
A Simple and Adaptive Learning Rate for FTRL in Online Learning with Minimax Regret of $\Theta(T^{2/3})$ and its Application to Best-of-Both-Worlds

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Follow-the-Regularized-Leader (FTRL) is a powerful framework for various on-
2 line learning problems. By designing its regularizer and learning rate to be adap-
3 tive to past observations, FTRL is known to work adaptively to various properties
4 of an underlying environment. However, most existing adaptive learning rates are
5 for online learning problems with a minimax regret of $\Theta(\sqrt{T})$ for the number of
6 rounds T , and there are only a few studies on adaptive learning rates for prob-
7 lems with a minimax regret of $\Theta(T^{2/3})$, which include several important prob-
8 lems dealing with indirect feedback. To address this limitation, we establish a new
9 adaptive learning rate framework for problems with a minimax regret of $\Theta(T^{2/3})$.
10 Our learning rate is designed by matching the stability, penalty, and bias terms
11 that naturally appear in regret upper bounds for problems with a minimax regret
12 of $\Theta(T^{2/3})$. As applications of this framework, we consider two major problems
13 dealing with indirect feedback: partial monitoring and graph bandits. We show that
14 FTRL with our learning rate and the Tsallis entropy regularizer improves existing
15 Best-of-Both-Worlds (BOBW) regret upper bounds, which achieve simultaneous
16 optimality in the stochastic and adversarial regimes. The resulting learning rate is
17 surprisingly simple compared to the existing learning rates for BOBW algorithms
18 for problems with a minimax regret of $\Theta(T^{2/3})$.

19 1 Introduction

20 Online learning is a problem setting in which a learner interacts with an environment for T rounds
21 with the goal of minimizing their cumulative loss. This framework includes many important online
22 decision-making problems, such as expert problems [21, 38, 57], multi-armed bandits [6, 8, 33],
23 linear bandits [1, 14], graph bandits [4, 42], and partial monitoring [9, 11].

24 For the sake of discussion in a general form, we consider the following *general online learning*
25 *framework*. In this framework, a learner is initially given a finite action set $\mathcal{A} = [k] := \{1, \dots, k\}$
26 and an observation set \mathcal{O} . At each round $t \in [T]$, the environment determines a loss function $\ell_t: \mathcal{A} \rightarrow$
27 $[0, 1]$, and the learner selects an action $A_t \in \mathcal{A}$ based on past observations without knowing ℓ_t . The
28 learner then suffers a loss $\ell_t(A_t)$ and observes a feedback $o_t \in \mathcal{O}$. The goal of the learner is to
29 minimize the (pseudo-)regret Reg_T , which is defined as the expectation of the difference between
30 the cumulative loss of the selected actions $(A_t)_{t=1}^T$ and that of an optimal action $a^* \in \mathcal{A}$ fixed in
31 hindsight. That is, $\text{Reg}_T = \mathbb{E}[\sum_{t=1}^T \ell_t(A_t) - \sum_{t=1}^T \ell_t(a^*)]$ for $a^* \in \arg \min_{a \in \mathcal{A}} \mathbb{E}[\sum_{t=1}^T \ell_t(a)]$.
32 For example in the multi-armed bandit problem, the observation is $o_t = \ell_t(A_t)$.

33 *Follow-the-Regularized-Leader (FTRL)* is a highly powerful framework for such online learning
 34 problems. In FTRL, a probability vector q_t over \mathcal{A} , which is used for determining action selection
 35 probability p_t so that $A_t \sim p_t$, is obtained by solving the following convex optimization problem:

$$q_t \in \arg \min_{q \in \mathcal{P}_k} \left\{ \sum_{s=1}^{t-1} \widehat{\ell}_s(q) + \beta_t \psi(q) \right\}, \quad (1)$$

36 where \mathcal{P}_k is the set of probability distributions over $\mathcal{A} = [k]$, $\widehat{\ell}_t: \mathcal{P}_k \rightarrow \mathbb{R}$ is an estimator of loss
 37 function ℓ_t , $\beta_t > 0$ is (a reciprocal of) learning rate at round t , and ψ is a convex regularizer. FTRL
 38 is known for its usefulness in various online learning problems [1, 4, 8, 27, 37]. Notably, FTRL can
 39 be viewed as a generalization of Online Gradient Descent [63] and the Hedge algorithm [21, 38, 57],
 40 and is closely related to Online Mirror Descent [36, 45].

41 The benefit of FTRL due to its generality is that one can design its regularizer ψ and learning rate
 42 $(\beta_t)_t$ so that it can perform adaptively to various properties of underlying loss functions. The *adaptive*
 43 *learning rate*, which exploits past observations, is often used to obtain such adaptivity. In order to
 44 see how it is designed, we consider the following stability–penalty decomposition, well-known in the
 45 literature [36, 45]:

$$\text{Reg}_T \lesssim \underbrace{\sum_{t=1}^T \frac{z_t}{\beta_t}}_{\text{stability term}} + \beta_1 h_1 + \underbrace{\sum_{t=2}^T (\beta_t - \beta_{t-1}) h_t}_{\text{penalty term}}. \quad (2)$$

46 Intuitively, the *stability* term arises from the regret when the difference in FTRL outputs, x_t and
 47 x_{t+1} , is large, and the *penalty* term is due to the strength of the regularizer. For example, in the Exp3
 48 algorithm for multi-armed bandits [8], h_t is the Shannon entropy of x_t or its upper bound, and z_t is
 49 the expectation of $(\nabla^2 \psi(x_t))^{-1}$ -norm of the importance-weighted estimator $\widehat{\ell}_t$ or its upper bound.

50 Adaptive learning rates have been designed so that it depends on the stability or penalty. For ex-
 51 ample, the well-known AdaGrad [19, 44] and the first-order algorithm [2] depend on stability com-
 52 ponents $(z_s)_{s=1}^{t-1}$ to determine β_t . More recently, there are learning rates that depend on penalty
 53 components $(h_s)_{s=1}^{t-1}$ [25, 54] and that depend on both stability and penalty components [26, 28, 55].

54 However, almost all adaptive learning rates developed so far have been limited to problems with a
 55 minimax regret of $\Theta(\sqrt{T})$, and there has been limited investigation into problems with a minimax re-
 56 gret of $\Theta(T^{2/3})$ [25, 54]. Such online learning are primarily related to indirect feedback and includes
 57 many important problems, such as partial monitoring [9, 34], graph bandits [4], dueling bandits [51],
 58 online ranking [12], bandits with switching costs [18], and bandits with paid observations [53].

59 **Contributions** To address this limitation, we establish a new learning rate framework for online
 60 learning with a minimax regret of $\Theta(T^{2/3})$. Henceforth, we will refer to problems with a minimax
 61 regret of $\Theta(T^{2/3})$ as *hard problems* to avoid repetition, abusing the terminology of partial monitor-
 62 ing. For hard problems, it is common to combine FTRL with *forced exploration* [4, 17, 34, 51]. In
 63 this study, we first observe that the regret of FTRL with forced exploration rate γ_t is roughly bounded
 64 as follows:

$$\text{Reg}_T \lesssim \underbrace{\sum_{t=1}^T \frac{z_t}{\beta_t \gamma_t}}_{\text{stability term}} + \beta_1 h_1 + \underbrace{\sum_{t=2}^T (\beta_t - \beta_{t-1}) h_t}_{\text{penalty term}} + \underbrace{\sum_{t=1}^T \gamma_t}_{\text{bias term}}. \quad (3)$$

65 Here, the third term, called the bias term, represents the regret incurred by forced exploration. In
 66 the aim of minimizing the RHS of (3), we will determine the exploration rate γ_t and learning rate
 67 β_t so that the above stability, penalty, and bias elements for each $t \in [T]$ are matched, where the
 68 resulting learning rate is called *Stability–Penalty–Bias matching learning rate (SPB-matching)*. This
 69 was inspired by the learning rate designed by matching the stability and penalty terms for problems
 70 with a minimax regret of $\Theta(\sqrt{T})$ [26]. Our learning rate is simultaneously adaptive to the stability
 71 component z_t and penalty component h_t , which have attracted attention in very recent years [26, 28,
 72 55]. The SPB-matching learning rate allows us to bound the RHS of (3) from above as follows:

73 **Theorem 1** (informal version of Theorem 6). *There exists learning rate $(\beta_t)_t$ and exploration rate*
 74 *$(\gamma_t)_t$ for which the RHS of (3) is bounded by $O\left(\left(\sum_{t=1}^T \sqrt{z_t h_t \log(\varepsilon T)}\right)^{2/3} + \left(\sqrt{z_{\max} h_{\max}}/\varepsilon\right)^{2/3}\right)$*
 75 *for any $\varepsilon \geq 1/T$, where $z_{\max} = \max_{t \in [T]} z_t$ and $h_{\max} = \max_{t \in [T]} h_t$.*

Table 1: Regret bounds for partial monitoring and graph bandits. The number of rounds is denoted as T , the number of actions as k , and the minimum suboptimality gap as Δ_{\min} . The variables $c_{\mathcal{G}}$ is defined in Section 5, D is a constant dependent on the outcome distribution. The graph complexity measures δ, δ^* , satisfying $\delta^* \leq \delta$ for graphs with no self-loops, are defined in Section 6, and $\tilde{\delta}^* \leq \delta$ is the fractional weak domination number [13]. AwSB is the abbreviation of the adversarial regime with a self-bounding constraint. MS-type means that the bound in AdvSB has a form similar to the bound established by Masoudian and Seldin [43].

Setting	Ref.	Stochastic	Adversarial	AwSB
Partial monitoring (with global observability)	[30]	$D \log T$	–	–
	[37]	–	$(c_{\mathcal{G}}T)^{2/3}(\log k)^{1/3}$	–
	[54]	$\frac{c_{\mathcal{G}}^2 \log T \log(kT)}{\Delta_{\min}^2}$	$(c_{\mathcal{G}}T)^{2/3}(\log T \log(kT))^{1/3}$	✓
	[56]	$\frac{c_{\mathcal{G}}^2 k \log T}{\Delta_{\min}^2}$	$(c_{\mathcal{G}}T)^{2/3}(\log T)^{1/3}$	✓
	Ours (Cor. 9)	$\frac{c_{\mathcal{G}}^2 \log k \log T}{\Delta_{\min}^2}$	$(c_{\mathcal{G}}T)^{2/3}(\log k)^{1/3}$	✓ (MS-type)
Graph bandits (with weak observability)	[4]	–	$(\delta \log k)^{1/3}T^{2/3}$	–
	[13]	–	$(\tilde{\delta}^* \log k)^{1/3}T^{2/3}$	–
	[25]	$\frac{\delta \log T \log(kT)}{\Delta_{\min}^2}$	$(\delta \log T \log(kT))^{1/3}T^{2/3}$	✓
	[15] ^a	$\frac{\delta \log k \log T}{\Delta_{\min}^2}$	$(\delta \log k)^{1/3}T^{2/3}$	✓
	Ours (Cor. 11)	$\frac{\delta^* \log k \log T}{\Delta_{\min}^2}$	$(\delta^* \log k)^{1/3}T^{2/3}$	✓ (MS-type)

^aThe bounds in [15] depend on δ , but their framework with the algorithm in [13] can achieve improved bounds replacing δ with $\tilde{\delta}^* \leq \delta$. The framework in [15] is a hierarchical reduction-based approach, rather than a direct FTRL method, discarding past observations as doubling-trick.

76 Within the general online learning framework, this theorem allows us to prove the following Best-
77 of-Both-Worlds (BOBW) guarantee [10, 58, 61], which achieves an $O(\log T)$ regret in the stochastic
78 regime and an $O(T^{2/3})$ regret in the adversarial regime simultaneously:

79 **Theorem 2** (informal version of Theorem 7). *Under some regularity conditions, an FTRL-based*
80 *algorithm with SPB-matching achieves $\text{Reg}_T \lesssim (z_{\max} h_{\max})^{1/3} T^{2/3}$ in the adversarial regime. In*
81 *the stochastic regime, if $\sqrt{z_t h_t} \leq \sqrt{\rho_1}(1 - q_{ta^*})$ holds for FTRL output $q_t \in \mathcal{P}_k$ and $\rho_1 > 0$ for all*
82 *$t \in [T]$, the same algorithm achieves $\text{Reg}_T \lesssim \rho_1 \log T / \Delta_{\min}^2$ for the minimum suboptimality gap Δ_{\min} .*

83 To assess the usefulness of the above result that holds for the general online learning framework,
84 this study focuses on two major hard problems: partial monitoring with global observability and
85 graph bandits with weak observability. We demonstrate that the assumptions in Theorem 2 are in-
86 deed satisfied for these problems by appropriately choosing the parameters in SPB-matching, thereby
87 improving the existing BOBW regret upper bounds in several respects. To obtain better bounds in
88 this analysis, we leverage the smallness of stability components z_t , which results from the forced
89 exploration. Additionally, SPB-matching is the first unified framework to achieve a BOBW guaran-
90 tee for hard online learning problems. Our learning rate is based on a surprisingly simple principle,
91 whereas existing learning rates for graph bandits and partial monitoring are extremely complicated
92 (see [25, Eq. (15)] and [54, Eq. (16)]). Due to its simplicity, we believe that SPB-matching will serve
93 as a foundation for building new BOBW algorithms for a variety of hard online learning problems.

94 Although omitted in Theorem 2, our approach achieves a refined regret bound devised by Masoudian
95 and Seldin [43] in the *adversarial regime with a self-bounding constraint* [61], which includes the
96 stochastic regime, adversarial regime, and the stochastic regime with adversarial corruptions [41] as
97 special cases. We call the refined bound *MS-type bound*, named after the author. The MS-type bound
98 maintains an ideal form even when $C = \Theta(T)$ or $\Delta_{\min} = \Theta(1/\sqrt{T})$ (see [43] for details), and our
99 bounds are the first MS-type bounds for hard problems. A comparison with existing regret bounds
100 is summarized in Table 1.

101 **2 Preliminaries**

102 **Notation** For a natural number $n \in \mathbb{N}$, we let $[n] = \{1, \dots, n\}$. For vector x , let x_i denote its i -th
 103 element and $\|x\|_p$ the ℓ_p -norm for $p \in [1, \infty]$. Let $\mathcal{P}_k = \{p \in [0, 1]^k : \|p\|_1 = 1\}$ be the $(k - 1)$ -
 104 dimensional probability simplex. The vector e_i is the i -th standard basis and $\mathbf{1}$ is the all-ones vector.
 105 Let $D_\psi(x, y)$ denote the Bregman divergence from y to x induced by a differentiable convex function
 106 ψ : $D_\psi(x, y) = \psi(x) - \psi(y) - \langle \nabla \psi(y), x - y \rangle$. To simplify the notation, we sometimes write $(a_t)_{t=1}^T$
 107 as $a_{1:T}$ and $f = O(g)$ as $f \lesssim g$. We regard function $f: \mathcal{A} = [k] \rightarrow \mathbb{R}$ as a k -dimensional vector.

108 **General online learning framework** To provide results that hold for a wide range of settings, we
 109 consider the following general online learning framework introduced in Section 1.

At each round $t \in [T] = \{1, \dots, T\}$:

1. The environment determines a loss vector $\ell_t: \mathcal{A} \rightarrow [0, 1]$;
2. The learner selects an action $A_t \in \mathcal{A}$ based on $p_t \in \mathcal{P}_k$ without knowing ℓ_t ;
3. The learner suffers a loss of $\ell_t(A_t) \in [0, 1]$ and observes a feedback $o_t \in \mathcal{O}$.

110 This framework includes many problems such as the expert problem, multi-armed bandits, graph
 111 bandits, partial monitoring as special cases.

112 **Stochastic, adversarial, and their intermediate regimes** Within the above general online frame-
 113 work, we study three different regimes for a sequence of loss functions $(\ell_t)_t$. In the stochastic regime,
 114 the sequence of loss functions is sampled from an unknown distribution \mathcal{D} in an i.i.d. manner. The
 115 suboptimality gap for action $a \in \mathcal{A}$ is given by $\Delta_a = \mathbb{E}_{\ell_t \sim \mathcal{D}}[\ell_t(a) - \ell_t(a^*)]$ and the minimum sub-
 116 optimality gap by $\Delta_{\min} = \min_{a \neq a^*} \Delta_a$. In the adversarial regime, the loss functions can be selected
 117 arbitrarily, possibly based on the past history up to round $t - 1$.

118 We also investigate, the adversarial regime with a self-bounding constraint [61], which is an inter-
 119 mediate regime between the stochastic and adversarial regimes.

