

---

# Meta-Learning Adversarial Bandit Algorithms

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study online meta-learning with bandit feedback, with the goal of improving  
2 performance across multiple tasks if they are similar according to some natural  
3 similarity measure. As the first to target the adversarial online-within-online  
4 partial-information setting, we design meta-algorithms that combine outer learners  
5 to simultaneously tune the initialization and other hyperparameters of an inner  
6 learner for two important cases: multi-armed bandits (MAB) and bandit linear  
7 optimization (BLO). For MAB, the meta-learners initialize and set hyperparameters  
8 of the Tsallis-entropy generalization of Exp3, with the task-averaged regret  
9 improving if the entropy of the optima-in-hindsight is small. For BLO, we learn  
10 to initialize and tune online mirror descent (OMD) with self-concordant barrier  
11 regularizers, showing that task-averaged regret varies directly with an action space-  
12 dependent measure they induce. Our guarantees rely on proving that unregularized  
13 follow-the-leader combined with two levels of low-dimensional hyperparameter  
14 tuning is enough to learn a sequence of affine functions of non-Lipschitz and  
15 sometimes non-convex Bregman divergences bounding the regret of OMD.

## 16 1 Introduction

17 Learning-to-learn [50] is an important area of research that studies how to improve the performance  
18 of a learning algorithm by *meta-learning* its parameters—e.g. initializations, step-sizes, and/or  
19 representations—across many similar tasks. The goal is to encode information from previous  
20 tasks in order to achieve better performance on future ones. Meta-learning has seen a great deal  
21 of experimental work [24, 48], practical impact [21, 29], and theoretical effort [11, 18, 22, 44, 20].  
22 One important setting is online-within-online meta-learning [19, 31], where the learner performs a  
23 sequence of tasks, each of which has a sequence of rounds. Past work has studied the *full-information*  
24 setting, where the loss for every arm is revealed after each round. This assumption is not realistic in  
25 many applications, e.g. recommender systems and experimental design, where often partial or *bandit*  
26 feedback—only the loss of the action taken—is revealed. Such feedback can be *stochastic*, e.g. the  
27 losses are i.i.d. from some distribution, or *adversarial*, i.e. chosen by an adversary. We establish  
28 the first formal guarantees for online-within-online meta-learning with adversarial bandit feedback.

29 As with past full-information meta-learning results, our goal when faced with a sequence of bandit  
30 tasks will be to achieve low regret *on average* across them. Specifically, our task-averaged regret  
31 should (a) be no worse than that of algorithms for the single-task setting, e.g. if the tasks are not very  
32 similar, and should (b) be much better on tasks that are closely related, e.g. if the same small set of  
33 arms do well on all of them. We show that a natural way to achieve both is to initialize and tune online  
34 mirror descent (OMD), an algorithm associated with a strictly convex regularizer whose hyperparam-  
35 eters have a significant impact on performance. Our approach works because it can learn the best  
36 hyperparameters in hindsight across tasks, which will recover OMD’s worst-case optimal performance  
37 if the tasks are dissimilar but will take advantage of more optimistic settings if they are related. As  
38 generalized distances, the regularizers also induce interpretable measures of similarity between tasks.

39 **1.1 Main contributions**

40 We design a meta-algorithm (Algorithm 1) for learning variants of OMD—specifically those with  
 41 entropic or self-concordant regularizers—that are used for adversarial bandits. This meta-algorithm  
 42 combines three *full-information* algorithms—follow-the-leader (FTL), exponentially weighted online  
 43 optimization (EWO), and multiplicative weights (MW)—to set the initialization, step-size, and  
 44 regularizer-specific parameters, respectively. It works by optimizing a sequence of functions that each  
 45 *upper-bound* the regret of OMD on a single task (Theorem 2.1), resulting in (a) interesting notions  
 46 of task-similarity because these functions depend on generalized notions of distances (Bregman  
 47 divergences) and (b) adaptivity, i.e not needing to know how similar the tasks are beforehand.

48 Our first application is to OMD with the Tsallis regularizer [3], a relative of Exp3 [6] that is optimal for  
 49 adversarial MAB. We bound the task-averaged regret by the Tsallis entropy of the *estimated* optima-  
 50 in-hindsight (Corollary 3.1), which we further extend to that of the *true* optima by assuming a gap  
 51 between the best and second-best arms (Corollary 3.2). Both results are the first known consequences  
 52 of the online learnability of Bregman divergences that are *non-convex* in their second arguments [31],  
 53 while the latter is obtained by showing that the loss estimators of a modified algorithm identify the opti-  
 54 mal arm w.h.p. As an example, our average  $m$ -round regret across  $T$  tasks under the gap assumption is

$$o_T(\text{poly}(m)) + 2 \min_{\beta \in (0,1]} \sqrt{H_\beta d^\beta m / \beta} + o(\sqrt{m}) \quad (1)$$

55 where  $d$  is the number of actions and  $H_\beta$  is the Tsallis entropy [51, 3] of the distribution of the optimal  
 56 actions ( $\beta = 1$  recovers the Shannon entropy).<sup>1</sup> This entropy is low if all tasks are solved by the same  
 57 few arms, making it a natural task-similarity notion. For example, if  $s \ll d$  are always optimal then  
 58  $H_\beta = \mathcal{O}(s)$ , so using  $\beta = 1/\log d$  in (1) yields an asymptotic task-averaged regret of  $\mathcal{O}(\sqrt{sm \log d})$ ,  
 59 dropping fast terms. For  $s = \mathcal{O}_d(1)$  this beats the minimax optimal rate of  $\Theta(\sqrt{dm})$  [5]. On the other  
 60 hand, since  $H_{1/2} = \mathcal{O}(\sqrt{d})$ , the same bound recovers this rate in the worst-case of dissimilar tasks.

61 Lastly, we adapt our meta-algorithm to the adversarial BLO problem by setting the regularizer  
 62 to be a self-concordant barrier function, as in Abernethy et al. [2]. Our bounds yield notions of  
 63 task-similarity that depend on the constraints of the action space, e.g. over the sphere the measure  
 64 is the closeness of the average of the estimated optima to the sphere’s surface (Corollary 4.1). We  
 65 also instantiate BLO on the bandit shortest-path problem (Corollary D.2) [49, 30].

66 **1.2 Related work**

67 While we are the first to consider meta-learning under adversarial bandit feedback, many have  
 68 studied meta-learning in various *stochastic* bandit settings [9, 34, 46, 47, 35, 13, 15, 40, 10]. The  
 69 latter three study stochastic bandits under various task-generation assumptions, e.g. Azizi et al. [10]  
 70 is in a batch-within-online setting where the optimal arms are adversarial. In contrast, we make no  
 71 distributional assumptions either within or without.

72 A setting that bears some similarity to online-within-online bandits is that of switching bandits [6],  
 73 and more generally online learning with dynamic comparators [4, 27, 38, 7, 53]. In such problems,  
 74 instead of using a static best arm as the comparator we use a piecewise constant sequence of arms,  
 75 with a limited number of arm switches. The key difference between such work and ours is our  
 76 assumption that task-boundaries are known; this makes the other setting more general. However,  
 77 while e.g. Exp3.S [6] can indeed be applied to online meta-learning, its guarantees are worse than  
 78 if we just repeatedly apply a base-learner such as Exp3 on each task. Furthermore, these approaches  
 79 usually quantify difficulty by the number of switches, whereas we focus on task-similarity.

80 There has been a variety of work on full-information online-within-online meta-learning [32, 12],  
 81 including tuning OMD [31, 19]. Doing so for bandit algorithms has many additional challenges,  
 82 including (1) their inherent and high-variance stochasticity, (2) the use of non-Lipschitz and even  
 83 unbounded regularizers, and (3) the lack of access to task-optima in order to adapt to deterministic,  
 84 algorithm-independent task-similarity measures. Theoretically our analysis draws on the average  
 85 regret-upper-bound analysis (ARUBA) framework [31], which observes that OMD can be tuned by  
 86 targeting its upper bounds, which are affine functions of Bregman divergences, and provide online  
 87 learning tools for doing so. Our core structural result shows that the distance generating functions  $\psi_\theta$   
 88 of these Bregman divergences can be tuned without interfering with meta-learning the initialization

---

<sup>1</sup>We use  $\mathcal{O}_n(\cdot)$  (and  $o_n(\cdot)$ ) to denote terms with constant and (sub-constant) dependence on  $n$ .

89 and step-size; tuning  $\theta$  is critical for adapting to settings such as that of a small set of optimal arms  
 90 in MAB. Doing so depends on several refinements of the original approach, including bounding the  
 91 task-averaged-regret via the spectral norm of  $\nabla^2\psi_\theta$  and expressing the loss of the meta-comparator  
 92 using only  $\psi_\theta$ , rather than via its Bregman divergence as in prior work. Finally, applying our structural  
 93 result requires setting-specific analysis, e.g. to show regularity w.r.t.  $\theta$  or to obtain MAB guarantees  
 94 in terms of the entropy of the true optimal arms. The latter is especially difficult, as Khodak et al. [31]  
 95 define task-similarity via full information upper bounds, and involves applying tools from the best-  
 96 arm-identification literature [1] to show that a constrained variant of Exp3 finds the optimal arm w.h.p.

## 97 2 Learning the regularizers of bandit algorithms

98 We consider the problem of meta-learning over bandit tasks  $t = 1, \dots, T$  over some fixed set  $\mathcal{K} \subset \mathbb{R}^d$ ,  
 99 a (possibly improper) subset of which is the action space  $\mathcal{A}$ . On each round  $i = 1, \dots, m$  of task  $t$  we  
 100 play action  $\mathbf{x}_{t,i} \in \mathcal{A}$  and receive feedback  $\ell_{t,i}(\mathbf{x}_{t,i})$  for some function  $\ell_{t,i} : \mathcal{A} \mapsto [-1, 1]$ . Note that  
 101 all functions we consider will be linear and so we will also write  $\ell_{t,i}(\mathbf{x}) = \langle \ell_{t,i}, \mathbf{x} \rangle$ . Additionally, we  
 102 assume the adversary is *oblivious within-task*, i.e. it chooses losses  $\ell_{t,1}, \dots, \ell_{t,m}$  at time  $t$ . We will  
 103 also denote  $\mathbf{x}(a)$  to be the  $a$ -th element of the vector  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathcal{K}^\circ$  to be the interior of  $\mathcal{K}$ ,  $\partial\mathcal{K}$  its bound-  
 104 ary, and  $\Delta$  to be the simplex on  $d$  elements. Finally, note that all proofs can be found in the Appendix.

105 In online learning, the goal on a single task  $t$  is to play actions  $\mathbf{x}_{t,1}, \dots, \mathbf{x}_{t,m}$  that minimize  
 106 the regret  $\sum_{i=1}^m \ell_{t,i}(\mathbf{x}_{t,i}) - \ell_{t,i}(\hat{\mathbf{x}}_t)$ , where  $\hat{\mathbf{x}}_t \in \arg \min_{\mathbf{x} \in \mathcal{K}} \sum_{i=1}^m \ell_{t,i}(\mathbf{x})$ . Lifting this to the  
 107 meta-learning setting, our goal as in past work [31, 19] will be to minimize the **task-averaged**  
 108 **regret**:  $\frac{1}{T} \sum_{t=1}^T \sum_{i=1}^m \ell_{t,i}(\mathbf{x}_{t,i}) - \ell_{t,i}(\hat{\mathbf{x}}_t)$ . In particular, we want to use multi-task data to improve  
 109 average performance as the number of tasks  $T \rightarrow \infty$ . For example, we wish to attain a task-averaged  
 110 regret bound of the form  $o_T(\text{poly}(m)) + \tilde{\mathcal{O}}(V\sqrt{m}) + o(\sqrt{m})$ , where  $V \in \mathbb{R}_{\geq 0}$  is a measure  
 111 of task-similarity that is small if the tasks are similar but still yields the worst-case single-task  
 112 performance— $\mathcal{O}(\sqrt{dm})$  for MAB and  $\mathcal{O}(d\sqrt{m})$  for BLO—if they are not.

### 113 2.1 Online mirror descent as a base-learner

114 In meta-learning we are commonly interested in learning a within-task algorithm or **base-learner**,  
 115 a parameterized method that we run on each task  $t$ . A popular approach is to learn the initialization  
 116 and other parameters of a gradient-based method such as gradient descent [24, 43, 36]. If the task  
 117 optima are close, the best initialization should perform well after only a few steps on a new task.  
 118 We take a similar approach applied to online mirror descent, a generalization of gradient descent  
 119 to non-Euclidean geometries [14]. Given a strictly convex **regularizer**  $\psi : \mathcal{K}^\circ \mapsto \mathbb{R}$ , step-size  $\eta > 0$ ,  
 120 and initialization  $\mathbf{x}_{t,1} \in \mathcal{K}^\circ$ , OMD has the iteration

$$\mathbf{x}_{t,i+1} = \arg \min_{\mathbf{x} \in \mathcal{K}} B(\mathbf{x} | \mathbf{x}_{t,i}) + \eta \sum_{j \leq i} \langle \nabla \ell_{t,j}(\mathbf{x}_{t,j}), \mathbf{x} \rangle \quad (2)$$

121 where  $B(\mathbf{x} | \mathbf{y}) = \psi(\mathbf{x}) - \psi(\mathbf{y}) - \langle \nabla \psi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$  is the **Bregman divergence** of  $\psi$ . OMD recovers  
 122 online gradient descent when  $\psi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$ , in which case  $B(\mathbf{x} | \mathbf{y}) = \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$ ; another  
 123 example is exponentiated gradient, for which  $\psi(\mathbf{p}) = \langle \mathbf{p}, \log \mathbf{p} \rangle$  is the negative Shannon entropy  
 124 on probability vectors  $\mathbf{p} \in \Delta$  and  $B$  is the KL-divergence [45]. An important property of  $B$  is that  
 125 the sum over functions  $B(\mathbf{x}_t | \cdot)$  is minimized at the mean  $\bar{\mathbf{x}}$  of the points  $\mathbf{x}_1, \dots, \mathbf{x}_T$ .

126 OMD on **loss estimators**  $\hat{\ell}_{t,i}$  constructed via partial feedback forms an important class of bandit  
 127 methods [6, 2, 3]. Their regularizers  $\psi$  are often non-Lipschitz, e.g. the negative entropy, or even  
 128 unbounded, e.g. the log-barrier. Thus full-information results for tuning OMD, e.g. by Khodak  
 129 et al. [31] and Denevi et al. [19], do not suffice. We do adapt the former’s approach of online  
 130 learning a sequence  $U_t(\mathbf{x}, \eta, \theta)$  of affine functions of Bregman divergences from initializations  $\mathbf{x}$   
 131 to known points in  $\mathcal{K}$ . We are interested in them because the regret of OMD w.r.t. a comparator  
 132  $\mathbf{y}$  is bounded by  $B(\mathbf{y} | \mathbf{x})/\eta + \mathcal{O}(\eta m)$  [45, 25]. In our case the comparator is based on the estimated  
 133 optimum  $\hat{\mathbf{x}}_t \in \arg \min_{\mathbf{x} \in \mathcal{K}} \langle \hat{\ell}_t, \mathbf{x} \rangle$ , where  $\hat{\ell}_t = \sum_{i=1}^m \hat{\ell}_{t,i}$ , resulting from running OMD on task  
 134  $t$  using initialization  $\mathbf{x} \in \mathcal{K}$  and hyperparameters  $\eta$  and  $\theta$ , which we denote  $\text{OMD}_{\eta,\theta}(\mathbf{x})$ . Unlike  
 135 full-information meta-learning, we use a parameter  $\varepsilon > 0$  to constrain this optimum to lie in a subset  
 136  $\mathcal{K}_\varepsilon \subset \mathcal{K}^\circ$ . Formally, we fix a point  $\mathbf{x}_1 \in \mathcal{K}^\circ$  to be the “center”—e.g.  $\mathbf{x}_1 = \mathbf{1}_d/d$  when  $\mathcal{K}$  is the  
 137  $d$ -simplex  $\Delta$ —and define the projection  $\mathbf{c}_\varepsilon(\mathbf{x}) = \mathbf{x}_1 + \frac{\mathbf{x} - \mathbf{x}_1}{1 + \varepsilon}$  mapping from  $\mathcal{K}$  to  $\mathcal{K}_\varepsilon$ . For example,  
 138  $\mathbf{c}_{\frac{\varepsilon}{1-\varepsilon}}(\mathbf{x}) = (1 - \varepsilon)\mathbf{x} + \varepsilon\mathbf{1}_d/d$  on the simplex. This projection allows us to handle regularizers  $\psi$

---

**Algorithm 1:** Tunes  $\text{OMD}_{\eta,\theta}$  with regularizer  $\psi_\theta : \mathcal{K}^o \mapsto \mathbb{R}$  and step-size  $\eta > 0$ , which when run over loss estimators  $\hat{\ell}_{t,1}, \dots, \hat{\ell}_{t,m}$ , yielding task-optima  $\hat{\mathbf{x}}_t = \arg \min_{\mathbf{x} \in \mathcal{K}} \sum_{i=1}^m \langle \hat{\ell}_{t,i}, \mathbf{x} \rangle$ .

---

**Input:** compact set  $\mathcal{K} \subset \mathbb{R}^d$ , initialization  $\mathbf{x}_1 \in \mathcal{K}$ , ordered subset  $\Theta_k \subset \mathbb{R}$  also used to index interval bounds  $\underline{\eta}, \bar{\eta} \in \mathbb{R}_{\geq 0}^k$  and hyperparameters  $\alpha \in \mathbb{R}_{\geq 0}^k$ , scalar hyperparameters

$\rho > 0$  and  $\lambda \geq 0$ , learners  $\text{OMD}_{\eta,\theta} : \mathcal{K} \mapsto \mathbb{R}^d$ , projections  $\mathbf{c}_\theta : \mathcal{K} \mapsto \mathcal{K}_\theta$

**for**  $\theta \in \Theta_k$  **do**

$\mathbf{w}_1(\theta) \leftarrow 1$  and  $\eta_1(\theta) \leftarrow \frac{\underline{\eta}(\theta) + \bar{\eta}(\theta)}{2}$  // initialize MW and EWOO

**for** task  $t = 1, \dots, T$  **do**

    sample  $\theta_t$  from  $\Theta_k$  w.p.  $\propto \exp(\mathbf{w}_t)$  // sample from MW distribution

$\hat{\mathbf{x}}_t \leftarrow \text{OMD}_{\eta_t(\theta_t), \theta_t}(\mathbf{c}_{\theta_t}(\mathbf{x}_t))$  // run bandit OMD within-task

$\mathbf{x}_{t+1} \leftarrow \frac{1}{t} \sum_{s=1}^t \hat{\mathbf{x}}_s$  // FTL update of initialization

**for**  $\theta \in \Theta_k$  **do**

$\eta_{t+1}(\theta) \leftarrow \frac{\int_{\underline{\eta}(\theta)}^{\bar{\eta}(\theta)} v \exp(-\alpha(\theta) \sum_{s=1}^t U_s^{(\rho)}(\mathbf{x}_s, v, \theta)) dv}{\int_{\underline{\eta}(\theta)}^{\bar{\eta}(\theta)} \exp(-\alpha(\theta) \sum_{s=1}^t U_s^{(\rho)}(\mathbf{x}_s, v, \theta)) dv}$  // EWOO step-size update

$\mathbf{w}_{t+1}(\theta) \leftarrow \mathbf{w}_t(\theta) - \lambda U_t(\mathbf{x}_t, \eta_t(\theta), \theta)$  // MW update of tuning parameter

---

139 that diverge near the boundary, but also introduces  $\varepsilon$ -dependent error terms. In the BLO case it also  
140 forces us to tune  $\varepsilon$  itself, as initializing too close to the boundary leads to unbounded regret while  
141 initializing too far away does not take advantage of task-similarity. Thus the general upper bounds of  
142 interest are the following functions of the initialization  $\mathbf{x}$ , the step-size  $\eta > 0$ , and a third parameter  
143  $\theta$  that is either  $\beta$  or  $\varepsilon$ , depending on the setting (MAB or BLO):

$$U_t(\mathbf{x}, \eta, \theta) = \frac{B_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}_t) \|\mathbf{x})}{\eta} + \eta g(\theta)m + f(\theta)m \quad (3)$$

144 Here  $B_\theta$  is the Bregman divergence of  $\psi_\theta$  while  $g(\theta) \geq 1$  and  $f(\theta) \geq 0$  are tunable constants. We  
145 overload  $\theta$  to be either  $\beta$  or  $\varepsilon$  for notational simplicity, as we will not tune them simultaneously; if  $\theta =$   
146  $\beta$  (for MAB) then  $\mathbf{c}_\theta(\mathbf{x}) = \mathbf{x}_1 + \frac{\mathbf{x} - \mathbf{x}_1}{1 + \varepsilon}$  for fixed  $\varepsilon$ , while if  $\theta = \varepsilon$  (for BLO) then  $B_\theta$  is the Bregman  
147 divergence of a fixed  $\psi$ . The reason to optimize this sequence of upper bounds  $U_t$  is because they di-  
148 rectly bound the task-averaged regret while being no worse than the worst-case single-task regret. Fur-  
149 thermore, an average over Bregman divergences is minimized at the average  $\hat{\mathbf{x}} = \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{x}}_t$ , where  
150 it attains the value  $\hat{V}_\theta^2 = \frac{1}{T} \sum_{t=1}^T \psi_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}_t)) - \psi_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}))$  (c.f. Claim A.1). We will show that this  
151 quantity leads to intuitive and interpretable notions of task-similarity in all the applications we study.

