

# A Unified Analysis of Nonstochastic Delayed Feedback for Combinatorial Semi-Bandits, Linear Bandits, and MDPs

**Lukas Zierahn**

LUKAS.ZIERAHN@GMAIL.COM

*Centrum Wiskunde & Informatica and Booking.com, The Netherlands*

**Dirk van der Hoeven**

DIRK@DIRKVANDERHOEVEN.COM

*Mathematical Institute, Leiden University, The Netherlands*

**Tal Lancewicki**

LANCEWICKI@MAIL.TAU.AC.IL

*Blavatnik School of Computer Science, Tel Aviv University, Israel*

**Aviv Rosenberg**

AVIVROS@GOOGLE.COM

*Google Research*

**Nicolò Cesa-Bianchi**

NICOLO.CESA-BIANCHI@UNIMI.IT

*Università degli Studi di Milano and Politecnico di Milano, Italy*

**Editor:** Ambuj Tewari

## Abstract

We derive a new analysis of Follow The Regularized Leader (FTRL) for online learning with delayed bandit feedback. By separating the cost of delayed feedback from that of bandit feedback, our analysis allows us to obtain new results in four important settings. We derive the first optimal (up to logarithmic factors) regret bounds for combinatorial semi-bandits with delay and adversarial Markov Decision Processes with delay (both known and unknown transition functions). Furthermore, we use our analysis to develop an efficient algorithm for linear bandits with delay achieving near-optimal regret bounds. In order to derive these results we show that FTRL remains stable across multiple rounds under mild assumptions on the regularizer.

**Keywords:** Online learning, bandit feedback, delayed feedback, Markov Decision Processes, combinatorial semi-bandits

## 1. Introduction

Delayed feedback is a phenomenon that cannot be avoided in many applications of online learning. For example, in digital advertisement a conversion event may happen with some delay after an ad is shown to a user. In healthcare, the effect of a drug on a patient may take some time before it becomes observable (Eick, 1988). A consequence of delayed feedback is that sequential decision makers have to act before knowing the effect of their previous actions, where the effect of multiple past actions may be potentially observed all at once. These challenges pertain not only to the algorithms, but also to the way they are analyzed, which is the reason why standard (non-delayed) proof techniques fail in the presence of delayed feedback.

Due to its fundamental nature in online learning, delayed feedback has been extensively studied in several different scenarios, including full-information feedback (Weinberger and Ordentlich, 2002; Joulani et al., 2013; Quanrud and Khashabi, 2015; Joulani et al., 2016; Flaspohler et al., 2021) and bandit feedback (Cesa-Bianchi et al., 2016; Thune et al., 2019; Bistritz et al., 2019; Zimmert

and Seldin, 2020; Ito et al., 2020a; Gyorgy and Joulani, 2021; Van der Hoeven and Cesa-Bianchi, 2022; Masoudian et al., 2022). In this work, we focus on the bandit feedback setting; that is, when the only way for the learner to know the effect of an action is to execute it. We develop a general framework for the analysis of delayed bandit feedback which we then apply to three important settings: combinatorial semi-bandits (which includes multi-armed bandits as a special case), linear bandits, and adversarial Markov Decision Processes (MDPs).

Our analysis, which is based on Follow The Regularized Leader (FTRL)—see, for example, (Orabona, 2019, Chapter 7), unifies previous analyses and sheds light on the impact of delayed bandit feedback in online learning. Our main insight is that one can separate the cost of delayed feedback and bandit feedback through a novel decomposition of the FTRL regret, which allows us to separately control these different regret components. This insight leads to new results in all of the settings we consider. We prove the first regret bounds for combinatorial semi-bandits with delays, which also turn out to be optimal for sufficiently large  $T$  (throughout the paper, by optimal we always mean optimal for sufficiently large  $T$ ). We provide the first optimal regret bounds for adversarial MDPs with delays and known transitions. For adversarial MDPs with delays and unknown transitions we provide the state-of-the-art results. Finally, we derive a computationally efficient algorithm for linear bandits, whose regret has an optimal dependence on delays.

