least one *point mass*: a point  $x \in \mathbb{R}$  such that  $\mathcal{P}(x) > 0$ .

**Corollary 4.5.** Consider the Noisy Lipschitz MAB problem such that the noise distribution  $\mathcal{P}$  has at least one point mass. Then the problem admits a confidence radius which is  $(\beta, c)$ -good for any given  $\beta > 0$  and a constant  $c = c(\beta, \mathcal{P})$ . The corresponding low-regret guarantees follow via Theorem 4.2.

Proof Sketch. Let  $S = \operatorname{argmax} \mathcal{P}(x)$  be the set of all points with the largest point mass  $p = \max_x \mathcal{P}(x)$ , and let  $q = \max_{x: \mathcal{P}(x) < p} \mathcal{P}(x)$  be the second largest point mass. Then  $n = \Theta(\log t)$  samples suffices to ensure that with high probability each node in S will get at least n(p+q)/2 hits whereas any other node will get less, which exactly locates all points in S. We use confidence radius  $r_t(v) = \Theta(i_{ph})(\frac{3}{4})^{n_t(v)}$ .

Third, we consider noise distributions with a "special region" which can be located using a few samples. This may be a more efficient way to estimate  $\mu(v)$  than using the standard Chernoff-style tail bounds. Moreover, in our examples  $\mathcal{P}$  may be heavy-tailed, so that Chernoff-style bounds do not hold.

**Corollary 4.6.** Consider the Noisy Lipschitz MAB problem with noise distribution  $\mathcal{P}$ . Suppose  $\mathcal{P}$  has a density f(x) which is symmetric around 0 and non-increasing for x > 0. Assume one of the following: (a) f(x) has a sharp peak:  $f(x) = \Theta(|x|^{-\alpha})$  for all small enough |x|, where  $\alpha \in (0, 1)$ .

<sup>&</sup>lt;sup>10</sup>Recall that throughout the paper the payoff distribution of each strategy x has support  $S(x) \subset [0, 1]$ . In this subsection, by a slight abuse of notation, we do not make this assumption.

(b) f(x) piecewise continuous on  $(0, \infty)$  with at least one jump.

Then for some constant  $c_{\mathcal{P}}$  that depends only on  $\mathcal{P}$  the problem admits a  $(\beta, c_{\mathcal{P}})$ -good confidence radius, where (a)  $\beta = 1 - \alpha$ , (b)  $\beta = 1$ . The corresponding low-regret guarantees follow via Theorem 4.2.

*Proof Sketch.* For part (a), note that for any x > 0 in a neighborhood of 0 we have  $\mathcal{P}[(-x,x)] = \Theta(x^{1-\alpha})$ . Therefore  $n = \Theta(x^{\alpha-1} \log t)$  samples suffices to separate with high probability any length-x sub-interval of (-x, x) from any length-x sub-interval of  $(2x, \infty)$ . It follows that using n samples we can approximate the mean reward up to  $\pm O(x)$ . Accordingly, we set  $r_t(v) = \Theta(i_{ph}/n_t(v))^{1/(1-\alpha)}$ .

For part (b), let  $x_0$  be the smallest positive point where density f has a jump. Then by continuity there exists some  $\epsilon > 0$  such that  $\inf_{x \in (x_0 - \epsilon, x_0)} f(x) > \sup_{x \in (x_0, x_0 + \epsilon)} f(x)$ . Therefore for any  $x < \epsilon$  using  $n = \Theta(\frac{1}{x} \log n)$  samples suffices to separate with high probability any length-x sub-interval of  $(0, x_0)$  from any length-x sub-interval of  $(x_0, \infty)$ . It follows that using n samples we can approximate the mean reward up to  $\pm O(x)$ . Accordingly, we set  $r_t(v) = \Theta(i_{ph}/n_t(v))$ .

#### 4.2 What if the maximal expected reward is 1?

We elaborate the algorithm from Section 2 so that it satisfies the guarantee (4) **and** performs much better if the maximal expected reward is 1.

**Definition 4.7.** Consider the Lipschitz MAB problem. Call an algorithm  $\beta$ -good if there exists an absolute constant  $c_0$  such that for any problem instance of c-zooming dimension d it has the properties (ab) in Theorem 4.2. Call a confidence radius  $\beta$ -good if it is  $(\beta, c_0)$ -good for some absolute constant  $c_0$ .