## 152 2.2 A meta-algorithm for tuning bandit algorithms

153 To learn these functions  $U_t(\mathbf{x}, \eta, \theta)$ —and thus to meta-learn  $\text{OMD}_{\eta,\theta}(\mathbf{x})$ —our meta-algorithm sets  
154  $\mathbf{x}$  to be the projection  $\mathbf{c}_\theta$  of the mean of the estimated optima—i.e. follow-the-leader (FTL) over  
155 the Bregman divergences in (3)—while simultaneously setting  $\eta$  via EWOO and  $\theta$  via discrete  
156 multiplicative weights (MW). We choose FTL, EWOO, and MW because each is well-suited to the  
157 way  $U_t$  depends on  $\mathbf{x}$ ,  $\eta$ , and  $\theta$ , respectively. First, the only effect of  $\mathbf{x}$  on  $U_t$  is via the Bregman  
158 divergence  $B_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}_t) \|\mathbf{x})$ , over which FTL attains logarithmic regret [31]. For  $\eta$ ,  $U_t$  is exp-concave  
159 on  $\eta > 0$  so long as the first term is nonzero, but it is also non-Lipschitz; the EWOO algorithm is  
160 one of the few methods with logarithmic regret on exp-concave losses without a dependence on the  
161 Lipschitz constant [26], and we ensure the first term is nonzero by *regularizing* the upper bounds as  
162 follows for some  $\rho > 0$  and  $D_\theta^2 = \max_{\mathbf{x}, \mathbf{y} \in \mathcal{K}_\theta} B_\theta(\mathbf{x} \|\mathbf{y})$ :

$$U_t^{(\rho)}(\mathbf{x}, \eta, \theta) = \frac{B_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}_t) \|\mathbf{x}) + \rho^2 D_\theta^2}{\eta} + \eta g(\theta)m + f(\theta)m \quad (4)$$

163 Note that this function is fully defined after obtaining  $\hat{\mathbf{x}}_t$  by running OMD on task  $t$ , which allows us  
164 to use full-information MW to tune  $\theta$  across the grid  $\Theta_k$ . Showing low regret w.r.t. any  $\theta \in \Theta \supset \Theta_k$   
165 then just requires sufficiently large  $k$  and Lipschitzness of  $U_t$  w.r.t.  $\theta$ . Combining all three algorithms  
166 together thus yields the guarantee in Theorem 2.1, which is our main structural result. It implies  
167 a generic approach for obtaining meta-learning algorithms by (1) bounding the task-averaged

168 regret by an average of functions of the form  $U_t$ , (2) applying the theorem to obtain a new bound  
 169  $o_T(1) + \min_{\theta, \eta} \frac{\hat{V}_\theta^2}{\eta} + \eta g(\theta)m + f(\theta)m$ , and (3) bounding the estimated task-similarity  $\hat{V}_\theta^2$  by an  
 170 interpretable quantity. Crucially, since we can choose any  $\eta > 0$ , the asymptotic regret is always  
 171 as good as the worst-case guarantee for running the base-learner separately on each task.

172 **Theorem 2.1** (c.f. Thm. A.1). *Suppose  $\mathbf{x}_1 = \arg \min_{\mathbf{x} \in \mathcal{K}} \psi_\theta(\mathbf{x}) \forall \theta$  and let  $D$ ,  $M$ ,  $F$ , and  $S$  be  
 173 maxima over  $\theta$  of  $D_\theta$ ,  $D_\theta \sqrt{g(\theta)m}$ ,  $f(\theta)$ , and  $\|\nabla^2 \psi_\theta\|_2$ , respectively. For each  $\rho \in (0, 1)$  we can set  
 174  $\underline{\eta}$ ,  $\bar{\eta}$ ,  $\alpha$ , and  $\lambda$  s.t. the expected average of the losses  $U_t(\mathbf{c}_{\theta_t}(\mathbf{x}_t), \eta_t(\theta_t), \theta_t)$  of Algorithm 1 is at most*

$$\min_{\theta \in \Theta, \eta > 0} \frac{\mathbb{E} \hat{V}_\theta^2}{\eta} + \eta g(\theta)m + f(\theta)m + \tilde{O} \left( \frac{M + Fm}{\sqrt{T}} + \frac{L_\eta}{k} + \frac{M}{\rho^2 T} + \min \left\{ \frac{\rho^2 D^2}{\eta}, \rho M \right\} + \frac{S}{\eta T} \right) \quad (5)$$

175 Here  $\hat{V}_\theta^2 = \frac{1}{T} \sum_{t=1}^T \psi_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}_t)) - \psi_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}))$  and  $L_\eta$  bounds the Lipschitz constant w.r.t.  $\theta$  at  
 176  $\hat{V}_\theta^2/\eta + \eta g(\theta)m + f(\theta)m$ . The same bound plus  $(M/\rho + Fm)\sqrt{\frac{1}{T} \log \frac{1}{\delta}}$  holds w.p.  $\geq 1 - \delta$ .

177 We keep details of the dependence on  $S$  and other constants as they are important in applying this  
 178 result, but in most cases setting  $\rho = \frac{1}{\sqrt[3]{T}}$  yields  $\tilde{O}(T^{\frac{3}{4}})$  regret. While a slow rate, the losses  $U_t$  are  
 179 non-Lipschitz and non-convex in-general, and learning them allows us to tune  $\theta$  over user-specified  
 180 intervals and  $\eta$  over all positive numbers, which will be crucial later. At the same time, this tuning  
 181 is what leads to the slow rate, as without tuning ( $k = 1$ ,  $L_\eta = 0$ ) the same  $\rho$  yields  $\tilde{O}(\sqrt{T})$  regret.  
 182 Lastly, while we focus on learning guarantees, we note that Algorithm 1 is reasonably efficient,  
 183 requiring a  $2k$  single-dimensional integrals per task; this is discussed in more detail in Section A.3.

### 184 3 Multi-armed bandits

185 We now turn to our first application: the multi-armed bandit problem, where at each round  $i$  of task  
 186  $t$  we take action  $a_{t,i} \in [d]$  and observe loss  $\ell_{t,i}(a_{t,i}) \in [0, 1]$ . As we are sampling actions from  
 187 distributions  $\mathbf{x} \in \mathcal{K} = \Delta$  on the  $k$ -simplex, the inner product  $\langle \ell_{t,i}, \mathbf{x}_{t,i} \rangle$  is the expected loss and the  
 188 optimal arm  $\hat{a}_t$  on task  $t$  can be encoded as a vector  $\hat{\mathbf{x}}_t$  s.t.  $\hat{\mathbf{x}}_t(a) = 1_{a=\hat{a}_t}$ .

189 We use as a base-learner a generalization of Exp3 that uses the negative Tsallis entropy  
 190  $\psi_\beta(\mathbf{p}) = \frac{1 - \sum_{a=1}^d \mathbf{p}^\beta(a)}{1 - \beta}$  for some  $\beta \in (0, 1]$  as the regularizer; this improves regret from Exp3's  
 191  $\mathcal{O}(\sqrt{dm \log d})$  to the optimal  $\mathcal{O}(\sqrt{dm})$  [3]. Note that  $-\psi_\beta$  is the Shannon entropy in the limit  
 192  $\beta \rightarrow 1$  and its Bregman divergence  $B_\beta(\mathbf{x}||\cdot)$  is non-convex in the second argument. As the  
 193 Tsallis entropy is non-Lipschitz at the simplex boundary, which is where the estimated and  
 194 true optima  $\hat{\mathbf{x}}_t$  and  $\tilde{\mathbf{x}}_t$  lie, we will project them using  $\mathbf{c}_{\frac{\varepsilon}{1-\varepsilon}}(\mathbf{x}) = (1 - \varepsilon)\mathbf{x} + \varepsilon \mathbf{1}_d/d$  to the set  
 195  $\mathcal{K}_{\frac{\varepsilon}{1-\varepsilon}} = \{\mathbf{x} \in \Delta : \min_a \mathbf{x}(a) \geq \varepsilon/d\}$ . We denote the resulting vectors using the superscript  $(\varepsilon)$ ,  
 196 e.g.  $\hat{\mathbf{x}}_t^{(\varepsilon)} = \mathbf{c}_{\frac{\varepsilon}{1-\varepsilon}}(\hat{\mathbf{x}}_t)$ , and also use  $\Delta^{(\varepsilon)} = \mathcal{K}_{\frac{\varepsilon}{1-\varepsilon}}$  to denote the constrained simplex. For MAB we  
 197 also study two base-learners: (1) **implicit exploration** and (2) **guaranteed exploration**. The former  
 198 uses low-variance loss *under*-estimators  $\hat{\ell}_{t,i}(a) = \frac{\ell_{t,i}(a) 1_{a_{t,i}=a}}{\mathbf{x}_{t,i}(a) + \gamma}$  for  $\gamma > 0$ , where  $\mathbf{x}_{t,i}(a)$  is the  
 199 probability of sampling  $a$  on task  $t$  round  $i$ , to enable high probability bounds [42]. On the other hand,  
 200 **guaranteed exploration** uses unbiased loss estimators (i.e.  $\gamma = 0$ ) but constrains the action space  
 201 to  $\Delta^{(\varepsilon)}$ , which we will use to adapt to a task-similarity determined by the *true* optima-in-hindsight.

#### 202 3.1 Adapting to low estimated entropy with high probability using implicit exploration

203 In our first setting, the base-learner runs  $\text{OMD}_{\eta_t, \beta_t}(\mathbf{x}_{t,1})$  on  $\gamma$ -regularized estimators with Tsallis  
 204 regularizer  $\psi_{\beta_t}$ , step-size  $\eta_t$ , and initialization  $\mathbf{x}_{t,1} \in \Delta^{(\varepsilon)}$ . Standard OMD analysis combined with  
 205 implicit exploration analysis [42] shows (43) that the task-averaged regret is bounded w.h.p. by

$$(\varepsilon + \gamma d)m + \tilde{O} \left( \frac{\sqrt{d}}{\gamma T} \right) + \frac{1}{T} \sum_{t=1}^T \frac{B_{\beta_t}(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x}_{t,1})}{\eta_t} + \frac{\eta_t d^{\beta_t} m}{\beta_t} \quad (6)$$

206 The summands have the desired form of  $U_t(\mathbf{x}_{t,1}, \eta_t, \beta_t)$ , so we can apply Theorem 2.1 to bound  
 207 their average by

$$\min_{\beta \in [\underline{\beta}, \bar{\beta}], \eta > 0} \frac{\hat{V}_\beta^2}{\eta} + \frac{\eta d^\beta m}{\beta} + \tilde{O} \left( \frac{L_\eta}{k} + \frac{\left(\frac{d}{\varepsilon}\right)^{2-\beta}}{\eta T} + \left( \rho + \frac{1}{\rho \sqrt{T}} + \frac{1}{\rho^2 T} \right) d \sqrt{m} \right) \quad (7)$$

208 where  $\hat{V}_\beta^2 = \frac{1}{T} \sum_{t=1}^T \psi_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)}) - \psi_\beta(\hat{\mathbf{x}}^{(\varepsilon)})$  is the average difference in Tsallis entropies between  
 209 the ( $\varepsilon$ -constrained) estimated optima  $\hat{\mathbf{x}}_t$  and their empirical distribution  $\hat{\mathbf{x}} = \frac{1}{T} \sum_{t=1}^T \hat{\mathbf{x}}_t$ , while  $L_\eta$   
 210 is the Lipschitz constant of  $\frac{\hat{V}_\beta^2}{\eta} + \frac{\eta d^\beta m}{\beta}$  w.r.t.  $\beta \in [\underline{\beta}, \bar{\beta}]$ . The specific instantiation of Algorithm 1  
 211 that (7) holds for is to do the following at each time  $t$ :

1. sample  $\beta_t$  via the MW distribution  $\propto \exp(\mathbf{w}_t)$  over the discretization  $\Theta_k$  of  $[\underline{\beta}, \bar{\beta}] \subset [0, 1]$
2. run  $\text{OMD}_{\eta_t, \beta_t}$  using the initialization  $\mathbf{x}_{t,1} = \frac{1}{t-1} \sum_{s<t} \hat{\mathbf{x}}_s^{(\varepsilon)} = \frac{\varepsilon}{d} \mathbf{1}_d + \frac{1-\varepsilon}{t-1} \sum_{s<t} \hat{\mathbf{x}}_s$  (FTL)
3. update EWOO at each  $\beta \in \Theta_k$  with loss  $\frac{B_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x}_{t,1}) + \rho^2 D_\beta^2}{\eta} + \frac{\eta d^\beta m}{\beta}$ , where  $D_\beta^2 = \frac{d^{1-\beta} - 1}{1-\beta}$
4. update  $\mathbf{p}_{t+1}$  using multiplicative weights with expert losses  $\frac{B_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x}_{t,1})}{\eta} + \frac{\eta d^\beta m}{\beta}$

(8)

212 The final guarantee for this procedure, given in full in Theorem B.1, follows by two properties of  
 213 the Tsallis entropy  $-\psi_\beta$ : (1) its Lipschitzness w.r.t.  $\beta \in [0, 1]$  (c.f. Lem B.1) and (2) the fact that  
 214  $\hat{V}_\beta^2$  is bounded by the entropy  $\hat{H}_\beta = -\psi_\beta(\hat{\mathbf{x}})$  of the empirical distribution of estimated optima (c.f.  
 215 Lem B.2), which yields our first notion of task-similarity: *multi-armed bandit tasks are similar if*  
 216 *the empirical distribution of their (estimated) optimal arms has low entropy.*

217 We exemplify the implications of Theorem B.1 in Corollary 3.1, where we consider three regimes  
 218 of the lower bound  $\underline{\beta}$  on the entropy parameter:  $\underline{\beta} = 1$ , i.e. always using Exp3;  $\underline{\beta} = 1/2$ , which  
 219 corresponds to the optimal worst-case setting [3]; and  $\underline{\beta} = 1/\log d$ , below which the OMD  
 220 regret-upper-bound always worsens (and so it does not make sense to try  $\beta < 1/\log d$ ).

221 **Corollary 3.1** (c.f. Cors. B.1, B.2, and B.3). *Suppose  $\bar{\beta} = 1$  and we set the initialization, step-size,*  
 222 *and entropy parameter of Tsallis OMD with implicit exploration via Algorithm 1 as in Theorem B.1.*

223 1. If  $\underline{\beta} = 1$  and  $T \geq \frac{d^2}{m}$  we can ensure  $\frac{1}{T} \sum_{t=1}^T \sum_{i=1}^m \ell_{t,i}(\mathbf{x}_{t,i}) - \ell_{t,i}(\hat{\mathbf{x}}_t) \leq 2\sqrt{\hat{H}_1 dm} + \tilde{\mathcal{O}}\left(\frac{d^{\frac{3}{2}} m^{\frac{3}{2}}}{\sqrt[3]{T}}\right)$  w.h.p.

224 2. If  $\underline{\beta} = \frac{1}{2}$  and  $T \geq \frac{d^{5/2}}{m}$  we can set  $k = \lceil \sqrt[4]{d\sqrt{T}} \rceil$  and ensure w.h.p. that task-averaged regret is

$$\min_{\beta \in [\frac{1}{2}, 1]} 2\sqrt{\hat{H}_\beta d^\beta m / \beta} + \tilde{\mathcal{O}}\left(\frac{d^{5/7} m^{5/7}}{T^{2/7}} + \frac{d\sqrt{m}}{\sqrt[4]{T}}\right) \quad (9)$$

225 3. If  $\underline{\beta} = \frac{1}{\log d}$  and  $T \geq \frac{d^3}{m}$  we can set  $k = \lceil \sqrt[4]{d\sqrt{T}} \rceil$  and ensure w.h.p. that task-averaged regret is

$$\min_{\beta \in (0, 1]} 2\sqrt{\hat{H}_\beta d^\beta m / \beta} + \tilde{\mathcal{O}}\left(\frac{d^{3/4} m^{3/4} + d\sqrt{m}}{\sqrt[4]{T}}\right) \quad (10)$$

226 In all three settings, as  $T \rightarrow \infty$  the regret scales directly with the entropy of the estimated  
 227 optima-in-hindsight, which is small if most tasks are estimated to be solved by one of a few arms  
 228 and large if all arms are used roughly equally. Corollary 3.1 demonstrates the importance of tuning  
 229  $\beta$ : even if tasks are dissimilar, we asymptotically recover the worst-case optimal guarantee  $\mathcal{O}(\sqrt{dm})$   
 230 in cases two and three because the entropy is at most  $\frac{d^{1-\beta}}{1-\beta}$ . On the other hand, if a constant  $s \ll d$   
 231 actions are always minimizers, i.e. the empirical distribution  $\hat{\mathbf{x}}$  is  $s$ -sparse, then the last bound (10)  
 232 implies that Algorithm 1 can achieve task-averaged regret  $o_T(md) + \mathcal{O}(\sqrt{sm \log d})$ . At the same  
 233 time, this tuning is costly, with the last two results having an extra  $\tilde{\mathcal{O}}\left(\frac{d\sqrt{m}}{\sqrt[4]{T}}\right)$  term because of it.  
 234 Furthermore, the bound of  $\underline{\beta} = \frac{1}{2}$  has a slightly better dependence on  $d$ ,  $m$ , and  $T$  compared to that  
 235 of  $\underline{\beta} = \frac{1}{\log d}$  due to the  $\left(\frac{d}{\varepsilon}\right)^{2-\beta}$  term in the bound (7) returned for MAB by our structural result.

236 We can compare the  $s$ -sparse result to Azizi et al. [10], who achieve task-averaged regret  
 237  $\tilde{\mathcal{O}}(m/\sqrt[3]{T} + \sqrt{sm \log T})$  for *stochastic* MAB. Despite our adversarial setting and no stipulations on  
 238 how tasks are related, our bounds are asymptotically comparable if the estimated and true optima are  
 239 roughly equivalent (ignoring their  $\mathcal{O}(\sqrt{\log T})$ -factor), as we also have  $\tilde{\mathcal{O}}(\sqrt{sm})$  average regret as  
 240  $T \rightarrow \infty$ . Their rate in the number of tasks is better, but at a cost of runtime exponential in  $s$ . Apart  
 241 from generality, we believe a great strength of our results is their adaptiveness; unlike Azizi et al.  
 242 [10], we do not need to know how many optimal arms there are to adapt to there being few of them.

243 **3.2 Adapting to the entropy of the true optima-in-hindsight using guaranteed exploration**

244 While the entropy of estimated optima-in-hindsight may be useful in some cases where we wish to  
 245 actually *compute* the task-similarity, it is otherwise generally more desirable to adapt to an intrinsic  
 246 and algorithm-independent measure, e.g. the entropy of the *true* optima-in-hindsight. However, doing  
 247 so is difficult without further assumptions, as the optima are both hard to identify and the measure  
 248 itself may not be fully defined in case of ties. Thus in this section we focus on the setting where we  
 249 have a nonzero performance gap  $\Delta > 0$  between the best and second-best arms:

250 **Assumption 3.1.** *For some  $\Delta > 0$  and all tasks  $t \in [T]$ ,  $\frac{1}{m} \sum_{i=1}^m \ell_{t,i}(a) - \ell_{t,i}(\hat{a}_t) \geq \Delta \forall a \neq \hat{a}_t$ .*

251 This assumption is common in the best-arm identification literature [28, 1], which we adapt to show  
 252 that the estimated optimal arms match the true optima, and thus so do their entropies. To do so, we  
 253 switch to *unbiased* loss estimators, i.e.  $\gamma = 0$ , and control their variance by lower-bounding the  
 254 probability of selecting an arm to be at least  $\frac{\varepsilon}{d}$ ; this can alternatively be expressed as running OMD  
 255 using the regularizer  $\psi_\beta + I_{\Delta(\varepsilon)}$ , where for any  $\mathcal{C} \subset \mathbb{R}^d$  the function  $I_{\mathcal{C}}(\mathbf{x}) = 0$  if  $\mathbf{x} \in \mathcal{C}$  and  $\infty$   
 256 otherwise. Guaranteed exploration allows us extend the analysis of Abbasi-Yadkori et al. [1] to show  
 257 that the estimated arm is optimal w.h.p.:

258 **Lemma 3.1** (c.f. Lem C.1). *Suppose for  $\varepsilon > 0$  and any  $\beta \in (0, 1]$  we run OMD on task  $t \in [T]$  with  
 259 regularizer  $\psi_\beta + I_{\Delta(\varepsilon)}$ . If  $m = \tilde{\Omega}(\frac{d}{\varepsilon \Delta^2})$  then  $\hat{\mathbf{x}}_t = \hat{\mathbf{x}}_t$  w.p.  $\geq 1 - d \exp(-\Omega(\varepsilon \Delta^2 m/d))$ .*

260 However, the constraint that the probabilities are at least  $\frac{\varepsilon}{d}$  does lead to  $\varepsilon m$  additional error on each  
 261 task, with the upper bound on the task-averaged expected regret becoming

$$\mathbb{E} \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^m \ell_{t,i}(a_{t,i}) - \ell_{t,i}(\hat{a}_t) \leq \varepsilon m + \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{E} B_{\beta_t}(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x}_{t,1})}{\eta_t} + \frac{\eta_t d^{\beta_t} m}{\beta_t} \quad (11)$$

262 Moreover, we will no longer set  $\varepsilon = o_T(1)$ , as this would require  $m$  to be *increasing* in  $T$  for the best-  
 263 arm identification result of Lemma C.1 to hold. Thus, unlike in the previous section, our results will  
 264 contain “fast” terms—terms in the task-averaged regret that are  $o(\sqrt{m})$  but not decreasing in  $T$  nor af-  
 265 fected by the task-similarity. They still allow us to circumvent the  $\Omega(\sqrt{dm})$  MAB lower bound if tasks  
 266 are similar, but the task-averaged regret will not converge to zero as  $T \rightarrow \infty$  if the tasks are identical.