We now formally introduce the setting of online learning with delayed bandit feedback studied in this paper. Online learning with delayed bandit feedback proceeds in rounds. In each round  $t \in [T]$  the learner chooses (possibly in a randomized manner) an action  $\mathbf{a}_t \in \mathcal{A} \subseteq \mathbb{R}^K$ , where  $\mathcal{A}$  is an action set of dimension  $K$ . The learner subsequently suffers loss  $\mathbf{a}_t^\top \ell_t$ , where  $\ell_t \in \mathbb{R}^K$  is bounded in some suitably chosen norm, and observes  $\{\mathcal{L}(\ell_\tau, \mathbf{a}_\tau) : \tau + d_\tau = t\}$ , where  $d_1, \dots, d_T$  is an unknown sequence of delays and  $\mathcal{L}$  is an application-specific (possibly randomized) feedback function, encoding which information about  $\ell_\tau$  the learner sees based on the action  $\mathbf{a}_\tau$ . For example, in the combinatorial semi-bandit setting the learner observes all loss components corresponding to the non-zero elements of the action, whereas in the linear bandit setting the learner only observes the scalar  $\mathbf{a}_\tau^\top \ell_\tau$ . We assume that delays  $d_1, \dots, d_T$  and losses  $\ell_1, \dots, \ell_T$  are both generated by an oblivious adversary.

## 1.1 Contributions

This work and the work it extends (Van der Hoeven et al., 2023) have the following main contributions:

**New analysis.** In section 3 we provide a novel analysis of FTRL under delayed bandit feedback. The main novelty is showing that we can decompose the regret into three main parts. The first part of the regret is standard, namely the pseudo-distance between the starting point of the algorithm and the optimal point in hindsight. The second part is the cost of delayed feedback. In our analysis, we show that the cost of delayed feedback is essentially the same as in the delayed full-information setting. The third part of the regret is the cost of bandit feedback, which is the same term that occurs in the standard analysis of FTRL for bandit feedback. A technical novelty is that we show that FTRL is stable across multiple rounds under some mild assumptions on the Hessian of the regularizer. In related work, Huang et al. (2023) provide an analysis of online mirror descent with delayed bandit feedback in several settings. However, their analysis does not lead to optimal bounds because it does not separate the cost of delayed and bandit feedback.

**Combinatorial semi-bandits with delayed feedback.** As far as we know we are the first to consider nonstochastic combinatorial semi-bandits under delayed feedback. In the combinatorial semi-bandit setting, we apply the newly gained insight from our analysis of FTRL to derive an optimal algorithm. We show that if  $\max_{a \in \mathcal{A}} \|a\|_1 \leq B$ , then the regret after  $T$  rounds is of order  $\sqrt{B(KT + BD) \log(K)}$ , where  $D = \sum_{t=1}^T d_t$  is the total delay after  $T$  rounds. In the worst case, the delay is constant (i.e.,  $d_t = d$  for all  $t$ ) and we provide a matching lower bound (up to logarithmic factors) showing that any learner must incur  $\Omega(\sqrt{BT(K + Bd)})$  regret.

**Linear bandits.** In the linear bandit setting, Ito et al. (2020a) provide an analysis of continuous exponential weights (Cover, 1991; Vovk, 1990; Littlestone and Warmuth, 1994) with delayed bandit feedback and constant delay  $d$  that obtains the optimal  $\tilde{O}(K\sqrt{T} + \sqrt{dT})$  regret bound. One drawback is that the per-round runtime of continuous exponential weights is prohibitively large, although it is polynomial in  $K$  and  $T$ . Building on Scribble (Abernethy et al., 2008), we derive an algorithm that achieves a slightly suboptimal  $\tilde{O}(K^{3/2}\sqrt{T} + \sqrt{D})$  regret, but with a much better per-round running time of order  $K^3$ , provided a self-concordant barrier for the decision set can be efficiently computed. Huang et al. (2023) show an algorithm with a similar running time, but with a worse regret bound of  $\tilde{O}(K^{3/2}\sqrt{T} + K^2\sqrt{D})$ .