120 **Definition 3.** Let $\Delta \in [0, 1]^k$ and $C \geq 0$. The environment is in an *adversarial regime with a*
 121 (Δ, C, T) self-bounding constraint if it holds for any algorithm that $\text{Reg}_T \geq \mathbb{E}[\sum_{t=1}^T \Delta_{A_t} - C]$.

122 From the definition, the stochastic and adversarial regimes are special cases of this regime. Addition-
 123 ally, the well-known stochastic regime with adversarial corruptions [41] also falls within this regime.
 124 For the adversarial regime with a self-bounding constraint, we assume that there exists a unique opti-
 125 mal action a^* . This assumption is standard in the literature of BOBW algorithms (e.g., [22, 39, 58]).

126 **3 SBP-matching: Simple and adaptive learning rate for hard problems**

127 This section designs a new learning rate framework for hard online learning problems.

128 **3.1 Objective function that adaptive learning rate aims to minimize**

129 In hard problems, the regret of FTRL with somewhat large exploration rate γ_t is known to be bounded
 130 in the following form [4, 25, 54]:

$$\text{Reg}_T \lesssim \sum_{t=1}^T \frac{z_t}{\beta_t \gamma_t} + \sum_{t=1}^T (\beta_t - \beta_{t-1}) h_t + \sum_{t=1}^T \gamma_t \tag{4}$$

131 for some stability component z_t and penalty component h_t , where we set $\beta_{T+1} = \beta_T$ and $\beta_0 = 0$
 132 for simplicity. Recall that the first term is the stability term, the second term is the penalty term, and
 133 the third term is the bias term, which arises from the forced exploration.

134 The goal when designing the adaptive learning rate is to minimize (4), under the constraints that
 135 $(\beta_t)_t$ is non-decreasing and β_t depends on $(z_{1:t}, h_{1:t})$ or $(z_{1:t-1}, h_{1:t})$. A naive way to choose γ_t to
 136 minimize (4) is to set $\gamma_t = \sqrt{z_t/\beta_t}$ so that the stability term and the bias term match. However, this
 137 choice does not work well in hard problems because to obtain a regret bound of (4), a lower bound
 138 of $\gamma_t \geq u_t/\beta_t$ for some $u_t > 0$ is needed. This lower bound is used to control the magnitude of the

139 loss estimator $\widehat{\ell}_t$.¹ Therefore, we consider exploration rate of $\gamma_t = \gamma'_t + u_t/\beta_t$ for $\gamma'_t = \sqrt{z_t/\beta_t}$ and
 140 some $u_t > 0$, where γ'_t is chosen so that the stability and bias terms are matched. With these choices,

$$\begin{aligned} \text{Eq. (4)} &\leq \sum_{t=1}^T \left(\frac{z_t}{\beta_t \gamma'_t} + (\beta_t - \beta_{t-1})h_t + \left(\gamma'_t + \frac{u_t}{\beta_t} \right) \right) \\ &= \sum_{t=1}^T \left(2\sqrt{\frac{z_t}{\beta_t}} + \frac{u_t}{\beta_t} + (\beta_t - \beta_{t-1})h_t \right) =: F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{1:T}). \end{aligned} \quad (5)$$

141 Note that the first two terms in F , $2\sqrt{z_t/\beta_t} + u_t/\beta_t$, come from the stability and bias terms and the
 142 last term, $(\beta_t - \beta_{t-1})h_t$, is the penalty term. In the following, we investigate adaptive learning rate
 143 $(\beta_t)_{t=1}^T$ that minimizes F in (5) instead of (4).

144 3.2 Stability–penalty–bias matching learning rate

145 We consider determining $(\beta_t)_t$ by matching the stability–bias terms and the penalty term as
 146 $2\sqrt{z_t/\beta_t} + u_t/\beta_t = (\beta_t - \beta_{t-1})h_t$. Assume that when choosing β_t , we have an access to \widehat{h}_t such
 147 that $h_t \leq \widehat{h}_t$. Then, inspired by the above matching, we consider the following two update rules:

$$\text{(Rule 1) } \beta_t = \beta_{t-1} + \frac{1}{\widehat{h}_t} \left(2\sqrt{\frac{z_t}{\beta_t}} + \frac{u_t}{\beta_t} \right), \quad \text{(Rule 2) } \beta_t = \beta_{t-1} + \frac{1}{\widehat{h}_t} \left(2\sqrt{\frac{z_{t-1}}{\beta_{t-1}}} + \frac{u_{t-1}}{\beta_{t-1}} \right). \quad (6)$$

148 We call these update rules *Stability–Penalty–Bias Matching (SPB-matching)*. These are designed by
 149 following the simple principle of matching the stability, penalty, and bias elements, and Rules 1 and
 150 2 differ only in the way indices are shifted. For the sake of convenience, we define G_1 and G_2 by

$$G_1(z_{1:T}, h_{1:T}) = \sum_{t=1}^T \frac{\sqrt{z_t}}{\left(\sum_{s=1}^t \sqrt{z_s/h_s} \right)^{1/3}}, \quad G_2(u_{1:T}, h_{1:T}) = \sum_{t=1}^T \frac{u_t}{\sqrt{\sum_{s=1}^t u_s/h_s}}. \quad (7)$$

151 Define $z_{\max} = \max_{t \in [T]} z_t$, $u_{\max} = \max_{t \in [T]} u_t$, and $h_{\max} = \max_{t \in [T]} h_t$. Then, using SPB-
 152 matching rules in (6), we can upper-bound F in terms of G_1 and G_2 as follows:

153 **Lemma 4.** Consider SPB-matching (6) and suppose that $h_t \leq \widehat{h}_t$ for all $t \in [T]$. Then, Rule 1
 154 achieves $F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{1:T}) \leq 3.2G_1(z_{1:T}, \widehat{h}_{1:T}) + 2G_2(u_{1:T}, \widehat{h}_{1:T})$ and Rule 2 achieves
 155 $F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{1:T}) \leq 4G_1(z_{1:T}, \widehat{h}_{2:T+1}) + 3G_2(u_{1:T}, \widehat{h}_{2:T+1}) + 10\sqrt{z_{\max}/\beta_1} + 5u_{\max}/\beta_1 +$
 156 $\beta_1 h_1$.

157 The proof of Lemma 4 can be found in Appendix B.1. One can see from the proof that the effect of
 158 using $\gamma_t = \sqrt{z_t/\beta_t} + u_t/\beta_t$ instead of $\gamma_t = \sqrt{z_t/\beta_t}$ only appears in G_2 , which has a less impact
 159 than G_1 when bounding F . We can further upper-bound G_1 as follows:

160 **Lemma 5.** Let $(z_t)_{t=1}^T \subseteq \mathbb{R}_{\geq 0}$ and $(h_t)_{t=1}^T \subseteq \mathbb{R}_{> 0}$ be any non-negative and positive se-
 161 quences, respectively. Let $\theta_0 > \theta_1 > \dots > \theta_J > \theta_{J+1} = 0$ and $\theta_0 \geq h_{\max}$ and de-
 162 fine $\mathcal{T}_j = \{t \in [T]: \theta_{j-1} \geq h_t > \theta_j\}$ for $j \in [J]$ and $\mathcal{T}_{J+1} = \{t \in [T]: \theta_J \geq h_t\}$. Then,
 163 $G_1(z_{1:T}, h_{1:T}) \leq \frac{3}{2} \sum_{j=1}^{J+1} (\sqrt{\theta_{j-1}} \sum_{t \in \mathcal{T}_j} \sqrt{z_t})^{2/3}$. This implies that for all $j \in \mathbb{N}$ it holds that

$$G_1(z_{1:T}, h_{1:T}) \leq \frac{3}{2} \min \left\{ \left(\sqrt{2J} \sum_{t=1}^T \sqrt{z_t h_t} \right)^{\frac{2}{3}} + \left(2^{-J/2} \sqrt{z_{\max} h_{\max}} \right)^{\frac{2}{3}} T^{\frac{2}{3}}, \left(\sum_{t=1}^T \sqrt{z_t h_{\max}} \right)^{\frac{2}{3}} \right\}.$$

164 Combining Lemmas 4 and 5 and the bound on G_2 in [26, Lemma 3], we obtain the following theorem.

¹This is particularly the case when we use the Shannon entropy or Tsallis entropy regularizers, which is a weaker regularization than the log-barrier regularizer.

Algorithm 1: Best-of-both-worlds framework based on FTRL with SPB-matching learning rate and Tsallis entropy for online learning with minimax regret of $\Theta(T^{2/3})$

- 1 **input:** action set \mathcal{A} , observation set \mathcal{O} , exponent of Tsallis entropy $\alpha, \beta_1, \bar{\beta}$
2 **for** $t = 1, 2, \dots$ **do**
3 Compute $q_t \in \mathcal{P}_k$ by (10) with a loss estimator \hat{y}_t .
4 Set $h_t = H_\alpha(q_t)$ and $z_t, u_t \geq 0$ defined for each problem.
5 Compute action selection probability p_t from q_t by (11).
6 Choose $A_t \in \mathcal{A}$ so that $\Pr[A_t = i | p_t] = p_{ti}$ and observe feedback $o_t \in \mathcal{O}$.
7 Compute loss estimator $\hat{\ell}_t$ based on p_t and o_t .
8 Compute β_{t+1} by Rule 2 of SPB-matching in (6) with $\hat{h}_{t+1} = h_t$.
-

165 **Theorem 6.** Let $(z_t)_{t=1}^T, (u_t)_{t=1}^T \subseteq \mathbb{R}_{\geq 0}$ and $(h_t)_{t=1}^T \subseteq \mathbb{R}_{> 0}$. Suppose that \hat{h}_t satisfies $h_t \leq \hat{h}_t$ for
166 all $t \in [T]$. Then, if β_t is given by Rule 1 in (6), then for all $\varepsilon \geq 1/T$ it holds that

$$F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{1:T}) \lesssim \min \left\{ \left(\sum_{t=1}^T \sqrt{z_t \hat{h}_t \log(\varepsilon T)} \right)^{\frac{2}{3}} + \left(\sqrt{z_{\max} \hat{h}_{\max}} / \varepsilon \right)^{\frac{2}{3}}, \left(\sum_{t=1}^T \sqrt{z_t \hat{h}_{\max}} \right)^{\frac{2}{3}} \right\} \\ + \min \left\{ \sqrt{\sum_{t=1}^T u_t \hat{h}_t \log(\varepsilon T)} + \sqrt{u_{\max} \hat{h}_{\max} / \varepsilon}, \sqrt{\sum_{t=1}^T u_t \hat{h}_{\max}} \right\}. \quad (8)$$

167 If β_t is given by Rule 2 in (6), then for all $\varepsilon \geq 1/T$ it holds that

$$F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{1:T}) \lesssim \min \left\{ \left(\sum_{t=1}^T \sqrt{z_t \hat{h}_{t+1} \log(\varepsilon T)} \right)^{\frac{2}{3}} + \left(\sqrt{z_{\max} \hat{h}_{\max}} / \varepsilon \right)^{\frac{2}{3}}, \left(\sum_{t=1}^T \sqrt{z_t \hat{h}_{\max}} \right)^{\frac{2}{3}} \right\} \\ + \min \left\{ \sqrt{\sum_{t=1}^T u_t \hat{h}_{t+1} \log(\varepsilon T)} + \sqrt{u_{\max} \hat{h}_{\max} / \varepsilon}, \sqrt{\sum_{t=1}^T u_t \hat{h}_{\max}} \right\} + \sqrt{\frac{z_{\max}}{\beta_1}} + \frac{u_{\max}}{\beta_1} + \beta_1 h_1. \quad (9)$$

168 Note that these bounds are for problems with a minimax regret of $\Theta(T^{2/3})$. Roughly speaking, our
169 bounds have an order of $\left(\sum_{t=1}^T \sqrt{z_t \hat{h}_{t+1} \log T} \right)^{1/3}$ and differ from the existing stability-penalty-
170 adaptive-type bounds of $\sqrt{z_t \hat{h}_{t+1} \log T}$ for problems with a minimax regret of $\Theta(\sqrt{T})$ [26, 55]. We
171 will see in the subsequent sections that our bounds are reasonable as they give nearly optimal regret
172 bounds in stochastic and adversarial regimes in partial monitoring and graph bandits.

173 4 Best-of-both-worlds framework for hard online learning problems

174 Using the SPB-matching learning rate established in Section 3, this section provides a BOBW algo-
175 rithm framework for hard online learning problems. We consider the following FTRL update:

$$q_t = \arg \min_{p \in \mathcal{P}_k} \left\{ \sum_{s=1}^{t-1} \langle \hat{\ell}_s, p \rangle + \beta_t (-H_\alpha(p)) + \bar{\beta} (-H_{\bar{\alpha}}(p)) \right\}, \quad \alpha \in (0, 1), \bar{\alpha} = 1 - \alpha, \quad (10)$$

176 where H_α is the α -Tsallis entropy defined as $H_\alpha(p) = \frac{1}{\alpha} \sum_{i=1}^k (p_i^\alpha - p_i)$, which satisfies $H_\alpha(p) \geq 0$
177 and $H_\alpha(e_i) = 0$. Based on this FTRL output q_t , we set $h_t = H_\alpha(q_t)$, which satisfies $h_1 = h_{\max}$.
178 Additionally, for q_t and some $p_0 \in \mathcal{P}_k$, we use the action selection probability $p_t \in \mathcal{P}_k$ defined by

$$p_t = (1 - \gamma_t) q_t + \gamma_t p_0 \quad \text{for} \quad \gamma_t = \gamma'_t + \frac{u_t}{\beta_t} = \sqrt{\frac{z_t}{\beta_t}} + \frac{u_t}{\beta_t}, \quad (11)$$

179 where β_1 is chosen so that $\gamma_t \in [0, 1/2]$. Let $\kappa = \sqrt{z_{\max}/\beta_1} + u_{\max}/\beta_1 + \beta_1 h_1 + \bar{\beta} \bar{h}$ and let $\mathbb{E}_t[\cdot]$
180 be the expectation given all observations before round t . Then the above procedure with Rule 2 of
181 SPB-matching in (6), summarized in Algorithm 1, achieves the following BOBW bound:

182 **Theorem 7.** Suppose that loss function ℓ_t satisfies $\|\ell_t\|_\infty \leq 1$ and the following three conditions
 183 (i)–(iii) are satisfied: (i) $\text{Reg}_T \leq \mathbb{E}[\sum_{t=1}^T \langle \widehat{\ell}_t, q_t - e_{a^*} \rangle + 2 \sum_{t=1}^T \gamma_t]$,

$$(ii) \mathbb{E}_t \left[\langle \widehat{\ell}_t, q_t - q_{t+1} \rangle - \beta_t D_{(-H_\alpha)}(q_{t+1}, q_t) \right] \lesssim \frac{z_t}{\beta_t \gamma_t}, \quad (iii) h_t \lesssim h_{t-1}. \quad (12)$$

184 Then, in the adversarial regime, Algorithm 1 achieves

$$\text{Reg}_T = O\left((z_{\max} h_1)^{1/3} T^{2/3} + \sqrt{u_{\max} h_1 T} + \kappa\right). \quad (13)$$

185 In the adversarial regime with a (Δ, C, T) -self-bounding constraint, further suppose that

$$\sqrt{z_t h_t} \leq \sqrt{\rho_1} \cdot (1 - q_{ta^*}) \quad \text{and} \quad u_t h_t \leq \rho_2 \cdot (1 - q_{ta^*}) \quad (14)$$

186 are satisfied for some $\rho_1, \rho_2 > 0$ for all $t \in [T]$. Then, the same algorithm achieves

$$\text{Reg}_T = O\left(\frac{\rho}{\Delta_{\min}^2} \log(T \Delta_{\min}^2) + \left(\frac{C^2 \rho}{\Delta_{\min}^2} \log\left(\frac{T \Delta_{\min}}{C}\right)\right)^{1/3} + \kappa'\right) \quad (15)$$

187 for $\rho = \max\{\rho_1, \rho_2\}$ and $\kappa' = \kappa + ((z_{\max} h_1)^{1/3} + \sqrt{u_{\max} h_1})(1/\Delta_{\min}^2 + C/\Delta_{\min})^{2/3}$ when $T \geq$
 188 $1/\Delta_{\min}^2 + C/\Delta_{\min} =: \tau$, and $\text{Reg}_T = O((z_{\max} h_1)^{1/3} \tau^{2/3} + \sqrt{u_{\max} h_1 \tau})$ when $T < \tau$.

189 The proof of Theorem 7 relies on Theorem 6 established in the last section and can be found in
 190 Appendix C. Note that the bound (15) becomes the bound for the stochastic regime when $C = 0$.