267 Nevertheless, the tuning-dependent component of the upper bounds in (11) has the appropriate  
 268 form for our structural result—in fact we can use the same meta-algorithm (8) as for implicit  
 269 exploration—and so we can again apply Theorem 2.1 to get a bound on the task-averaged regret  
 270 in terms of the average difference  $\hat{V}_\beta^2 = \frac{1}{T} \sum_{t=1}^T \psi_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)}) - \psi_\beta(\hat{\mathbf{x}}^{(\varepsilon)})$  of the entropies of the  
 271  $\varepsilon$ -constrained estimated task-optima  $\hat{\mathbf{x}}_t^{(\varepsilon)}$  and their mean  $\hat{\mathbf{x}}^{(\varepsilon)}$ . The easiest way to apply Lemma C.1  
 272 to bound  $\hat{V}_\beta^2$  in terms of  $H_\beta = \frac{1}{T} \sum_{t=1}^T \psi_\beta(\hat{\mathbf{x}}_t) - \psi_\beta(\hat{\mathbf{x}})$  is via union bound on all  $T$  tasks to show  
 273 that  $\hat{\mathbf{x}}_t = \hat{\mathbf{x}}_t \forall t$  w.p.  $\geq 1 - dT \exp(-\Omega(\varepsilon \Delta^2 m/d))$ ; however, setting a constant failure probability  
 274 leads to  $m$  growing, albeit only logarithmically, in  $T$ . Instead, by analyzing the worst-case best-arm  
 275 identification probabilities, we show in Lemma C.2 that the expectation of  $\hat{V}_\beta^2$  is bounded by  
 276  $H_\beta + 3\beta \frac{(d/\varepsilon)^{1-\beta} - 1}{1-\beta} \exp\left(-\frac{3\varepsilon \Delta^2 m}{28d}\right)$  without resorting to  $m = \omega_T(1)$ . Assuming  $m \geq \frac{75d}{\varepsilon \Delta^2} \log \frac{d}{\varepsilon \Delta^2}$   
 277 is enough (68) to bound the second term by  $\frac{56}{dm}$ . Then the final result (c.f. Thm. C.1) bounds the  
 278 expected task-averaged regret as follows (ignoring terms that become  $o_T(1)$  after setting  $\rho$  and  $k$ ):

$$\varepsilon m + \min_{\beta \in [\underline{\beta}, \bar{\beta}], \eta > 0} \frac{h_\beta(\Delta)}{\eta} + \frac{\eta d^\beta m}{\beta} \quad \text{for } h_\beta(\Delta) = \begin{cases} H_\beta + \frac{56}{md} & \text{if } m \geq \frac{75d}{\varepsilon \Delta^2} \log \frac{d}{\varepsilon \Delta^2} \\ \frac{d^{1-\beta} - 1}{1-\beta} & \text{otherwise} \end{cases} \quad (12)$$

279 If the gap  $\Delta$  is known and sufficiently large, then we can set  $\varepsilon = \Theta(\frac{d}{\Delta^2 m})$  to obtain an asymptotic  
 280 task-averaged regret that scales only with the entropy  $H_\beta$  and a fast term that is logarithmic in  $m$ :

281 **Corollary 3.2** (c.f. Cor. C.3). *Suppose we set the initialization, step-size, and entropy parameter of  
 282 Tsallis OMD with guaranteed exploration via Algorithm 1 as specified in Theorem C.1. If  $[\underline{\beta}, \bar{\beta}] =$   
 283  $[\frac{1}{\log d}, 1]$  and  $m \geq \frac{75d}{\Delta^2} \log \frac{d}{\Delta^2}$ , then setting  $\varepsilon = \tilde{\Theta}(\frac{d}{\Delta^2 m})$ ,  $\rho = \frac{1}{\sqrt[3]{d} \sqrt[3]{mT}}$ , and  $k = \lceil \sqrt[3]{d^2 m T} \rceil$   
 284 ensures that the expected task-averaged regret is at most*

$$\min_{\beta \in (0,1]} 2\sqrt{H_\beta d^\beta m / \beta} + \tilde{\mathcal{O}} \left( \frac{d}{\Delta^2} + \frac{d^{\frac{4}{3}} m^{\frac{2}{3}}}{\sqrt[3]{T}} + \frac{d^{\frac{5}{3}} m^{\frac{5}{6}}}{T^{\frac{2}{3}}} + \frac{d \Delta^4 m^3}{T} \right) \quad (13)$$

285 Knowing the gap  $\Delta$  is a strong assumption, as ideally we could set  $\varepsilon$  without it. Note that if  $\varepsilon = \Omega(\frac{1}{m^p})$   
 286 for some  $p \in (0, 1)$  then the condition  $m \geq \frac{75d}{\varepsilon\Delta^2} \log \frac{d}{\varepsilon\Delta^2}$  only fails if  $m \leq \text{poly}(\frac{1}{\Delta})$ , i.e. for gap  
 287 decreasing in  $m$ . We can use this together with the fact that minimizing over  $\eta$  and  $\beta$  in our bound  
 288 allows us to replace them with any value, even a gap-dependent one, to derive a gap-independent  
 289 setting of  $\varepsilon$  that ensures a task-similarity-adaptive bound when  $\Delta$  is not too small and falls back to  
 290 the worst-case optimal guarantee otherwise. Specifically, for indicator  $\iota_\Delta = 1_{m \geq \frac{75d}{\varepsilon\Delta^2} \log \frac{d}{\varepsilon\Delta^2}}$ , setting

291  $\eta = \Theta\left(\sqrt{\frac{h_\beta(\Delta)}{d^\beta m/\beta}}\right)$  in (12) and using  $\beta = \frac{1}{2}$  if the condition  $\iota_\Delta$  fails yields asymptotic regret at most

$$\varepsilon m + \min_{\beta \in (0,1]} \mathcal{O}\left(\iota_\Delta \sqrt{\frac{H_\beta d^\beta m}{\beta}} + (1 - \iota_\Delta) \sqrt{dm}\right) \leq \varepsilon m + \tilde{\mathcal{O}}\left(\min\left\{\min_{\beta \in (0,1]} \sqrt{\frac{H_\beta d^\beta m}{\beta}} + \frac{d}{\Delta\sqrt{\varepsilon}}, \sqrt{dm}\right\}\right) \quad (14)$$

292 Thus setting  $\varepsilon = \Theta(\sqrt{d}/m^{\frac{3}{5}})$  yields the desired dependence on the entropy  $H_\beta$  and a fast term in  $m$ :

293 **Corollary 3.3** (c.f. Cor. C.4). *In the setting of Corollary 3.2 but with  $m = \Omega(d^{\frac{3}{4}})$  and unknown  $\Delta$ ,*  
 294 *using  $\varepsilon = \Theta(\sqrt{d}/m^{\frac{3}{5}})$  ensures expected task-averaged regret at most*

$$\min\left\{\min_{\beta \in (0,1]} 2\sqrt{H_\beta d^\beta m/\beta} + \tilde{\mathcal{O}}\left(\frac{d^{\frac{3}{4}} \sqrt[3]{m}}{\Delta}\right), 8\sqrt{dm}\right\} + \tilde{\mathcal{O}}\left(\frac{d^{\frac{4}{3}} m^{\frac{2}{3}}}{\sqrt[3]{T}} + \frac{d^{\frac{5}{3}} m^{\frac{5}{6}}}{T^{\frac{2}{3}}} + \frac{d^2 m^{\frac{7}{3}}}{T}\right) \quad (15)$$

295 While not logarithmic, the gap-dependent term is still  $o(\sqrt{m})$ , and moreover the asymptotic regret is  
 296 no worse than the worst-case optimal  $\mathcal{O}(\sqrt{dm})$ . Note that the latter is only needed if  $\Delta = o(1/\sqrt[5]{m})$ .

297 The main improvement in this section is in using the entropy of the true optima, which can be much  
 298 smaller than that of the estimated optima if there are a few good arms but large noise. Our use of  
 299 the gap assumption for this seems difficult to avoid for this notion of task-similarity. We can also  
 300 compare to Corollary 3.1 (10), which did not require  $\Delta > 0$  and had no fast terms but had a worse  
 301 rate in  $T$ ; in contrast, the  $\mathcal{O}(\frac{1}{\sqrt[3]{T}})$  rates above match that of the closest stochastic bandit result [10].

302 As before, for  $s \ll d$  “good” arms we obtain  $\mathcal{O}(\sqrt{sm \log d})$  asymptotic regret, assuming the gap  
 303 is not too small. Finally, we can also compare to the classic shifting regret bound for Exp3.S [6],  
 304 which translated to task-averaged regret is  $\mathcal{O}(\sqrt{dm \log(dmT)})$ . This is worse than even running  
 305 OMD separately on each task, albeit under weaker assumptions (not knowing task boundaries). It  
 306 also cannot take advantage of repeated optimal arms, e.g. the case of  $s \ll d$  good arms.

## 307 4 Bandit linear optimization

308 Our last application is bandit linear optimization, in which at task  $t$  round  $i$  we play  $\mathbf{x}_{t,i} \in \mathcal{K}$  in  
 309 some convex  $\mathcal{K} \subset \mathbb{R}^d$  and observe loss  $\langle \ell_{t,i}, \mathbf{x}_{t,i} \rangle \in [-1, 1]$ . We will again use a variant of mirror  
 310 descent, using a **self-concordant barrier** for  $\psi$  and the specialized loss estimators of Abernethy  
 311 et al. [2, Alg. 1]. More information on such regularizers can be found in the literature on interior  
 312 point methods [41]. We pick this class of algorithms because of their optimal dependence on the  
 313 number of rounds and their applicability to any convex domain  $\mathcal{K}$  via specific barriers  $\psi$ , which will  
 314 yield interesting notions of task-similarity. Our ability to handle non-smooth regularizers via the  
 315 structural result (Thm. 2.1) is even more important here, as barriers are infinite at the boundaries.  
 316 Indeed, we will *not* learn a  $\beta$  parameterizing the regularizer and instead focus on tuning a boundary  
 317 offset  $\varepsilon > 0$ . Here we make use of notation from Section 2, where  $\mathbf{c}_\varepsilon$  maps points in  $\mathcal{K}$  to a subset  
 318  $\mathcal{K}_\varepsilon$  defined by the Minkowski function (c.f. Def. D.1) centered at  $\mathbf{x}_1 = \arg \min_{\mathbf{x} \in \mathcal{K}} \psi(\mathbf{x})$ .

319 From Abernethy et al. [2] we have an upper bound on the expected task-averaged regret of their  
 320 algorithm run from initializations  $\mathbf{x}_{t,1} \in \mathcal{K}^\circ$  with step-sizes  $\eta_t > 0$  and offsets  $\varepsilon_t > 0$ :

$$\mathbb{E} \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^m \langle \ell_{t,i}, \mathbf{x}_{t,i} - \hat{\mathbf{x}}_t \rangle \leq \frac{1}{T} \sum_{t=1}^T \frac{\mathbb{E} B(\mathbf{c}_{\varepsilon_t}(\hat{\mathbf{x}}_t) | | \mathbf{x}_{t,1})}{\eta_t} + (32\eta_t d^2 + \varepsilon_t) m \quad (16)$$

321 We can show (86) that  $D_\varepsilon^2 = \max_{\mathbf{x}, \mathbf{y} \in \mathcal{K}_\varepsilon} B(\mathbf{x} | | \mathbf{y}) \leq \frac{9\nu^{\frac{3}{2}} K \sqrt{S_1}}{\varepsilon}$ , where  $\nu$  is the self-concordance  
 322 constant of  $\psi$  and  $S_1 = \|\nabla^2 \psi(\mathbf{x}_1)\|_2$  is the spectral norm of its Hessian at the center  $\mathbf{x}_1$  of  $\mathcal{K}$ .  
 323 Restricting to tuning  $\varepsilon \in [\frac{1}{m}, 1]$ —which is enough to obtain constant task-averaged regret above if  
 324 the estimated optima  $\hat{\mathbf{x}}_t$  are identical—we can now apply Algorithm 1 via the following instantiation:

1. sample  $\varepsilon_t$  via the MW distribution  $\propto \exp(\mathbf{w}_t)$  over the discretization  $\Theta_k$  of  $[\frac{1}{m}, 1]$
2. run  $\text{OMD}_{\eta_t, \varepsilon_t}$  using the initialization  $\mathbf{x}_{t,1} = \frac{1}{t-1} \sum_{s < t} \mathbf{c}_{\varepsilon_t}(\hat{\mathbf{x}}_s) = \mathbf{x}_1 + \frac{\sum_{s < t} \hat{\mathbf{x}}_s - \mathbf{x}_1}{(1+\varepsilon_t)(t-1)}$  (FTL)
3. update EWOO at each  $\varepsilon \in \Theta_k$  with loss  $\frac{B(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \|\mathbf{x}_{t,1}) + \rho^2 D_\varepsilon^2}{\eta} + 32\eta d^2$  for  $D_\varepsilon^2 = \frac{9\nu^{\frac{3}{2}} K \sqrt{S_1}}{\varepsilon}$
4. update  $\mathbf{p}_{t+1}$  using multiplicative weights with expert losses  $\frac{B(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \|\mathbf{x}_{t,1})}{\eta} + \varepsilon m$

(17)

325 Note the similarity to the MAB case (8), with the difference being the upper bound passed to EWOO  
 326 and MW. Our structural result bounds the expected task-averaged regret as follows (c.f. Thm. D.1):

$$\mathbb{E} \min_{\varepsilon \in [\frac{1}{m}, 1], \eta > 0} \frac{\hat{V}_\varepsilon^2}{\eta} + (32\eta d^2 + \varepsilon)m + \tilde{O} \left( \frac{\frac{m^2}{T} + \frac{1}{k} + m}{\eta} + m \min \left\{ \frac{\rho^2}{\eta}, d\rho \right\} + \frac{dm}{\rho} \sqrt{\frac{\log k}{T} + \frac{dm}{\rho^2 T}} \right) \quad (18)$$

327 For  $\rho = o_T(1)$  and  $k = \omega_T(1)$  this becomes  $o_T(\text{poly}(m)) + \mathbb{E} \min_{\varepsilon \in [\frac{1}{m}, 1], \eta > 0} \frac{\hat{V}_\varepsilon^2}{\eta} + 32\eta d^2 m + \varepsilon m$ ,  
 328 where  $\hat{V}_\varepsilon^2 = \frac{1}{T} \sum_{t=1}^T \psi(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) - \psi(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t)))$ . Then by tuning  $\eta$  we get an asymptotic ( $T \rightarrow \infty$ )  
 329 regret of  $4d\hat{V}_\varepsilon \sqrt{2m} + \varepsilon m$  for any  $\varepsilon \in [\frac{1}{m}, 1]$ . Our analysis removes the explicit dependence on  $\sqrt{\nu}$   
 330 that appears in the single-task regret [2]; as an example,  $\nu$  equals the number of inequalities defining  
 331 a polytope  $\mathcal{K}$ , as in the bandit shortest-path application below.

332 The remaining challenge is to interpret  $\hat{V}_\varepsilon^2$ , which as we did for MAB we do via specific examples,  
 333 in this case concrete action domains  $\mathcal{K}$ . Our first example is for BLO over the unit sphere  $\mathcal{K} = \{\mathbf{x} \in$   
 334  $\mathbb{R}^d : \|\mathbf{x}\|_2 \leq 1\}$  using the appropriate log-barrier regularizer  $\psi(\mathbf{x}) = -\log(1 - \|\mathbf{x}\|_2^2)$ :

335 **Corollary 4.1** (c.f. Cor. D.1). *For BLO on the sphere, Algorithm 1 has expected task-averaged regret*

$$\tilde{O} \left( \frac{dm^{\frac{3}{2}}}{T^{\frac{3}{4}}} + \frac{dm}{\sqrt[4]{T}} \right) + \min_{\varepsilon \in [\frac{1}{m}, 1]} 4d \sqrt{2m \log \left( 1 + \frac{1 - \mathbb{E} \|\hat{\mathbf{x}}\|_2^2}{2\varepsilon + \varepsilon^2} \right)} + \varepsilon m \quad (19)$$

336 The bound above is decreasing in  $\mathbb{E} \|\hat{\mathbf{x}}\|_2^2$ , the expected squared norm of the average of the estimated  
 337 optima  $\hat{\mathbf{x}}_t$ . We thus say that *bandit linear optimization tasks over the sphere are similar if the norm*  
 338 *of the empirical mean of their (estimated) optima is large*. This makes intuitive sense: if the tasks'  
 339 optima are uniformly distributed, we should expect  $\mathbb{E} \|\hat{\mathbf{x}}\|_2^2$  to be small, even decreasing in  $d$ . On the  
 340 other hand, in the degenerate case where the estimated optima  $\hat{\mathbf{x}}_t$  are the same across all tasks  $t \in [T]$ ,  
 341 we have  $\mathbb{E} \|\hat{\mathbf{x}}\|_2^2 = 1$ , so the asymptotic task-averaged regret is 1 because we can use  $\varepsilon = \frac{1}{m}$ . Perhaps  
 342 slightly more realistically, if it is  $\frac{1}{m^p}$ -away from 1 for some power  $p \geq \frac{1}{2}$  then setting  $\varepsilon = \frac{1}{\sqrt{m}}$  can  
 343 remove the logarithmic dependence on  $m$ . These two regimes illustrate the importance of tuning  $\varepsilon$ .

344 Motivated by the bandit shortest-path problem [49, 30] and described in full in Section D.3, our last  
 345 application specializes Theorem D.1 to polytopes. There, the induced task-similarity is a sum across  
 346 polytope boundaries, with each summand the logarithm of a quotient of arithmetic and geometric  
 347 means aggregating how close the task optima are to the boundary being considered. When all  
 348 distances across tasks are the same, the two means are the same and so the log of their quotient is  
 349 zero, making that summand zero. Thus, the task-averaged regret improves if the optima for different  
 350 tasks are at similar distances from different boundaries of the polytope.

## 351 5 Conclusion and limitations

352 We develop and apply a meta-algorithm for learning to initialize and tune bandit algorithms, obtaining  
 353 task-averaged regret guarantees for both multi-armed and linear bandits that depend on natural,  
 354 setting-specific notions of task similarity. For MAB, we meta-learn the initialization, step-size,  
 355 and entropy parameter of Tsallis-entropic OMD and show good performance if the entropy of the  
 356 optimal arms is small. For BLO, we use OMD with self-concordant regularizers and meta-learn  
 357 the initialization, step-size, and boundary-offset, yielding interesting domain-specific task-similarity  
 358 measures. Some natural directions for future work involve overcoming some limitations of our results:  
 359 can we adapt to a notion of task-similarity that depends on the true optima without assuming a gap  
 360 for MAB, or at all for BLO? Alternatively, can we design meta-learning algorithms that adapt to both  
 361 stochastic and adversarial bandits, i.e. a ‘‘best-of-both-worlds’’ guarantee? Beyond this, one could  
 362 explore other partial information settings, such as contextual bandits or bandit convex optimization.

## References

- 363
- 364 [1] Yasin Abbasi-Yadkori, Peter Bartlett, Victor Gabillon, Alan Malek, and Michal Valko. Best  
365 of both worlds: Stochastic & adversarial best-arm identification. In *Proceedings of the 31st*  
366 *Conference On Learning Theory*, 2018.
- 367 [2] Jacob Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient  
368 algorithm for bandit linear optimization. In *Proceedings of the International Conference on*  
369 *Computational Learning Theory*, 2008.
- 370 [3] Jacob Abernethy, Chansoo Lee, and Ambuj Tewari. Fighting bandits with a new kind of  
371 smoothness. In *Advances in Neural Information Processing Systems*, 2015.
- 372 [4] Oren Anava and Zohar S. Karnin. Multi-armed bandits: Competing with optimal sequences. In  
373 *Advances in Neural Information Processing Systems*, 2016.
- 374 [5] Jean-Yves Audibert, Sébastien Bubeck, and Gábor Lugosi. Minimax policies for combinatorial  
375 prediction games. In *Proceedings of the International Conference on Computational Learning*  
376 *Theory*, 2011.
- 377 [6] Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic  
378 multiarmed bandit problem. *SIAM Journal of Computing*, 32:48–77, 2002.
- 379 [7] Peter Auer, Pratik Gajane, and Ronald Ortner. Adaptively tracking the best bandit arm with an  
380 unknown number of distribution changes. In *Proceedings of the 32nd Annual Conference on*  
381 *Learning Theory*, 2019.
- 382 [8] Baruch Awerbuch and Robert D. Kleinberg. Adaptive routing with end-to-end feedback:  
383 Distributed learning and geometric approaches. In *Proceedings of the Thirty-Sixth Annual ACM*  
384 *Symposium on Theory of Computing*, 2004.
- 385 [9] Mohammad Gheshlaghi Azar, Alessandro Lazaric, Emma Brunskill, et al. Sequential transfer  
386 in multi-armed bandit with finite set of models. In *Advances in Neural Information Processing*  
387 *Systems*, 2013.
- 388 [10] MohammadJavad Azizi, Thang Duong, Yasin Abbasi-Yadkori, András György, Claire Vernade,  
389 and Mohammad Ghavamzadeh. Non-stationary bandits and meta-learning with a small set of  
390 optimal arms. arXiv, 2022.
- 391 [11] Maria-Florina Balcan, Avrim Blum, and Santosh Vempala. Efficient representations for lifelong  
392 learning and autoencoding. In *Proceedings of the 28th Annual Conference on Learning Theory*,  
393 2015.
- 394 [12] Maria-Florina Balcan, Mikhail Khodak, Dravyansh Sharma, and Ameet Talwalkar. Learning-to-  
395 learn non-convex piecewise-Lipschitz functions. In *Advances in Neural Information Processing*  
396 *Systems*, 2021.
- 397 [13] Soumya Basu, Branislav Kveton, Manzil Zaheer, and Csaba Szepesvári. No regrets for learning  
398 the prior in bandits. In *Advances in Neural Information Processing Systems*, 2021.
- 399 [14] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for  
400 convex optimization. *Operations Research Letters*, 31:167–175, 2003.
- 401 [15] Leonardo Cella, Alessandro Lazaric, and Massimiliano Pontil. Meta-learning with stochastic  
402 linear bandits. In *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- 403 [16] Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, Learning, and Games*. Cambridge  
404 University Press, 2006.
- 405 [17] Varsha Dani, Thomas Hayes, and Sham Kakade. The price of bandit information for online  
406 optimization. In *Advances in Neural Information Processing Systems*, 2008.
- 407 [18] Giulia Denevi, Carlo Ciliberto, Dimitris Stamos, and Massimiliano Pontil. Learning to learning  
408 around a common mean. In *Advances in Neural Information Processing Systems*, 2018.