**Theorem 4.8.** Consider the standard Lipschitz MAB problem. There is an algorithm A which is 2-good in general, and 1-good when the maximal expected reward is 1.

The key ingredient here is a refined version of the confidence radius which is much sharper than (2) when the sample average is close to 1. For phase  $i_{ph}$ , we define

$$r_t(v) := \frac{\alpha}{1 + n_t(v)} + \sqrt{\alpha} \ \frac{1 - \mu_t(v)}{1 + n_t(v)} \text{ for some } \alpha = \Theta(i_{\text{ph}}).$$

$$(20)$$

In order to analyze (20) we need to establish the following Chernoff-style bound which, to the best of our knowledge, has not appeared in the literature:

**Lemma 4.9.** Consider n i.i.d. random variables  $X_1 \dots X_n$  on [0,1]. Let  $\mu$  be their mean, and let X be their average. Then for any  $\alpha > 0$  the following holds:

$$\Pr\left[\left|X-\mu\right| < r(\alpha,X) < 3\,r(\alpha,\mu)\,\right] > 1 - e^{-\Omega(\alpha)}, \text{ where } r(\alpha,x) = \frac{\alpha}{n} + \sqrt{\frac{\alpha x}{n}}.$$

*Proof.* We will use two well-known Chernoff Bounds which we state below (e.g. see p. 64 of [20]):

(CB1) 
$$\Pr[|X - \mu| > \delta\mu] < 2e^{-\mu n\delta^2/3}$$
 for any  $\delta \in (0, 1)$ .

(CB2)  $\Pr[X > a] < 2^{-an}$  for any  $a > 6\mu$ .

First, suppose  $\mu \geq \frac{\alpha}{6n}$ . Apply (CB1) with  $\delta = \frac{1}{2}\sqrt{\frac{\alpha}{6\mu n}}$ . Thus with probability at least  $1 - e^{-\Omega(\alpha)}$  we have  $|X - \mu| < \delta \mu \leq \mu/2$ . Moreover, plugging in the value for  $\delta$ ,

$$|X - \mu| < \frac{1}{2}\sqrt{\alpha\mu/n} \le \sqrt{\alpha X/n} \le r(\alpha, X) < 1.5 r(\alpha, \mu).$$

Now suppose  $\mu < \frac{\alpha}{6n}$ . Then using (CB2) with  $a = \frac{\alpha}{n}$ , we obtain that with probability at least  $1 - 2^{-\Omega(\alpha)}$  we have  $X < \frac{\alpha}{n}$ , and therefore

$$|X - \mu| < \frac{\alpha}{n} < r(\alpha, X) < (1 + \sqrt{2}) \frac{\alpha}{n} < 3r(\alpha, \mu). \quad \Box$$

**Proof of Theorem 4.8:** Let us fix a strategy v and time t. Let us use Lemma 4.9 with  $n = n_t(v)$  and  $\alpha = \Theta(i_{ph})$  as in (20), setting each random variable  $X_i$  equal to 1 minus the reward from the *i*-th time strategy v is played in the current phase. Then  $\mu = \mu(v)$  and  $X = \mu_t(v)$ , so the Lemma says that

$$\Pr\left[|\mu_t(v) - \mu(v)| < r_t(v) < 3\left(\frac{\alpha}{n_t(v)} + \sqrt{\frac{\alpha(1 - \mu(v))}{n_t(v)}}\right)\right] > 1 - 2^{\Omega(\alpha)}.$$
(21)

Note that (20) is indeed a confidence radius with respect to  $\mu_t(v)$ : property (i) in Definition 4.1 holds by (21), and it is easy to check that property (ii) holds, too. It is easy to see that (20) is a 2-good confidence radius. It remains to show that it is 1-good when the maximal reward is 1; this is where we use the upper bound on  $r_t(v)$  in (21). It suffices to prove the following claim:

If the maximal reward is 1 and  $\Delta(v) \leq 4r_t(v)$  then  $n_t(v) \leq O(\log t) \Delta(v)^{-1}$ .