191 5 Case study (1): Partial monitoring with global observability

192 This section provides a new BOBW algorithm for globally observable partial monitoring games.

193 5.1 Problem setting and some concepts in partial monitoring

194 **Partial monitoring games** A Partial Monitoring (PM) game $\mathcal{G} = (\mathcal{L}, \Phi)$ consists of a loss matrix
 195 $\mathcal{L} \in [0, 1]^{k \times d}$ and feedback matrix $\Phi \in \Sigma^{k \times d}$, where k and d are the number of actions and out-
 196 comes, respectively, and Σ is the set of feedback symbols. The game unfolds over T rounds between
 197 the learner and the environment. Before the game starts, the learner is given \mathcal{L} and Φ . At each round
 198 $t \in [T]$, the environment picks an outcome $x_t \in [d]$, and then the learner chooses an action $A_t \in [k]$
 199 without knowing x_t . Then the learner incurs an unobserved loss $\mathcal{L}_{A_t x_t}$ and only observes a feed-
 200 back symbol $\sigma_t := \Phi_{A_t x_t}$. This framework can be indeed expressed as the general online learning
 201 framework in Section 2, by setting $\mathcal{O} = \Sigma$, $\ell_t(a) = \mathcal{L}_{a x_t} = e_a^\top \mathcal{L} e_{x_t}$ and $o_t = \sigma_t = \Phi_{A_t x_t}$.

202 We next introduce fundamental concepts for PM games. Based on the loss matrix \mathcal{L} , we can
 203 decompose all distributions over outcomes. For each action $a \in [k]$, the cell of action a , den-
 204 oted as \mathcal{C}_a , is the set of probability distributions over $[d]$ for which action a is optimal. That is,
 205 $\mathcal{C}_a = \{u \in \mathcal{P}_d : \max_{b \in [k]} (\ell_a - \ell_b)^\top u \leq 0\}$, where $\ell_a \in \mathbb{R}^d$ is the a -th row of \mathcal{L} .

206 To avoid the heavy notions and concepts of PM, we assume that the PM game has no duplicate actions
 207 $a \neq b$ such that $\ell_a = \ell_b$ and its all actions are *Pareto optimal*; that is, $\dim(\mathcal{C}_a) = d - 1$ for all $a \in [k]$.
 208 The discussion of the effect of this assumption can be found *e.g.*, in [34, 37].

209 **Observability and loss estimation** Two Pareto optimal actions a and b are *neighbors* if $\dim(\mathcal{C}_a \cap$
 210 $\mathcal{C}_b) = d - 2$. Then, this neighborhood relations defines *globally observable games*, for which the
 211 minimax regret of $\Theta(T^{2/3})$ is known in the literature [9, 34]. Two neighbouring actions a and b are
 212 *globally observable* if there exists a function $w_{e(a,b)} : [k] \times \Sigma \rightarrow \mathbb{R}$ satisfying

$$\sum_{c=1}^k w_{e(a,b)}(c, \Phi_{cx}) = \mathcal{L}_{ax} - \mathcal{L}_{bx} \quad \text{for all } x \in [d], \quad (16)$$

213 where $e(a, b) = \{a, b\}$. A PM game is said to be globally observable if all neighboring actions are
 214 globally observable. To the end, we assume that \mathcal{G} is globally observable.²

²Another representative class of PM is locally observable games, for which we can achieve a minimax regret of $\Theta(\sqrt{T})$. See [9, 36, 37] for local observability and [54, 55] for BOBW algorithms for it.

215 Based on the neighborhood relations, we can estimate the loss *difference* between actions, instead of
 216 estimating the loss itself. The *in-tree* is the edges of a directed tree with vertices $[k]$ and let $\mathcal{T} \subseteq$
 217 $[k] \times [k]$ be an in-tree over the set of actions induced by the neighborhood relations with an arbitrarily
 218 chosen root $r \in [k]$. Then, we can estimate the loss differences between Pareto optimal actions as
 219 follows. Let $G(a, \sigma)_b = \sum_{e \in \text{path}_{\mathcal{T}}(b)} w_e(a, \sigma)$ for $a \in [k]$, where $\text{path}_{\mathcal{T}}(b)$ is the set of edges from
 220 $b \in [k]$ to the root r on \mathcal{T} . Then, it is known that this G satisfies that for any Pareto optimal actions a
 221 and b , $\sum_{c=1}^k (G(c, \Phi_{cx})_b - G(b, \Phi_{cx})_c) = \mathcal{L}_{ax} - \mathcal{L}_{bx}$ for all $x \in [d]$ (e.g., [37, Lemma 4]). From this
 222 fact, one can see that we can use $\hat{y}_t = G(A_t, \Phi_{A_t x_t}) / p_{t A_t} \in \mathbb{R}^k$ as the loss (difference) estimator,
 223 following the standard construction of the importance-weighted estimator [8, 36]. In fact, \hat{y}_t satisfies
 224 $\mathbb{E}_{A_t \sim p_t} [\hat{y}_{ta} - \hat{y}_{tb}] = \sum_{c=1}^k (G(c, \sigma_t)_a - G(c, \sigma_t)_b) = \mathcal{L}_{ax} - \mathcal{L}_{bx}$. We let $c_G = \max\{1, k \|G\|_{\infty}\}$
 225 be a game-dependent constant, where $\|G\|_{\infty} = \max_{a \in [k], \sigma \in \Sigma} |G(a, \sigma)|$.

226 5.2 Algorithm and regret upper bounds

227 Here, we present a new BOBW algorithm based on Algorithm 1. We use the following parameters
 228 for Algorithm 1. We use the loss (difference) estimator of $\hat{\ell}_t = \hat{y}_t$. We set p_0 in (11) to $p_0 = \mathbf{1}/k$.
 229 For $\tilde{I}_t \in \arg \max_{i \in [k]} q_{ti}$ and $q_{t*} = \min\{q_{t\tilde{I}_t}, 1 - q_{t\tilde{I}_t}\}$, let

$$\beta_1 \geq \frac{64c_G^2}{1-\alpha}, \bar{\beta} = \frac{32c_G\sqrt{k}}{(1-\alpha)^2\sqrt{\beta_1}}, z_t = \frac{4c_G^2}{1-\alpha} \left(\sum_{i \neq \tilde{I}_t} q_{ti}^{2-\alpha} + q_{t*}^{2-\alpha} \right), u_t = \frac{8c_G}{1-\alpha} q_{t*}^{1-\alpha}. \quad (17)$$

230 Note that $z_{\max} = \frac{4c_G^2}{1-\alpha}$, $u_{\max} = \frac{8c_G}{1-\alpha}$, and $h_{\max} = h_1 = \frac{1}{\alpha} k^{1-\alpha}$. Then, we can prove the following:

231 **Theorem 8.** *In globally observable partial monitoring, for any $\alpha \in (0, 1)$, Algorithm 1 with (17)*
 232 *satisfies the assumptions of Theorem 7 with $\rho_1 = \Theta\left(\frac{c_G^2 k^{1-\alpha}}{\alpha(1-\alpha)}\right)$ and $\rho_2 = \Theta\left(\frac{c_G k^{1-\alpha}}{\alpha(1-\alpha)}\right)$.*

233 The proof of Theorem 8 is given in Appendix E. Setting $\alpha = 1 - 1/(\log k)$ gives the following:

234 **Corollary 9.** *In globally observable partial monitoring with $T \geq \tau$, Algorithm 1 with (17) for*
 235 *$\alpha = 1 - 1/(\log k)$ achieves $\text{Reg}_T = O((c_G T)^{2/3} (\log k)^{1/3} + \kappa)$ in the adversarial regime and*

$$\text{Reg}_T = O\left(\frac{c_G^2 \log k}{\Delta_{\min}^2} \log(T \Delta_{\min}^2) + \left(\frac{C^2 c_G^2 \log k}{\Delta_{\min}^2} \log\left(\frac{T \Delta_{\min}}{C}\right)\right)^{1/3} + \kappa'\right) \quad (18)$$

236 *in the adversarial regime with a (Δ, C, T) -self-bounding constraint.*

237 This regret upper bound is better than the bound in [54, 56] in both stochastic and adversarial regimes,
 238 notably by a factor of $\log T$ or k in the stochastic regime. The bound for the adversarial regime with
 239 a (Δ, C, T) -self-bounding constraint is the first MS-type bound in PM.

240 6 Case study (2): Graph bandits with weak observability

241 This section presents a new BOBW algorithm for weakly observable graph bandits.

242 6.1 Problem setting and some concepts in graph bandits

243 **Problem setting** In the graph bandit problem, the learner is given a directed feedback graph $G =$
 244 (V, E) with $V = [k]$ and $E \subseteq V \times V$. For each $i \in V$, let $N^{\text{in}}(i) = \{j \in V : (j, i) \in E\}$ and
 245 $N^{\text{out}}(i) = \{j \in V : (i, j) \in E\}$ be the in-neighborhood and out-neighborhood of vertex $i \in V$,
 246 respectively. The game proceeds as the general online learning framework provided in Section 2,
 247 with action set $\mathcal{A} = V$, loss function $\ell_t: V \rightarrow [0, 1]$, and observation $o_t = \{\ell_t(j) : j \in N^{\text{out}}(I_t)\}$.

248 **Observability and domination number** Similar to partial monitoring, the minimax regret of
 249 graph bandits is characterized by the properties of the feedback graph G [4]. A graph G is *ob-*
 250 *servable* if it contains no self-loops, $N^{\text{in}}(i) \neq \emptyset$ for all $i \in V$. A graph G is *strongly observable* if
 251 $i \in N^{\text{in}}(i)$ or $V \setminus \{i\} \subseteq N^{\text{in}}(i)$ for all $i \in V$. Then, a graph G is *weakly observable* if it is observable
 252 but not strongly observable.³ The minimax regret of the weakly observable is known to be $\Theta(T^{2/3})$.

³Similar to the locally observable games of partial monitoring, we can achieve an $O(\sqrt{T})$ regret for graph bandits with strong observability. See e.g., [4] for details.

253 The weak domination number characterizes precisely the minimax regret. The *weakly dominating*
 254 *set* $D \subseteq V$ is a set of vertices such that $\{i \in V : i \notin N^{\text{out}}(i)\} \subseteq \bigcup_{i \in D} N^{\text{out}}(i)$. Then, the *weak*
 255 *domination number* $\delta(G)$ of graph G is the size of the smallest weakly dominating set. For weakly
 256 observable G , the minimax regret of $\Theta(\delta^{1/3}T^{2/3})$ is known [4]. Instead, our bound depends on the
 257 *fractional domination number* $\delta^*(G)$, defined by the optimal value of the following linear program:

$$\text{minimize } \sum_{i \in V} x_i \quad \text{subject to } \sum_{i \in N^{\text{in}}(j)} x_i \geq 1 \quad \forall j \in V, \quad 0 \leq x_i \leq 1 \quad \forall i \in V. \quad (19)$$

258 We use $(x_i^*)_{i \in V}$ to denote the optimal solution of (19) and define its normalized version $u \in \mathcal{P}_k$
 259 by $u_i = x_i^* / \sum_{j \in V} x_j^*$. The advantage of using the fractional domination number mainly lies in its
 260 computational complexity; further details are provided in Appendix F.1.

261 6.2 Algorithm and regret analysis

262 Here, we present a new BOBW algorithm based on Algorithm 1. We use the following parameters
 263 for Algorithm 1. We use the estimator $\hat{\ell}_t \in \mathbb{R}^k$ defined by $\hat{\ell}_{ti} = \frac{\ell_{ti}}{P_{ti}} \mathbb{1}[i \in N^{\text{out}}(I_t)]$ for $P_{ti} =$
 264 $\sum_{j \in N^{\text{in}}(i)} p_{tj}$, which is unbiased and has been employed in the literature [4, 13]. We set p_0 in (11)
 265 to $p_0 = u$. For $\tilde{I}_t \in \arg \max_{i \in [k]} q_{ti}$ and $q_{t*} = \min\{q_{t\tilde{I}_t}, 1 - q_{t\tilde{I}_t}\}$, let

$$\beta_1 \geq \frac{64\delta^*}{1-\alpha}, \quad \bar{\beta} = \frac{32\sqrt{k}\delta^*}{(1-\alpha)^2\sqrt{\beta_1}}, \quad z_t = \frac{4\delta^*}{1-\alpha} \left(\sum_{i \in V \setminus \{\tilde{I}_t\}} q_{ti}^{2-\alpha} + q_{t*}^{2-\alpha} \right), \quad u_t = \frac{8\delta^*}{1-\alpha} q_{t*}^{1-\alpha}. \quad (20)$$

266 Note that $z_{\max} = \frac{4\delta^*}{1-\alpha}$, $u_{\max} = \frac{8\delta^*}{1-\alpha}$, and $h_{\max} = h_1 = \frac{1}{\alpha}k^{1-\alpha}$. Then, we can prove the following:

267 **Theorem 10.** *In the weakly observable graph bandit problem, for any $\alpha \in (0, 1)$, Algorithm 1*
 268 *with (20) satisfies the assumptions of Theorem 7 with $\rho_1 = \rho_2 = \Theta\left(\frac{\delta^*k^{1-\alpha}}{\alpha(1-\alpha)}\right)$.*

269 The proof of Theorem 10 is given in Appendix F. Setting $\alpha = 1 - 1/(\log k)$ gives the following:

270 **Corollary 11.** *In weakly observable graph bandits with $T \geq \max\{\delta^*(\log k)^2, \tau\}$, Algorithm 1 with*
 271 *(20) for $\alpha = 1 - 1/(\log k)$ achieves $\text{Reg}_T = O(\delta^{*1/3}T^{2/3}(\log k)^{1/3} + \kappa)$ in adversarial regime and*

$$\text{Reg}_T = O\left(\frac{\delta^* \log k}{\Delta_{\min}^2} \log(T\Delta_{\min}^2) + \left(\frac{C^2\delta^* \log k}{\Delta_{\min}^2} \log\left(\frac{T\Delta_{\min}}{C}\right)\right)^{1/3} + \kappa'\right) \quad (21)$$

272 *in the adversarial regime with a (Δ, C, T) -self-bounding constraint.*

273 Our bound is the first BOBW FTRL-based algorithm with the $O(\log T)$ bound in the stochastic
 274 regime, improving the existing best FTRL-based algorithm in [25]. Compared to the reduction-based
 275 approach in [15], the dependences on T are the same. However, our bound unfortunately depends on
 276 the fractional domination number δ^* instead of the weak domination number δ , which can be smaller
 277 than δ^* . Roughly speaking, this comes from the use of Tsallis entropy instead of Shannon entropy
 278 employed for the existing BOBW bound [25]. The technical challenges of making our bound depend
 279 on δ instead of δ^* or the weak fractional domination number $\tilde{\delta}^*$ are further discussed in Appendix F.3.
 280 Still, we believe that our algorithm can perform better since the reduction-based algorithm discards
 281 past observations as the doubling trick. Furthermore, the bound for the adversarial regime with a
 282 (Δ, C, T) -self-bounding constraint is the first MS-type bound in weakly observable graph bandits.

283 7 Conclusion and future work

284 In this work, we investigated hard online learning problems, that is online learning with a minimax
 285 regret of $\Theta(T^{2/3})$, and established a simple and adaptive learning rate framework called stability-
 286 penalty-bias matching (SPB-matching). We showed that FTRL with this framework and the Tsallis
 287 entropy regularization improves the existing BOBW regret bounds based on FTRL for two typical
 288 hard problems, partial monitoring with global observability and graph bandits with weak observabil-
 289 ity. Interestingly, the optimal exponent of Tsallis entropy in both settings is $1 - 1/(\log k)$, suggest-
 290 ing the reasonableness of using Shannon entropy in existing algorithms for partial monitoring [37]
 291 and graph bandits [4]. Our learning rate is surprisingly simple compared to existing ones for hard
 292 problems [25, 54]. Hence, it is important future work to investigate whether this simplicity can be
 293 leveraged to apply SPB-matching to other hard problems, such as bandits with switching costs [18]
 294 or with paid observations [53] and dueling bandits with Borda winner [51].