- 409 [19] Giulia Denevi, Carlo Ciliberto, Riccardo Grazi, and Massimiliano Pontil. Online-within-online  
410 meta-learning. In *Advances in Neural Information Processing Systems*, 2019.
- 411 [20] Simon S. Du, Wei Hu, Sham M. Kakade, Jason D. Lee, and Qi Lei. Few-shot learning via  
412 learning the representation, provably. In *Proceedings of the 9th International Conference on*  
413 *Learning Representations*, 2021.
- 414 [21] Yan Duan, Marcin Andrychowicz, Bradly Stadie, Jonathan Ho, Jonas Schneider, Ilya Sutskever,  
415 Pieter Abbeel, and Wojciech Zaremba. One-shot imitation learning. In *Advances in Neural*  
416 *Information Processing Systems*, 2017.
- 417 [22] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. On the convergence theory of gradient-  
418 based model-agnostic meta-learning algorithms. In *Proceedings of the 23rd International*  
419 *Conference on Artificial Intelligence and Statistics*, 2020.
- 420 [23] Xiequan Fan, Ion Grama, and Quansheng Liu. Hoeffding’s inequality for supermartingales.  
421 *Stochastic Processes and their Applications*, 122(10):3545–3559, 2012.
- 422 [24] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adap-  
423 tation of deep networks. In *Proceedings of the 34th International Conference on Machine*  
424 *Learning*, 2017.
- 425 [25] Elad Hazan. Introduction to online convex optimization. In *Foundations and Trends in*  
426 *Optimization*, volume 2, pages 157–325. now Publishers Inc., 2015.
- 427 [26] Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex  
428 optimization. *Machine Learning*, 69:169–192, 2007.
- 429 [27] Ali Jadbabaie, Alexander Rakhlin, and Shahin Shahrampour. Online optimization: Competing  
430 with dynamic comparators. In *Proceedings of the 18th International Conference on Artificial*  
431 *Intelligence and Statistics*, 2015.
- 432 [28] Kevin Jamieson and Ameet Talwalkar. Non-stochastic best arm identification and hyperparame-  
433 ter optimization. In *Proceedings of the 18th International Conference on Artificial Intelligence*  
434 *and Statistics*, 2015.
- 435 [29] Yihan Jiang, Jakub Konečný, Keith Rush, and Sreeram Kannan. Improving federated learning  
436 personalization via model agnostic meta learning. arXiv, 2019.
- 437 [30] Adam Kalai and Santosh Vempala. Efficient algorithms for online decision problems. *Journal*  
438 *of Computer and System Sciences*, 71:291–307, 2005.
- 439 [31] Mikhail Khodak, Maria-Florina Balcan, and Ameet Talwalkar. Adaptive gradient-based meta-  
440 learning methods. In *Advances in Neural Information Processing Systems*, 2019.
- 441 [32] Mikhail Khodak, Renbo Tu, Tian Li, Liam Li, Maria-Florina Balcan, Virginia Smith, and  
442 Ameet Talwalkar. Federated hyperparameter tuning: Challenges, baselines, and connections to  
443 weight-sharing. In *Advances in Neural Information Processing Systems*, 2021.
- 444 [33] Mikhail Khodak, Maria-Florina Balcan, Ameet Talwalkar, and Sergei Vassilvitskii. Learning  
445 predictions for algorithms with predictions. In *Advances in Neural Information Processing*  
446 *Systems*, 2022.
- 447 [34] Branislav Kveton, Martin Mladenov, Chih-Wei Hsu, Manzil Zaheer, Csaba Szepesvári, and  
448 Craig Boutilier. Meta-learning bandit policies by gradient ascent. arXiv, 2020.
- 449 [35] Branislav Kveton, Mikhail Konobeev, Manzil Zaheer, Chih-Wei Hsu, Martin Mladenov, Craig  
450 Boutilier, and Csaba Szepesvári. Meta-Thompson sampling. In *Proceedings of the 38th*  
451 *International Conference on Machine Learning*, 2021.
- 452 [36] Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-SGD: Learning to learning quickly  
453 for few-shot learning. arXiv, 2017.
- 454 [37] Haipeng Luo. CSCI 699 Lecture 13. 2017. URL [https://haipeng-luo.net/  
455 courses/CSCI699/lecture13.pdf](https://haipeng-luo.net/courses/CSCI699/lecture13.pdf).

- 456 [38] Haipeng Luo, Chen-Yu Wei, Alekh Agarwal, and John Langford. Efficient contextual bandits  
457 in non-stationary worlds. In *Proceedings of the 31st Annual Conference on Learning Theory*,  
458 2018.
- 459 [39] Michael Mitzenmacher and Sergei Vassilvitskii. Algorithms with predictions. In Tim Rough-  
460 garden, editor, *Beyond the Worst-Case Analysis of Algorithms*. Cambridge University Press,  
461 Cambridge, UK, 2021.
- 462 [40] Ahmadreza Moradipari, Mohammad Ghavamzadeh, Taha Rajabzadeh, Christos Thrampoulidis,  
463 and Mahnoosh Alizadeh. Multi-environment meta-learning in stochastic linear bandits. In  
464 *Proceedings of the 2022 IEEE International Symposium on Information Theory*, 2022.
- 465 [41] Yurii Nesterov and Arkadii Nemirovskii. *Interior-Point Polynomial Algorithms in Convex*  
466 *Programming*. SIAM Studies in Applied and Numerical Mathematics, 1994.
- 467 [42] Gergely Neu. Explore no more: Improved high-probability regret bounds for non-stochastic  
468 bandits. In *Advances in Neural Information Processing Systems*, 2015.
- 469 [43] Alex Nichol, Joshua Achiam, and John Schulman. On first-order meta-learning algorithms.  
470 arXiv, 2018.
- 471 [44] Nikunj Saunshi, Yi Zhang, Mikhail Khodak, and Sanjeev Arora. A sample complexity separation  
472 between non-convex and convex meta-learning. In *Proceedings of the 37th International*  
473 *Conference on Machine Learning*, 2020.
- 474 [45] Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends*  
475 *in Machine Learning*, 4(2):107–194, 2011.
- 476 [46] Amr Sharaf and Hal Daumé III. Meta-learning effective exploration strategies for contextual  
477 bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021.
- 478 [47] Max Simchowitz, Christopher Tosh, Akshay Krishnamurthy, Daniel J. Hsu, Thodoris Lykouris,  
479 Miro Dudik, and Robert E. Schapire. Bayesian decision-making under misspecified priors with  
480 applications to meta-learning. In *Advances in Neural Information Processing Systems*, 2021.
- 481 [48] Jake Snell, Kevin Swersky, and Richard S. Zemel. Prototypical networks for few-shot learning.  
482 In *Advances in Neural Information Processing Systems*, 2017.
- 483 [49] Eiji Takimoto and Manfred K. Warmuth. Path kernels and multiplicative updates. *Journal of*  
484 *Machine Learning Research*, 4:773–818, 2003.
- 485 [50] Sebastian Thrun and Lorien Pratt. *Learning to Learn*. Springer Science & Business Media,  
486 1998.
- 487 [51] Constantino Tsallis. Possible generalization of Boltzmann-Gibbs statistics. *Journal of Statistical*  
488 *Physics*, 52:479–487, 1988.
- 489 [52] Takuya Yamano. Some properties of q-logarithm and q-exponential functions in Tsallis statistics.  
490 *Physica A: Statistical Mechanics and its Applications*, 305:486–496, 2002.
- 491 [53] Peng Zhao, Guanghui Wang, Lijun Zhang, and Zhi-Hua Zhou. Bandit convex optimization in  
492 non-stationary environments. *Journal of Machine Learning Research*, 22(125):1–45, 2021.

493 **A Structural results**

494 **A.1 Properties of the Bregman divergence**

495 **Lemma A.1.** *Let  $\psi : \mathcal{C} \mapsto \mathbb{R}$  be a strictly convex function with  $\max_{\mathbf{x} \in \mathcal{C}} \|\nabla^2 \psi(\mathbf{x})\|_2 \leq S$  over*  
 496 *a convex set  $\mathcal{C} \subset \mathbb{R}^d$  over size  $\max_{\mathbf{x} \in \mathcal{C}} \|\mathbf{x}\|_2 \leq K$ , and let  $B(\cdot|\cdot)$  be the Bregman divergence*  
 497 *generated by  $\psi$ . Then for any points  $\mathbf{x}_1, \dots, \mathbf{x}_T \in \mathcal{C}$  the actions  $\mathbf{y}_1 = \arg \min_{\mathbf{x} \in \mathcal{C}} \psi(\mathbf{x})$  and*  
 498  *$\mathbf{y}_t = \frac{1}{t-1} \sum_{s < t} \mathbf{x}_s$  have regret*

$$\sum_{t=1}^T B(\mathbf{x}_t|\mathbf{y}_t) - B(\mathbf{x}_t|\mathbf{y}_{T+1}) \leq \sum_{t=1}^T \frac{8SK^2}{2t-1} \leq 8SK^2(1 + \log T) \quad (20)$$

499 *Proof.* Note that

$$\nabla_{\mathbf{y}} B(\mathbf{x}|\mathbf{y}) = -\nabla \psi(\mathbf{y}) - \nabla_{\mathbf{y}} \langle \nabla \psi(\mathbf{y}), \mathbf{x} \rangle + \nabla_{\mathbf{y}} \langle \nabla \psi(\mathbf{y}), \mathbf{y} \rangle = \text{diag}(\nabla^2 \psi(\mathbf{y}))(\mathbf{y} - \mathbf{x}) \quad (21)$$

500 so  $B(\mathbf{x}_t|\mathbf{y}_t)$  is  $2SK$ -Lipschitz w.r.t. the Euclidean norm. Applying Khodak et al. [31, Prop. B.1]  
 501 yields the result (note that its assumption of strong convexity of the regularizer can be replaced with  
 502 strict convexity without changing the proof or result).  $\square$

503 **Claim A.1.** *Let  $\psi : \mathcal{K} \mapsto \mathbb{R}$  be a strictly-convex function with Bregman divergence  $B(\cdot|\cdot)$  over a*  
 504 *convex set  $\mathcal{K} \subset \mathbb{R}^d$  containing points  $\mathbf{x}_1, \dots, \mathbf{x}_T$ . Then their mean  $\bar{\mathbf{x}} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$  satisfies*

$$\sum_{t=1}^T B(\mathbf{x}_t|\bar{\mathbf{x}}) = \sum_{t=1}^T \psi(\mathbf{x}_t) - \psi(\bar{\mathbf{x}}) \quad (22)$$

*Proof.*

$$\begin{aligned} \sum_{t=1}^T B(\mathbf{x}_t|\bar{\mathbf{x}}) &= \sum_{t=1}^T \psi(\mathbf{x}_t) - \psi(\bar{\mathbf{x}}) - \langle \nabla \psi(\bar{\mathbf{x}}), \mathbf{x}_t - \bar{\mathbf{x}} \rangle \\ &= \sum_{t=1}^T \psi(\mathbf{x}_t) - \psi(\bar{\mathbf{x}}) - \langle \nabla \psi(\bar{\mathbf{x}}), \sum_{t=1}^T \mathbf{x}_t - T\bar{\mathbf{x}} \rangle = \sum_{t=1}^T \psi(\mathbf{x}_t) - \psi(\bar{\mathbf{x}}) \end{aligned} \quad (23)$$

505  $\square$

506 **A.2 Tuning the step-size**

507 **Lemma A.2.** *Let  $\ell_1, \dots, \ell_T : \mathbb{R}_{>0} \mapsto \mathbb{R}_{>0}$  be a sequence of functions of form  $\ell_t(x) = \frac{B_t^2}{x} + G^2 x$*   
 508 *for adversarially chosen  $B_t \in [0, D]$  and some  $G > 0$ . Then for any  $\rho \geq 0$ , the actions of*  
 509 *EWOO [26, Fig. 4] with parameter  $\frac{2\rho^2}{DG}$  run on the modified losses  $\frac{B_t^2 + \rho^2 D^2}{x} + G^2 x$  over the domain*  
 510  *$\left[ \frac{\rho D}{G}, \frac{D}{G} \sqrt{1 + \rho^2} \right]$  achieves regret w.r.t. any  $x > 0$  of*

$$\sum_{t=1}^T \ell_t(x) - \ell_t(x) \leq \min \left\{ \frac{\rho^2 D^2}{x}, \rho DG \right\} T + \frac{DG(1 + \log(T+1))}{2\rho^2} \quad (24)$$

511 *Proof.* By Khodak et al. [31, Prop. C.1] the modified functions are  $\frac{2\rho^2}{DG}$ -exp-concave. Then Khodak  
 512 et al. [31, Cor. C.2] with  $B_t$  set to  $\frac{B_t}{G}$ ,  $D$  to  $\frac{D}{G}$ ,  $\alpha_t = G^2$ , and  $\varepsilon = \frac{\rho D}{G}$  yields the result.  $\square$

513 **Lemma A.3.** For  $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_T \in \partial\mathcal{K}$  consider a sequence of functions of form

$$U_t(\mathbf{x}, \eta) = \frac{B(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \|\mathbf{x})}{\eta} + \eta G^2 m \quad (25)$$

514 where  $B$  is the Bregman divergence of a strictly convex d.g.f.  $\psi : \mathcal{K}^\circ \mapsto \mathbb{R}$  and where  $\mathbf{x}_1 =$   
515  $\arg \min_{\mathbf{x} \in \mathcal{K}} \psi(\mathbf{x})$  defines the projection  $\mathbf{c}_\varepsilon(\mathbf{x}) = \mathbf{x}_1 + \frac{\mathbf{x} - \mathbf{x}_1}{1 + \varepsilon}$  for some  $\varepsilon > 0$ . Suppose we play  
516  $\mathbf{x}_{t+1} \leftarrow \mathbf{c}_\varepsilon\left(\frac{1}{t} \sum_{s=1}^t \hat{\mathbf{x}}_s\right)$  and set  $\eta_t$  using the actions of EWOO [26, Fig. 4] with parameter  $\frac{2\rho^2}{DG}$  for  
517 some  $\rho, D_\varepsilon > 0$  s.t.  $B(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \|\mathbf{x}) \leq D_\varepsilon^2 \forall \mathbf{x} \in \mathcal{K}_\varepsilon$  on the functions  $\frac{B(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \|\mathbf{x}) + \rho^2 D_\varepsilon^2}{\eta} + \eta G^2 m$  over  
518 the domain  $\left[\frac{\rho D_\varepsilon}{G\sqrt{m}}, \frac{D_\varepsilon}{G} \sqrt{\frac{1+\rho^2}{m}}\right]$ , with  $\eta_1$  being at the midpoint of the domain. Then  $U_t(\mathbf{x}_t, \eta_t) \leq$   
519  $D_\varepsilon G \sqrt{m} \left(\frac{1}{\rho} + \sqrt{1 + \rho^2}\right) \forall t \in [T]$  and

$$\begin{aligned} \sum_{t=1}^T U_t(\mathbf{x}_t, \eta_t) &\leq \min_{\eta > 0, \mathbf{x} \in \mathcal{K}} \sum_{t=1}^T \frac{B(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \|\mathbf{x})}{\eta} + \eta G^2 m \\ &\quad + \min \left\{ \frac{\rho^2 D_\varepsilon^2}{\eta}, \rho D_\varepsilon G \right\} T + \frac{D_\varepsilon G (1 + \log(T+1))}{2\rho^2} + \frac{8S_\varepsilon K^2 (1 + \log T)}{\eta} \end{aligned} \quad (26)$$

520 for  $K = \max_{\mathbf{x} \in \mathcal{K}} \|\mathbf{x}\|_2$  and  $S_\varepsilon = \max_{\mathbf{x} \in \mathcal{K}_\varepsilon} \|\nabla^2 \psi(\mathbf{x})\|_2$ .

521 *Proof.* The first claim follows by directly substituting the worst-case values of  $\eta$  into  $U_t(\mathbf{x}, \eta)$ . For  
522 the second, apply Lemma A.2 followed by Lemma A.1:

$$\begin{aligned} &\sum_{t=1}^T U_t(\mathbf{x}_t, \eta_t) \\ &= \sum_{t=1}^T \frac{B(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \|\mathbf{x}_t)}{\eta_t} + \eta_t G^2 m \\ &\leq \min_{\eta > 0} \min \left\{ \frac{\rho^2 D_\varepsilon^2}{\eta}, \rho D_\varepsilon G \right\} T + \frac{D_\varepsilon G (1 + \log(T+1))}{2\rho^2} + \sum_{t=1}^T \frac{B(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \|\mathbf{x})}{\eta} + \eta G^2 m \\ &\leq \min_{\eta > 0} \min \left\{ \frac{\rho^2 D_\varepsilon^2}{\eta}, \rho D_\varepsilon G \right\} T + \frac{D_\varepsilon G (1 + \log(T+1))}{2\rho^2} + \frac{8S_\varepsilon K^2 (1 + \log T)}{\eta} \\ &\quad + \min_{\mathbf{x} \in \mathcal{K}_\varepsilon} \sum_{t=1}^T \frac{B(\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \|\mathbf{x})}{\eta} + \eta G^2 m \end{aligned} \quad (27)$$

523 Conclude by noting that the sum of Bregman divergence to  $\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t)$  is minimized on their convex hull,  
524 a subset of  $\mathcal{K}_\varepsilon$ .  $\square$

### 525 A.3 Computational and space complexity

526 Algorithm 1 implicitly maintains a separate copy of FTL for each hyperparameter in the continuous  
527 space of EWOO and the grid  $\Theta_k$  over the domain of  $\theta$ , but explicitly just needs to average the estimated  
528 task-optima  $\hat{\mathbf{x}}_t$ ; this is due to the mean-as-minimizer property of Bregman divergences and the linear-  
529 ity of  $\mathbf{c}_\varepsilon$ . Thus the memory it uses is  $\mathcal{O}(d+k)$ , where  $k$  is size of the discretization of  $\Theta$  and should be  
530 viewed as sublinear in  $T$ , e.g. for MAB with implicit exploration and BLO  $k = \mathcal{O}(\sqrt[4]{d}\sqrt{T})$ . Computa-  
531 tionally, at each timestep  $t$  and for each grid point we must compute two single-dimensional integrals;  
532 the integrands are sums of upper bounds that just need to be incremented once per round, leading to a  
533 total per-iteration complexity of  $\mathcal{O}(k)$  (ignoring the running of OMD). Although outside the scope of  
534 this work, it may be possible to avoid integration by tuning  $\eta$  with MW as well, rather than EWOO,  
535 but likely at the cost of worse regret because it would not take advantage of the exp-concavity of  $U_t^{(\rho)}$ .

536 **A.4 Main structural result**

537 **Theorem A.1.** Consider a family of strictly convex functions  $\psi_\theta : \mathcal{K}^\circ \mapsto \mathbb{R}$  parameterized by  $\theta$  lying  
538 in an interval  $\Theta \subset \mathbb{R}$  of radius  $R_\Theta$  that are all minimized at the same  $\mathbf{x}_1 \in \mathcal{K}^\circ$ , and for  $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_T \in$   
539  $\partial\mathcal{K}$  consider a sequence of functions of form  $U_t(\mathbf{x}, \eta, \theta)$  (3), as well as the associated regularized up-  
540 per bounds  $U_t^{(\rho)}$  (4). Define the maximum divergence  $D = \max_{\theta \in \Theta} D_\theta$ , radius  $K = \max_{\mathbf{x} \in \mathcal{K}} \|\mathbf{x}\|_2$ ,  
541 and  $L_\eta$  the Lipschitz constant w.r.t.  $\theta \in \Theta$  of  $\frac{\hat{V}_\theta^2}{\eta} + \eta g(\theta)m + f(\theta)m$ . Then Algorithm 1 with  $\Theta_k \subset \Theta$   
542 the uniform discretization of  $\Theta$  s.t.  $\max_{\theta \in \Theta} \min_{\theta' \in \Theta_k} |\theta - \theta'| \leq \frac{R_\Theta}{k}$ ,  $\rho \in (0, 1)$ ,  $\underline{\eta}(\theta) = \frac{\rho D_\theta}{\sqrt{g(\theta)m}}$ ,  
543  $\bar{\eta}(\theta) = D_\theta \sqrt{\frac{1+\rho^2}{g(\theta)m}}$ ,  $\alpha(\theta) = \frac{2\rho^2}{D_\theta \sqrt{g(\theta)m}}$ , and  $\lambda = \left( M \left( \frac{1}{\rho} + \sqrt{1+\rho^2} \right) + Fm \right)^{-1} \sqrt{\frac{\log k}{2T}}$  leads to  
544 a sequence  $(\mathbf{x}_t, \eta_t(\theta_t), \theta_t)$  s.t.  $\mathbb{E} \sum_{t=1}^T U_t(\mathbf{x}_t, \eta_t(\theta_t), \theta_t)$  is bounded by

$$\begin{aligned} & \mathbb{E} \min_{\theta \in \Theta, \eta > 0} \frac{8SK^2(1 + \log T)}{\eta} + \left( \frac{\hat{V}_\theta^2}{\eta} + \eta g(\theta)m + f(\theta)m + \frac{L_\eta R_\Theta}{k} + \min \left\{ \frac{\rho^2 D^2}{\eta}, \rho M \right\} \right) T \\ & + \left( \frac{4M}{\rho} + Fm \right) \sqrt{T \log k} + \frac{M(1 + \log(T+1))}{2\rho^2} \end{aligned} \quad (28)$$

545 and  $\sum_{t=1}^T U_t(\mathbf{x}_t, \eta_t(\theta_t), \theta_t)$  is bounded w.p.  $\geq 1 - \delta \mathbf{1}_{k>1}$  by

$$\begin{aligned} & \min_{\theta \in \Theta, \eta > 0} \frac{8SK^2(1 + \log T)}{\eta} + \left( \frac{\hat{V}_\theta^2}{\eta} + \eta g(\theta)m + f(\theta)m + \frac{L_\eta R_\Theta}{k} + \min \left\{ \frac{\rho^2 D^2}{\eta}, \rho M \right\} \right) T \\ & + \left( \frac{4M}{\rho} + Fm \right) \left( \sqrt{T \log k} + \mathbf{1}_{k>1} \sqrt{\frac{T}{2} \log \frac{1}{\delta}} \right) + \frac{M(1 + \log(T+1))}{2\rho^2} \end{aligned} \quad (29)$$

546 *Proof.* In the following proof, we first consider online learning  $U_t(\cdot, \cdot, \theta)$  for fixed  $\theta \in \Theta_k$ . To tune  
547  $\eta$ , we online learn the one-dimensional losses  $B_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}_t) | | \mathbf{c}_\theta(\mathbf{x}_t)) / \eta + \eta g(\theta)$ , where  $\mathbf{c}_\theta(\hat{\mathbf{x}}_t)$  is the  
548  $(\eta_t(\theta)$ -independent) action of FTL at time  $t$ . As discussed, the corresponding regularized losses  $U_t^{(\rho)}$   
549 are exp-concave, and so running EWOO yields  $\tilde{\mathcal{O}}(M/\rho^2 + \min\{\rho^2 D^2/\eta, \rho M\}T)$  regret w.r.t. the  
550 original sequence [31, Cor. C.2]. At the same time, we show that FTL has logarithmic regret on the  
551 sequence  $B_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}_t) | | \cdot)$  that scales with the spectral norm  $S$  of  $\nabla^2 \psi_\theta$  (c.f. Lem. A.1), and that the  
552 average loss of the optimal comparator is  $\hat{V}_\theta^2$  (c.f. Claim A.1). Thus, since we only care about a fixed  
553 comparator  $\eta$ , dividing by  $\eta T$  yields the first and last terms (5). We run a copy of these algorithms  
554 for each  $\theta \in \Theta_k$ ; since their losses are bounded by  $\tilde{\mathcal{O}}(M/\rho + Fm)$ , textbook results for MW yield  
555  $\mathcal{O}(\sqrt{T \log k})$  regret w.r.t.  $\theta \in \Theta_k$ , which we then extend to  $\Theta \supset \Theta_k$  using  $L_\eta$ -Lipschitzness.