**Adversarial Markov Decision Processes.** Delayed feedback in adversarial (finite-horizon and episodic) MDPs was first studied by Lancewicki et al. (2022a). Under full-information feedback, where the agent observes the entire cost function at the end of the episode, they achieve the optimal regret bound:  $\tilde{O}(H\sqrt{T + D})$ , where  $T$  is the number of episodes and  $H$  is the horizon. However, with bandit feedback (where the only observed costs are those along the agent’s trajectory), their regret bound is of order  $T^{2/3} + D^{2/3}$ . The current state-of-the-art guarantees under delayed bandit feedback are by Jin et al. (2022) and Lancewicki et al. (2023) who achieve a regret bound of  $\tilde{O}(H\sqrt{SAT} + H(HSA)^{1/4}\sqrt{D})$  and  $\tilde{O}(H^2\sqrt{SAT} + H^3\sqrt{D})$  in the known transition setting, and a regret bound of  $\tilde{O}(H^2S\sqrt{AT} + H(HSA)^{1/4}\sqrt{D})$  and  $\tilde{O}(H^3S\sqrt{AT} + H^3\sqrt{D})$  in the unknown transition setting, respectively. Here,  $S$  is the number of states in the MDP and  $A$  the number of actions. However, there is still a gap compared to the lower bound of Lancewicki et al. (2022a). Remarkably, the application of our FTRL analysis to adversarial MDPs allows us to close this gap and achieve the first optimal regret bound of  $\tilde{O}(H\sqrt{SAT} + H\sqrt{D})$  for the case of known transitions. Moreover, our bound of  $\tilde{O}(H^2S\sqrt{AT} + H\sqrt{D})$  for unknown transitions, achieves the first optimal regret in the delay term and matches the best known regret bound (even for the standard non-delayed setting) in the other term.

## 1.2 Additional related work

**Combinatorial semi-bandits with delayed feedback.** Stochastic combinatorial semi-bandits have first been introduced by Gai et al. (2012) but featured an undesirable dependency on the reciprocal of the square of the smallest gap between arms, which was improved by Chen et al. (2013) by removing the square. The first matching upper and lower bounds are due to Kveton et al. (2015) by using an upper confidence bound (UCB) based approach, though bounds using Thompson sampling are also known Wen et al. (2015). A special case of the stochastic combinatorial semi-bandit setting with delayed feedback, namely stochastic multi-armed bandits with delayed feedback, has been studied in many different variations (Dudik et al., 2011; Agarwal and Duchi, 2012; Pike-Burke et al., 2018; Zhou et al., 2019; Gael et al., 2020; Lancewicki et al., 2021; Cohen et al., 2021).

In the nonstochastic combinatorial semi-bandit setting there have been several results. Adversarial online path-finding problems, a special case of semi-bandits, has been studied by György et al. (2007) achieving an sub-optimal upper bound, an optimal upper bound for  $m$ -sets is due to Kale et al. (2010) and Uchiya et al. (2010). The optimal bound for semi-bandits in general, which we recover in the non-delayed setting, is due to Audibert et al. (2014). Even though we are the first to study combinatorial semi-bandits with delayed feedback, a special case, namely multi-armed bandits with delayed feedback, is well understood. Neu et al. (2010, 2014) were among the first ones to study the impact of delayed feedback in the nonstochastic setting. Subsequently, Cesa-Bianchi et al. (2019) proved a  $\Omega(\sqrt{KT} + \sqrt{dT \log(K)})$  lower bound when  $d_t = d$  for all  $t$ . The matching upper bound was provided by Zimmert and Seldin (2020), but nearly matching upper bounds also exist (Thune et al., 2019; Bistritz et al., 2019; Gyorgy and Joulani, 2021; Van der Hoeven and Cesa-Bianchi, 2022). Conversely, (special cases of) combinatorial semi-bandits without delay have also received considerable attention (György et al., 2007; Kale et al., 2010; Uchiya et al., 2010; Cesa-Bianchi and Lugosi, 2012; Audibert et al., 2014; Combes et al., 2015; Lattimore et al., 2018; Zimmert et al., 2019).

**Adversarial Markov Decision Processes.** There is a rich literature on regret minimization in MDPs with non-delayed feedback (Even-Dar et al., 2009; Jaksch et al., 2010; Zimin and Neu, 2013; Dick et al., 2014; Rosenberg and Mansour, 2019b,a, 2021; Jin et al., 2020; Shani et al., 2020; Luo et al., 2021). Under delayed feedback, apart from the literature mentioned earlier, Dai et al. (2022) recently presented a Follow-The-Perturbed-Leader approach that can also handle delayed feedback in adversarial MDPs. However, their regret bound is slightly weaker than that of Jin et al. (2022) mentioned earlier. Finally, a different line of work (Katsikopoulos and Engelbrecht, 2003; Walsh et al., 2009) considers delays in observing the current state, which is inherently different than our setting—for a thorough discussion on the differences between the models we refer the reader to Lancewicki et al. (2022a). A stochastic version of MDPs with delayed feedback has been studied by (Howson et al., 2023a).