295 References

- 296 [1] Jacob D Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient
297 algorithm for bandit linear optimization. In *The 21st Annual Conference on Learning Theory*,
298 pages 263–274, 2008.
- 299 [2] Jacob D. Abernethy, Elad Hazan, and Alexander Rakhlin. Interior-point methods for full-
300 information and bandit online learning. *IEEE Transactions on Information Theory*, 58(7):
301 4164–4175, 2012.
- 302 [3] Jacob D Abernethy, Chansoo Lee, and Ambuj Tewari. Fighting bandits with a new kind of
303 smoothness. In *Advances in Neural Information Processing Systems*, volume 28, pages 2197–
304 2205, 2015.
- 305 [4] Noga Alon, Nicolò Cesa-Bianchi, Ofer Dekel, and Tomer Koren. Online learning with feed-
306 back graphs: Beyond bandits. In *Proceedings of The 28th Conference on Learning Theory*,
307 volume 40, pages 23–35, 2015.
- 308 [5] Jean-Yves Audibert and Sébastien Bubeck. Minimax policies for adversarial and stochastic
309 bandits. In *Conference on Learning Theory*, volume 7, pages 1–122, 2009.
- 310 [6] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Ma-
311 chine Learning Research*, 3(Nov):397–422, 2002.
- 312 [7] Peter Auer and Chao-Kai Chiang. An algorithm with nearly optimal pseudo-regret for both
313 stochastic and adversarial bandits. In *29th Annual Conference on Learning Theory*, volume 49,
314 pages 116–120, 2016.
- 315 [8] Peter Auer, Nicolás Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic
316 multiarmed bandit problem. *SIAM Journal on Computing*, 32(1):48–77, 2002.
- 317 [9] Gábor Bartók, Dávid Pál, and Csaba Szepesvári. Minimax regret of finite partial-monitoring
318 games in stochastic environments. In *Proceedings of the 24th Annual Conference on Learning
319 Theory*, volume 19, pages 133–154, 2011.
- 320 [10] Sébastien Bubeck and Aleksandrs Slivkins. The best of both worlds: Stochastic and adversarial
321 bandits. In *Proceedings of the 25th Annual Conference on Learning Theory*, volume 23, pages
322 42.1–42.23, 2012.
- 323 [11] Nicolò Cesa-Bianchi, Gábor Lugosi, and Gilles Stoltz. Regret minimization under partial mon-
324 itoring. *Mathematics of Operations Research*, 31(3):562–580, 2006.
- 325 [12] Sougata Chaudhuri and Ambuj Tewari. Online learning to rank with top-k feedback. *Journal
326 of Machine Learning Research*, 18(103):1–50, 2017.
- 327 [13] Houshuang Chen, zengfeng Huang, Shuai Li, and Chihao Zhang. Understanding bandits with
328 graph feedback. In *Advances in Neural Information Processing Systems*, volume 34, pages
329 24659–24669, 2021.
- 330 [14] Varsha Dani, Thomas P. Hayes, and Sham M Kakade. Stochastic linear optimization under
331 bandit feedback. In *The 21st Annual Conference on Learning Theory*, volume 2, pages 355–
332 366, 2008.
- 333 [15] Chris Dann, Chen-Yu Wei, and Julian Zimmert. A blackbox approach to best of both worlds in
334 bandits and beyond. In *Proceedings of Thirty Sixth Conference on Learning Theory*, volume
335 195, pages 5503–5570, 2023.
- 336 [16] Steven de Rooij, Tim van Erven, Peter D. Grünwald, and Wouter M. Koolen. Follow the leader
337 if you can, hedge if you must. *Journal of Machine Learning Research*, 15(37):1281–1316,
338 2014.
- 339 [17] Ofer Dekel, Ambuj Tewari, and Raman Arora. Online bandit learning against an adaptive
340 adversary: from regret to policy regret. In *Proceedings of the 29th International Conference
341 on Machine Learning*, pages 1747–1754, 2012.

- 342 [18] Ofer Dekel, Jian Ding, Tomer Koren, and Yuval Peres. Bandits with switching costs: $T^{2/3}$
343 regret. In *Proceedings of the Forty-Sixth Annual ACM Symposium on Theory of Computing*,
344 pages 459–467, 2014.
- 345 [19] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning
346 and stochastic optimization. *Journal of Machine Learning Research*, 12(61):2121–2159, 2011.
- 347 [20] Liad Erez and Tomer Koren. Towards best-of-all-worlds online learning with feedback graphs.
348 In *Advances in Neural Information Processing Systems*, volume 34, pages 28511–28521, 2021.
- 349 [21] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and
350 an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- 351 [22] Pierre Gaillard, Gilles Stoltz, and Tim van Erven. A second-order bound with excess losses. In
352 *Proceedings of The 27th Conference on Learning Theory*, volume 35, pages 176–196, 2014.
- 353 [23] Pratik Gajane and Tanguy Urvoy. Utility-based dueling bandits as a partial monitoring game.
354 *arXiv preprint arXiv:1507.02750*, 2015.
- 355 [24] Shinji Ito. Parameter-free multi-armed bandit algorithms with hybrid data-dependent regret
356 bounds. In *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134, pages
357 2552–2583, 2021.
- 358 [25] Shinji Ito, Taira Tsuchiya, and Junya Honda. Nearly optimal best-of-both-worlds algorithms for
359 online learning with feedback graphs. In *Advances in Neural Information Processing Systems*,
360 volume 35, pages 28631–28643, 2022.
- 361 [26] Shinji Ito, Taira Tsuchiya, and Junya Honda. Adaptive learning rate for follow-the-regularized-
362 leader: Competitive analysis and best-of-both-worlds. *arXiv preprint arXiv:2403.00715*, 2024.
- 363 [27] Tiancheng Jin, Longbo Huang, and Haipeng Luo. The best of both worlds: stochastic and adver-
364 sarial episodic MDPs with unknown transition. In *Advances in Neural Information Processing*
365 *Systems*, volume 34, pages 20491–20502, 2021.
- 366 [28] Tiancheng Jin, Junyan Liu, and Haipeng Luo. Improved best-of-both-worlds guarantees for
367 multi-armed bandits: FTRL with general regularizers and multiple optimal arms. In *Advances*
368 *in Neural Information Processing Systems*, volume 36, pages 30918–30978, 2023.
- 369 [29] Robert Kleinberg and Tom Leighton. The value of knowing a demand curve: Bounds on re-
370 gret for online posted-price auctions. In *the 44th Annual IEEE Symposium on Foundations of*
371 *Computer Science*, pages 594–605, 2003.
- 372 [30] Junpei Komiyama, Junya Honda, and Hiroshi Nakagawa. Regret lower bound and optimal
373 algorithm in finite stochastic partial monitoring. In *Advances in Neural Information Processing*
374 *Systems*, volume 28, pages 1792–1800, 2015.
- 375 [31] Fang Kong, Yichi Zhou, and Shuai Li. Simultaneously learning stochastic and adversarial
376 bandits with general graph feedback. In *Proceedings of the 39th International Conference on*
377 *Machine Learning*, volume 162, pages 11473–11482, 2022.
- 378 [32] Joon Kwon and Vianney Perchet. Gains and losses are fundamentally different in regret mini-
379 mization: The sparse case. *Journal of Machine Learning Research*, 17(227):1–32, 2016.
- 380 [33] T. L. Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in*
381 *Applied Mathematics*, 6(1):4–22, 1985.
- 382 [34] Tor Lattimore and Csaba Szepesvári. Cleaning up the neighborhood: A full classification for
383 adversarial partial monitoring. In *Proceedings of the 30th International Conference on Algo-*
384 *rithmic Learning Theory*, volume 98, pages 529–556, 2019.
- 385 [35] Tor Lattimore and Csaba Szepesvári. An information-theoretic approach to minimax regret
386 in partial monitoring. In *the 32nd Annual Conference on Learning Theory*, volume 99, pages
387 2111–2139, 2019.
- 388 [36] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.

- 389 [37] Tor Lattimore and Csaba Szepesvári. Exploration by optimisation in partial monitoring. In
390 *Proceedings of Thirty Third Conference on Learning Theory*, volume 125, pages 2488–2515,
391 2020.
- 392 [38] Nick Littlestone and Manfred K Warmuth. The weighted majority algorithm. *Information and*
393 *computation*, 108(2):212–261, 1994.
- 394 [39] Haipeng Luo and Robert E. Schapire. Achieving all with no parameters: AdaNormalHedge. In
395 *Proceedings of The 28th Conference on Learning Theory*, volume 40, pages 1286–1304, 2015.
- 396 [40] Haipeng Luo, Chen-Yu Wei, and Kai Zheng. Efficient online portfolio with logarithmic regret.
397 In *Advances in Neural Information Processing Systems*, volume 31, pages 8235–8245, 2018.
- 398 [41] Thodoris Lykouris, Vahab Mirrokni, and Renato Paes Leme. Stochastic bandits robust to ad-
399 versarial corruptions. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory*
400 *of Computing*, pages 114–122, 2018.
- 401 [42] Shie Mannor and Ohad Shamir. From bandits to experts: On the value of side-observations. In
402 *Advances in Neural Information Processing Systems*, volume 24, pages 684–692, 2011.
- 403 [43] Saeed Masoudian and Yevgeny Seldin. Improved analysis of the Tsallis-INF algorithm in
404 stochastically constrained adversarial bandits and stochastic bandits with adversarial corrup-
405 tions. In *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134, pages
406 3330–3350, 2021.
- 407 [44] H. Brendan McMahan and Matthew J. Streeter. Adaptive bound optimization for online convex
408 optimization. In *The 23rd Conference on Learning Theory*, pages 244–256, 2010.
- 409 [45] Francesco Orabona. A modern introduction to online learning. *arXiv preprint*
410 *arXiv:1912.13213*, 2019.
- 411 [46] Aldo Pacchiano, Christoph Dann, and Claudio Gentile. Best of both worlds model selection.
412 In *Advances in Neural Information Processing Systems*, volume 35, pages 1883–1895, 2022.
- 413 [47] Antonio Piccolboni and Christian Schindelhauer. Discrete prediction games with arbitrary feed-
414 back and loss (extended abstract). In *Computational Learning Theory*, pages 208–223, 2001.
- 415 [48] Chloé Rouyer and Yevgeny Seldin. Tsallis-INF for decoupled exploration and exploitation in
416 multi-armed bandits. In *Proceedings of Thirty Third Conference on Learning Theory*, volume
417 125, pages 3227–3249, 2020.
- 418 [49] Chloé Rouyer, Dirk van der Hoeven, Nicolò Cesa-Bianchi, and Yevgeny Seldin. A near-optimal
419 best-of-both-worlds algorithm for online learning with feedback graphs. In *Advances in Neural*
420 *Information Processing Systems*, volume 35, pages 35035–35048, 2022.
- 421 [50] Aldo Rustichini. Minimizing regret: The general case. *Games and Economic Behavior*, 29(1):
422 224–243, 1999.
- 423 [51] Aadirupa Saha, Tomer Koren, and Yishay Mansour. Adversarial dueling bandits. In *Proceed-*
424 *ings of the 38th International Conference on Machine Learning*, volume 139, pages 9235–9244,
425 2021.
- 426 [52] Yevgeny Seldin and Aleksandrs Slivkins. One practical algorithm for both stochastic and ad-
427 versarial bandits. In *Proceedings of the 31st International Conference on Machine Learning*,
428 volume 32, pages 1287–1295, 2014.
- 429 [53] Yevgeny Seldin, Peter Bartlett, Koby Crammer, and Yasin Abbasi-Yadkori. Prediction with
430 limited advice and multiarmed bandits with paid observations. In *Proceedings of the 31st*
431 *International Conference on Machine Learning*, volume 32, pages 280–287, 2014.
- 432 [54] Taira Tsuchiya, Shinji Ito, and Junya Honda. Best-of-both-worlds algorithms for partial moni-
433 toring. In *Proceedings of The 34th International Conference on Algorithmic Learning Theory*,
434 pages 1484–1515, 2023.

- 435 [55] Taira Tsuchiya, Shinji Ito, and Junya Honda. Stability-penalty-adaptive follow-the-regularized-
436 leader: Sparsity, game-dependency, and best-of-both-worlds. In *Advances in Neural Informa-*
437 *tion Processing Systems*, volume 36, 2023.
- 438 [56] Taira Tsuchiya, Shinji Ito, and Junya Honda. Exploration by optimization with hybrid regu-
439 larizers: Logarithmic regret with adversarial robustness in partial monitoring. *arXiv preprint*
440 *arXiv:2402.08321*, 2024.
- 441 [57] Vladimir Vovk. Aggregating strategies. In *Proceedings of the Third Annual Workshop on*
442 *Computational Learning Theory*, pages 371–383, 1990.
- 443 [58] Chen-Yu Wei and Haipeng Luo. More adaptive algorithms for adversarial bandits. In *Proceed-*
444 *ings of the 31st Conference On Learning Theory*, volume 75, pages 1263–1291, 2018.
- 445 [59] Julian Zimmert and Tor Lattimore. Connections between mirror descent, thompson sampling
446 and the information ratio. In *Advances in Neural Information Processing Systems*, pages 11973–
447 11982, 2019.
- 448 [60] Julian Zimmert and Yevgeny Seldin. An optimal algorithm for stochastic and adversarial ban-
449 dits. In *Proceedings of the Twenty-Second International Conference on Artificial Intelligence*
450 *and Statistics*, volume 89, pages 467–475, 2019.
- 451 [61] Julian Zimmert and Yevgeny Seldin. Tsallis-INF: An optimal algorithm for stochastic and
452 adversarial bandits. *Journal of Machine Learning Research*, 22(28):1–49, 2021.
- 453 [62] Julian Zimmert, Haipeng Luo, and Chen-Yu Wei. Beating stochastic and adversarial semi-
454 bandits optimally and simultaneously. In *Proceedings of the 36th International Conference on*
455 *Machine Learning*, volume 97, pages 7683–7692, 2019.
- 456 [63] Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent.
457 In *the Twentieth International Conference on Machine Learning*, pages 928–935, 2003.

458 A Additional related work

459 **Best-of-both-worlds algorithms** The study of BOBW algorithms was initiated by Bubeck and
460 Slivkins [10], who focused on multi-armed bandits. The motivation arises from the difficulty of
461 determining in advance whether the underlying environment is stochastic or adversarial in real-world
462 problems. Since then, BOBW algorithms have been extensively studied [7, 16, 22, 40, 46, 52], and
463 recently, FTRL is the common approach for developing BOBW algorithms [24, 28, 60, 62]. One
464 reason is by appropriately designing the learning rate and regularizer of FTRL, we can prove a BOBW
465 guarantee for various problem settings. Another reason is that FTRL-based approaches not only
466 perform well in both stochastic and adversarial regimes but also achieve favorable regret bounds in
467 the adversarial regime with a self-bounding constraint, intermediate settings including stochastically
468 constrained adversarial regime [58] and stochastic regime with adversarial corruptions [41]. This
469 intermediate regime is particularly useful, considering that real-world problems often lie between
470 purely stochastic and purely adversarial regimes.

471 This study is closely related to FTRL with the Tsallis entropy regularization. Tsallis entropy in online
472 learning was introduced in [3, 5], and its significance for BOBW algorithms was established in [61].
473 In the multi-armed bandit problem, using the exponent of Tsallis entropy $\alpha = 1/2$ provides optimal
474 upper bounds, up to logarithmic factors, in both stochastic and adversarial regimes [61]. However,
475 in the graph bandits, where the dependence on k is critical or in decoupled settings, optimal upper
476 bounds can be achieved with $\alpha \neq 1/2$ [26, 32, 48, 59]. In this work, we demonstrate that using the
477 exponent of $\alpha = 1 - 1/(\log k)$ for the number of actions k results in favorable regret bounds, as
478 shown in Corollaries 9 and 11.

479 **Partial monitoring** Partial monitoring [11, 47, 50] is a very general online decision-making frame-
480 work and includes a wide range of problems such as multi-armed bandits, (utility-based) dueling
481 bandits [23], online ranking [12], and dynamic pricing [29]. The characterization of the minimax
482 regret in partial monitoring has been progressively understood through various studies. It is known
483 that all partial monitoring games can be classified into trivial, easy, hard, and hopeless games, where
484 their minimax regrets are 0 , $\Theta(\sqrt{T})$, $\Theta(T^{2/3})$ and $\Omega(T)$. For comprehensive literature, refer to [9]
485 and the improved results presented in [34, 35]. The games for which we can achieve a regret bound
486 of $O(T^{2/3})$ correspond to globally observable games.

487 There is limited research on BOBW algorithms for partial monitoring with global observability [54,
488 56]. The existing bounds exhibit suboptimal dependencies on k and T , particularly in the stochastic
489 regime, which comes from the use of the Shannon entropy or the log-barrier regularization. By
490 employing Tsallis entropy, our algorithm is the first to achieve ideal dependencies on both k and T .
491 It remains uncertain whether our upper bound in the stochastic regime is optimal with respect to
492 variables other than T . While there is an asymptotic lower bound for the stochastic regime [30], its
493 coefficient is expressed as a complex optimization problem. Investigating this lower bound further is
494 important future work.