556 Formally, we have that

$$\begin{aligned} & \mathbb{E} \sum_{t=1}^T U_t(\mathbf{x}_t, \eta_t(\theta_t), \theta_t) \\ & = \mathbb{E} \sum_{t=1}^T \frac{B_{\theta_t}(\mathbf{c}_{\theta_t}(\hat{\mathbf{x}}_t) | | \mathbf{x}_t)}{\eta_t(\theta_t)} + \eta_t(\theta_t)g(\theta_t)m + f(\theta_t)m \\ & \leq \left( M \left( \frac{1}{\rho} + \sqrt{2} \right) + Fm \right) \sqrt{2T \log k} + \mathbb{E} \min_{\theta \in \Theta_k} \sum_{t=1}^T \frac{B_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}_t) | | \mathbf{x}_t)}{\eta_t(\theta)} + \eta_t(\theta)g(\theta)m + f(\theta)m \\ & \leq \left( \frac{4M}{\rho} + Fm \right) \sqrt{T \log k} + \mathbb{E} \min_{\theta \in \Theta_k, \eta > 0, \mathbf{x} \in \mathcal{K}} \sum_{t=1}^T \frac{B_\theta(\mathbf{c}_\theta(\hat{\mathbf{x}}_t) | | \mathbf{x})}{\eta} + \eta g(\theta)m + f(\theta)m \\ & \quad + \min \left\{ \frac{\rho^2 D_\theta^2}{\eta}, \rho D_\theta \sqrt{g(\theta)m} \right\} T + \frac{D_\theta \sqrt{g(\theta)m}(1 + \log(T+1))}{2\rho^2} + \frac{8SK^2(1 + \log T)}{\eta} \end{aligned} \quad (30)$$

557 where the first inequality is the regret of multiplicative weights with step-size  $\lambda$  [45, Cor. 2.14] and  
 558 the second is by applying Lemma A.3 for each  $\theta$ . We then simplify and apply the definition of  $\hat{V}_\theta^2$   
 559 via Claim A.1 and conclude by applying Lipschitzness w.r.t.  $\theta$ :

$$\begin{aligned}
 & \mathbb{E} \sum_{t=1}^T U_t(\mathbf{x}_t, \eta_t(\theta_t), \theta_t) \\
 & \leq \left( \frac{4M}{\rho} + Fm \right) \sqrt{T \log k} + \mathbb{E} \min_{\theta \in \Theta_k, \eta > 0} \frac{\hat{V}_\theta^2 T}{\eta} + \eta g(\theta) m T + f(\theta) m T \\
 & \quad + \min \left\{ \frac{\rho^2 D^2}{\eta}, \rho M \right\} T + \frac{M(1 + \log(T+1))}{2\rho^2} + \frac{8SK^2(1 + \log T)}{\eta} \\
 & \leq \mathbb{E} \min_{\theta \in \Theta, \eta > 0} \frac{8SK^2(1 + \log T)}{\eta} + \left( \frac{\hat{V}_\theta^2}{\eta} + \eta g(\theta) m + f(\theta) m + \frac{L_\eta R_\Theta}{k} + \min \left\{ \frac{\rho^2 D^2}{\eta}, \rho M \right\} \right) T \\
 & \quad + \left( \frac{4M}{\rho} + Fm \right) \sqrt{T \log k} + \frac{M(1 + \log(T+1))}{2\rho^2}
 \end{aligned} \tag{31}$$

560 The w.h.p. guarantee follows by Cesa-Bianchi and Lugosi [16, Lem. 4.1].  $\square$

## 561 B Implicit exploration

### 562 B.1 Properties of the Tsallis entropy

563 **Lemma B.1.** For any  $\varepsilon \in (0, 1]$  and  $\mathbf{x} \in \Delta$  s.t.  $\mathbf{x}(a) \geq \frac{\varepsilon}{d} \forall a \in [d]$  the  $\beta$ -Tsallis entropy  
 564  $H_\beta(\mathbf{x}) = -\frac{1 - \sum_{a=1}^d \mathbf{x}^\beta(a)}{1 - \beta}$  is  $d \log \frac{d}{\varepsilon}$ -Lipschitz w.r.t.  $\beta \in [0, 1]$ .

565 *Proof.* Let  $\log_\beta x = \frac{x^{1-\beta} - 1}{1-\beta}$  be the  $\beta$ -logarithm function and note that by Yamano [52, Equation 6]  
 566 we have  $\log_\beta x - \log x = (1 - \beta)(\partial_b \log_\beta x + \log_\beta x \log x) \geq 0 \forall \beta \in [0, 1]$ . Then we have for  
 567  $\beta \in [0, 1)$  that

$$\begin{aligned}
 |\partial_\beta H_\beta(\mathbf{x})| &= \left| \frac{-H_\beta(\mathbf{x}) - \sum_{a=1}^d \mathbf{x}^\beta(a) \log \mathbf{x}(a)}{1 - \beta} \right| \\
 &= \frac{1}{1 - \beta} \left| \sum_{a=1}^d \mathbf{x}^\beta(a) (\log_\beta \mathbf{x}(a) - \log \mathbf{x}(a)) \right| \\
 &= \frac{1}{1 - \beta} \sum_{a=1}^d \mathbf{x}^\beta(a) (\log_\beta \mathbf{x}(a) - \log \mathbf{x}(a)) \\
 &\leq \frac{1}{1 - \beta} \left( \sum_{a=1}^d \mathbf{x}(a) \right)^\beta \left( \sum_{a=1}^d (\log_\beta \mathbf{x}(a) - \log \mathbf{x}(a))^{\frac{1}{1-\beta}} \right)^{1-\beta} \\
 &\leq \frac{1}{1 - \beta} \sum_{a=1}^d \log_\beta \mathbf{x}(a) - \log \mathbf{x}(a) \leq \frac{d}{1 - \beta} (\log_\beta \frac{d}{\varepsilon} - \log \frac{d}{\varepsilon}) \leq -d \log \frac{d}{\varepsilon}
 \end{aligned} \tag{32}$$

568 where the fourth inequality follows by Hölder's inequality, the fifth by subadditivity of  $x^a$  for  
 569  $a \in (0, 1]$ , the sixth by the fact that  $\partial_x (\log_\beta x - \log x) = x^{-\beta} - 1/x \leq 0 \forall \beta, x \in [0, 1]$ , and  
 570 the last line by substituting  $\beta = 0$  since  $\partial_\beta \left( \frac{\log_\beta x - \log x}{1 - \beta} \right) = \frac{2(x - x^\beta) - (1 - \beta)(x^\beta + x) \log x}{x^\beta (1 - \beta)^3} \leq 0 \forall \beta \in$   
 571  $[0, 1), x \in (0, 1/d]$ . For  $\beta = 1$ , applying L'Hôpital's rule yields

$$\lim_{\beta \rightarrow 1} \partial_\beta H_\beta(\mathbf{x}) = -\frac{1}{2} \lim_{\beta \rightarrow 1} \sum_{a=1}^d \mathbf{x}^\beta(a) \log^2 \mathbf{x}(a) (1 - (1 - \beta) \log \mathbf{x}(a)) = -\frac{1}{2} \sum_{a=1}^d \mathbf{x}(a) \log^2 \mathbf{x}(a) \tag{33}$$

572 which is bounded on  $[-2d/e^2, 0]$ .  $\square$

573 **Lemma B.2.** Consider  $\mathbf{x}_1, \dots, \mathbf{x}_T \in \Delta$  s.t.  $\mathbf{x}_t(a_t) = 1$  for some  $a_t \in [d]$ , and let  $\bar{\mathbf{x}} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$   
574 be their average. For any  $\varepsilon \in (0, 1]$  and  $\beta \in (0, 1]$  we have that for every  $t \in [T]$

$$H_\beta(\bar{\mathbf{x}}^{(\varepsilon)}) - H_\beta(\mathbf{x}_t^{(\varepsilon)}) \leq H_\beta(\bar{\mathbf{x}}) \quad (34)$$

575 where recall that  $\mathbf{x}^{(\varepsilon)} = \mathbf{c}_{\frac{\varepsilon}{1-\varepsilon}}(\mathbf{x}) = \mathbf{1}_d/d + (1-\varepsilon)(\mathbf{x} - \mathbf{1}_d/d) = (1-\varepsilon)\mathbf{x} + \frac{\varepsilon}{d}\mathbf{1}_d$ .

576 *Proof.* Assume w.l.o.g. that  $\bar{\mathbf{x}}(1) \leq \bar{\mathbf{x}}(2) \leq \dots \leq \bar{\mathbf{x}}(d)$  and  $a_t = 1$ , so that  $\mathbf{x}_t^{(\varepsilon)} = \mathbf{e}_1^{(\varepsilon)}$ . We take the  
577 derivative

$$\begin{aligned} & \partial_\varepsilon H_\beta \left( (1-\varepsilon)\bar{\mathbf{x}} + \frac{\varepsilon}{d}\mathbf{1}_d \right) - \partial_\varepsilon H_\beta \left( \mathbf{e}_1^{(\varepsilon)} \right) \\ &= \frac{d}{1-\beta} \sum_{a=1}^{d-1} \left( \frac{1}{((1-\varepsilon)\bar{\mathbf{x}}(a) + \varepsilon/d)^{1-\beta}} - \frac{1}{(\varepsilon/d)^{1-\beta}} \right) \\ &+ \frac{d}{1-\beta} \sum_{a=1}^{d-1} \left( \frac{1}{((1-\varepsilon) + \varepsilon/d)^{1-\beta}} - \frac{1}{((1-\varepsilon)\bar{\mathbf{x}}(d) + \varepsilon/d)^{1-\beta}} \right) \\ &+ \frac{d^2}{1-\beta} \sum_{a=1}^{d-1} \bar{\mathbf{x}}(a) \left( \frac{1}{((1-\varepsilon)\bar{\mathbf{x}}(d) + \varepsilon/d)^{1-\beta}} - \frac{1}{((1-\varepsilon)\bar{\mathbf{x}}(a) + \varepsilon/d)^{1-\beta}} \right) \end{aligned} \quad (35)$$

578 By the assumption that  $\bar{\mathbf{x}}(a)$  is non-decreasing in  $a$ , each of the summands above become non-positive.  
579 So for  $\varepsilon \in (0, 1]$  the derivative is non-positive, and for  $\varepsilon \rightarrow 0^+$  it goes to  $-\infty$ . Thus the l.h.s. of the  
580 bound is monotonically non-increasing in  $\varepsilon$  for all  $\varepsilon \in [0, 1]$ . The result then follows from the fact  
581 that for  $\varepsilon = 0$  we have  $H_\beta((1-\varepsilon)\bar{\mathbf{x}} + \frac{\varepsilon}{d}\mathbf{1}_d) - H_\beta(\mathbf{e}_1^{(\varepsilon)}) = H_\beta(\bar{\mathbf{x}})$ .  $\square$

## 582 B.2 Implicit exploration bounds

583 **Lemma B.3.** Suppose we play  $\text{OMD}_{\beta, \eta}$  with regularizer  $\psi_\beta$  the negative Tsallis entropy and initial-  
584 ization  $\mathbf{x}_1 \in \Delta$  on the sequence of linear loss functions  $\ell_1, \dots, \ell_T \in [0, 1]^d$ . Then for any  $\mathbf{x} \in \Delta$  we  
585 have

$$\sum_{t=1}^T \langle \ell_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \frac{B_\beta(\mathbf{x}|\mathbf{x}_1)}{\eta} + \frac{\eta}{\beta} \sum_{a=1}^d \mathbf{x}_t^{2-\beta}(a) \ell_t^2(a) \quad (36)$$

586 *Proof.* Note that the following proof follows parts of the course notes by Luo [37], which we  
587 reproduce for completeness. The OMD update at each step  $t$  involves the following two steps: set  
588  $\mathbf{y}_{t+1} \in \Delta$  s.t.  $\nabla \phi_\beta(\mathbf{y}_{t+1}) = \nabla \phi_\beta(\mathbf{x}_t) - \eta \ell_t$  and then set  $\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \Delta} B_\beta(\mathbf{x}, \mathbf{y}_{t+1})$  [25,  
589 Algorithm 14]. Note that by Hazan [25, Equation 5.3] and nonnegativity of the Bregman divergence  
590 we have

$$\sum_{t=1}^T \langle \ell_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \frac{B_\beta(\mathbf{x}|\mathbf{x}_1)}{\eta} + \frac{1}{\eta} \sum_{t=1}^T B_\beta(\mathbf{x}_t|\mathbf{y}_{t+1}) \quad (37)$$

591 To bound the second term, note that when  $\phi_\beta$  is the negative Tsallis entropy we have

$$\begin{aligned} & B_\beta(\mathbf{x}_t|\mathbf{y}_{t+1}) \\ &= \frac{1}{1-\beta} \sum_{a=1}^d \left( \mathbf{y}_{t+1}^\beta(a) - \mathbf{x}_t^\beta(a) + \frac{\beta}{\mathbf{y}_{t+1}^{1-\beta}(a)} (\mathbf{x}_t(a) - \mathbf{y}_{t+1}(a)) \right) \\ &= \frac{1}{1-\beta} \sum_{a=1}^d \left( (1-\beta)\mathbf{y}_{t+1}^\beta(a) - \mathbf{x}_t^\beta(a) + \beta \left( \frac{1}{\mathbf{x}_t^{1-\beta}(a)} + \frac{1-\beta}{\beta} \eta \ell_t(a) \right) \mathbf{x}_t(a) \right) \\ &= \sum_{a=1}^d \left( \mathbf{y}_{t+1}^\beta(a) - \mathbf{x}_t^\beta(a) + \eta \mathbf{x}_t(a) \ell_t(a) \right) \end{aligned} \quad (38)$$

592 Plugging the following result, which follows from  $(1+x)^\alpha \leq 1 + \alpha x + \alpha(\alpha-1)x^2 \forall x \geq 0, \alpha < 0$ ,  
 593 into the above yields the desired bound.

$$\begin{aligned} \mathbf{y}_{t+1}^\beta(a) &= \mathbf{x}_t^\beta(a) \left( \frac{\mathbf{y}_{t+1}^{\beta-1}(a)}{\mathbf{x}_t^{\beta-1}(a)} \right)^{\frac{\beta}{\beta-1}} = \mathbf{x}_t^\beta(a) \left( 1 + \frac{1-\beta}{\beta} \eta \mathbf{x}_t^{1-\beta}(a) \ell_t(a) \right)^{\frac{\beta}{\beta-1}} \\ &\leq \mathbf{x}_t^\beta(a) \left( 1 - \eta \mathbf{x}_t^{1-\beta}(a) \ell_t(a) + \frac{\eta^2}{\beta} \mathbf{x}_t^{2-2\beta}(a) \ell_t(a)^2 \right) \quad (39) \\ &= \mathbf{x}_t^\beta(a) - \eta \mathbf{x}_t(a) \ell_t(a) + \frac{\eta^2}{\beta} \mathbf{x}_t^{2-\beta}(a) \ell_t(a)^2 \end{aligned}$$

594

□

595 **Theorem B.1.** In Algorithm 1, let  $OMD_{\eta,\beta}$  be online mirror descent with the Tsallis entropy regular-  
 596 izer  $\psi_\beta$  over  $\gamma$ -offset loss estimators,  $\Theta_k$  is a subset of  $[\underline{\beta}, \bar{\beta}] \subset [\frac{1}{\log d}, 1]$ , and

$$U_t(\mathbf{x}, \eta, \beta) = \frac{B_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x})}{\eta} + \frac{\eta d^\beta m}{\beta} \quad (40)$$

597 where  $\hat{\mathbf{x}}_t^{(\varepsilon)} = (1-\varepsilon)\hat{\mathbf{x}}_t + \varepsilon \mathbf{1}_d/d$ . Note that  $U_t^{(\rho)}(\mathbf{x}, \eta, \beta) = U_t(\mathbf{x}, \eta, \beta) + \frac{\rho^2(d^{1-\beta}-1)}{\eta(1-\beta)}$ . Then there  
 598 exists settings of  $\underline{\eta}, \bar{\eta}, \alpha, \lambda$  s.t. for all  $\varepsilon, \rho, \gamma \in (0, 1)$  we have w.p.  $\geq 1 - \delta$  that

$$\begin{aligned} &\sum_{t=1}^T \sum_{i=1}^m \ell_{t,i}(a_{t,i}) - \ell_{t,i}(\hat{a}_t) \\ &\leq (\varepsilon + \gamma d)mT + \frac{2 + \sqrt{\frac{d \log d}{\varepsilon m}}}{\gamma} \log \frac{5}{\delta} + \frac{8d\sqrt{m}}{\rho} \left( 1_{k>1} \sqrt{T \log \frac{5k}{\delta}} + \frac{1 + \log(T+1)}{16\rho} \right) \\ &\quad + \min_{\beta \in [\underline{\beta}, \bar{\beta}], \eta > 0} \frac{8 \left(\frac{d}{\varepsilon}\right)^{2-\beta} (1 + \log T)}{\eta} + \left( \frac{\hat{H}_\beta}{\eta} + \frac{\eta d^\beta m}{\beta} + \frac{L_\eta(\bar{\beta} - \beta)}{2k} + d \min \left\{ \frac{\rho^2}{2\eta}, \rho\sqrt{m} \right\} \right) T \quad (41) \end{aligned}$$

599 for  $L_\eta = \left( \frac{\log \frac{d}{\varepsilon}}{\eta} + \eta m \log^2 d \right) d$ .