**Linear bandits.** Early work in the non-delayed linear bandit setting suffered from suboptimal results in terms of  $T$  (McMahan and Blum, 2004; Awerbuch and Kleinberg, 2004; Dani and Hayes, 2006). Dani et al. (2007); Abernethy et al. (2008) were the first to prove a regret bound with optimal scaling in  $T$ . Subsequent works by (Bubeck and Eldan, 2015; Hazan and Karnin, 2016; Ito et al., 2020b; Zimmert and Lattimore, 2022) obtained the optimal  $O(K\sqrt{T})$  regret bound. Stochastic linear bandits with delayed feedback has been studied by (Vernade et al., 2020; Howson et al., 2023b).

## 2. Preliminaries

We denote by  $\hat{\ell}_t \in \mathbb{R}^K$  the estimate of the loss  $\ell_t$  in round  $t$ . We will define a loss estimator for each application separately. We focus on Follow The Regularized Leader (FTRL) and define the FTRL prediction given a sum of losses  $\mathbf{L}$  (or estimated losses  $\hat{\mathbf{L}}$ ) as follows,

$$\mathbf{w}_t(\mathbf{L}) = \arg \min_{\mathbf{v} \in \mathcal{W}} \mathbf{L}^\top \mathbf{v} + R_t(\mathbf{v}),$$

where  $\mathcal{W} \subseteq \mathbb{R}^K$  is a compact closed convex set,  $R_t$  is a twice-differentiable strongly convex function. Note that the domain  $\mathcal{W}$  and the action set  $\mathcal{A}$  do not necessarily coincide, as is the case of

combinatorial semi-bands for example, where  $\mathcal{W}$  is the convex hull of  $\mathcal{A}$ , i.e.  $\mathcal{W} = \text{Conv}(\mathcal{A})$ . Similarly,  $\mathbf{a}_t$  and  $\mathbf{w}_t(\hat{\mathbf{L}}_t)$  do not necessarily coincide. We will specify the relationship between  $\mathbf{a}_t$  and  $\mathbf{w}_t(\hat{\mathbf{L}}_t)$  in each application.

If we define  $\tilde{o}_t = \{\tau : \tau + d_\tau < t\}$  as the set of all losses available at the beginning of round  $t$ , then FTRL predicts  $\mathbf{w}_t(\hat{\mathbf{L}}_t)$ , where  $\hat{\mathbf{L}}_t = \sum_{\tau \in \tilde{o}_t} \hat{\ell}_\tau$ . We then use the notation  $[N] = \{1, \dots, N\}$  and define  $\tilde{m}_t = [t-1] \setminus \tilde{o}_t$  to be the set of indices of losses that have not been observed at the start of round  $t$  due to delay. As a simplifying assumption, we assume that  $d_{\max} = \max_{t \in [T]} d_t \geq 1$  which is known to the learner. This assumption is without loss of generality, as we may employ the standard doubling trick to overcome the need to know this parameter (Bistritz et al., 2019; Lancewicki et al., 2022a), see also Appendix E.

**Additional notations.** We denote by  $\mathbf{w}_t(\mathbf{L}, i)$  the  $i$ -th element of the vector  $\mathbf{w}_t(\mathbf{L})$ . We define a filtration of all random events observed by the learner up to round  $t$  as  $\mathcal{F}_t = \{(\tau, \mathbf{a}_\tau, \mathcal{L}(\ell_\tau, \mathbf{a}_\tau)) : \tau + d_\tau < t\}$  and we use  $\mathbb{E}_t[\cdot] = \mathbb{E}[\cdot | \mathcal{F}_t]$ . For a twice-differentiable function  $\phi$  such that  $\nabla^2 \phi(\mathbf{v}) \succ 0\mathbf{I}$  for all  $\mathbf{v} \in \mathcal{W}$  we denote by  $\|\mathbf{L}\|_{\phi, \mathbf{v}} = \sqrt{\mathbf{L}^\top (\nabla^2 \phi(\mathbf{v}))^{-1} \mathbf{L}}$  and by  $\|\mathbf{L}\|_{\phi, \mathbf{v}}^* = \sqrt{\mathbf{L}^\top \nabla^2 \phi(\mathbf{v}) \mathbf{L}}$ . The Dikin ellipsoid with radius  $r$  around  $\mathbf{v}$  induced by  $\phi$  is defined as  $\mathcal{D}_\phi(\mathbf{v}, r) = \{\mathbf{x} \in \mathcal{W} : \|\mathbf{x} - \mathbf{v}\|_{\phi, \mathbf{v}}^* \leq r\}$ . The notation  $\tilde{O}(\cdot)$  hides poly-logarithmic factors, whereas  $\lesssim$  denotes inequalities that hide constant factors.