495 **Graph bandits** The study on the graph bandit problem, which is also known as online learning
496 with feedback graphs, was initiated by [42]. This problem includes several important problems such
497 as the expert setting, multi-armed bandits, and label-efficient prediction. For example, considering
498 a feedback graph with only self-loops, one can see that this corresponds to the multi-armed bandit
499 problem. One of the most seminal studies on the graph bandit problem is by Alon et al. [4], who
500 elucidated how the structure of the feedback graph influences its minimax regret. They demonstrated
501 that the minimax regret is characterized by the observability of the feedback graph, introducing the
502 notions of weakly observable graphs and strongly observable graphs. Of particular relevance to this
503 study is the minimax regret of $\tilde{O}(\delta T^{2/3})$ for weakly observable graphs, where δ is the weak dom-
504 ination number and $\tilde{O}(\cdot)$ ignores logarithmic factors. Recently, this upper bound was improved to
505 $\tilde{O}(\delta^* T^{2/3})$ by replacing the weak domination number with the fractional weak domination num-
506 ber $\tilde{\delta}^*$ [13].

507 There are several BOBW algorithms for graph bandits [15, 20, 25, 31, 49]. However, only a few
508 of these studies consider the weakly observable setting [15, 25, 31]. The existing results based on
509 FTRL rely on the domination number rather than the weak domination number [31] or exhibit poor
510 dependence on T [25, 31], and the best regret bound of them still exhibited a dependence on T of

511 $(\log T)^2$ [25]. Our algorithm is the first FTRL-based algorithm in the weakly observable setting that
 512 achieves an $O(\log T)$ stochastic bound.

513 **B Proofs for SPB-matching learning rate (Section 3)**

514 **B.1 Proof of Lemma 4**

515 *Proof of Lemma 4.* We first consider Rule 1 in (6). The learning rate β_t is lower-bounded as

$$\beta_t^{3/2} \geq \beta_t^{1/2} \left(\beta_{t-1} + \frac{2}{\widehat{h}_t} \sqrt{z_t} \right) \geq \beta_{t-1}^{3/2} + \frac{2\sqrt{z_t}}{\widehat{h}_t} \geq 2 \sum_{s=1}^t \frac{\sqrt{z_s}}{\widehat{h}_s}, \quad (22)$$

516 where the first inequality follows from the definition of β_t in (6) and the second inequality from the
 517 fact that $(\beta_t)_t$ is non-decreasing. We also have

$$\beta_t^2 \geq \beta_t \left(\beta_{t-1} + \frac{1}{\widehat{h}_t} \frac{u_t}{\beta_t} \right) \geq \beta_{t-1}^{3/2} + \frac{u_t}{\widehat{h}_t} \geq \sum_{s=1}^t \frac{u_s}{\widehat{h}_s}. \quad (23)$$

518 Using the last two lower bounds on β_t , we can bound F in (5) as

$$\begin{aligned} F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{1:T}) &\leq \sum_{t=1}^T \left(2\sqrt{\frac{z_t}{\beta_t}} + \frac{u_t}{\beta_t} + (\beta_t - \beta_{t-1})\widehat{h}_t \right) \\ &\leq \sum_{t=1}^T \left(4\sqrt{\frac{z_t}{\beta_t}} + 2\frac{u_t}{\beta_t} \right) \\ &\leq 4 \sum_{t=1}^T \sqrt{\frac{z_t}{\left(2 \sum_{s=1}^t \sqrt{z_s/\widehat{h}_s} \right)^{1/3}}} + 2 \sum_{t=1}^T \frac{u_t}{\sqrt{\sum_{s=1}^t u_t/\widehat{h}_t}} \\ &= 3.2G_1(z_{1:T}, \widehat{h}_{1:T}) + 2G_2(u_{1:T}, \widehat{h}_{1:T}), \end{aligned} \quad (24)$$

519 where the second inequality follows from the definition of β_t in (6) and the third inequality from
 520 (22) and (23). This completes the proof of the first statement in Lemma 4.

521 We next consider Rule 2 in (6). In this case, we can bound F as follows:

$$\begin{aligned} F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{1:T}) &\leq 2\sqrt{\frac{z_1}{\beta_1}} + \frac{u_1}{\beta_1} + \beta_1 h_1 + \sum_{t=2}^T \left(2\sqrt{\frac{z_t}{\beta_t}} + \frac{u_t}{\beta_t} + (\beta_t - \beta_{t-1})\widehat{h}_t \right) \\ &= 2\sqrt{\frac{z_1}{\beta_1}} + \frac{u_1}{\beta_1} + \beta_1 h_1 + \sum_{t=2}^T \left(2\sqrt{\frac{z_t}{\beta_t}} + \frac{u_t}{\beta_t} + 2\sqrt{\frac{z_{t-1}}{\beta_{t-1}}} + \frac{u_{t-1}}{\beta_{t-1}} \right) \\ &\leq \beta_1 h_1 + \sum_{t=1}^T \left(4\sqrt{\frac{z_t}{\beta_t}} + 2\frac{u_t}{\beta_t} \right), \end{aligned} \quad (25)$$

522 where the equality follows from (6).

523 We then first consider bounding $\sum_{t=1}^T \sqrt{z_t/\beta_t}$. We can lower-bound $\beta_t^{3/2}$ as

$$\beta_t^{3/2} \geq \beta_t^{1/2} \left(\beta_{t-1} + \frac{2}{\widehat{h}_t} \sqrt{z_{t-1}} \right) \geq \beta_{t-1}^{3/2} + \frac{2\sqrt{z_{t-1}}}{\widehat{h}_t} \geq \beta_1^{3/2} + 2 \sum_{s=2}^t \frac{\sqrt{z_{s-1}}}{\widehat{h}_s} =: (\beta_t^{(1)})^{3/2}, \quad (26)$$

524 where we define

$$\beta_t^{(1)} = \left(\beta_1^{3/2} + 2 \sum_{s=2}^t \frac{\sqrt{z_{s-1}}}{\widehat{h}_s} \right)^{2/3} = \left(\beta_1^{3/2} + 2 \sum_{s=1}^{t-1} \frac{\sqrt{z_s}}{\widehat{h}_{s+1}} \right)^{2/3} \leq \beta_t. \quad (27)$$

525 In the following, we will upper-bound $\sum_{t=1}^T \sqrt{z_t/\beta_t} \leq \sum_{t=1}^T \sqrt{z_t/\beta_t^{(1)}}$. Let $c = (1+\delta)^2$ for $\delta > 0$
 526 and we then define $\mathcal{S} = \{t \in [T]: \beta_{t+1}^{(1)} \leq c^2 \beta_t^{(1)}\}$ and $\mathcal{S}^c = [T] \setminus \mathcal{S} = \{t \in [T]: \beta_{t+1}^{(1)} >$
 527 $c^2 \beta_t^{(1)}\}$. From these definitions, we have

$$\sum_{t \in \mathcal{S}^c} \sqrt{\frac{z_t}{\beta_t^{(1)}}} \leq \sum_{t \in \mathcal{S}^c} \sqrt{\frac{z_{\max}}{\beta_t^{(1)}}} \leq \sum_{s=0}^{\infty} \left(\frac{1}{c}\right)^s \sqrt{\frac{z_{\max}}{\beta_1}} \leq \frac{1}{1-1/c} \sqrt{\frac{z_{\max}}{\beta_1}}. \quad (28)$$

528 Hence, using the last inequality, we obtain

$$\begin{aligned} \sum_{t=1}^T \sqrt{\frac{z_t}{\beta_t}} &\leq \sum_{t \in \mathcal{S}} \sqrt{\frac{z_t}{\beta_t^{(1)}}} + \sum_{t \in \mathcal{S}^c} \sqrt{\frac{z_t}{\beta_t^{(1)}}} \\ &\leq c \sum_{t \in \mathcal{S}} \sqrt{\frac{z_t}{\beta_{t+1}^{(1)}}} + \frac{1}{1-1/c} \sqrt{\frac{z_{\max}}{\beta_1}} \\ &\leq c \sum_{t \in \mathcal{S}} \sqrt{\frac{z_t}{\left(2 \sum_{s=1}^t \sqrt{z_s/\hat{h}_{s+1}}\right)^{2/3}}} + \frac{1}{1-1/c} \sqrt{\frac{z_{\max}}{\beta_1}} \\ &= \frac{c}{2^{1/3}} G_1(z_{1:T}, \hat{h}_{2:T+1}) + \frac{c}{c-1} \sqrt{\frac{z_{\max}}{\beta_1}}, \end{aligned} \quad (29)$$

529 where the third inequality follows from the definition of $\beta^{(1)}$ in (26).

530 We next bound $\sum_{t=1}^T u_t/\beta_t$. We can lower-bound β_t^2 as

$$\beta_t^2 \geq \beta_t \left(\beta_{t-1} + \frac{1}{\hat{h}_t} \frac{u_{t-1}}{\beta_{t-1}} \right) \geq \beta_{t-1}^2 + \frac{u_{t-1}}{\hat{h}_t} \geq \beta_1^2 + \sum_{s=2}^t \frac{u_{s-1}}{\hat{h}_s} =: (\beta_t^{(2)})^2, \quad (30)$$

531 where we define

$$\beta_t^{(2)} = \sqrt{\beta_1^2 + \sum_{s=2}^t \frac{u_{s-1}}{\hat{h}_s}} = \sqrt{\beta_1^2 + \sum_{s=1}^{t-1} \frac{u_s}{\hat{h}_{s+1}}} \leq \beta_t. \quad (31)$$

532 In the following, we will upper-bound $\sum_{t=1}^T u_t/\beta_t \leq \sum_{t=1}^T u_t/\beta_t^{(2)}$. Let us define $\mathcal{T} =$
 533 $\{t \in [T]: \beta_{t+1}^{(2)} \leq c\beta_t^{(2)}\}$ and $\mathcal{T}^c = [T] \setminus \mathcal{T} = \{t \in [T]: \beta_{t+1}^{(2)} > c\beta_t^{(2)}\}$. From these definitions,
 534 we have

$$\sum_{t \in \mathcal{T}^c} \frac{u_t}{\beta_t^{(2)}} \leq \sum_{t \in \mathcal{T}^c} \frac{u_{\max}}{\beta_t^{(2)}} \leq \sum_{s=0}^{\infty} \left(\frac{1}{c}\right)^s \frac{u_{\max}}{\beta_1} \leq \frac{1}{1-1/c} \frac{u_{\max}}{\beta_1}. \quad (32)$$

535 Hence, using the last inequality, we obtain

$$\begin{aligned} \sum_{t=1}^T \frac{u_t}{\beta_t} &\leq \sum_{t \in \mathcal{T}} \frac{u_t}{\beta_t^{(2)}} + \sum_{t \in \mathcal{T}^c} \frac{u_t}{\beta_t^{(2)}} \\ &\leq c \sum_{t \in \mathcal{T}} \frac{u_t}{\beta_{t+1}^{(2)}} + \frac{1}{1-1/c} \frac{u_{\max}}{\beta_1} \\ &\leq c \sum_{t \in \mathcal{T}} \frac{u_t}{\sqrt{\sum_{s=1}^t u_s/\hat{h}_{s+1}}} + \frac{1}{1-1/c} \frac{u_{\max}}{\beta_1} \\ &= c G_2(u_{1:T}, \hat{h}_{2:T+1}) + \frac{c}{c-1} \frac{z_{\max}}{\beta_1}. \end{aligned} \quad (33)$$

536 Finally, combining (25) with (29) and (33), we obtain

$$\begin{aligned} F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{1:T}) &\leq 3.2c G_1(z_{1:T}, \hat{h}_{2:T+1}) + 2c G_2(u_{1:T}, \hat{h}_{2:T+1}) \\ &\quad + \frac{c}{c-1} \left(2\sqrt{\frac{z_{\max}}{\beta_1}} + \frac{u_{\max}}{\beta_1} \right) + \beta_1 h_1. \end{aligned} \quad (34)$$

537 Setting $c = 1.25$ completes the proof. \square

538 **B.2 Proof of Lemma 5**

539 Before proving Lemma 5, we prepare the following lemma, a variant of [45, Lemma 4.13].

540 **Lemma 12.** *Let $\mathcal{T} \subseteq [T] = \{1, \dots, T\}$ and $(x_t)_{t \in \mathcal{T}}$ be a non-negative sequence. Then,*

$$\sum_{t \in \mathcal{T}} \frac{x_t}{\left(\sum_{s \in [t] \cap \mathcal{T}} x_s\right)^{1/3}} \leq \frac{3}{2} \left(\sum_{t \in \mathcal{T}} x_t\right)^{2/3}. \quad (35)$$

541 *Proof.* Let $S_t = \sum_{s \in [t] \cup \mathcal{T}} x_s$. Then,

$$\frac{x_t}{\left(\sum_{s \in [t] \cap \mathcal{T}} x_s\right)^{1/3}} = \frac{x_t}{S_t^{1/3}} = \int_{S_{t-1}}^{S_t} S_t^{-1/3} dz \leq \int_{S_{t-1}}^{S_t} z^{-1/3} dz = \frac{3}{2} \left(S_t^{2/3} - S_{t-1}^{2/3}\right). \quad (36)$$

542 Summing up the last inequality over \mathcal{T} , we obtain

$$\sum_{t \in \mathcal{T}} \frac{x_t}{\left(\sum_{s \in [t] \cap \mathcal{T}} x_s\right)^{1/3}} = \frac{3}{2} \sum_{t \in \mathcal{T}} \left(S_t^{2/3} - S_{t-1}^{2/3}\right) \leq \frac{3}{2} S_T^{2/3}, \quad (37)$$

543 where the last inequality follows from the telescoping argument with the assumption that $x_t \geq 0$. \square

544 *Proof of Lemma 5.* We upper-bound G_1 as follows:

$$\begin{aligned} G_1(z_{1:T}, h_{1:T}) &= \sum_{t=1}^T \frac{\sqrt{z_t}}{\left(\sum_{s=1}^t \sqrt{z_s}/h_s\right)^{1/3}} = \sum_{j=1}^{J+1} \sum_{t \in \mathcal{T}_j} \frac{\sqrt{z_t}}{\left(\sum_{s=1}^t \sqrt{z_s}/h_s\right)^{1/3}} \\ &\leq \sum_{j=1}^{J+1} \sum_{t \in \mathcal{T}_j} \frac{\sqrt{z_t}}{\left(\sum_{s \in \mathcal{T}_j \cap [t]} \sqrt{z_s}/h_s\right)^{1/3}} \leq \sum_{j=1}^{J+1} \sum_{t \in \mathcal{T}_j} \frac{\sqrt{z_t}}{\left(\sum_{s \in \mathcal{T}_j \cap [t]} \sqrt{z_s}/\theta_{j-1}\right)^{1/3}} \\ &= \sum_{j=1}^{J+1} \theta_{j-1}^{1/3} \sum_{t \in \mathcal{T}_j} \frac{\sqrt{z_t}}{\left(\sum_{s \in \mathcal{T}_j \cap [t]} \sqrt{z_s}\right)^{1/3}} \leq \frac{3}{2} \sum_{j=1}^{J+1} \left(\sqrt{\theta_{j-1}} \sum_{t \in \mathcal{T}_j} \sqrt{z_t}\right)^{2/3}, \quad (38) \end{aligned}$$

545 where the last inequality follows from Lemma 12. This completes the proof of the first statement in
546 Lemma 5. Setting $J = 0$ and $\theta_0 = h_{\max}$ in (38) yields that

$$G_1(z_{1:T}, h_{1:T}) \leq \frac{3}{2} \left(\sum_{t=1}^T \sqrt{z_t} h_{\max}\right)^{2/3}. \quad (39)$$

547 Setting $\theta_j = 2^{-j} h_{\max}$ for $j \in \{0\} \cup [J]$ in (38) also gives

$$\begin{aligned} G_1(z_{1:T}, h_{1:T}) &\leq \frac{3}{2} \sum_{j=1}^{J+1} \left(\sqrt{\theta_{j-1}} \sum_{t \in \mathcal{T}_j} \sqrt{z_t}\right)^{2/3} \\ &\leq \frac{3}{2} \sum_{j=1}^J \left(\sqrt{\frac{\theta_{j-1}}{\theta_j}} \sum_{t \in \mathcal{T}_j} \sqrt{z_t} h_t\right)^{2/3} + \frac{3}{2} \left(\sqrt{\theta_J} \sum_{t \in \mathcal{T}_J} \sqrt{z_t}\right)^{2/3} \\ &= \frac{3}{2} \sum_{j=1}^J \left(\sqrt{2} \sum_{t \in \mathcal{T}_j} \sqrt{z_t} h_t\right)^{2/3} + \frac{3}{2} \left(2^{-J/2} \sum_{t \in \mathcal{T}_J} \sqrt{z_t} h_{\max}\right)^{2/3} \\ &\leq \frac{3}{2} \left(\sqrt{2J} \sum_{j=1}^J \sum_{t \in \mathcal{T}_j} \sqrt{z_t} h_t\right)^{2/3} + \frac{3}{2} \left(2^{-J/2} \sum_{t \in \mathcal{T}_J} \sqrt{z_t} h_{\max}\right)^{2/3} \end{aligned}$$