600 *Proof.* In this setting we have  $g(\beta) = d^\beta/\beta$ ,  $f(\beta) = 0$ ,  $D_\beta^2 = \frac{d^{1-\beta}-1}{1-\beta}$ ,  $D \leq \sqrt{d/2}$ ,  $M = d\sqrt{m}$ ,  
 601  $F = 0$ ,  $S = (d/\varepsilon)^{2-\beta}$ , and  $K = 1$ . We have that

$$\begin{aligned} &\sum_{t=1}^T \sum_{i=1}^m \ell_{t,i}(a_{t,i}) - \ell_{t,i}(\hat{a}_t) \\ &= \sum_{t=1}^T \sum_{i=1}^m \langle \hat{\ell}_{t,i}, \mathbf{x}_{t,i} \rangle - \ell_{t,i}(\hat{a}_t) + \gamma \sum_{a=1}^d \hat{\ell}_{t,i}(a) \\ &\leq \sum_{t=1}^T \frac{B_{\beta_t}(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x}_{t,1})}{\eta_t} + \sum_{i=1}^m \langle \hat{\ell}_{t,i}, \hat{\mathbf{x}}_t^{(\varepsilon)} \rangle - \ell_{t,i}(\hat{a}_t) + \frac{\eta_t}{\beta_t} \sum_{a=1}^d \mathbf{x}_{t,i}^{2-\beta_t}(a) \hat{\ell}_{t,i}^2(a) + \gamma \sum_{a=1}^d \hat{\ell}_{t,i}(a) \\ &\leq \varepsilon mT + \sum_{t=1}^T \frac{B_{\beta_t}(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x}_{t,1})}{\eta_t} + \sum_{i=1}^m \langle \hat{\ell}_{t,i}, \hat{\mathbf{x}}_t^{(\varepsilon)} \rangle - \langle \ell_{t,i}, \hat{\mathbf{x}}_t^{(\varepsilon)} \rangle \\ &\quad + \sum_{t=1}^T \frac{\eta_t}{\beta_t} \sum_{i=1}^m \sum_{a=1}^d \mathbf{x}_{t,i}^{1-\beta_t}(a) \hat{\ell}_{t,i}(a) + \gamma \sum_{a=1}^d \hat{\ell}_{t,i}(a) \quad (42) \end{aligned}$$

602 where the equality follows similarly to Luo [37] since  $\langle \hat{\ell}_{t,i}, \mathbf{x}_{t,i} \rangle = \ell_{t,i}(a_{t,i}) - \gamma \sum_{a=1}^d \hat{\ell}_{t,i}(a)$ , the  
 603 first inequality follows by Lemma B.3 and the second by Hölder's inequality and the definitions of

604  $\hat{\ell}_{t,i}$  and  $\hat{\mathbf{x}}_{t,i}^{(\varepsilon)}$ . We next apply the optimality of  $\hat{a}_t$  for  $\sum_{i=1}^m \hat{\ell}_{t,i}$  to get

$$\begin{aligned}
& \sum_{t=1}^T \sum_{i=1}^m \ell_{t,i}(a_{t,i}) - \ell_{t,i}(\hat{a}_t) \\
& \leq \varepsilon mT + \sum_{t=1}^T \frac{B_{\beta_t}(\hat{\mathbf{x}}_t^{(\varepsilon)} \|\mathbf{x}_{t,1})}{\eta_t} + (1 - \varepsilon) \sum_{i=1}^m \hat{\ell}_{t,i}(\hat{a}_t) - \ell_{t,i}(\hat{a}_t) + \frac{\varepsilon}{d} \sum_{a=1}^d \hat{\ell}_{t,i}(a) - \ell_{t,i}(a) \\
& \quad + \sum_{t=1}^T \frac{\eta_t}{\beta_t} \sum_{i=1}^m \sum_{a=1}^d \mathbf{x}_{t,i}^{1-\beta_t}(a) \hat{\ell}_{t,i}(a) + \gamma \sum_{a=1}^d \hat{\ell}_{t,i}(a) \\
& \leq \varepsilon mT + \frac{1 + \frac{\varepsilon}{d} + \frac{\bar{\eta}}{\beta} + \gamma}{2\gamma} \log \frac{5}{\delta} + \sum_{t=1}^T \frac{B_{\beta_t}(\hat{\mathbf{x}}_t^{(\varepsilon)} \|\mathbf{x}_{t,1})}{\eta_t} \\
& \quad + \sum_{t=1}^T \frac{\eta_t}{\beta_t} \sum_{i=1}^m \sum_{a=1}^d \mathbf{x}_{t,i}^{1-\beta_t}(a) \ell_{t,i}(a) + \gamma \sum_{a=1}^d \ell_{t,i}(a) \\
& \leq \varepsilon mT + \frac{2 + \sqrt{\frac{d \log d}{em}}}{\gamma} \log \frac{5}{\delta} + \gamma dmT + \sum_{t=1}^T \frac{B_{\beta_t}(\hat{\mathbf{x}}_t^{(\varepsilon)} \|\mathbf{x}_{t,1})}{\eta_t} + \frac{\eta_t d^{\beta_t} m}{\beta_t}
\end{aligned} \tag{43}$$

605 where the second inequality follows by Neu [42, Lemma 1] applied to each of the last four terms  
606 and the fifth by the definition of  $\ell_{t,i}$  and using  $\max_{\beta \in [\frac{1}{\log d}, 1]} \bar{\eta}(\beta) \leq \sqrt{\frac{d}{em \log d}}$ . Substituting into  
607 Theorem A.1 and simplifying yields the result except with  $\hat{V}_\beta^2 = \frac{1}{T} \sum_{t=1}^T \psi_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)}) - \psi_\beta(\hat{\mathbf{x}}^{(\varepsilon)})$  in  
608 place of  $\hat{H}_\beta$ , but the former is bounded by the latter by Lemma B.2.  $\square$

609 **Corollary B.1.** *Let  $\underline{\beta} = \bar{\beta} = 1$ . Then w.h.p. we can ensure task-averaged regret at most*

$$2\sqrt{\hat{H}_1 dm} + \tilde{\mathcal{O}} \left( \frac{d\sqrt{m} + d^{\frac{3}{2}} m^{\frac{3}{2}}}{\sqrt[3]{T}} \right) \tag{44}$$

610 so long as  $mT \geq d^2$  or alternatively ensure

$$\min \left\{ 2\sqrt{\hat{H}_1 dm} + \tilde{\mathcal{O}} \left( \frac{d^{\frac{3}{4}} m^{\frac{3}{4}} + d\sqrt{m}}{\sqrt[4]{T}} \right), 2\sqrt{dm \log d} + \tilde{\mathcal{O}} \left( \frac{d^{\frac{3}{2}} \sqrt{m}}{\sqrt{T}} \right) \right\} \tag{45}$$

611 so long as  $mT \geq d$ .

612 *Proof.* Applying Theorem B.1, simplifying, and dividing by  $T$  yields task-averaged regret at most

$$\begin{aligned}
& (\varepsilon + \gamma d)m + \frac{2 + \sqrt{\frac{d \log d}{em}}}{\gamma T} \log \frac{5}{\delta} + \left( \frac{1 + \log(T+1)}{2\rho^2 T} + \min \left\{ \frac{\rho^2}{\eta\sqrt{m}}, \rho \right\} \right) d\sqrt{m} \\
& \quad + \min_{\eta > 0} \frac{8d(1 + \log T)}{\varepsilon\eta T} + \left( \frac{\hat{H}_1}{\eta} + \eta dm \right)
\end{aligned} \tag{46}$$

613 Set  $\gamma = \frac{1}{\sqrt{dmT}}$ . Then set  $\varepsilon = \sqrt[3]{\frac{d^2}{mT}}$  and  $\rho = \frac{1}{\sqrt[3]{T}}$ , and use  $\eta = \sqrt{\frac{\hat{H}_1}{dm}} + \frac{1}{\sqrt[3]{dmT}}$  to get the first result.

614 Otherwise, set  $\varepsilon = \sqrt{\frac{d}{mT}}$  and  $\rho = \frac{1}{\sqrt[4]{T}}$ , and use the better of  $\eta = \sqrt{\frac{\hat{H}_1}{dm}} + \frac{1}{\sqrt[4]{dmT}}$  and  $\eta = \sqrt{\frac{\log d}{dm}}$   
615 to get the second.  $\square$

616 **Corollary B.2.** *Let  $\underline{\beta} = \frac{1}{2}$  and  $\bar{\beta} = 1$  and assume  $mT \geq d^{\frac{5}{2}}$ . Then w.h.p. we can ensure*  
617 *task-averaged regret at most*

$$\min_{\beta \in [\frac{1}{2}, 1]} 2\sqrt{\hat{H}_\beta d^\beta m / \beta} + \tilde{\mathcal{O}} \left( \frac{d^{\frac{5}{2}} m^{\frac{5}{2}}}{T^{\frac{2}{7}}} + \frac{d\sqrt{m}}{\sqrt[4]{T}} \right) \tag{47}$$

618 using  $k = \lceil \sqrt[4]{d\sqrt{T}} \rceil$ .

619 *Proof.* Applying Theorem B.1, simplifying, and dividing by  $T$  yields task-averaged regret at most

$$\begin{aligned}
& (\varepsilon + \gamma d)m + \frac{2 + \sqrt{\frac{d \log d}{em}}}{\gamma T} \log \frac{5}{\delta} + \frac{8d\sqrt{m}}{\rho} \left( \sqrt{\frac{\log \frac{5k}{\delta}}{T}} + \frac{1 + \log(T+1)}{16\rho T} \right) \\
& + \min_{\beta \in [\underline{\beta}, \bar{\beta}], \eta > 0} \frac{8d^{\frac{3}{2}}(1 + \log T)}{\varepsilon^{\frac{3}{2}}\eta T} + \left( \frac{\hat{H}_\beta}{\eta} + \frac{\eta d^\beta m}{\beta} + \frac{d}{4k} \left( \frac{\log \frac{d}{\varepsilon}}{\eta} + \eta m \log^2 d \right) + \rho d\sqrt{m} \right)
\end{aligned} \tag{48}$$

620 Set  $\gamma = \frac{1}{\sqrt{dmT}}$ ,  $\varepsilon = \frac{d^{\frac{5}{7}}}{(mT)^{\frac{2}{7}}}$ ,  $\rho = \frac{1}{\sqrt[4]{T}}$ , and use  $\eta = \sqrt{\frac{\beta \hat{H}_\beta}{md^\beta}} + \frac{1}{(dmT)^{\frac{2}{7}}}$  to get the result.  $\square$

621 **Corollary B.3.** Let  $\underline{\beta} = \frac{1}{\log d}$  and  $\bar{\beta} = 1$  and assume  $mT \geq d^3$ . Then w.h.p. we can ensure  
622 task-averaged regret at most

$$\min_{\beta \in (0,1]} 2\sqrt{\hat{H}_\beta d^\beta m / \beta} + \tilde{O} \left( \frac{d^{\frac{3}{4}} m^{\frac{3}{4}} + d\sqrt{m}}{\sqrt[4]{T}} \right) \tag{49}$$

623 using  $k = \lceil \sqrt[4]{d\sqrt{T}} \rceil$ .

624 *Proof.* Applying Theorem B.1, dividing by  $T$ , and simplifying yields

$$\begin{aligned}
& (\varepsilon + \gamma d)m + \frac{2 + \sqrt{\frac{d \log d}{em}}}{\gamma T} \log \frac{5}{\delta} + \frac{8d\sqrt{m}}{\rho} \left( \sqrt{\frac{\log \frac{5k}{\delta}}{T}} + \frac{1 + \log(T+1)}{16\rho T} \right) \\
& + \min_{\beta \in [\underline{\beta}, \bar{\beta}], \eta > 0} \frac{8d^2(1 + \log T)}{\varepsilon^2\eta T} + \left( \frac{\hat{H}_\beta}{\eta} + \frac{\eta d^\beta m}{\beta} + \frac{d}{2k} \left( \frac{\log \frac{d}{\varepsilon}}{\eta} + \eta \log^2 d \right) + \rho d\sqrt{m} \right)
\end{aligned} \tag{50}$$

625 Note that  $\hat{H}_\beta$  and  $\frac{d^\beta}{\beta}$  are both decreasing on  $\beta < \frac{1}{\log d}$ , so  $\beta$  in the chosen interval is optimal over all  
626  $\beta \in (0, 1]$ . Set  $\gamma = \frac{1}{\sqrt{dmT}}$ ,  $\varepsilon = \frac{d^{\frac{3}{4}}}{\sqrt[4]{mT}}$ ,  $\rho = \frac{1}{\sqrt[4]{T}}$ , and use  $\eta = \sqrt{\frac{\beta \hat{H}_\beta}{md^\beta}} + \frac{1}{\sqrt[4]{dmT}}$  to get the result.  $\square$

## 627 C Guaranteed exploration

### 628 C.1 Best-arm identification

629 **Lemma C.1.** Suppose for  $\varepsilon > 0$  we run OMD on task  $t \in [T]$  with initialization  $\mathbf{x}_{t,1} \in \Delta^{(\varepsilon)}$ ,  
630 regularizer  $\psi_{\beta_t} + I_{\Delta^{(\varepsilon)}}$  for some  $\beta_t \in (0, 1]$ , and unbiased loss estimators ( $\gamma = 0$ ). If Assumption 3.1  
631 holds and  $m > \frac{28d \log d}{3\varepsilon \Delta^2}$  then  $\hat{\mathbf{x}}_t = \hat{\mathbf{x}}_t$  w.p.  $\geq 1 - d\kappa$ , where  $\kappa = \exp\left(-\frac{3\varepsilon \Delta^2 m}{28d}\right)$ .

632 *Proof.* We extend the proof by Abbasi-Yadkori et al. [1, Appendices B and F] to arbitrary lower  
633 bounds  $\varepsilon/d$  on the probability. First, since  $0 \leq \hat{\ell}_{t,i}(a) \leq \frac{d}{\varepsilon} \ell_{t,i}(a)$  we have that

$$-\frac{d}{\varepsilon} \leq -1 \leq -\ell_{t,i}(a) \leq \hat{\ell}_{t,i}(a) - \ell_{t,i}(a) \leq \left( \frac{d}{\varepsilon} - 1 \right) \ell_{t,i}(a) \leq \frac{d}{\varepsilon} \tag{51}$$

634 and so  $|\hat{\ell}_{t,i}(a) - \ell_{t,i}(a)| \leq \frac{d}{\varepsilon}$ . Therefore since the variance of the estimated losses is a scaled  
635 Bernoulli we have that

$$\text{Var}(\hat{\ell}_{t,i}(a) - \ell_{t,i}(a)) = \text{Var}(\hat{\ell}_{t,i}(a)) = \mathbf{x}_{t,i}(a)(1 - \mathbf{x}_{t,i}(a)) \left( \frac{\ell_{t,i}(a)}{\mathbf{x}_{t,i}(a)} \right)^2 \leq \frac{\ell_{t,i}^2(a)}{\mathbf{x}_{t,i}(a)} \leq \frac{d}{\varepsilon} \tag{52}$$

636 We can thus apply a martingale concentration inequality of Fan et al. [23, Corollary 2.1] to the  
 637 martingale difference sequence (MDS)  $\frac{\varepsilon}{d}(\hat{\ell}_{t,i}(a) - \ell_{t,i}(a)) \in [-\frac{\varepsilon}{d}, 1]$  to obtain

$$\begin{aligned}
 \Pr\left(\sum_{i=1}^m \hat{\ell}_{t,i}(a) - \ell_{t,i}(a) \geq \frac{m\Delta_a}{2}\right) &= \Pr\left(\frac{\varepsilon}{d} \sum_{i=1}^m \hat{\ell}_{t,i}(a) - \ell_{t,i}(a) \geq \frac{\varepsilon m\Delta_a}{2d}\right) \\
 &\leq \Pr\left(\max_{j \in [m]} \frac{\varepsilon}{d} \sum_{i=j}^m \hat{\ell}_{t,i}(a) - \ell_{t,i}(a) \geq \frac{\varepsilon m\Delta_a}{2d}\right) \\
 &\leq \exp\left(-\frac{2\left(\frac{\varepsilon m\Delta_a}{2d}\right)^2}{\min\{m(1 + \varepsilon/d)^2, 4(\varepsilon m/d + \frac{\varepsilon m\Delta_a}{6})\}}\right) \quad (53) \\
 &\leq \exp\left(-\frac{2\left(\frac{\varepsilon m\Delta_a}{2d}\right)^2}{4(\varepsilon m/d + \frac{\varepsilon m\Delta_a}{6})}\right) \\
 &= \exp\left(-\frac{3\varepsilon m\Delta_a^2}{4d(6 + \Delta_a)}\right) \\
 &\leq \exp\left(-\frac{3\varepsilon m\Delta_a^2}{28d}\right)
 \end{aligned}$$

638 where  $\Delta_a = \frac{1}{m} |\sum_{i=1}^m \ell_{t,i}(a) - \min_{a' \neq a} \sum_{i=1}^m \ell_{t,i}(a')|$  is the per-arm loss gap in the last step we  
 639 apply  $\Delta_a \leq 1$ . For the symmetric MDS  $-\frac{\varepsilon}{d} \leq \ell_{t,i}(a) - \hat{\ell}_{t,i}(a) \leq 1$  we have

$$\begin{aligned}
 \Pr\left(\sum_{i=1}^m \hat{\ell}_{t,i}(a) - \ell_{t,i}(a) \leq -\frac{m\Delta_a}{2}\right) &= \Pr\left(\sum_{i=1}^m \ell_{t,i}(a) - \hat{\ell}_{t,i}(a) \geq \frac{m\Delta_a}{2}\right) \\
 &\leq \exp\left(-\frac{2\left(\frac{m\Delta_a}{2}\right)^2}{4\left(\frac{dm}{\varepsilon} + \frac{m\Delta_a}{6}\right)}\right) \quad (54) \\
 &\leq \exp\left(-\frac{3\varepsilon m\Delta_a^2/d}{4(6 + \varepsilon\Delta_a/d)}\right) \\
 &\leq \exp\left(-\frac{3\varepsilon m\Delta_a^2}{28d}\right)
 \end{aligned}$$

640 We can then conclude that

$$\begin{aligned}
 &\Pr(\hat{\mathbf{x}}_t \neq \mathring{\mathbf{x}}_t) \\
 &\leq \Pr\left(\exists a \neq \hat{a}_t : \sum_{i=1}^m \hat{\ell}_{t,i}(a) \leq \sum_{i=1}^m \ell_{t,i}(\hat{a}_t)\right) \\
 &\leq \Pr\left(\sum_{i=1}^m \hat{\ell}_{t,i}(\hat{a}_t) \geq \sum_{i=1}^m \ell_{t,i}(\hat{a}_t) + \frac{m\Delta_{\hat{a}_t}}{2} \vee \exists a \neq \hat{a}_t : \sum_{i=1}^m \hat{\ell}_{t,i}(a) \leq \sum_{i=1}^m \ell_{t,i}(a) - \frac{m\Delta_a}{2}\right) \\
 &\leq \Pr\left(\sum_{i=1}^m \hat{\ell}_{t,i}(\hat{a}_t) \geq \sum_{i=1}^m \ell_{t,i}(\hat{a}_t) + \frac{m\Delta_{\hat{a}_t}}{2}\right) + \sum_{a \neq \hat{a}_t} \Pr\left(\sum_{i=1}^m \hat{\ell}_{t,i}(a) \leq \sum_{i=1}^m \ell_{t,i}(a) - \frac{m\Delta_a}{2}\right) \\
 &\leq \exp\left(-\frac{3\varepsilon m\Delta_{\hat{a}_t}^2}{28d}\right) + \sum_{a \neq \hat{a}_t} \exp\left(-\frac{3\varepsilon m\Delta_a^2}{28d}\right) \\
 &\leq d \exp\left(-\frac{3\varepsilon m\Delta^2}{28d}\right) \quad (55)
 \end{aligned}$$

641 where the second-to-last line follows by substituting the bounds (53) and (54) into the left and right  
 642 terms, respectively.  $\square$

643 **Lemma C.2.** *Suppose on each task  $t \in [T]$  we run OMD as in Lemma C.1. Then for any  $\beta \in (0, 1]$*   
 644 *we have  $\frac{1}{T} \mathbb{E} \sum_{t=1}^T \psi_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)}) - \psi_\beta(\hat{\mathbf{x}}^{(\varepsilon)}) \leq -\psi_\beta(\hat{\mathbf{x}}) + \frac{3d\kappa\beta}{1-\beta} \left( \left(\frac{d}{\varepsilon}\right)^{1-\beta} - 1 \right)$ .*

645 *Proof.* We consider the expected divergence of the best initialization under the worst-case distribution  
 646 of best arm estimation, which satisfies Lemma C.1 and (55). We have by Claim A.1 and the mean-as-  
 647 minimizer property of Bregman divergences that

$$\begin{aligned}
\frac{1}{T} \mathbb{E} \sum_{t=1}^T \psi_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)}) - \psi_\beta(\hat{\mathbf{x}}^{(\varepsilon)}) &= \mathbb{E} \min_{\mathbf{x} \in \Delta^{(\varepsilon)}} \frac{1}{T} \sum_{t=1}^T B_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x}) \\
&\leq \min_{\mathbf{x} \in \Delta^{(\varepsilon)}} \mathbb{E} \frac{1}{T} \sum_{t=1}^T B_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x}) \\
&= \min_{\mathbf{x} \in \Delta^{(\varepsilon)}} \frac{1}{T} \sum_{t=1}^T \sum_{a=1}^d \mathbb{P}(a = \hat{a}_t) B_\beta(\mathbf{e}_a^{(\varepsilon)} || \mathbf{x}) \\
&\leq \max_{\substack{\mathbf{p}_t \in \Delta, \forall t \in [T] \\ \mathbf{p}_t(a) \leq 2\kappa, \forall t \in [T], a \neq \hat{a}_t \\ 1-d\kappa \leq \mathbf{p}_t(a), \forall t \in [T], a = \hat{a}_t}} \min_{\mathbf{x} \in \Delta^{(\varepsilon)}} \frac{1}{T} \sum_{t=1}^T \sum_{a=1}^d \mathbf{p}_t(a) B_\beta(\mathbf{e}_a^{(\varepsilon)} || \mathbf{x})
\end{aligned} \tag{56}$$

648 To simplify the last expression, we define  $\bar{\mathbf{p}} = \frac{1}{T} \sum_{t=1}^T \mathbf{p}_t$  and again apply the (weighted) mean-as-  
 649 minimizer property, followed by Claim A.1:

$$\begin{aligned}
\min_{\mathbf{x} \in \Delta^{(\varepsilon)}} \frac{1}{T} \sum_{t=1}^T \sum_{a=1}^d \mathbf{p}_t(a) B_\beta(\mathbf{e}_a^{(\varepsilon)} || \mathbf{x}) &= \min_{\mathbf{x} \in \Delta^{(\varepsilon)}} \sum_{a=1}^d \bar{\mathbf{p}}(a) B_\beta(\mathbf{e}_a^{(\varepsilon)} || \mathbf{x}) = \sum_{a=1}^d B_\beta(\mathbf{e}_a^{(\varepsilon)} || \bar{\mathbf{p}}^{(\varepsilon)}) \\
&= \psi_\beta(\mathbf{e}_1^{(\varepsilon)}) - \psi_\beta(\bar{\mathbf{p}}^{(\varepsilon)})
\end{aligned} \tag{57}$$

650 By substituting into the previous inequality, we can bound the expected divergence for the worst-case  
 651  $\mathbf{p}_t$  as follows:

$$\begin{aligned}
\frac{1}{T} \mathbb{E} \sum_{t=1}^T \psi_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)}) - \psi_\beta(\hat{\mathbf{x}}^{(\varepsilon)}) &\leq \psi_\beta(\mathbf{e}_1^{(\varepsilon)}) + \max_{\substack{\mathbf{p}_t \in \Delta, \forall t \in [T] \\ \mathbf{p}_t(a) \leq 2\kappa, \forall t \in [T], a \neq \hat{a}_t \\ 1-d\kappa \leq \mathbf{p}_t(a), \forall t \in [T], a = \hat{a}_t}} -\psi_\beta(\bar{\mathbf{p}}^{(\varepsilon)}) \\
&\leq \psi_\beta(\mathbf{e}_1^{(\varepsilon)}) + \max_{\substack{\sum_{t=1}^T \sum_{a=1}^d \mathbf{p}_t(a) = T \\ \sum_{t=1}^T \mathbf{p}_t(a) \geq (1-d\kappa)\hat{\mathbf{x}}(a)T, \forall a \\ \sum_{t=1}^T \mathbf{p}_t(a) \leq (2\kappa(1-\hat{\mathbf{x}}(a))T + \hat{\mathbf{x}}(a)T), \forall a}} -\psi_\beta(\bar{\mathbf{p}}^{(\varepsilon)}) \\
&= \psi_\beta(\mathbf{e}_1^{(\varepsilon)}) - \min_{\substack{\bar{\mathbf{p}} \in \Delta \\ \bar{\mathbf{p}}(a) \geq (1-d\kappa)\hat{\mathbf{x}}(a), \forall a \\ \bar{\mathbf{p}}(a) \leq 2\kappa + (1-2\kappa)\hat{\mathbf{x}}(a), \forall a}} \psi_\beta(\bar{\mathbf{p}}^{(\varepsilon)})
\end{aligned} \tag{58}$$