**Changing domains.** Some settings require changing domains. In the MDP setting with unknown transitions, the domain is related to the estimate of the transition function and as we update and become more confident in our estimates, we may wish to shrink the domain. We overload the notation slightly and define

$$\mathbf{w}_t(\mathbf{L}) = \arg \min_{\mathbf{v} \in \mathcal{W}_t} \mathbf{L}^\top \mathbf{v} + R_t(\mathbf{v}), \quad (1)$$

where we require all  $\mathcal{W}_t$  to be compact closed convex sets. Our analysis requires that if we observe the feedback from round  $\tau$  in timestep  $t$ , then the corresponding iterate of FTRL,  $\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)$ , must be in the same Dikin ellipsoid as the current iterate  $\mathbf{w}_t(\hat{\mathbf{L}}_t)$ . To ensure that condition holds the domains of timestep  $\tau$  and timestep  $t$  have to agree. If that is not the case, we have to skip round  $\tau$ , which means trivially bounding the regret of round  $\tau$  with an appropriate constant value (like the length of the episode in the MDP setting) and not building a loss estimator using the information of round  $\tau$ . We define  $\Lambda \subseteq [T]$  to be the set of rounds that we skip and  $\bar{\Lambda} = [T] \setminus \Lambda$  be the rounds that we do not skip. Since we chose not to use the loss estimators of skipped rounds, we intersect the set of observed losses and the set of missing losses at time  $t$  with the rounds that we did not skip:  $o_t = \tilde{o}_t \cap \bar{\Lambda}$ ,  $m_t = \tilde{m}_t \cap \bar{\Lambda}$ . When we observe the loss of round  $\tau$ , we know if we have changed the domain since  $\tau$  and thus  $o_t$  is well defined and non-random given the history  $\mathcal{F}_t$ . The same is not true for  $m_t$ , which can depend on future rounds. This is not a problem for the algorithms considered here, as  $m_t$  is a quantity only used in the analysis and for tuning the learning rates, where  $|\tilde{m}_t|$  can be used as an upper bound for  $|m_t|$ . The constraints that must be fulfilled to use changing domains are formalised in Assumption 1.

**Assumption 1** For all  $t \in [T]$  we assume that  $o_t$  is non-random given the history  $\mathcal{F}_t$  and that  $\mathcal{W}_t = \mathcal{W}_\tau$  for all  $\tau \in o_t \setminus o_{t-1}$ . We also assume that  $\mathcal{W}_t \subseteq \mathcal{W}_{t-1}$  is a compact convex set such that  $\mathcal{W}_T$  is non-empty.

If the domain is constant and no rounds are skipped then Assumption 1 reduces to the standard assumption that  $\mathcal{W}$  is compact, convex, and non-empty as in that case  $o_t = \tilde{o}_t$  and  $m_t = \tilde{m}_t$ .

In the remainder of the paper we use the following notation for cumulative loss estimates:

$$\hat{\mathbf{L}}_t = \sum_{\tau \in o_t} \hat{\ell}_\tau, \quad \hat{\mathbf{L}}_t^m = \hat{\mathbf{L}}_t + \sum_{\tau \in m_t} \hat{\ell}_\tau, \quad \hat{\mathbf{L}}_t^* = \sum_{\tau \in [t]} \hat{\ell}_\tau.$$

Note that  $\hat{\mathbf{L}}_t^* = \hat{\mathbf{L}}_t^m + \hat{\ell}_t$  and that  $\mathbf{w}_t(\hat{\mathbf{L}}_t^m)$  is equivalent to FTRL in the non-delayed setting. We also make the following regularity assumptions on the regularizer  $R_t$ .