(Hölder's inequality)

$$\leq \frac{3}{2} \left(\sqrt{2J} \sum_{t=1}^T \sqrt{z_t h_t} \right)^{2/3} + \frac{3}{2} \left(2^{-J/2} \sqrt{z_{\max} h_{\max}} \right)^{2/3} T^{2/3}, \quad (40)$$

548 where the second inequality follows from $(x+y)^{2/3} \leq x^{2/3} + y^{2/3}$ for $x, y \geq 0$. Combining the
549 last inequality and (39) completes the proof of the second statement in Lemma 5. \square

550 **C Proof for best-of-both-worlds analysis in general online learning** 551 **framework (Theorem 7, Section 4)**

552 This section provides the proof of Theorem 7.

553 *Proof.* From Assumption (i), the regret is bounded as

$$\text{Reg}_T \leq \mathbb{E} \left[\sum_{t=1}^T \langle \hat{\ell}_t, q_t - e_{a^*} \rangle + 2 \sum_{t=1}^T \gamma_t \right]. \quad (41)$$

554 From the standard FTRL analysis in [36, Exercise 28.12], we obtain

$$\sum_{t=1}^T \langle \hat{\ell}_t, q_t - e_{a^*} \rangle \leq \sum_{t=1}^T \left(\langle \hat{\ell}_t, q_t - q_{t+1} \rangle - \beta_t D_{(-H_\alpha)}(q_{t+1}, q_t) + (\beta_t - \beta_{t-1}) h_t \right) + \bar{\beta} \bar{h}. \quad (42)$$

555 Combining the last two inequalities, we obtain

$$\begin{aligned} \text{Reg}_T &\leq \mathbb{E} \left[\sum_{t=1}^T \left(\langle \hat{\ell}_t, q_t - q_{t+1} \rangle - \beta_t D_{(-H_\alpha)}(q_{t+1}, q_t) + (\beta_t - \beta_{t-1}) h_t + 2\gamma_t \right) + \bar{\beta} \bar{h} \right] \\ &\lesssim \mathbb{E} \left[\sum_{t=1}^T \left(\frac{z_t}{\beta_t \gamma'_t} + (\beta_t - \beta_{t-1}) h_t + \gamma_t \right) + \bar{\beta} \bar{h} \right] \quad (\text{Assumption (ii) in (12)}) \\ &\lesssim \mathbb{E} \left[\sum_{t=1}^T \left(\frac{z_t}{\beta_t \gamma'_t} + (\beta_t - \beta_{t-1}) h_t + \gamma'_t + \frac{u_t}{\beta_t} \right) + \bar{\beta} \bar{h} \right] \quad (\text{definition of } \gamma_t \text{ in (11)}) \\ &\lesssim \mathbb{E} \left[\sum_{t=1}^T \left(\sqrt{\frac{z_t}{\beta_t}} + \frac{u_t}{\beta_t} + (\beta_t - \beta_{t-1}) h_{t-1} \right) + \bar{\beta} \bar{h} \right] \quad (\text{definition of } \gamma'_t \text{ and Assumption (iii)}) \\ &\lesssim \mathbb{E}[F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{0:T-1})] + \bar{\beta} \bar{h}, \quad (43) \end{aligned}$$

556 where the last inequality follows from (5). Now, since β_t follows Rule 2 in (6) with $\hat{h}_t = h_{t-1}$,
557 Eq. (9) in Theorem 6 gives

$$F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{0:T-1}) \lesssim \left(\sum_{t=1}^T \sqrt{z_t h_1} \right)^{\frac{2}{3}} + \sqrt{\sum_{t=1}^T u_t h_1} + \sqrt{\frac{z_{\max}}{\beta_1} + \frac{u_{\max}}{\beta_1}} + \beta_1 h_1, \quad (44)$$

$$\begin{aligned} F(\beta_{1:T}, z_{1:T}, u_{1:T}, h_{0:T-1}) &\lesssim \inf_{\varepsilon \geq 1/T} \left\{ \left(\sum_{t=1}^T \sqrt{z_t h_t \log(\varepsilon T)} \right)^{\frac{2}{3}} + \left(\frac{\sqrt{z_{\max} h_1}}{\varepsilon} \right)^{\frac{2}{3}} \right. \\ &\quad \left. + \sqrt{\sum_{t=1}^T u_t h_t \log(\varepsilon T) + \frac{u_{\max} h_1}{\varepsilon}} \right\} + \sqrt{\frac{z_{\max}}{\beta_1} + \frac{u_{\max}}{\beta_1}} + \beta_1 h_1. \quad (45) \end{aligned}$$

558 Hence, in the adversarial regime, combining (43) and (44) gives

$$\text{Reg}_T \lesssim \mathbb{E} \left[\left(\sum_{t=1}^T \sqrt{z_t h_1} \right)^{2/3} + \sqrt{\sum_{t=1}^T u_t h_1} \right] + \kappa \leq (z_{\max} h_1)^{1/3} T^{2/3} + \sqrt{u_{\max} h_1 T} + \kappa, \quad (46)$$

559 where we recall that $\kappa = \sqrt{z_{\max}/\beta_1} + u_{\max}/\beta_1 + \beta_1 h_1 + \bar{\beta} \bar{h}$. This completes the proof of (13).

560 We next consider the adversarial regime with a (Δ, C, T) -self-bounding constraint. For any $\varepsilon \geq 1/T$,
 561 combining (43) and (45) gives

$$\begin{aligned} \text{Reg}_T &\lesssim \mathbb{E} \left[\left(\sum_{t=1}^T \sqrt{z_t h_t \log(\varepsilon T)} \right)^{\frac{2}{3}} + \sqrt{\sum_{t=1}^T u_t h_t \log(\varepsilon T)} + \left(\frac{\sqrt{z_{\max} h_1}}{\varepsilon} \right)^{\frac{2}{3}} + \sqrt{\frac{u_{\max} h_1}{\varepsilon}} + \kappa \right] \\ &\leq \left(\mathbb{E} \left[\sum_{t=1}^T \sqrt{z_t h_t} \right] \sqrt{\log(\varepsilon T)} \right)^{\frac{2}{3}} + \sqrt{\mathbb{E} \left[\sum_{t=1}^T u_t h_t \right] \log(\varepsilon T)} + \left(\frac{\sqrt{z_{\max} h_1}}{\varepsilon} \right)^{\frac{2}{3}} + \sqrt{\frac{u_{\max} h_1}{\varepsilon}} + \kappa, \end{aligned} \quad (47)$$

562 where the last inequality follows from Jensen's inequality. Now, using the assumption (14) and defin-
 563 ing $Q(a^*) = \mathbb{E}[\sum_{t=1}^T (1 - q_{ta^*})] \in [0, T]$, we have

$$\mathbb{E} \left[\sum_{t=1}^T \sqrt{z_t h_t} \right] \leq \sqrt{\rho_1} \mathbb{E} \left[\sum_{t=1}^T (1 - q_{ta^*}) \right] = \sqrt{\rho_1} Q(a^*), \quad (48)$$

$$\mathbb{E} \left[\sum_{t=1}^T u_t h_t \right] \leq \rho_2 \mathbb{E} \left[\sum_{t=1}^T (1 - q_{ta^*}) \right] = \rho_2 Q(a^*). \quad (49)$$

564 Since we consider the adversarial regime with a (Δ, C, T) -self-bounding constraint, the regret is
 565 lower-bounded as

$$\begin{aligned} \text{Reg}_T &\geq \mathbb{E} \left[\sum_{t=1}^T \langle \Delta, p \rangle \right] - C \geq \frac{1}{2} \mathbb{E} \left[\sum_{t=1}^T \langle \Delta, q \rangle \right] - C \\ &\geq \frac{1}{2} \Delta_{\min} \mathbb{E} \left[\sum_{t=1}^T (1 - q_{ta^*}) \right] - C = \frac{1}{2} \Delta_{\min} Q(a^*) - C, \end{aligned} \quad (50)$$

566 where the second inequality follows from $p = (1 - \gamma_t)q_t + \gamma_t p_0 \geq q_t/2$. Hence, combining (47)
 567 with (48), (49) and (50), we can bound the regret for any $\lambda \in (0, 1]$ as follows:

$$\begin{aligned} \text{Reg}_T &= (1 + \lambda) \text{Reg}_T - \lambda \text{Reg}_T \\ &\lesssim (1 + \lambda) \left(\sqrt{\rho_1} Q(a^*) \sqrt{\log(\varepsilon T)} \right)^{2/3} - \frac{\lambda}{4} \Delta_{\min} Q(a^*) + (1 + \lambda) \sqrt{\rho_2 Q(a^*) \log(\varepsilon T)} - \frac{\lambda}{4} \Delta_{\min} Q(a^*) \\ &\quad + (1 + \lambda) \left(\left(\frac{\sqrt{z_{\max} h_1}}{\varepsilon} \right)^{2/3} + \sqrt{\frac{u_{\max} h_1}{\varepsilon}} + \kappa \right) + \lambda C \\ &\lesssim \frac{(1 + \lambda)^3}{\lambda^2} \frac{\rho_1 \log(\varepsilon T)}{\Delta_{\min}^2} + \frac{(1 + \lambda)^2}{\lambda} \frac{\rho_2 \log(\varepsilon T)}{\Delta_{\min}} + \left(\frac{\sqrt{z_{\max} h_1}}{\varepsilon} \right)^{2/3} + \sqrt{\frac{u_{\max} h_1}{\varepsilon}} + \kappa + \lambda C \\ &\lesssim \frac{\rho_1 \log(\varepsilon T)}{\Delta_{\min}^2} + \frac{\rho_2 \log(\varepsilon T)}{\Delta_{\min}} + \frac{1}{\lambda^2} \left(\frac{\rho_1 \log(\varepsilon T)}{\Delta_{\min}^2} + \frac{\rho_2 \log(\varepsilon T)}{\Delta_{\min}} \right) + \left(\frac{\sqrt{z_{\max} h_1}}{\varepsilon} \right)^{2/3} + \sqrt{\frac{u_{\max} h_1}{\varepsilon}} + \kappa + \lambda C \\ &\lesssim \frac{\rho \log(\varepsilon T)}{\Delta_{\min}^2} + \frac{1}{\lambda^2} \frac{\rho \log(\varepsilon T)}{\Delta_{\min}^2} + \left(\frac{\sqrt{z_{\max} h_1}}{\varepsilon} \right)^{2/3} + \sqrt{\frac{u_{\max} h_1}{\varepsilon}} + \kappa + \lambda C, \end{aligned} \quad (51)$$

568 where in the first inequality we used (47) with (48), (49), (50), and Jensen's inequality, in the second
 569 inequality we used $ax^2 - bx^3 \leq 4a^3/(27b^2)$ for $a \geq 0, b > 0$ and $x \geq 0$ and $ax - bx^2 \leq$
 570 $a^2/(4b)$ for $a \geq 0, b > 0$ and $x \geq 0$ and in the third inequality we used $\lambda \in (0, 1]$. Setting
 571 $\lambda = \Theta((\rho \log(\varepsilon T)/C)^{1/3})$ in the last inequality, we obtain

$$\text{Reg}_T \lesssim \frac{\rho \log(\varepsilon T)}{\Delta_{\min}^2} + \left(\frac{C^2 \rho \log(\varepsilon T)}{\Delta_{\min}^2} \right)^{1/3} + \left(\frac{\sqrt{z_{\max} h_1}}{\varepsilon} \right)^{2/3} + \sqrt{\frac{u_{\max} h_1}{\varepsilon}} + \kappa.$$

572 Finally, when $T \geq \tau = 1/\Delta_{\min}^2 + C/\Delta_{\min}$, setting

$$\varepsilon = \frac{1}{\rho^2/\Delta_{\min}^2 + C\rho/\Delta_{\min}} \geq \frac{1}{T} \quad (52)$$

573 yields that

$$\begin{aligned}
\text{Reg}_T &\lesssim \frac{\rho}{\Delta_{\min}^2} \log_+ \left(\frac{T}{1/\Delta_{\min}^2 + C/\Delta_{\min}} \right) + \left(\frac{C^2 \rho}{\Delta_{\min}^2} \log_+ \left(\frac{T}{1/\Delta_{\min}^2 + C/\Delta_{\min}} \right) \right)^{1/3} \\
&\quad + (z_{\max} h_1)^{1/3} \left(\frac{1}{\Delta_{\min}^2} + \frac{C}{\Delta_{\min}} \right)^{2/3} + \sqrt{u_{\max} h_1} \sqrt{\frac{1}{\Delta_{\min}^2} + \frac{C}{\Delta_{\min}}} + \kappa \\
&\lesssim \frac{\rho}{\Delta_{\min}^2} \log_+ (T \Delta_{\min}^2) + \left(\frac{C^2 \rho}{\Delta_{\min}^2} \log_+ \left(\frac{T \Delta_{\min}}{C} \right) \right)^{1/3} \\
&\quad + \left((z_{\max} h_1)^{1/3} + \sqrt{u_{\max} h_1} \right) \left(\frac{1}{\Delta_{\min}^2} + \frac{C}{\Delta_{\min}} \right)^{2/3} + \kappa, \tag{53}
\end{aligned}$$

574 which completes the proof. \square

575 D Auxiliary lemmas

576 This section provides auxiliary lemmas useful for proving the BOBW guarantee.

577 **Lemma 13.** *Let $\alpha \in (0, 1)$ and $i^* \in [k]$. Then, the α -Tsallis entropy H_α is bounded from above as*

$$H_\alpha(q) = \frac{1}{\alpha} \sum_{i=1}^k (q_i^\alpha - q_i) \leq \frac{1}{\alpha} (k-1)^\alpha (1 - q_{i^*})^\alpha \tag{54}$$

578 for any $q \in \mathcal{P}_k$.

579 *Proof.* From Jensen's inequality and the fact that $x \mapsto x^\alpha$ is concave for $\alpha \in (0, 1)$,

$$\begin{aligned}
\sum_{i=1}^k (q_i^\alpha - q_i) &\leq \sum_{i \neq i^*} q_i^\alpha = (k-1) \sum_{i \neq i^*} \frac{1}{k-1} q_i^\alpha \leq (k-1) \left(\frac{1}{k-1} \sum_{i \neq i^*} q_i \right)^\alpha \\
&= (k-1)^{1-\alpha} \left(\sum_{i \neq i^*} q_i \right)^\alpha = (k-1)^{1-\alpha} (1 - q_{i^*})^\alpha, \tag{55}
\end{aligned}$$

580 which completes the proof. \square

581 **Lemma 14** ([26, Lemma 10]). *Let $q \in \mathcal{P}_k$ and $\tilde{I} \in \arg \max_{i \in [k]} q_i$. For $\ell \in \mathbb{R}^k$, if $|\ell_i| \leq$*

582 $\frac{1-\alpha}{4} \frac{1}{\min\{q_{\tilde{I}}, 1-q_{\tilde{I}}\}^{1-\alpha}}$ for all $i \in [k]$, it holds that

$$\max_{p \in \mathcal{P}_k} \{ \langle \ell, q - p \rangle - D_{(-H_\alpha)}(p, q) \} \leq \frac{4}{1-\alpha} \left(\sum_{i \neq \tilde{I}} q_i^{2-\alpha} \ell_i^2 + \min\{q_{\tilde{I}}, 1-q_{\tilde{I}}\}^{2-\alpha} \ell_{\tilde{I}}^2 \right). \tag{56}$$

583 **Lemma 15** ([26, Lemmas 11 and 12]). *Let $L \in \mathbb{R}^k$ and $\ell \in \mathbb{R}^k$ and suppose that $q, r \in \mathcal{P}_k$ are*

584 given by

$$\begin{aligned}
q &\in \arg \min_{p \in \mathcal{P}_k} \{ \langle L, p \rangle + \beta(-H_\alpha(p)) + \bar{\beta}(-H_{\bar{\alpha}}(p)) \} \\
r &\in \arg \min_{p \in \mathcal{P}_k} \{ \langle L + \ell, p \rangle + \beta'(-H_\alpha(p)) + \bar{\beta}(-H_{\bar{\alpha}}(p)) \} \tag{57}
\end{aligned}$$