652 We use the shorthand  $h(\mathbf{x}) = \psi_\beta((1-\varepsilon)\mathbf{x} + \frac{\varepsilon}{d}\mathbf{1}_d)$ . We have

$$\begin{aligned}
-\partial_{\mathbf{x}(a)}(\psi_\beta(\mathbf{x})) &= \partial_{\mathbf{x}(a)} \left( \frac{1}{(1-\beta)} \left( \sum_{b=1}^d \mathbf{x}(b)^\beta - 1 \right) \right) \\
&= \partial_{\mathbf{x}(a)} \left( \frac{1}{(1-\beta)} \left( \sum_{b=1}^d \mathbf{x}(b)^\beta + \beta d^{1-\beta} \left( 1 - \sum_{b=1}^d \mathbf{x}(b) \right) - 1 \right) \right) \\
&= \frac{\beta}{1-\beta} \cdot (\mathbf{x}(a)^{\beta-1} - d^{1-\beta})
\end{aligned} \tag{59}$$

653 and therefore

$$\begin{aligned}
\|\nabla h(\mathbf{x})\|_\infty &= \max_{a=1,\dots,d} \left| \partial_{\mathbf{x}(a)} \psi_\beta \left( (1-\varepsilon)\mathbf{x} + \frac{\varepsilon}{d} \mathbf{1}_d \right) \right| \\
&\leq \frac{\beta}{1-\beta} \max_{a=1,\dots,d} \left| ((1-\varepsilon)\mathbf{x}(a) + \varepsilon/d)^{\beta-1} - d^{1-\beta} \right| \\
&\leq \frac{\beta}{1-\beta} \left( \left( \frac{d}{\varepsilon} \right)^{1-\beta} - 1 \right) = \beta \log_\beta \left( \frac{d}{\varepsilon} \right)
\end{aligned} \tag{60}$$

654 Finally, by convexity of  $h$  we have

$$\begin{aligned}
\min_{\substack{\bar{\mathbf{p}} \in \Delta \\ \bar{\mathbf{p}}(a) \geq (1-d\kappa)\hat{\mathbf{x}}(a), \forall a \\ \bar{\mathbf{p}}(a) \leq 2\kappa + (1-2\kappa)\hat{\mathbf{x}}(a), \forall a}} h(\bar{\mathbf{p}}) &\geq h(\hat{\mathbf{x}}) - \|\nabla h(\hat{\mathbf{x}})\|_\infty \max_{\substack{\bar{\mathbf{p}} \in \Delta \\ \bar{\mathbf{p}}(a) \geq (1-d\kappa)\hat{\mathbf{x}}(a), \forall a \\ \bar{\mathbf{p}}(a) \leq 2\kappa + (1-2\kappa)\hat{\mathbf{x}}(a), \forall a}} \|\bar{\mathbf{p}} - \hat{\mathbf{x}}\|_1 \\
&\geq h(\hat{\mathbf{x}}) - 3d\kappa \|\nabla h(\hat{\mathbf{x}})\|_\infty \\
&\geq h(\hat{\mathbf{x}}) - 3d\kappa\beta \log_\beta \left( \frac{d}{\varepsilon} \right)
\end{aligned} \tag{61}$$

655 so we can substitute into (58) to get

$$\frac{1}{T} \mathbb{E} \sum_{t=1}^T \psi_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)}) - \psi_\beta(\hat{\mathbf{x}}^{(\varepsilon)}) \leq -\psi_\beta(\hat{\mathbf{x}}^{(\varepsilon)}) + \frac{3d\kappa\beta}{1-\beta} \left( \left( \frac{d}{\varepsilon} \right)^{1-\beta} - 1 \right) \tag{62}$$

656 Applying Lemma B.2 completes the proof.  $\square$

## 657 C.2 Guaranteed exploration bounds

658 **Lemma C.3.** Suppose we play  $OMD_{\beta,\eta}$  with initialization  $\mathbf{x}_1 \in \Delta^{(\varepsilon)}$ , regularizer  $\psi_\beta + I_{\Delta^{(\varepsilon)}}$  for some  
659  $\beta \in (0, 1]$ , and unbiased loss estimators ( $\gamma = 0$ ) on the sequence of loss functions  $\ell_1, \dots, \ell_T \in [0, 1]^d$ .  
660 Then for any  $\hat{a} \in [d]$  we have expected regret

$$\mathbb{E} \sum_{t=1}^T \ell_t(a_t) - \ell_t(\hat{a}) \leq \frac{\mathbb{E} B_\beta(\hat{\mathbf{x}}^{(\varepsilon)} \mid \mathbf{x}_1)}{\eta} + \frac{\eta d^\beta m}{\beta} + \varepsilon m \tag{63}$$

661 for  $\hat{\mathbf{x}}$  the estimated optimum of the loss estimators  $\hat{\ell}_1, \dots, \hat{\ell}_T$ .

*Proof.*

$$\begin{aligned}
\mathbb{E} \sum_{t=1}^T \ell_t(a_t) - \ell_t(\hat{a}) &= \mathbb{E} \sum_{t=1}^T \ell_t(a_t) - \langle \ell_t, \hat{\mathbf{x}} \rangle \\
&\leq \mathbb{E} \sum_{t=1}^T \ell_t(a_t) - \langle \ell_t, \hat{\mathbf{x}}^{(\varepsilon)} \rangle + \varepsilon m \\
&= \mathbb{E} \sum_{t=1}^m \hat{\ell}_t(a_t) - \langle \hat{\ell}_t, \hat{\mathbf{x}}^{(\varepsilon)} \rangle + \varepsilon m \\
&\leq \mathbb{E} \sum_{t=1}^m \hat{\ell}_t(a_t) - \langle \hat{\ell}_t, \hat{\mathbf{x}}^{(\varepsilon)} \rangle + \varepsilon m \\
&\leq \mathbb{E} \left( \frac{B_\beta(\hat{\mathbf{x}}^{(\varepsilon)} \mid \mathbf{x}_1)}{\eta} + \frac{\eta}{\beta} \sum_{t=1}^T \sum_{a=1}^d \hat{\ell}_t^2(a) \mathbf{x}_t^{2-\beta}(a) \right) + \varepsilon m \\
&\leq \frac{\mathbb{E} B_\beta(\hat{\mathbf{x}}^{(\varepsilon)} \mid \mathbf{x}_1)}{\eta} + \frac{\eta d^\beta m}{\beta} + \varepsilon m
\end{aligned} \tag{64}$$

662 where the second inequality follows by optimality of  $\hat{\mathbf{x}}$  for the estimated losses  $\hat{\ell}_t$ , the third by  
663 Lemma B.3 constrained to  $\Delta^{(\varepsilon)}$ , and the fourth similarly to Theorem B.1 (note both are also  
664 effectively shown in Luo [37]).  $\square$

665 **Theorem C.1.** In Algorithm 1, let  $\text{OMD}_{\eta,\beta}$  be online mirror descent with the regularizer  $\psi_\beta + I_{\Delta(\varepsilon)}$   
666 over unbiased ( $\gamma = 0$ ) loss estimators,  $\Theta_k$  is a subset of  $[\underline{\beta}, \bar{\beta}] \subset [\frac{1}{\log d}, 1]$ , and

$$U_t(\mathbf{x}, \eta, \beta) = \frac{B_\beta(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x})}{\eta} + \frac{\eta d^\beta m}{\beta} \quad (65)$$

667 where  $\hat{\mathbf{x}}_t^{(\varepsilon)} = (1 - \varepsilon)\hat{\mathbf{x}}_t + \varepsilon \mathbf{1}_d/d$ . Note that  $U_t^{(\rho)}(\mathbf{x}, \eta, \beta) = U_t(\mathbf{x}, \eta, \beta) + \frac{\rho^2(d^{1-\beta}-1)}{\eta(1-\beta)}$ . Then under  
668 Assumption 3.1 there exists settings of  $\underline{\eta}, \bar{\eta}, \alpha, \lambda$  s.t. for all  $\varepsilon, \rho \in (0, 1)$  we have that

$$\begin{aligned} & \mathbb{E} \frac{1}{T} \sum_{t=1}^T \sum_{i=1}^m \ell_{t,i}(a_{t,i}) - \ell_{t,i}(\hat{a}_t) \\ & \leq \varepsilon m + \frac{8d\sqrt{m}}{\rho} \left( \mathbf{1}_{k>1} \sqrt{\frac{\log k}{T}} + \frac{1 + \log(T+1)}{16\rho T} \right) \\ & \quad + \min_{\beta \in [\underline{\beta}, \bar{\beta}], \eta > 0} \frac{8\left(\frac{d}{\varepsilon}\right)^{2-\beta} (1 + \log T)}{\eta T} + \frac{h_\beta(\Delta)}{\eta} + \frac{\eta d^\beta m}{\beta} + \frac{L_\eta(\bar{\beta} - \underline{\beta})}{2k} + d \min \left\{ \frac{\rho^2}{2\eta}, \rho\sqrt{m} \right\} \end{aligned} \quad (66)$$

669 for  $L_\eta = \left( \frac{\log \frac{d}{\varepsilon}}{\eta} + \eta m \log^2 d \right) d$  and  $h_\beta(\Delta) = (H_\beta + \frac{56}{dm})\iota_\Delta + \frac{d^{1-\beta}-1}{1-\beta}(1 - \iota_\Delta)$  for  $\iota_\Delta =$   
670  $\mathbf{1}_{m \geq \frac{75d}{\varepsilon\Delta^2} \log \frac{d}{\varepsilon\Delta^2}}$ .

671 *Proof.* By Lemma C.3 we have

$$\mathbb{E} \sum_{t=1}^T \sum_{i=1}^m \ell_{t,i}(a_{t,i}) - \ell_{t,i}(\hat{a}_t) \leq \varepsilon m T + \mathbb{E} \sum_{t=1}^T \frac{B_{\beta_t}(\hat{\mathbf{x}}_t^{(\varepsilon)} || \mathbf{x}_{t,1})}{\eta_t} + \frac{\eta_t d^{\beta_t} m}{\beta_t} \quad (67)$$

672 Since we have the same environment-dependent quantities as in Theorem B.1, we can substitute the  
673 above bound into Theorem A.1 and then apply the Lemma C.2 bound

$$\begin{aligned} \mathbb{E} \hat{V}_\beta^2 & \leq H_\beta + \frac{3d\kappa\beta}{1-\beta} \left( \left( \frac{d}{\varepsilon} \right)^{1-\beta} - 1 \right) \leq H_\beta + \frac{3d^2}{\varepsilon} \exp \left( -\frac{3\varepsilon\Delta^2 m}{28d} \right) \\ & = H_\beta + \frac{3\varepsilon\Delta^2}{d^2} \exp \left( 4 \log \frac{d}{\varepsilon\Delta^2} - \frac{3\varepsilon\Delta^2 m}{28d} \right) \\ & \leq H_\beta + \frac{3\varepsilon\Delta^2/d^2}{\frac{3\varepsilon\Delta^2 m}{28d} - 4 \log \frac{d}{\varepsilon\Delta^2}} \\ & \leq H_\beta + \frac{56}{dm} \end{aligned} \quad (68)$$

674 where the last line follows by assuming  $m \geq \frac{75d}{\varepsilon\Delta^2} \log \frac{d}{\varepsilon\Delta^2}$ . If this condition does not hold, then we  
675 apply the default bound of  $\mathbb{E} \hat{V}_\beta^2 \leq \frac{1}{T} \sum_{t=1}^T \psi_\beta(\hat{\mathbf{x}}_t) - \psi_\beta(\hat{\mathbf{x}}) \leq \frac{d^{1-\beta}-1}{1-\beta}$ .  $\square$

676 **Corollary C.1.** Let  $\underline{\beta} = \bar{\beta} = 1$ . Then for known  $\Delta$  and assuming  $m \geq \frac{75d}{\Delta^2} \log \frac{d}{\Delta^2}$  we can ensure  
677 expected task-averaged regret at most

$$2\sqrt{H_1 dm + 56} + \frac{75d}{\Delta^2} W \left( \frac{m}{75} \right) + \tilde{O} \left( \frac{d^{\frac{3}{2}} m^{\frac{3}{4}}}{\sqrt{T}} + \frac{d\Delta^2 m^2}{T} \right) \quad (69)$$

678 where  $W$  is the Lambert  $W$ -function, while for unknown  $\Delta$  we can ensure expected task-averaged  
679 regret at most

$$2\sqrt{H_1 dm + 56} + \frac{3}{\Delta} \sqrt[3]{50dm \log d \log \frac{d^2 m^2}{150\Delta^6 \log d}} + \tilde{O} \left( \frac{d^{\frac{3}{2}} m^{\frac{3}{4}}}{\sqrt{T}} + \frac{d^{\frac{4}{3}} m^{\frac{5}{3}}}{T} \right) \quad (70)$$

680 so long as  $m^2 \geq 150d \log d$ .

681 *Proof.* Applying Theorem C.1 and simplifying yields

$$\varepsilon m + \frac{8d\sqrt{m}(1 + \log(T + 1))}{16\rho^2 T} + \min_{\eta > 0} \frac{8d(1 + \log T)}{\varepsilon \eta T} + \frac{h_1(\Delta)}{\eta} + \eta dm + \frac{d\rho^2}{2\eta} \quad (71)$$

682 Then substitute  $\eta = \sqrt{\frac{h_1(\Delta)}{dm}}$  and set  $\rho = \sqrt[4]{\frac{1}{dT\sqrt{m}}}$  and  $\varepsilon = \frac{75d}{\Delta^2 m} W\left(\frac{m}{75}\right)$  (for known  $\Delta$ ) or

683  $\varepsilon = \sqrt[3]{\frac{150d \log d}{m^2}}$  (otherwise).  $\square$

684 **Corollary C.2.** Let  $\beta = \frac{1}{2}$  and  $\bar{\beta} = 1$ . Then for known  $\Delta$  and assuming  $m \geq \frac{75d}{\Delta^2} \log \frac{d}{\Delta^2}$  we can  
685 ensure task-averaged regret at most

$$\min_{\beta \in [\frac{1}{2}, 1]} 2\sqrt{(H_\beta m + 56/d)d^\beta/\beta} + \frac{75d}{\Delta^2} W\left(\frac{m}{75}\right) + \tilde{\mathcal{O}}\left(\frac{d^{\frac{4}{3}} m^{\frac{2}{3}}}{\sqrt[3]{T}} + \frac{d^{\frac{5}{3}} m^{\frac{5}{6}}}{T^{\frac{2}{3}}} + \frac{d\Delta^3 m^{\frac{5}{2}}}{T}\right) \quad (72)$$

686 using  $k = \lceil \sqrt[3]{d^2 m T} \rceil$ , while for unknown  $\Delta$  we can ensure expected task-averaged regret at most

$$\min_{\beta \in [\frac{1}{2}, 1]} 2\sqrt{(H_\beta m + 56/d)d^\beta/\beta} + \frac{3}{\Delta} \sqrt[3]{50d^2 m \log \frac{dm^2}{150\Delta^6}} + \tilde{\mathcal{O}}\left(\frac{d^{\frac{4}{3}} m^{\frac{2}{3}}}{\sqrt[3]{T}} + \frac{d^{\frac{5}{3}} m^{\frac{5}{6}}}{T^{\frac{2}{3}}} + \frac{d^{\frac{3}{2}} m^2}{T}\right) \quad (73)$$

687 so long as  $m \geq 5d\sqrt{6}$ .

688 *Proof.* Applying Theorem C.1 and simplifying yields

$$\begin{aligned} \varepsilon m + \frac{8d\sqrt{m}}{\rho} \left( \sqrt{\frac{\log k}{T}} + \frac{1 + \log(T + 1)}{16\rho T} \right) \\ + \min_{\beta \in [\underline{\beta}, \bar{\beta}], \eta > 0} \frac{8d^{\frac{3}{2}}(1 + \log T)}{\varepsilon^{\frac{3}{2}} \eta T} + \frac{h_\beta(\Delta)}{\eta} + \frac{\eta d^\beta m}{\beta} + \frac{d}{4k} \left( \frac{\log \frac{d}{\varepsilon}}{\eta} + \eta m \log^2 d \right) + \frac{d\rho^2}{2\eta} \end{aligned} \quad (74)$$

689 Then substitute  $\eta = \sqrt{\frac{h_\beta(\Delta)}{d^\beta m/\beta}}$  and set  $\rho = \sqrt[3]{\frac{1}{d\sqrt{mT}}}$  and  $\varepsilon = \frac{75d}{\Delta^2 m} W\left(\frac{m}{75}\right)$  (for known  $\Delta$ ) or

690  $\varepsilon = \sqrt[3]{\frac{150d^2}{m^2}}$  (otherwise).  $\square$

691 **Corollary C.3.** Let  $\underline{\beta} = \frac{1}{\log d}$  and  $\bar{\beta} = 1$ . Then for known  $\Delta$  and assuming  $m \geq \frac{75d}{\Delta^2} \log \frac{d}{\Delta^2}$  we can  
692 ensure task-averaged regret at most

$$\min_{\beta \in (0, 1]} 2\sqrt{(H_\beta m + 56/d)d^\beta/\beta} + \frac{75d}{\Delta^2} W\left(\frac{m}{75}\right) + \tilde{\mathcal{O}}\left(\frac{d^{\frac{4}{3}} m^{\frac{2}{3}}}{\sqrt[3]{T}} + \frac{d^{\frac{5}{3}} m^{\frac{5}{6}}}{T^{\frac{2}{3}}} + \frac{d\Delta^4 m^3}{T}\right) \quad (75)$$

693 using  $k = \lceil \sqrt[3]{d^2 m T} \rceil$ , while for unknown  $\Delta$  we can ensure expected task-averaged regret at most

$$\min_{\beta \in (0, 1]} 2\sqrt{(H_\beta m + 56/d)d^\beta/\beta} + \frac{3}{\Delta} \sqrt[3]{50d^2 m \log \frac{dm^2}{150\Delta^6}} + \tilde{\mathcal{O}}\left(\frac{d^{\frac{4}{3}} m^{\frac{2}{3}}}{\sqrt[3]{T}} + \frac{d^{\frac{5}{3}} m^{\frac{5}{6}}}{T^{\frac{2}{3}}} + \frac{d^{\frac{5}{3}} m^{\frac{7}{3}}}{T}\right) \quad (76)$$

694 so long as  $m \geq 5d\sqrt{6}$ .

695 *Proof.* Applying Theorem C.1 and simplifying yields

$$\begin{aligned} \varepsilon m + \frac{8d\sqrt{m}}{\rho} \left( \sqrt{\frac{\log k}{T}} + \frac{1 + \log(T + 1)}{16\rho T} \right) \\ + \min_{\beta \in [\underline{\beta}, \bar{\beta}], \eta > 0} \frac{8d^2(1 + \log T)}{\varepsilon^2 \eta T} + \frac{h_\beta(\Delta)}{\eta} + \frac{\eta d^\beta m}{\beta} + \frac{d}{2k} \left( \frac{\log \frac{d}{\varepsilon}}{\eta} + \eta m \log^2 d \right) + \frac{d\rho^2}{2\eta} \end{aligned} \quad (77)$$

696 Then substitute  $\eta = \sqrt{\frac{h_\beta(\Delta)}{d^\beta m/\beta}}$  and set  $\rho = \sqrt[3]{\frac{1}{d\sqrt{mT}}}$  and  $\varepsilon = \frac{75d}{\Delta^2 m} W\left(\frac{m}{75}\right)$  (for known  $\Delta$ ) or

697  $\varepsilon = \sqrt[3]{\frac{150d^2}{m^2}}$  (otherwise).  $\square$

698 **Corollary C.4.** Let  $\underline{\beta} = \frac{1}{\log d}$  and  $\bar{\beta} = 1$ . Then for unknown  $\Delta$  and assuming  $m \geq \max\{d^{\frac{3}{4}}, 56\}$   
699 we can ensure task-averaged regret at most

$$\min_{\beta \in (0,1]} \min \left\{ 8\sqrt{dm}, 2\sqrt{\left(H_{\beta}m + \frac{56}{d}\right) \frac{d^{\beta}}{\beta}} + \frac{21d^{\frac{3}{4}}\sqrt{m}}{\Delta} \sqrt{3 \log \frac{dm}{\Delta^2}} \right\} + \tilde{O} \left( \frac{d^{\frac{4}{3}}m^{\frac{2}{3}}}{\sqrt[3]{T}} + \frac{d^{\frac{5}{3}}m^{\frac{5}{6}}}{T^{\frac{2}{3}}} + \frac{d^2m^{\frac{7}{3}}}{T} \right) \quad (78)$$

700 using  $k = \lceil \sqrt[3]{d^2mT} \rceil$ .