**Assumption 2** *Let  $R_t$  be the regularizer associated with equation (1) and let  $\kappa > 0$ . Suppose that for all  $t \in [T]$*

- (a)  $4\nabla^2 R_t(\mathbf{v}) \succeq \nabla^2 R_t(\mathbf{v}') \succeq \frac{1}{4}\nabla^2 R_t(\mathbf{v})$  for all  $\mathbf{v} \in \mathcal{W}_t$  and  $\mathbf{v}' \in \mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2\sqrt{\kappa}})$ .
- (b)  $\kappa \left( \nabla R_t(\mathbf{w}_t(\hat{\mathbf{L}}_t)) - \nabla R_{t+\delta}(\mathbf{w}_t(\hat{\mathbf{L}}_t)) \right)^\top \mathbf{x} \leq \frac{\sqrt{\kappa}}{32} \sqrt{\mathbf{x}^\top \nabla^2 R_{t+\delta}(\mathbf{w}_t(\hat{\mathbf{L}}_t)) \mathbf{x}}$  for all  $\mathbf{x} \in \mathbb{R}^d$  and all  $\delta \in [d_{\max}]$ .
- (c)  $R_t(\mathbf{v}) \leq R_{t'}(\mathbf{v})$  and  $\nabla^2 R_t(\mathbf{v}) \preceq \nabla^2 R_{t'}(\mathbf{v})$  for all  $\mathbf{v} \in \mathcal{W}_t$  all  $t \leq t'$ .

Assumption 2(a) allows us to relate the Hessian of the regularizer at different iterates of FTRL, which is crucial in our analysis. Since essentially all regularizers we use in this paper are approximately self-concordant, assumption 2(a) is almost automatically satisfied (Nemirovski, 2004), see also equation 18. Assumption (b) tells us that the regularizer should not change too much between rounds and is used show that the different iterates of FTRL are close to each other. As we will see, assumption (b) can be verified for most standard regularizers given that the learning rate does not change too much between rounds. Assumption 2(c) is a technical assumption and is satisfied by almost all standard regularizers, including those that we use in this paper.

### 3. Analysis

In this section we establish general results that are then applied to combinatorial bandits, MDPs, and linear bandits in the next sections. First, we give a broad overview of the proof ideas and then prove the statements rigorously.

#### 3.1 Overview

We build on the analysis of Flaspohler et al. (2021) for delayed feedback in the full-information setting, where they observe that delayed feedback can be interpreted as poor hints in the sense of optimistic online learning (Rakhlin and Sridharan, 2013). Taking this idea one step further, we analyze what would happen had the algorithm received slightly different hints, and subsequently bound the change between different instances of FTRL.

Suppose for a moment that the domain  $\mathcal{W}_t = \mathcal{W}$  is constant, we are not skipping any rounds  $\Lambda = \emptyset$ , and that our loss estimates satisfy  $\mathbb{E}[\hat{\ell}_t | \mathcal{F}_t] = \ell_t + \mathbf{b}_t$ , where  $\mathbf{b}_t$  is the estimator's bias. Let

$\mathbf{u} \in \mathbb{R}^K$  be any comparator. Our analysis relies on the following decomposition of the regret

$$\begin{aligned} \sum_{t=1}^T \mathbb{E} [(\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{u})^\top \ell_t] &= \underbrace{\sum_{t=1}^T -\mathbb{E} [(\mathbf{w}_t(\hat{\mathbf{L}}_t^m) - \mathbf{u})^\top \mathbf{b}_t]}_{\text{bias}} + \underbrace{\sum_{t=1}^T \mathbb{E} [(\mathbf{w}_t(\hat{\mathbf{L}}_t^*) - \mathbf{u})^\top \hat{\ell}_t]}_{\text{cheating regret}} \quad (2) \\ &+ \sum_{t=1}^T \left( \underbrace{\mathbb{E} [(\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{w}_t(\hat{\mathbf{L}}_t^m))^\top \ell_t]}_{H_1} + \underbrace{\mathbb{E} [(\mathbf{w}_t(\hat{\mathbf{L}}_t^m) - \mathbf{w}_t(\hat{\mathbf{L}}_t^*))^\top \hat{\ell}_t]}_{H_2} \right). \end{aligned}$$

If  $\hat{\ell}_t$  is an unbiased estimator of the loss then  $\mathbf{b}_t = \mathbf{0}$ , which implies that the bias term of the decomposition is also 0. The cheating regret can be found in different forms in online learning—see, for example, the proof of (Shalev-Shwartz, 2012, Lemma 2.3) or (Gyorgy and Joulani, 2021, Equation 4)—and can be bounded using the standard be-the-leader lemma (Lemma 18 in Appendix F), see also (Joulani et al., 2020, Theorem 3). Now we focus on the second line of Equation (2). Typically,  $H_1$  and  $H_2$  are analysed simultaneously and referred to as "drift", for example, see (Gyorgy and Joulani, 2021). We split the drift into  $H_1$  and  $H_2$  because we want to analyze the cost of delay and the cost of bandit feedback separately.