585 for the Tsallis entropy H_α and $H_{\bar{\alpha}}$, $0 < \beta \leq \beta'$. Suppose also that

$$\|\ell\|_\infty \leq \max \left\{ \frac{1 - (\sqrt{2})^{\alpha-1}}{2} q_*^{\alpha-1} \beta, \frac{1 - (\sqrt{2})^{\bar{\alpha}-1}}{2} q_*^{\bar{\alpha}-1} \bar{\beta} \right\}, \tag{58}$$

$$0 \leq \beta' - \beta \leq \max \left\{ \left(1 - (\sqrt{2})^{\alpha-1} \right) \beta, \frac{1 - (\sqrt{2})^{\bar{\alpha}-1}}{\sqrt{2}} q_*^{\bar{\alpha}-\alpha} \bar{\beta} \right\}. \tag{59}$$

586 Then, it holds that $H_\alpha(r) \leq 2H_\alpha(q)$.

587 **E Proof for partial monitoring (Theorem 8, Section 5)**

588 This section provides the proof of Theorem 8.

589 *Proof of Theorem 8.* It suffices to prove that assumptions in Theorem 7 are satisfied. We first verify
 590 Assumptions (i)–(iii) in (12). Let us start from checking Assumption (i). From the definition of the
 591 loss difference estimator \hat{y}_t , the regret is bounded as

$$\begin{aligned}
 \text{Reg}_T &= \mathbb{E} \left[\sum_{t=1}^T (\mathcal{L}_{A_t x_t} - \mathcal{L}_{a^* x_t}) \right] = \mathbb{E} \left[\sum_{t=1}^T \langle p_t - e_{a^*}, \mathcal{L} e_{x_t} \rangle \right] \\
 &= \mathbb{E} \left[\sum_{t=1}^T \langle q_t - e_{a^*}, \mathcal{L} e_{x_t} \rangle + \sum_{t=1}^T \gamma_t \left\langle \frac{1}{k} \mathbf{1} - q_t, \mathcal{L} e_{x_t} \right\rangle \right] \\
 &\leq \mathbb{E} \left[\sum_{t=1}^T \langle q_t - e_{a^*}, \mathcal{L} e_{x_t} \rangle + \sum_{t=1}^T \gamma_t \right] = \mathbb{E} \left[\sum_{t=1}^T \sum_{a=1}^k q_{ta} (\mathcal{L}_{a x_t} - \mathcal{L}_{a^* x_t}) + \sum_{t=1}^T \gamma_t \right] \\
 &= \mathbb{E} \left[\sum_{t=1}^T \sum_{a=1}^k q_{ta} (\hat{y}_{ta} - \hat{y}_{ta^*}) + \sum_{t=1}^T \gamma_t \right] = \mathbb{E} \left[\sum_{t=1}^T \langle q_t - e_{a^*}, \hat{y}_t \rangle + \sum_{t=1}^T \gamma_t \right], \quad (60)
 \end{aligned}$$

592 where the inequality holds since $\mathcal{L} \in [0, 1]^{k \times d}$. This implies that Assumption (i) is indeed satisfied.

593 We next check Assumption (ii) in (12). For any $b \in [k]$ we have

$$\left| \frac{\hat{y}_{tb}}{\beta_t} \right| = \left| \frac{G(A_t, \sigma_t)_b}{\beta_t p_{tA_t}} \right| \leq \frac{|G(A_t, \sigma_t)_b| k}{\beta_t \gamma_t} \leq \frac{c_{\mathcal{G}}}{\beta_t \gamma_t} \leq \frac{c_{\mathcal{G}}}{u_t} = \frac{1 - \alpha}{8} \frac{1}{(\min\{q_{t\bar{i}_t}, 1 - q_{t\bar{i}_t}\})^{1-\alpha}}, \quad (61)$$

594 where the third inequality follows from $\gamma_t \geq u_t/\beta_t$ and the last equality follows from the definition
 595 of u_t in (17). Hence, from Lemma 14 the LHS of Assumption (ii) is bounded as

$$\begin{aligned}
 &\mathbb{E}_t [\langle \hat{y}_t, q_t - q_{t+1} \rangle - \beta_t D_{(-H_\alpha)}(q_{t+1}, q_t)] = \beta_t \mathbb{E}_t \left[\left\langle \frac{\hat{y}_t}{\beta_t}, q_t - q_{t+1} \right\rangle - D_{(-H_\alpha)}(q_{t+1}, q_t) \right] \\
 &\leq \mathbb{E}_t \left[\frac{4}{\beta_t (1 - \alpha)} \left(\sum_{i \neq \bar{i}_t} q_{ti}^{2-\alpha} \hat{y}_{ti}^2 + (\min\{q_{t\bar{i}_t}, 1 - q_{t\bar{i}_t}\})^{2-\alpha} \hat{y}_{t\bar{i}_t}^2 \right) \right] \\
 &= \frac{4}{\beta_t (1 - \alpha)} \left(\sum_{i \neq \bar{i}_t} q_{ti}^{2-\alpha} \mathbb{E}_t [\hat{y}_{ti}^2] + q_{t*}^{2-\alpha} \mathbb{E}_t [\hat{y}_{t\bar{i}_t}^2] \right). \quad (62)
 \end{aligned}$$

596 Since the variance of \hat{y}_t is bounded from above as

$$\mathbb{E}_t [\hat{y}_{ti}^2] = \sum_{a=1}^k p_{ta} \frac{G(a, \sigma_t)_i^2}{p_{ta}^2} \leq \sum_{a=1}^k \frac{k \|G\|_\infty^2}{\gamma_t} = \frac{c_{\mathcal{G}}^2}{\gamma_t} \quad (63)$$

597 for any $i \in [k]$, the LHS of Assumption (ii) is further bounded as

$$\mathbb{E}_t [\langle \hat{y}_t, q_t - q_{t+1} \rangle - \beta_t D_{\psi_t}(q_{t+1}, q_t)] \leq \frac{4c_{\mathcal{G}}^2}{\beta_t \gamma_t (1 - \alpha)} \left(\sum_{i \neq \bar{i}_t} q_{ti}^{2-\alpha} + q_{t*}^{2-\alpha} \right) = \frac{z_t}{\beta_t \gamma_t} \leq \frac{z_t}{\beta_t \gamma_t'}, \quad (64)$$

598 which implies that Assumption (ii) in (12) is satisfied.

599 Next, we will prove $h_{t+1} \lesssim h_t$ of Assumption (iii) in (12). To prove this, we will check the condition
 600 in Lemma 15. For any $a \in [k]$,

$$|\hat{y}_{ta}| \leq \frac{\|G\|_\infty}{p_{tA_t}} \leq \frac{k \|G\|_\infty}{\gamma_t} \leq \frac{c_{\mathcal{G}} \beta_t}{u_t} \leq \frac{1 - \alpha}{8} \frac{\beta_t}{q_{t*}^{1-\alpha}} \leq \frac{1 - (\sqrt{2})^{\alpha-1}}{2} \frac{\beta_t}{q_{t*}^{1-\alpha}}, \quad (65)$$

601 where the second inequality follows from $p_{ta} \geq \gamma_t/k$, the third inequality from $\gamma_t \geq u_t/\beta_t$, and the
 602 last inequality from the fact that $(1 - x)/4 \leq 1 - (\sqrt{2})^{x-1}$ for $x \in [0, 1]$. Thus, the condition (58)
 603 is satisfied.

604 We next check the condition (59). Recalling $q_{t^*} = \min\{q_{t\tilde{I}_t}, 1 - q_{t\tilde{I}_t}\}$, the parameters z_t and u_t
605 satisfy

$$\sqrt{z_t} = \frac{2c_G}{\sqrt{1-\alpha}} \sqrt{\sum_{i \neq \tilde{I}_t} q_{ti}^{2-\alpha} + q_{t^*}^{2-\alpha}} \leq \frac{2\sqrt{k}c_G}{\sqrt{1-\alpha}} q_{t^*}^{1-\frac{1}{2}\alpha}, \quad u_t = \frac{8c_G}{1-\alpha} q_{t^*}^{1-\alpha}, \quad (66)$$

606 where the inequality follows from $q_{ti} \leq q_{t^*}$ for $i \neq \tilde{I}_t$. The penalty component h_t is lower-bounded
607 as

$$h_t = H_\alpha(q_t) = \frac{1}{\alpha} \sum_{i=1}^k (q_{ti}^\alpha - q_{ti}) \geq \frac{1 - (1/2)^{1-\alpha}}{\alpha} q_{t^*}^\alpha \geq \frac{1-\alpha}{4\alpha} q_{t^*}^\alpha, \quad (67)$$

608 where the last inequality in (67) follows from $1 - (1/2)^{1-x} \geq (1-x)/4$ for $x \leq 0$, and the first
609 inequality can be proven as follows: when $q_{t\tilde{I}_t} \leq 1/2$, it holds that $\sum_{i=1}^k (q_{ti}^\alpha - q_{ti}) \geq q_{t\tilde{I}_t}^\alpha - q_{t\tilde{I}_t} =$
610 $q_{t\tilde{I}_t}^\alpha (1 - q_{t\tilde{I}_t}^{1-\alpha}) \geq q_{t\tilde{I}_t}^\alpha (1 - (1/2)^{1-\alpha}) = q_{t^*}^\alpha (1 - (1/2)^{1-\alpha})$, and when $q_{t\tilde{I}_t} > 1/2$, it holds that
611 $\sum_{i=1}^k (q_{ti}^\alpha - q_{ti}) \geq \sum_{i=1}^k q_{ti}^\alpha - 1 \geq \sum_{i \neq \tilde{I}_t} q_{ti}^\alpha + (1/2)^\alpha - 1 \geq (\sum_{i \neq \tilde{I}_t} q_{ti})^\alpha + (1/2)^\alpha - 1 =$
612 $(1 - q_{t\tilde{I}_t})^\alpha + (1/2)^\alpha - 1 = q_{t^*}^\alpha + (1/2)^\alpha - 1 \geq q_{t^*}^\alpha (1 - (1/2)^{1-\alpha})$. Using the bounds on z_t , u_t ,
613 and h_t in (66) and (67), we have

$$\begin{aligned} \beta_{t+1} - \beta_t &= \frac{1}{\tilde{h}_{t+1}} \left(2\sqrt{\frac{z_t}{\beta_t}} + \frac{u_t}{\beta_t} \right) = \frac{2}{h_t} \sqrt{\frac{z_t}{\beta_t}} + \frac{1}{h_t} \frac{u_t}{\beta_t} \\ &\leq \frac{16\alpha c_G \sqrt{k}}{\sqrt{\beta_1} (1-\alpha)^{3/2}} q_{t^*}^{1-\frac{3}{2}\alpha} + \frac{32\alpha c_G}{\sqrt{\beta_1} (1-\alpha)^2} q_{t^*}^{1-2\alpha} \\ &\leq \alpha \bar{\beta} q_{t^*}^{1-\frac{3}{2}\alpha} + \alpha \bar{\beta} q_{t^*}^{1-2\alpha} \\ &\leq 2(1-\bar{\alpha}) \bar{\beta} q_{t^*}^{\bar{\alpha}-\alpha} \leq 2 \frac{1 - (\sqrt{2})^{\bar{\alpha}-1}}{\sqrt{2}} \bar{\beta} q_{t^*}^{\bar{\alpha}-\alpha}, \end{aligned} \quad (68)$$

614 where the first inequality follows from (66), (67), and the fact that $\beta_t \geq \beta_1 \geq 1$, the second inequality
615 from the definition of $\bar{\beta}$ in (17), the third inequality from $\min\{1 - \frac{3}{2}\alpha, 1 - 2\alpha\} \geq \bar{\alpha} - \alpha$ since
616 $\bar{\alpha} = 1 - \alpha$, and the last inequality from $1 - x \leq (1 - (\sqrt{2})^{x-1})/\sqrt{2}$ for $x \leq 1$. Therefore, the
617 condition (59) is satisfied. Hence, from Lemma 15, we have $h_{t+1} = H_\alpha(q_{t+1}) \leq 2H_\alpha(q_t) = 2h_t$,
618 which implies that Assumption (iii) in (12) is satisfied.

619 Finally, we check the assumption (14) in Theorem 7. We first consider the first inequality in (14).
620 From the definition of z_t and the fact that $q_{ti} \leq q_{t\tilde{I}_t}$ for $i \neq \tilde{I}_t$, the stability component z_t is bounded
621 as

$$\begin{aligned} z_t &= \frac{4c_G^2}{1-\alpha} \left\{ \sum_{i \neq \tilde{I}_t} q_{ti}^{2-\alpha} + (\min\{q_{t\tilde{I}_t}, 1 - q_{t\tilde{I}_t}\})^{2-\alpha} \right\} \\ &\leq \frac{4c_G^2}{1-\alpha} \left\{ \sum_{i \neq \tilde{I}_t} q_{ti}^{2-\alpha} + \left(\sum_{i \neq \tilde{I}_t} q_{ti} \right)^{2-\alpha} \right\} \\ &\leq \frac{8c_G^2}{1-\alpha} \left(\sum_{i \neq \tilde{I}_t} q_{ti} \right)^{2-\alpha} \leq \frac{8c_G^2}{1-\alpha} \left(\sum_{i \neq a^*} q_{ti} \right)^{2-\alpha} = \frac{8c_G^2}{1-\alpha} (1 - q_{ta^*})^{2-\alpha}, \end{aligned} \quad (69)$$

622 where the second inequality holds from the inequality $x^a + y^a \leq (x+y)^a$ for $x, y \geq 0$ and $a \in [0, 1]$,
623 and the third inequality from $q_{ti} \leq q_{t\tilde{I}_t}$ for $i \neq \tilde{I}_t$. From Lemma 13, we also obtain that

$$h_t = H_\alpha(q_t) \leq \frac{1}{\alpha} (k-1)^{1-\alpha} (1 - q_{ta^*})^\alpha. \quad (70)$$

624 Hence, combining this with (69), we obtain

$$z_t h_t \leq \frac{8c_G^2}{1-\alpha} (1 - q_{ta^*})^{2-\alpha} \cdot \frac{1}{\alpha} (k-1)^{1-\alpha} (1 - q_{ta^*})^\alpha = \underbrace{\frac{8c_G^2 (k-1)^{1-\alpha}}{\alpha(1-\alpha)}}_{=\rho_1} (1 - q_{ta^*})^2. \quad (71)$$

625 We next consider the second inequality in (14). We can bound u_t from above as

$$\begin{aligned} u_t &= \frac{8c_{\mathcal{G}}}{1-\alpha} (\min\{q_{t\bar{I}_t}, 1 - q_{t\bar{I}_t}\})^{1-\alpha} \leq \frac{8c_{\mathcal{G}}}{1-\alpha} \left(\sum_{i \neq \bar{I}_t} q_{ti} \right)^{1-\alpha} \\ &\leq \frac{8c_{\mathcal{G}}}{1-\alpha} \left(\sum_{i \neq a^*} q_{ti} \right)^{1-\alpha} = \frac{8c_{\mathcal{G}}}{1-\alpha} (1 - q_{ta^*})^{1-\alpha}, \end{aligned} \quad (72)$$

626 where the second inequality follows from $q_{t\bar{I}_t} \geq q_{ti}$ for all $i \in [k]$. Hence, combining the last two
627 inequality and (70),

$$u_t h_t \leq \underbrace{\frac{4c_{\mathcal{G}}(k-1)^{1-\alpha}}{\alpha(1-\alpha)}}_{=\rho_2} (1 - q_{ta^*}). \quad (73)$$

628 Hence, the assumption (14) is satisfied with above ρ_1 and ρ_2 , and thus we have completed the proof.
629 \square

630 F Proof for graph bandits (Theorem 10, Section 6)

631 This section provides the missing detail of Section 6.

632 F.1 Fractional domination number

633 Before introducing the fractional domination number, we define the domination number $\tilde{\delta} \leq \delta$. A
634 *dominating set* $D \subseteq V$ is a set of vertices such that $V \subseteq \bigcup_{i \in D} N^{\text{out}}(i)$. The *domination number*
635 $\tilde{\delta}(G)$ of graph G is the size of the smallest dominating set. From the definition, the domination
636 number $\tilde{\delta}$ can also be written as the optimal value of the following optimization problem:

$$\text{minimize } \sum_{i \in V} x_i \quad \text{subject to } \sum_{i \in N^{\text{in}}(j)} x_i \geq 1 \quad \forall j \in V, \quad x_i \in \{0, 1\} \quad \forall i \in V, \quad (74)$$

637 where $x_i \in \{0, 1\}$ a binary variable indicating whether vertex i is in the dominating set ($x_i = 1$) or
638 not ($x_i = 0$).