701 *Proof.* Applying Theorem C.1 and simplifying yields

$$\begin{aligned} \varepsilon m + \frac{8d\sqrt{m}}{\rho} \left( \sqrt{\frac{\log k}{T}} + \frac{1 + \log(T+1)}{16\rho T} \right) \\ + \min_{\beta \in [\underline{\beta}, \bar{\beta}], \eta > 0} \frac{8d^2(1 + \log T)}{\varepsilon^2 \eta T} + \frac{h_{\beta}(\Delta)}{\eta} + \frac{\eta d^{\beta} m}{\beta} + \frac{d}{2k} \left( \frac{\log \frac{d}{\varepsilon}}{\eta} + \eta m \log^2 d \right) + \frac{d\rho^2}{2\eta} \end{aligned} \quad (79)$$

702 Then substitute  $\eta = \sqrt{\frac{h_{\beta}(\Delta)}{d^{\beta}m/\beta}}$  and set  $\rho = \sqrt[3]{\frac{1}{d\sqrt{mT}}}$  and  $\varepsilon = \frac{\sqrt{d}}{\sqrt[3]{m^2}}$ .  $\square$

## 703 D Online learning with self-concordant barrier regularizers

### 704 D.1 General results

705 **Lemma D.1.** Let  $\mathcal{K} \subset \mathbb{R}^d$  be a convex set and  $\psi : \mathcal{K}^{\circ} \mapsto \mathbb{R}^d$  be a self-concordant barrier. Suppose  
706  $\ell_1, \dots, \ell_T$  are a sequence of loss functions satisfying  $|\langle \ell_t, \mathbf{x} \rangle| \leq 1 \forall \mathbf{x} \in \mathcal{K}$ . Then if we run OMD  
707 with step-size  $\eta > 0$  as in Abernethy et al. [2, Alg. 1] on the sequence of estimators  $\hat{\ell}_t$  our estimated  
708 regret w.r.t. any  $\mathbf{x} \in \mathcal{K}_{\varepsilon}$  for  $\varepsilon > 0$  will satisfy

$$\sum_{t=1}^T \langle \hat{\ell}_t, \mathbf{x}_t - \mathbf{x} \rangle \leq \frac{B(\mathbf{x}|\mathbf{x}_1)}{\eta} + 32d^2\eta T \quad (80)$$

709 *Proof.* The result follows from Abernethy et al. [2] by stopping the derivation on the second inequality  
710 below Equation 10.  $\square$

711 **Definition D.1.** For any convex set  $\mathcal{K}$  and any point  $\mathbf{y} \in \mathcal{K}$ ,  $\pi_{\mathbf{y}}(\mathbf{x}) = \inf_{t \geq 0, \mathbf{y} + \frac{\mathbf{x}-\mathbf{y}}{t} \in \mathcal{K}}$   $t$  is the

712 **Minkowski function with pole  $\mathbf{y}$ .**

713 **Lemma D.2.** For any  $\mathbf{x} \in \mathcal{K} \subset \mathbb{R}^d$  and  $\psi : \mathcal{K}^{\circ} \mapsto \mathbb{R}$  a  $\nu$ -self-concordant regularizer with minimum  
714  $\mathbf{x}_1 \in \mathcal{K}^{\circ}$ , the quantity  $\psi(\mathbf{c}_{\varepsilon}(\mathbf{x}))$  is  $\nu\sqrt{2}$ -Lipschitz w.r.t.  $\varepsilon \in [0, 1]$ .

715 *Proof.* Consider any  $\varepsilon, \varepsilon' \in [0, 1]$  s.t.  $\varepsilon' - \varepsilon \in (0, \frac{1}{2}]$  Note that for  $t = \frac{\varepsilon' - \varepsilon}{1 + \varepsilon}$  we have

$$\mathbf{c}_{\varepsilon'}(\mathbf{x}) + \frac{\mathbf{c}_{\varepsilon'}(\mathbf{x}) - \mathbf{c}_{\varepsilon}(\mathbf{x})}{t} = \mathbf{x}_1 + \frac{\mathbf{x} - \mathbf{x}_1}{1 + \varepsilon'} + \frac{\mathbf{x}_1 + \frac{\mathbf{x} - \mathbf{x}_1}{1 + \varepsilon} - \mathbf{x}_1 - \frac{\mathbf{x} - \mathbf{x}_1}{1 + \varepsilon'}}{t} = \mathbf{x} \in \mathcal{K} \quad (81)$$

716 so  $\pi_{\mathbf{c}_{\varepsilon'}(\mathbf{x})}(\mathbf{c}_{\varepsilon}(\mathbf{x})) \leq \frac{\varepsilon' - \varepsilon}{1 + \varepsilon} \leq \varepsilon' - \varepsilon$ . Therefore by Nesterov and Nemirovskii [41, Prop. 2.3.2] we  
717 have

$$\psi(\mathbf{c}_{\varepsilon}(\mathbf{x})) - \psi(\mathbf{c}_{\varepsilon'}(\mathbf{x})) \leq \nu \log \left( \frac{1}{1 - \pi_{\mathbf{c}_{\varepsilon'}(\mathbf{x})}(\mathbf{c}_{\varepsilon}(\mathbf{x}))} \right) \leq \nu \log \left( \frac{1}{1 + \varepsilon - \varepsilon'} \right) \leq \nu(\varepsilon' - \varepsilon)\sqrt{2} \quad (82)$$

718 where for the last inequality we used  $-\log(1 - x) \leq x\sqrt{2}$  for  $x \in [0, \frac{1}{2}]$ . The case of  $\varepsilon' - \varepsilon \in (0, 1]$   
719 follows by considering  $\varepsilon'' = \frac{\varepsilon' + \varepsilon}{2}$  and applying the above twice.  $\square$

720 **Theorem D.1.** In Algorithm 1, let  $\text{OMD}_{\eta,\varepsilon}$  be online mirror descent over loss estimators specified in  
721 Abernethy et al. [2] with a  $\nu$ -self-concordant barrier regularizer  $\psi : \mathcal{K}^\circ \mapsto \mathbb{R}$  that satisfies  $\nu \geq 1$   
722 and  $\|\nabla^2\psi(\mathbf{x}_1)\|_2 = S_1 \geq 1$ . Let  $\Theta_k$  be a subset of  $[\frac{1}{m}, 1]$  and

$$U_t(\mathbf{x}, \eta, \varepsilon) = \frac{B(\mathbf{c}_\varepsilon(\hat{\mathbf{x}})|\mathbf{x})}{\eta} + 32\eta d^2 + \varepsilon m \quad (83)$$

723 Note that  $U_t^{(\rho)}(\mathbf{x}, \eta, \varepsilon) = U_t(\mathbf{x}, \eta, \varepsilon) + \frac{9\nu^{\frac{3}{2}}\rho^2 K m \sqrt{S_1}}{\eta}$ . Then there exists settings of  $\underline{\eta}, \bar{\eta}, \alpha, \lambda$  s.t. for  
724 all  $\varepsilon, \rho \in (0, 1)$  we have expected task averaged regret at most

$$\begin{aligned} & \mathbb{E} \min_{\varepsilon \in [\frac{1}{m}, 1], \eta > 0} \frac{512\nu^2 K^2 S_1 m^2 (1 + \log T)}{\eta} + \left( \frac{\hat{V}_\varepsilon^2}{\eta} + 32\eta d^2 m + \varepsilon m + \frac{\nu\sqrt{2}/\eta + m}{k} \right) T \\ & + 3\nu^{\frac{3}{4}} m \min \left\{ \frac{3\rho^2 \nu^{\frac{3}{4}} K \sqrt{S_1}}{\eta}, 4d\rho\sqrt{2K\sqrt{S_1}} \right\} T \\ & + \frac{7dm}{\rho} \sqrt{2K\sqrt{\nu^3 S_1}} \left( 7\sqrt{T \log k} + \frac{1 + \log(T+1)}{\rho} \right) \end{aligned} \quad (84)$$

725 *Proof.* Let  $\underline{\varepsilon} = \frac{1}{m}$ . For any  $\varepsilon \in [\underline{\varepsilon}, 1]$  and  $\mathbf{x} \in \mathcal{K}$  we have  $\pi_{\mathbf{x}_1}(\mathbf{c}_\varepsilon(\mathbf{x})) \leq \frac{1}{1+\varepsilon}$ , so by Nesterov and  
726 Nemirovskii [41, Prop. 2.3.2] we have

$$\|\nabla^2\psi(\mathbf{c}_\varepsilon(\mathbf{x}))\|_2 \leq \left( \frac{1 + 3\nu}{1 - \pi_{\mathbf{x}_1}(\mathbf{c}_\varepsilon(\mathbf{x}))} \right)^2 \|\nabla^2\psi(\mathbf{x}_1)\|_2 \leq \frac{64\nu^2 S_1}{\varepsilon^2} \quad (85)$$

727 Thus  $S = \max_{\mathbf{x}, \mathbf{y} \in \mathcal{K}, \varepsilon \in [\underline{\varepsilon}, 1]} \|\nabla^2\psi(\mathbf{c}_\varepsilon(\mathbf{x}))\|_2 = \frac{64\nu^2 S_1}{\underline{\varepsilon}^2}$  and also

$$\begin{aligned} D_\varepsilon^2 &= \max_{\mathbf{x}, \mathbf{y} \in \mathcal{K}} B(\mathbf{c}_\varepsilon(\mathbf{x})|\mathbf{c}_\varepsilon(\mathbf{y})) \\ &= \max_{\mathbf{x}, \mathbf{y} \in \mathcal{K}} \psi(\mathbf{c}_\varepsilon(\mathbf{x})) - \psi(\mathbf{c}_\varepsilon(\mathbf{y})) - \langle \nabla\psi(\mathbf{c}_\varepsilon(\mathbf{y})), \mathbf{x} - \mathbf{y} \rangle \\ &\leq \max_{\mathbf{x}, \mathbf{y} \in \mathcal{K}} \nu \log \left( \frac{1}{1 - \pi_{\mathbf{x}_1}(\mathbf{c}_\varepsilon(\mathbf{x}))} \right) + \sqrt{\nu \|\nabla^2\psi(\mathbf{c}_\varepsilon(\mathbf{y}))\|_2} \|\mathbf{x} - \mathbf{y}\|_2 \\ &\leq \nu \log \frac{2}{\varepsilon} + \frac{8\nu^{\frac{3}{2}} K \sqrt{S_1}}{\varepsilon} \\ &\leq \frac{9\nu^{\frac{3}{2}} K \sqrt{S_1}}{\varepsilon} \end{aligned} \quad (86)$$

728 where the first inequality follows by Nesterov and Nemirovskii [41, Prop. 2.3.2] and the definition  
729 of a self-concordant barrier [2, Def. 5]. In addition, we have  $g(\varepsilon) = 32d^2$ ,  $f(\varepsilon) = \varepsilon$ ,  $M =$   
730  $12d\sqrt{2Km}/\underline{\varepsilon}\sqrt[4]{\nu^3 S_1}$ , and  $F = 1$ . We have

$$\begin{aligned} \mathbb{E} \sum_{t=1}^T \sum_{i=1}^m \langle \ell_{t,i}, \mathbf{x}_{t,i} - \hat{\mathbf{x}}_t \rangle &\leq \mathbb{E} \sum_{t=1}^T \varepsilon_t m + \sum_{i=1}^m \langle \ell_{t,i}, \mathbf{x}_{t,i} - \mathbf{c}_{\varepsilon_t}(\hat{\mathbf{x}}_t) \rangle \\ &\leq \mathbb{E} \sum_{t=1}^T \varepsilon_t m + \sum_{i=1}^m \langle \hat{\ell}_{t,i}, \mathbf{x}_{t,i} - \mathbf{c}_{\varepsilon_t}(\hat{\mathbf{x}}_t) \rangle \\ &\leq \mathbb{E} \sum_{t=1}^T \varepsilon_t m + \sum_{i=1}^m \langle \hat{\ell}_{t,i}, \mathbf{x}_{t,i} - \mathbf{c}_{\varepsilon_t}(\hat{\mathbf{x}}_t) \rangle \\ &\leq \sum_{t=1}^T \frac{\mathbb{E} B(\mathbf{c}_{\varepsilon_t}(\hat{\mathbf{x}}_t)|\mathbf{x}_{t,1})}{\eta_t} + (32\eta_t d^2 + \varepsilon_t) m \end{aligned} \quad (87)$$

731 where the first inequality follows by Abernethy et al. [2, Lem. 8], the second by Abernethy et al. [2,  
732 Lem. 3], the third by optimality of  $\hat{\mathbf{x}}_t$ , and the fourth by Lemma D.1. Substituting into Theorem A.1  
733 and simplifying yields the result.  $\square$

734 **D.2 Specialization to the unit sphere**

735 **Corollary D.1.** *Let  $\mathcal{K}$  be the unit sphere with the self-concordant barrier  $\psi(\mathbf{x}) = -\log(1 - \|\mathbf{x}\|_2^2)$ .*  
 736 *Then Algorithm 1 attains expected task-averaged regret bounded by*

$$\tilde{\mathcal{O}} \left( \frac{dm^{\frac{3}{2}}}{T^{\frac{3}{4}}} + \frac{dm}{\sqrt[4]{T}} \right) + \min_{\varepsilon \in [\frac{1}{m}, 1]} 4d \sqrt{2m \log \left( 1 + \frac{1 - \mathbb{E}\|\hat{\mathbf{x}}\|_2^2}{2\varepsilon + \varepsilon^2} \right)} + \varepsilon m \quad (88)$$

737 using  $k = \lceil \sqrt{T} \rceil$ .

738 *Proof.* Using the fact the  $\nu = 1$  and  $K = S_1 = 2$ , we apply Theorem D.1 and simplify to obtain

$$\mathbb{E} \min_{\varepsilon \in [\frac{1}{m}, 1], \eta > 0} \frac{\hat{V}_\varepsilon^2}{\eta} + 32\eta d^2 m + \varepsilon m + \tilde{\mathcal{O}} \left( \frac{m^2}{\eta T} + \frac{1}{\eta k} + \frac{m}{k} + m \min \left\{ \frac{\rho^2}{\eta}, d\rho \right\} + \frac{dm}{\rho\sqrt{T}} + \frac{dm}{\rho^2 T} \right) \quad (89)$$

739 Then substitute  $\eta = \frac{\hat{V}_\varepsilon}{4\sqrt{2dm}} + \frac{\sqrt{m}}{d\sqrt[4]{T}}$ , set  $\rho = \frac{1}{\sqrt[4]{T}}$ , and note that

$$\begin{aligned} \mathbb{E} \hat{V}_\varepsilon &= \mathbb{E} \sqrt{\log \left( \frac{1 - \|\mathbf{c}_\varepsilon(\hat{\mathbf{x}})\|_2^2}{\sqrt[T]{\prod_{t=1}^T 1 - \|\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t)\|_2^2}} \right)} = \mathbb{E} \sqrt{\log \left( \frac{1 - (1 + \varepsilon)^{-2} \|\hat{\mathbf{x}}\|_2^2}{1 - (1 + \varepsilon)^{-2}} \right)} \\ &\leq \sqrt{\log \left( 1 + \frac{1 - \mathbb{E}\|\hat{\mathbf{x}}\|_2^2}{2\varepsilon + \varepsilon^2} \right)} \end{aligned} \quad (90)$$

740 where we use the fact that  $\|\hat{\mathbf{x}}_t\|_2 = 1$  and the inequality is Jensen's.  $\square$

741 **D.3 Specialization to polytopes, specifically the bandit online shortest-path problem**

742 As a last application, we apply our meta-BLO result to the shortest-path problem in online  
 743 optimization [49, 30]. In its bandit variant [8, 17], at each step  $i = 1, \dots, m$  the player must choose  
 744 a path  $p_i$  from a fixed source  $u \in V$  to a fixed sink  $v \in V$  in a directed graph  $G(V, E)$ . At the same  
 745 time the adversary chooses edge weights  $\ell_i \in \mathbb{R}^{|E|}$  and the player suffers the sum  $\sum_{e \in p_t} \ell_i(e)$  of  
 746 the weights in their chosen path  $p_t$ . This can be relaxed as BLO over vectors  $\mathbf{x}$  in a set  $\mathcal{K} \subset [0, 1]^{|E|}$   
 747 defined by a set  $\mathcal{C}$  of  $\mathcal{O}(|E|)$  linear constraints  $\langle \mathbf{a}, \mathbf{x} \rangle \leq b$  enforcing flows from  $u$  to  $v$ ;  $u$  to  
 748  $v$  paths can be sampled from any  $\mathbf{x} \in \mathcal{K}$  in an unbiased manner [2, Proposition 1]. In the single-task  
 749 case the BLO method of Abernethy et al. [2] has  $\mathcal{O}(|E|^{\frac{3}{2}} \sqrt{m})$  regret on this problem.

750 In the multi-task case consider a sequence of  $t = 1, \dots, T$  shortest path instances, each with  $m$   
 751 adversarial edge loss vectors  $\ell_{t,i}$ . The goal is to minimize average regret across instances. This setup  
 752 may be viewed as learning a prediction of the optimal path, as in the algorithms with predictions  
 753 paradigm in beyond-worst-case-analysis [39]; in particular, we have incorporated predictions into  
 754 the algorithm of Abernethy et al. [2] via the meta-initialization approach and now present the  
 755 learning-theoretic result for an end-to-end guarantee [33].

756 **Corollary D.2** (c.f. Cor. D.3). *For multi-task bandit online shortest path, Algorithm 1 with regularizer*  
 757  *$\psi(\mathbf{x}) = -\sum_{\mathbf{a}, b \in \mathcal{C}} \log(b - \langle \mathbf{a}, \mathbf{x} \rangle)$  attains the following expected average regret across instances*

$$\tilde{\mathcal{O}} \left( \frac{|E|^4 m^{\frac{3}{2}}}{T^{\frac{3}{4}}} + \frac{|E|^{\frac{5}{2}} m^{\frac{5}{6}}}{\sqrt[4]{T}} \right) + \min_{\varepsilon \in [\frac{1}{m}, 1]} 4|E|\mathbb{E} \sqrt{2m \sum_{\mathbf{a}, b \in \mathcal{C}} \log \left( \frac{\frac{1}{T} \sum_{t=1}^T b - \langle \mathbf{a}, \mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \rangle}{\sqrt[T]{\prod_{t=1}^T b - \langle \mathbf{a}, \mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \rangle}} \right)} + \varepsilon m \quad (91)$$

758 Here the asymptotic regret scales with the sum across all constraints  $\mathbf{a}, b \in \mathcal{C}$  of the log of the ratio  
 759 between the arithmetic and geometric means across tasks of the distances  $b - \langle \mathbf{a}, \mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \rangle$  from the  
 760 estimated optimum flow  $\mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t)$  to the constraint boundary. As it is difficult to separate the effect of the  
 761 offset  $\varepsilon$ , we do not state an explicit task-similarity measure like in our previous settings. Nevertheless,  
 762 since the arithmetic and geometric means are equal exactly when all entries are equal—and otherwise  
 763 the former is larger—the bound does show that regret is small when the estimated optimal flows  
 764  $\hat{\mathbf{x}}_t$  for each task are at similar distances from the constraints. Indeed, just as on the sphere, if the  
 765 estimated optima are all the same then setting  $\varepsilon = \frac{1}{m}$  again yields constant averaged regret.

766 **Corollary D.3.** Let  $\mathcal{K} = \{\mathbf{x} \in [0, 1]^{|E|} : \langle \mathbf{a}, \mathbf{x} \rangle \leq b \forall (\mathbf{a}, b) \in \mathcal{C}\}$  be the set of flows from  $u$  to  $v$   
767 on a graph  $G(V, E)$ , where  $\mathcal{C} \subset \mathbb{R}^{|E|} \times \mathbb{R}$  is a set of  $\mathcal{O}(|E|)$  linear constraints. Suppose we see  $T$   
768 instances of the bandit online shortest path problem with  $m$  timesteps each. Then sampling from  
769 probability distributions over paths from  $u$  to  $v$  returned by running Algorithm 1 with regularizer  
770  $\psi(\mathbf{x}) = -\sum_{\mathbf{a}, b \in \mathcal{C}} \log(b - \langle \mathbf{a}, \mathbf{x} \rangle)$  attains the following expected average regret across instances

$$\tilde{\mathcal{O}} \left( \frac{|E|^4 m^{\frac{3}{2}}}{T^{\frac{3}{4}}} + \frac{|E|^{\frac{5}{2}} m^{\frac{5}{6}}}{\sqrt[4]{T}} \right) + \min_{\varepsilon \in [\frac{1}{m}, 1]} 4|E| \mathbb{E} \sqrt{2m \sum_{\mathbf{a}, b \in \mathcal{C}} \log \left( \frac{\frac{1}{T} \sum_{t=1}^T b - \langle \mathbf{a}, \mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \rangle}{\sqrt[4]{\prod_{t=1}^T b - \langle \mathbf{a}, \mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \rangle}} \right)} + \varepsilon m \quad (92)$$

771 using  $k = \lceil \sqrt{T} \rceil$ .

772 *Proof.* Using the fact that  $d = |E|$ ,  $\nu = \mathcal{O}(|E|)$ ,  $K = \sqrt{|E|}$ , and  $S_1 \leq \sum_{\mathbf{a}, b \in \mathcal{C}} \frac{\|\mathbf{a}\mathbf{a}^T\|_2}{(\langle \mathbf{a}, \mathbf{1}_{|E|} / |E| \rangle - b)^2} =$   
773  $\mathcal{O}(|E|^3)$ , we apply Theorem D.1 and simplify to obtain

$$\mathbb{E} \min_{\varepsilon \in [\frac{1}{m}, 1], \eta > 0} \frac{\hat{V}_\varepsilon^2}{\eta} + 32\eta |E|^2 m + \varepsilon m$$

$$+ \tilde{\mathcal{O}} \left( \frac{|E|^6 m^2}{\eta T} + \frac{|E|}{\eta k} + \frac{m}{k} + m \min \left\{ \frac{\rho^2 |E|^{\frac{7}{2}}}{\eta}, \rho |E|^{\frac{11}{4}} \right\} + \frac{|E|^{\frac{11}{4}} m}{\rho} \left( \frac{1}{\sqrt{T}} + \frac{1}{\rho T} \right) \right) \quad (93)$$

774 Then substitute  $\eta = \frac{\hat{V}_\varepsilon}{4\sqrt{2dm}} + \frac{|E|^2 \sqrt{m}}{\sqrt[4]{T}}$ , set  $\rho = \sqrt[4]{\frac{|E|}{T}} \sqrt[6]{m}$ , and note that

$$\hat{V}_\varepsilon^2 = \sum_{\mathbf{a}, b \in \mathcal{C}} \log \left( \frac{b - \langle \mathbf{a}, \mathbf{c}_\varepsilon(\hat{\mathbf{x}}) \rangle}{\sqrt[4]{\prod_{t=1}^T b - \langle \mathbf{a}, \mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \rangle}} \right) = \sum_{\mathbf{a}, b \in \mathcal{C}} \log \left( \frac{\frac{1}{T} \sum_{t=1}^T b - \langle \mathbf{a}, \mathbf{c}_\varepsilon(\hat{\mathbf{x}}) \rangle}{\sqrt[4]{\prod_{t=1}^T b - \langle \mathbf{a}, \mathbf{c}_\varepsilon(\hat{\mathbf{x}}_t) \rangle}} \right) \quad (94)$$

775 □