Indeed, let  $n = n_t(v)$  and  $\Delta = \Delta(v)$ , and suppose that the maximal reward is 1 and  $\Delta(v) \leq 4r_t(v)$ . Then by (21) we have  $\Delta \leq 4r_t(v) \leq \frac{\alpha}{n} + \sqrt{\alpha\Delta/n}$  for some  $\alpha = O(\log t)$ . Now there are two cases. If  $\frac{\alpha}{n} < \Delta/2$  then  $\sqrt{\alpha\Delta/n} \geq \Delta - \frac{\alpha}{n} > \Delta(v)/2$ , which implies the desired inequality. Else we simply have  $n \leq O(\alpha/\Delta)$ . Claim proved.

### **4.3** Example: expected reward = distance to the target

We consider a version of the Lipschitz MAB problem where the expected reward of a given strategy is equal to its distance to some *target set* which is not revealed to the algorithm.

**Definition 4.10.** The Target MAB problem on a metric space (L, X) with a *target set*  $S \subset X$  is the standard Lipschitz MAB problem on (L, X) with payoff function  $\mu(u) = 1 - L(u, S)$ .

*Remark.* It is a well-known fact that  $L(u, v) \ge L(u, S) - L(v, S)$  for any  $u, v \in X$  and any set  $S \subset X$ . Therefore the payoff function  $\mu$  in Definition 4.10 is Lipschitz on (L, X).

Note that in the Target MAB problem the maximal reward is 1, so we can take advantage of the zooming algorithm  $\mathcal{A}$  from Theorem 4.8. Recall that  $R_{\mathcal{A}}(t) \leq \tilde{O}(c t^{1-1/(1+d)})$  where d is the c-zooming dimension. In this example zooming dimension is about covering B(S, r) with sets of diameter  $\Theta(r)$ : it is the smallest d such that for each r > 0 the ball B(S, r) can be covered with  $c r^{-d}$  sets of diameter  $\leq r/8$ .

Let us refine this bound for metric spaces of finite doubling dimension. In particular, we show that for a finite target set the zooming algorithm from Theorem 4.8 achieves *poly-logarithmic* regret.

**Theorem 4.11.** Consider the Target MAB problem on a metric space of finite doubling dimension  $d^*$ . Let A be the zooming algorithm from Theorem 4.8. Then

$$R_{\mathcal{A}}(t) \le (c \, 2^{O(d^*)} \log^2 t) \ t^{1-1/(1+d)} \ \text{for all } t,$$
(22)

where d is the c-covering dimension of the target set S.

*Proof.* By Theorem 4.8 it suffices to prove that the K-zooming dimension of the pair  $(L, \mu)$  is at most d, for some  $K = c 2^{O(d^*)}$ . In other words, it suffices to cover the set  $S_{\delta} = \{u \in Y : \Delta(u) \leq \delta\}$  with  $K \delta^{-d}$  sets of diameter  $\leq \delta/16$ , for any given  $\delta > 0$ .

Fix  $\delta > 0$  and note that  $\Delta(u) = L(u, S)$ . Note that set S can be covered with  $c \, \delta^{-d}$  sets  $\{C_i\}_i$  of diameter  $\leq \delta$ . It follows that the set  $S_{\delta}$  can be covered with  $r^{-d}$  sets  $\{B(C_i, r)\}_i$  of diameter  $\leq 3r$ . Moreover, each set  $B(C_i, r)$  can be covered with  $2^{O(d^*)}$  of sets of diameter  $\leq \delta/16$ .

*Remarks.* This theorem is useful when  $d < d^*$ , i.e. when the target set is a low-dimensional subset of the metric space. Recall that the zooming algorithm is self-tuning: it does not need to know  $d^*$  and d, and in fact it does not even need to know that it is presented with an instance of the Target MAB problem!

We note in passing that it is very easy to extend Theorem 4.11 to a setting in which the strategy set Y is a proper subset of the metric space (L, X) and does not contain the target set S. If L(Y, S) = 0 then the guarantee (22) holds as is. If L(Y, S) > 0 then the following guarantee holds:

$$R_{\mathcal{A}}(t) \le (c \, 2^{O(d^*)} \log^2 t) t^{1-1/(2+d)}$$
 for all  $t$ ,

where d is the c-covering dimension of the set B(S, r), r = L(Y, S).