639 Then, one can see that the fractional domination number δ^* is defined as the optimal value of the
640 following optimization problem, in which the variables $(x_i)_{i \in V}$ are allowed to take values in $[0, 1]$
641 instead of $\{0, 1\}$:

$$\text{minimize } \sum_{i \in V} x_i \quad \text{subject to } \sum_{i \in N^{\text{in}}(j)} x_i \geq 1 \quad \forall j \in V, \quad 0 \leq x_i \leq 1 \quad \forall i \in V, \quad (75)$$

642 which is the linear program provided in (19). From the definitions, the fractional domination number
643 is less than or equal to the domination number, $\delta^* \leq \tilde{\delta}$. Another advantage of using δ^* instead of $\tilde{\delta}$ is
644 that the fractional domination number δ^* can be computed in polynomial time, while the computation
645 of the domination number $\tilde{\delta}$ is NP-hard. See [13] for more benefits of using the fractional version of
646 the (weak) domination number.

647 F.2 Proof of Theorem 10

648 Here, we provide the proof of Theorem 10.

649 *Proof.* It suffices to prove that assumptions in Theorem 7 are satisfied. We first verify Assumptions
650 (i)–(iii) in (12). We start from checking Assumption (i). The regret is bounded as

$$\begin{aligned} \text{Reg}_T &= \mathbb{E} \left[\sum_{t=1}^T \ell_t(A_t) - \sum_{t=1}^T \ell_t(a^*) \right] = \mathbb{E} \left[\sum_{t=1}^T \langle \ell_t, p_t - e_{a^*} \rangle \right] = \mathbb{E} \left[\sum_{t=1}^T \langle \ell_t, q_t - e_{a^*} \rangle + \sum_{t=1}^T \langle \ell_t, p_t - q_t \rangle \right] \\ &= \mathbb{E} \left[\sum_{t=1}^T \langle \ell_t, q_t - e_{a^*} \rangle + \sum_{t=1}^T \gamma_t \langle \ell_t, q_t - u \rangle \right] \leq \mathbb{E} \left[\sum_{t=1}^T \langle \hat{\ell}_t, q_t - e_{a^*} \rangle + \sum_{t=1}^T \gamma_t \right], \end{aligned} \quad (76)$$

651 where the third equality follows from the definition of γ_t . This implies that Assumption (i) is indeed
 652 satisfied.

653 We next check Assumption (ii) in (12). Now, recalling the definition of the fractional domination
 654 number and the optimal value x^* of (19), and $u_i = x_i^*/\sum_{j \in V} x_j^*$, we have

$$\sum_{j \in N^{\text{in}}(i)} u_j = \frac{\sum_{j \in N^{\text{in}}(i)} x_j^*}{\sum_{i \in V} x_i^*} \geq \frac{1}{\sum_{i \in V} x_i^*} = \frac{1}{\delta^*}, \quad (77)$$

655 where the inequality follows from the first constraint in (19). Hence, combining this with the defini-
 656 tion of $p_t = (1 - \gamma_t)q_t + \gamma_t u$, we can lower-bound P_{ti} as

$$P_{ti} = \sum_{j \in N^{\text{in}}(i)} p_{tj} \geq \gamma_t \sum_{j \in N^{\text{in}}(i)} u_j \geq \frac{\gamma_t}{\delta^*} \quad \text{for all } i \in V. \quad (78)$$

657 This lower bound yields that for any $i \in V$

$$\left| \frac{\widehat{\ell}_{ti}}{\beta_t} \right| = \frac{\ell_{ti}}{\beta_t P_{ti}} \leq \frac{\delta^*}{\beta_t \gamma_t} = \frac{\delta^*}{u_t} = \frac{1 - \alpha}{8} \frac{1}{(\min\{q_{t\bar{I}_t}, 1 - q_{t\bar{I}_t}\})^{1-\alpha}}. \quad (79)$$

658 Hence, from Lemma 14 we obtain

$$\begin{aligned} & \mathbb{E}_t \left[\left\langle \widehat{\ell}_t, q_t - q_{t+1} \right\rangle - \beta_t D_{(-H_\alpha)}(q_{t+1}, q_t) \right] = \beta_t \mathbb{E}_t \left[\left\langle \frac{\widehat{\ell}_t}{\beta_t}, q_t - q_{t+1} \right\rangle - D_{(-H_\alpha)}(q_{t+1}, q_t) \right] \\ & \leq \mathbb{E}_t \left[\frac{4}{\beta_t(1-\alpha)} \left(\sum_{i \in V \setminus \{\bar{I}_t\}} q_{ti}^{2-\alpha} \widehat{\ell}_{ti}^2 + (\min\{q_{t\bar{I}_t}, 1 - q_{t\bar{I}_t}\})^{2-\alpha} \widehat{\ell}_{i\bar{I}_t}^2 \right) \right] \\ & = \frac{4}{\beta_t(1-\alpha)} \left(\sum_{i \in V \setminus \{\bar{I}_t\}} q_{ti}^{2-\alpha} \mathbb{E}_t \left[\widehat{\ell}_{ti}^2 \right] + q_{t*}^{2-\alpha} \mathbb{E}_t \left[\widehat{\ell}_{i\bar{I}_t}^2 \right] \right). \end{aligned} \quad (80)$$

659 Then, by using the lower bound of P_t in (78), for any $i \in V$ the variance of the loss estimator $\widehat{\ell}_{ti}$ is
 660 bounded as

$$\mathbb{E}_t \left[\widehat{\ell}_{ti}^2 \right] = \sum_{j=1}^k p_{tj} \frac{\ell_{ti}^2}{P_{ti}^2} \mathbb{1}[i \in N^{\text{out}}(j)] = \frac{\ell_{ti}^2}{P_{ti}^2} \sum_{j \in V: i \in N^{\text{out}}(j)} p_{tj} = \frac{\ell_{ti}^2}{P_{ti}} \leq \frac{\delta^*}{\gamma_t}. \quad (81)$$

661 Hence, combining (80) with (81), we obtain

$$\mathbb{E}_t \left[\left\langle \widehat{y}_t, q_t - q_{t+1} \right\rangle - \beta_t D_{\psi_t}(q_{t+1}, q_t) \right] \leq \frac{4\delta^*}{\beta_t \gamma_t (1-\alpha)} \left(\sum_{i \in V \setminus \{\bar{I}_t\}} q_{ti}^{2-\alpha} + q_{t*}^{2-\alpha} \right) = \frac{z_t}{\beta_t \gamma_t} \leq \frac{z_t}{\beta_t \gamma_t}, \quad (82)$$

662 which implies that Assumption (ii) in (12) is satisfied.

663 Next, we will prove $h_{t+1} \lesssim h_t$ of Assumption (iii) in (12). To prove this, we will check the condition
 664 in Lemma 15. For any $i \in V$,

$$|\widehat{\ell}_{ti}| \leq \frac{1}{P_{ti}} \leq \frac{\delta^*}{\gamma_t} \leq \frac{\delta^* \beta_t}{u_t} = \frac{1 - \alpha}{8} \frac{\beta_t}{q_{t*}^{1-\alpha}} \leq \frac{1 - (\sqrt{2})^{\alpha-1}}{2} \frac{\beta_t}{q_{t*}^{1-\alpha}}, \quad (83)$$

665 where the second inequality follows from (78), the third inequality from $\gamma_t \geq u_t/\beta_t$, and the last
 666 inequality from the fact that $(1-x)/4 \leq 1 - (\sqrt{2})^{x-1}$ for $x \in [0, 1]$. Thus, the condition (58) is
 667 satisfied.

668 We next check the condition (59). Recalling $q_{t*} = \min\{q_{t\bar{I}_t}, 1 - q_{t\bar{I}_t}\}$, we observe that the parameters
 669 z_t and u_t satisfy

$$\sqrt{z_t} = \sqrt{\frac{4\delta^*}{1-\alpha} \left(\sum_{i \in V \setminus \{\bar{I}_t\}} q_{ti}^{2-\alpha} + q_{t*}^{2-\alpha} \right)} \leq \frac{2\sqrt{k\delta^*}}{\sqrt{1-\alpha}} q_{t*}^{1-\frac{1}{2}\alpha}, \quad u_t = \frac{8\delta^*}{1-\alpha} q_{t*}^{1-\alpha}, \quad (84)$$

670 where the last inequality follows from $q_{ti} \leq q_{t^*}$ for $i \neq \tilde{I}_t$. We can also lower-bound h_t as

$$h_t = H_\alpha(q_t) = \frac{1}{\alpha} \sum_{i=1}^k (q_{ti}^\alpha - q_{ti}) \geq \frac{1 - (1/2)^{1-\alpha}}{\alpha} q_{t^*}^\alpha \geq \frac{1-\alpha}{4\alpha} q_{t^*}^\alpha, \quad (85)$$

671 which can be proven by the same manner as in (67). Hence, using the upper bounds on z_t , u_t , and
672 h_t in (84) and (85), we have

$$\begin{aligned} \beta_{t+1} - \beta_t &= \frac{1}{\hat{h}_{t+1}} \left(2\sqrt{\frac{z_t}{\beta_t}} + \frac{u_t}{\beta_t} \right) = \frac{2}{h_t} \sqrt{\frac{z_t}{\beta_t}} + \frac{1}{h_t} \frac{u_t}{\beta_t} \\ &\leq \frac{16\alpha\sqrt{k}\delta^*}{\sqrt{\beta_1}(1-\alpha)^{3/2}} q_{t^*}^{1-\frac{3}{2}\alpha} + \frac{32\alpha\delta^*}{\sqrt{\beta_1}(1-\alpha)^2} q_{t^*}^{1-2\alpha} \\ &\leq \alpha\bar{\beta}q_{t^*}^{1-\frac{3}{2}\alpha} + \alpha\bar{\beta}q_{t^*}^{1-2\alpha} \\ &\leq 2(1-\bar{\alpha})\bar{\beta}q_{t^*}^{\bar{\alpha}-\alpha} \leq 2\frac{1-(\sqrt{2})^{\bar{\alpha}-1}}{\sqrt{2}}\bar{\beta}q_{t^*}^{\bar{\alpha}-\alpha}, \end{aligned} \quad (86)$$

673 where the first inequality follows from (84), (85), and $\beta_t \geq \beta_1 \geq 1$, the second inequality from
674 the definition of $\bar{\beta}$, the third inequality from $\min\{1 - \frac{3}{2}\alpha, 1 - 2\alpha\} \geq \bar{\alpha} - \alpha$ since $\bar{\alpha} = 1 - \alpha$,
675 and the last inequality from $1 - x \leq (1 - (\sqrt{2})^{x-1})/\sqrt{2}$ for $x \leq 1$. Thus the condition (59) is
676 satisfied. Therefore, from Lemma 15, we have $h_{t+1} = H_\alpha(q_{t+1}) \leq 2H_\alpha(q_t) = 2h_t$, which implies
677 that Assumption (iii) in (12) is satisfied.

678 Finally, we check the assumption (14) in Theorem 7. We first consider the first inequality in (14).
679 From the definition of z_t and the fact that $q_{ti} \leq q_{t\tilde{I}_t}$ for $i \neq \tilde{I}_t$, we get

$$\begin{aligned} z_t &= \frac{4\delta^*}{1-\alpha} \left\{ \sum_{i \in V \setminus \{\tilde{I}_t\}} q_{ti}^{2-\alpha} + (\min\{q_{t\tilde{I}_t}, 1 - q_{t\tilde{I}_t}\})^{2-\alpha} \right\} \\ &\leq \frac{4\delta^*}{1-\alpha} \left\{ \sum_{i \in V \setminus \{\tilde{I}_t\}} q_{ti}^{2-\alpha} + \left(\sum_{i \neq \tilde{I}_t} q_{ti} \right)^{2-\alpha} \right\} \\ &\leq \frac{8\delta^*}{1-\alpha} \left(\sum_{i \in V \setminus \{\tilde{I}_t\}} q_{ti} \right)^{2-\alpha} \leq \frac{8\delta^*}{1-\alpha} \left(\sum_{i \neq a^*} q_{ti} \right)^{2-\alpha} = \frac{8\delta^*}{1-\alpha} (1 - q_{ta^*})^{2-\alpha}, \end{aligned} \quad (87)$$

680 where the second inequality holds from the inequality $x^a + y^a \leq (x+y)^a$ for $x, y \geq 0$ and $a \in [0, 1]$,
681 and the third inequality from $q_{ti} \leq q_{t\tilde{I}_t}$. Hence, combining this with (87) with the upper bound on
682 h_t in (70), we obtain

$$z_t h_t \leq \frac{8\delta^*}{1-\alpha} (1 - q_{ta^*})^{2-\alpha} \cdot \frac{1}{\alpha} (k-1)^{1-\alpha} (1 - q_{ta^*})^\alpha = \underbrace{\frac{8\delta^*(k-1)^{1-\alpha}}{\alpha(1-\alpha)}}_{=\rho_1} (1 - q_{ta^*})^2. \quad (88)$$

683 We next consider the second inequality in (14). We can bound u_t from above as

$$\begin{aligned} u_t &= \frac{8\delta^*}{1-\alpha} (\min\{q_{t\tilde{I}_t}, 1 - q_{t\tilde{I}_t}\})^{1-\alpha} \leq \frac{8\delta^*}{1-\alpha} \left(\sum_{i \neq \tilde{I}_t} q_{ti} \right)^{1-\alpha} \\ &\leq \frac{8\delta^*}{1-\alpha} \left(\sum_{i \neq a^*} q_{ti} \right)^{1-\alpha} = \frac{8\delta^*}{1-\alpha} (1 - q_{ta^*})^{1-\alpha}, \end{aligned} \quad (89)$$

684 where the second inequality follows from $q_{t\tilde{I}_t} \geq q_{ti}$ for all $i \neq \tilde{I}_t$. Hence, combining the last
685 inequality with (70),

$$u_t h_t \leq \underbrace{\frac{4\delta^*(k-1)^{1-\alpha}}{\alpha(1-\alpha)}}_{=\rho_2} (1 - q_{ta^*}). \quad (90)$$

686 Hence, the assumption (14) is satisfied with above ρ_1 and ρ_2 , and thus we have completed the proof.
 687 □

688 **F.3 Technical challenges to derive best-of-both-worlds bounds depending on (fractional)** 689 **weak domination number**

690 Here, we discuss the technical challenges of making our upper bound in Theorem 10 depend on the
 691 weak domination number δ instead of the fractional domination number δ^* or the weak fractional
 692 domination number $\tilde{\delta}^* \leq \delta$.

693 First, we need to use Tsallis entropy to derive a regret upper bound with a stochastic bound of $\log T$.
 694 While we can prove a BOBW bound if we use the Shannon entropy regularizer [25], the bound in the
 695 stochastic regime is $O((\log T)^2)$, which is not desirable. Hence, a possible
 696 approach is to use the log-barrier regularizer or the Tsallis entropy. The log-barrier regularizer has
 697 a penalty term of $\Omega(k)$ due to the strength of its regularization, and the regret upper bound in the
 698 final adversarial regime is $\Omega(k^{1/3})$, which can be much larger than $\delta^{1/3}$. Therefore, the most hopeful
 699 solution would be to use Tsallis entropy with an appropriate exponent $\alpha \simeq 1$, where we note that the
 700 Tsallis entropy with $\alpha \rightarrow 1$ corresponds to the Shannon entropy.

701 Recalling the definition of the weak domination number in Section 6, we can see that the weak dom-
 702 ination set dominates only vertices without self-loop $U = \{i \in V : i \notin N^{\text{out}}(i)\}$. Thus, to achieve
 703 a BOBW bound that depends on the weak domination number, vertices with self-loop and those
 704 without self-loop should be treated separately by decomposing the stability term as follows:

$$\begin{aligned} & \langle \hat{\ell}_t, q_t - q_{t+1} \rangle - \beta_t D_{(-H_\alpha)}(q_{t+1}, q_t) \\ &= \sum_{i \in U} \left(\hat{\ell}_{ti}(q_{ti} - q_{t+1,i}) - \beta_t d(q_{t+1,i}, q_{t,i}) \right) + \sum_{i \in V \setminus U} \left(\hat{\ell}_{ti}(q_{ti} - q_{t+1,i}) - \beta_t d(q_{t+1,i}, q_{t,i}) \right), \end{aligned}$$

705 where $d(p, q)$ is the Bregman divergence induced by the real-valued convex function $x \mapsto -\frac{1}{\alpha}(x^\alpha -$
 706 $x)$. However, if we use this approach, we cannot use Lemma 14, which is useful to prove an upper
 707 bound with $(1 - q_{ta^*})$ (see (14)). This is because this lemma exploits the fact that q and r are
 708 probability vectors. This prevents us from deriving an upper bound with an $O(\log T)$ stochastic
 709 bound depending on the weak domination number.