#### 4.4 The Lipschitz MAB problem under relaxed assumptions

The analysis in Section 2 does not require all the assumptions in the Lipschitz MAB problem. In fact, it never uses the triangle inequality, and applies the Lipschitz condition (1) only if (essentially) one of the two strategies in (1) is optimal. Let us formulate our results under the properly relaxed assumptions. In what follows, the *zooming algorithm* will refer to the algorithm in Theorem 4.8.

**Theorem 4.12.** Consider a version of the Lipschitz MAB problem on (L, X) in which the similarity metric is not required to satisfy triangle inequality,<sup>11</sup> and the Lipschitz condition (1) is replaced by

$$(\forall u \in X) \quad \Delta(u) \le L(u, v^*) \quad \text{for some } v^* = \operatorname*{argmax}_{v \in X} \mu(v)$$
 (23)

More generally, if such node  $v^*$  does not exist, assume

$$(\forall \epsilon > 0) \quad (\exists v^* \in X) \quad (\forall u \in X) \quad \Delta(u) \le L(u, v^*) + \epsilon.$$
(24)

Then the guarantees for the zooming algorithm in Theorem 4.8 still hold.

We apply this theorem to a generalization of the Target MAB problem in which  $\mu(u) = f(L(u, S))$ for some known non-decreasing *shape function*  $f : [0, 1] \rightarrow [0, 1]$ . Let us define a quasi-distance  $L_f$  by  $L_f(u, v) = f(L(u, v)) - f(0)$ . It is easy to see that  $L_f$  satisfies (23). Indeed, fix any  $u^* \in S$ . Then

$$\Delta(u) = f(L(u, S)) - f(0) \le f(L(u, u^*)) - f(0) = L_f(u, u^*) \qquad (\forall u \in X).$$

Thus we can use the zooming algorithm on the quasi-distance  $L_f$  and enjoy the guarantees in Theorem 4.8. Below we refine these guarantees for several examples.

Our goal here is to provide clean illustrative statements rather than cover the most general setting to which our refined guarantees apply. Therefore we start with the most concrete example which we formulate as a theorem, and follow up with some extensions which we list without a proof.

<sup>&</sup>lt;sup>11</sup>Formally, we require L to be a symmetric function  $X \times X \to [0, \infty]$  such that L(x, x) = 0 for all  $x \in X$ . We call such function a *quasi-distance* on X.

**Theorem 4.13.** Consider the Target MAB problem on a metric space (L, X) of finite doubling dimension  $d^*$ , with shape function  $f(x) = x^{1/\alpha}$ ,  $\alpha > 0$ . Let A be the zooming algorithm on  $(L_f, X)$ . Then

$$R_{\mathcal{A}}(t) \le (c \, 2^{O(d^*)} \log^2 t) \ t^{1-1/(1+\alpha d)} \quad \text{for all } t,$$
(25)

where d is the c-covering dimension of the target set S.

*Proof.* Consider the pair  $(L_f, \mu)$ . Since the maximal reward is 1, by Theorem 4.8 it suffices to prove that for some  $c^* = c 2^{O(d^*)}$  the  $c^*$ -zooming dimension of this pair is at most  $\alpha d$ . Specifically, for each  $\delta > 0$  we need to cover the set  $S_{\delta} = \{u \in X : \Delta(u) \leq \delta\}$  with  $c^* \delta^{-\alpha d}$  sets of  $L_f$ -diameter at most  $\delta/16$ .

Indeed, since  $\Delta(u) = L_f(u, S)$ , for each  $u \in S_{\delta}$  we have  $L(u, S) \leq \delta^{\alpha}$ . Thus  $S_{\delta} \subset B(S, \delta^{\alpha})$ . Let  $\epsilon = 16^{-\alpha}$ . As in the proof of Theorem 4.11, we can show that  $S_{\delta}$  can be covered by  $c \epsilon^{-O(d^*)} \delta^{-\alpha d}$  sets of diameter  $\epsilon \delta^{\alpha}$ . Each of these sets has  $L_f$ -diameter at most  $f(\epsilon \delta^{\alpha})$ , which is at most  $\delta/16$ .

*Remarks.* This theorem includes Theorem 4.3 as a special case f(x) = x. Like the latter, this theorem is useful when the target set is a low-dimensional subset of the metric space.

We consider extensions to more general shape functions and to strategy sets which do not contain S:

• Suppose the shape function f satisfies the following constraints for some constants  $\alpha \ge \alpha^* > 0$ :

 $\forall x \in (0,1] \quad g(x) \geq x^{1/\alpha} \text{ and } g(x) \geq 2^{1/\alpha^*} g(\tfrac{x}{2}),$ 

where g(x) = f(x) - f(0). Then for  $\beta = 1 + 1_{\{f(0)>0\}}$  we have

 $R_A(t) \le (c \, 2^{O(\alpha^* d^*/\alpha)} \log^2 t) t^{1-1/(\beta+\alpha d)}$  for all t.

Consider the setting in which the strategy set Y is a proper subset of the metric space (L, X) and does not contain the target set S. If L(u<sup>\*</sup>, S) = 0 for some u<sup>\*</sup> ∈ Y then the guarantee (25) holds as is. In general, if we restrict the shape function to f(x) = c + x<sup>1/α</sup>, α ∈ (0, 1] then

$$R_{\mathcal{A}}(t) \le (c \, 2^{O(d^*)} \log^2 t) t^{1-1/(2+d)}$$
 for all  $t$ ,

where d is the c-covering dimension of the set  $S^* = B(S, r)$ , r = L(Y, S). Moreover, one can prove similar guarantees with  $d^*$  being the doubling dimension of an open neighborhood of  $S^*$ , rather than that of the entire metric space

#### 4.5 Heavy-tailed reward distributions

Consider the Lipschitz MAB problem and let  $X_n(v)$  be the reward from the *n*-th trial of strategy v. The current problem formulation restricts  $X_n(v)$  to support [0,1]. In this section we remove this restriction. In fact, it suffices to assume that  $X_n(v)$  is an independent random variable with mean  $\mu(v) \in [0,1]$  and a uniformly bounded bounded absolute third moment. Note that different trials of the same strategy  $\{X_n(v) : n \in \mathbb{N}\}$  do not need to be identically distributed.

**Theorem 4.14.** Consider the standard Lipschitz MAB problem. Let  $X_n(v)$  be the reward from the *n*-th trial of strategy v. Assume that each  $X_n(v)$  is an independent random variable with mean  $\mu(v) \in [0, 1]$  and furthermore that  $E\left[|X_n(v)|^3\right] < \rho$  for some constant  $\rho$ . Then there is an algorithm  $\mathcal{A}$  such that if for some c the problem instance has c-zooming dimension d then

$$R_{\mathcal{A}}(t) \le a(t) t^{1-1/(3d+6)}$$
 for all t, where  $a(t) = O(c\rho \log t)^{1/(3d+6)}$ . (26)

The proof relies on the non-uniform Berry-Esseen theorem (e.g. see [21] for a nice survey) which we use to obtain a tail inequality similar to Claim 2.2: for any  $\alpha > 0$ 

$$\Pr[|\mu_t(v) - \mu(v)| > r_t(v)] < O(t^{-3\alpha}), \text{ where } r_t(v) = \Theta(t^{\alpha}) / \sqrt{n_t(v)}.$$
(27)

However, this inequality gives much higher failure probability than Claim 2.2; in particular, we cannot take a union bound over all active strategies. Accordingly, we need a more refined version of Theorem 4.2 which is parameterized by the failure probability in (27). In the analysis, instead of the failure events when the phase is not clean (see Definition 2.1) we need to consider the  $\rho$ -failure events when the tail bound from (27) is violated by some strategy v such that  $\Delta(v) > \rho$ . Then using the technique from Section 2 we can upper-bound  $R_A(T)$  in terms of T, d,  $\rho$  and  $\alpha$  and choose the optimal values for  $\rho$  and  $\alpha$ .

## References

- [1] Rajeev Agrawal. The continuum-armed bandit problem. SIAM J. Control and Optimization, 33(6), 1995.
- [2] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. Machine Learning, 47(2-3):235–256, 2002. Preliminary version in 15th ICML, 1998.
- [3] Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. SIAM J. Comput., 32(1):48–77, 2002. Preliminary version in 36th IEEE FOCS, 1995.
- [4] Peter Auer, Ronald Ortner, and Csaba Szepesvári. Improved Rates for the Stochastic Continuum-Armed Bandit Problem. In 20th Conference on Learning Theory (COLT), pages 454–468, 2007.
- [5] Baruch Awerbuch and Robert Kleinberg. Online linear optimization and adaptive routing. *Journal of Computer and System Sciences*, 74(1):97–114, February 2008. Preliminary version appeared in STOC 2004.
- [6] Jeffrey Banks and Rangarajan Sundaram. Denumerable-armed bandits. *Econometrica*, 60(5):1071–1096, 1992.
- [7] Donald Berry and Bert Fristedt. Bandit problems: sequential allocation of experiments. Chapman&Hall, 1985.
- [8] Nicolò Cesa-Bianchi and Gábor Lugosi. Prediction, learning, and games. Cambridge University Press, 2006.
- [9] Eric Cope. Regret and convergence bounds for immediate-reward reinforcement learning with continuous action spaces, 2004. Unpublished manuscript.
- [10] Varsha Dani, Thomas Hayes, and Sham M. Kakade. The Price of Bandit Information for Online Optimization. In 21st Advances in Neural Information Processing Systems (NIPS), 2007.
- [11] Varsha Dani and Thomas P. Hayes. Robbing the bandit: Less regret in online geometric optimization against an adaptive adversary. In 16th ACM-SIAM Symp. on Discrete Algorithms (SODA), pages 937–943, 2006.
- [12] Abraham D. Flaxman, Adam Tauman Kalai, and H. Brendan McMahan. Online convex optimization in the bandit setting: Gradient descent without a gradient. In *Proceedings of the 16th ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 385–394, 2005.
- [13] J. C. Gittins and D. M. Jones. A dynamic allocation index for the sequential design of experiments. In J. Gani et al., editor, *Progress in Statistics*, pages 241–266. North-Holland, 1974.
- [14] Sham M. Kakade, Adam T. Kalai, and Katrina Ligett. Playing Games with Approximation Algorithms. In 39th ACM Symp. on Theory of Computing (STOC), 2007.
- [15] Robert Kleinberg. Nearly tight bounds for the continuum-armed bandit problem. In *18th Advances in Neural Information Processing Systems (NIPS)*, 2004. Full version appeared in the author's thesis (MIT, 1995).
- [16] Robert Kleinberg. Online Decision Problems with Large Strategy Sets. PhD thesis, MIT, Boston, MA, 2005.
- [17] Robert Kleinberg and Aleksandrs Slivkins. Multi-Armed Bandits in Metric Spaces: A dichotomy between logarithmic and square-root regret. Under submission, 2008.
- [18] Robert Kleinberg, Aleksandrs Slivkins, and Eli Upfal. Multi-Armed Bandits in Metric Spaces. In 40th ACM Symp. on Theory of Computing (STOC), pages 681–690, 2008.
- [19] H. Brendan McMahan and Avrim Blum. Online geometric optimization in the bandit setting against an adaptive adversary. In *Proceedings of the 17th Annual Conference on Learning Theory (COLT)*, volume 3120 of *Lecture Notes in Computer Science*, pages 109–123. Springer Verlag, 2004.
- [20] Michael Mitzenmacher and Eli Upfal. *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge University Press, 2005.
- [21] K. Neammanee. On the constant in the nonuniform version of the Berry-Esseen theorem. Intl. J. of Mathematics and Mathematical Sciences, 2005:12:1951–1967, 2005.
- [22] Sandeep Pandey, Deepak Agarwal, Deepayan Chakrabarti, and Vanja Josifovski. Bandits for Taxonomies: A Model-based Approach. In SIAM Intl. Conf. on Data Mining (SDM), 2007.
- [23] Sandeep Pandey, Deepayan Chakrabarti, and Deepak Agarwal. Multi-armed Bandit Problems with Dependent Arms. In 24th Intl. Conf. on Machine Learning (ICML), 2007.