$H_1$  can be interpreted as capturing the influence of the missing observations.  $H_2$  captures the influence of knowing the loss estimated one step in advance against running a non-delayed version of FTRL. To bound  $H_1$  and  $H_2$  we will use the same tools. First we need to relate the differences between the predictions of the different FTRL instances to the losses used in computing the different FTRL iterates. Lemma 5 states that if  $\mathbf{w}_t(\mathbf{L}'), \mathbf{w}_t(\mathbf{L}) \in \mathcal{D}_R(\mathbf{v}, \frac{1}{2\sqrt{\kappa}})$  for some  $\kappa > 0$ , some  $\mathbf{v} \in \mathcal{W}$ , and some  $\mathbf{L}, \mathbf{L}' \in \mathbb{R}^K$ , and the regularizer is sufficiently nice, then  $\|\mathbf{w}_t(\mathbf{L}) - \mathbf{w}_t(\mathbf{L}')\|_{R_t, \mathbf{v}}^* \leq 8\|\mathbf{L}' - \mathbf{L}\|_{R_t, \mathbf{v}}$ . In order to apply this result for different  $\mathbf{w}_t(\cdot)$ , we require them to lie in the same Dikin ellipsoid, and Lemma 6 establishes machinery to allow us to determine when that is the case. Specifically, if  $\mathbf{L}' = \mathbf{L} + \sum_{\tau \in z} \hat{\ell}_\tau$  for some finite set  $z$  and  $\sum_{\tau \in z} \kappa \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_t(\mathbf{L})} \leq \frac{1}{32}$ , then  $\mathbf{w}_{t'}(\mathbf{L}') \in \mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$ . We apply the last result in Lemma 7 to establish that  $\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau) \in \mathcal{D}_{R_t}(\mathbf{w}_t(\hat{\mathbf{L}}_t), \frac{1}{2\sqrt{\kappa}})$  for all  $\tau \in m_t$ , which in turn allows us to conclude that  $\mathbf{w}_t(\hat{\mathbf{L}}_t^m), \mathbf{w}_t(\hat{\mathbf{L}}_t^*) \in \mathcal{D}_{R_t}(\mathbf{w}_t(\hat{\mathbf{L}}_t), \frac{1}{2\sqrt{\kappa}})$ . Thus, we can repeatedly apply  $\|\mathbf{w}_t(\mathbf{L}) - \mathbf{w}_t(\mathbf{L}')\|_{R_t, \mathbf{v}}^* \leq 8\|\mathbf{L}' - \mathbf{L}\|_{R_t, \mathbf{v}}$  due to Lemma 5, which leads to Lemma 3.

**Lemma 3** Suppose that  $\mathbb{E}[\hat{\ell}_t | \mathcal{F}_t] = \ell_t + \mathbf{b}_t$  and suppose that Assumption 1 and Assumption 2 hold. Let  $t \in [T]$  and  $\tau \in m_t \cup \{t\}$ . Suppose that  $\|\ell_\tau\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq \alpha_t$  and  $\mathbb{E}[\|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}^2] \leq \beta_t^2$ . Suppose that for all  $t, t' \in [T]$   $\sqrt{\kappa} \|\hat{\ell}_t\|_{R_{t'}, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq \frac{1}{128d_{\max}}$ . Then for all  $\mathbf{u} \in \mathcal{W}_T$ ,

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{u})^\top \ell_t \right] &\leq \mathbb{E} \left[ \sum_{t \in \Lambda} (\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{u})^\top \ell_t \right] + R_T(\mathbf{u}) - \min_{v \in \mathcal{W}_1} R_1(v) + \sum_{t \in \bar{\Lambda}} 8\beta_t^2 \\ &- \sum_{t \in \bar{\Lambda}} \mathbb{E} [(\mathbf{w}_t(\hat{\mathbf{L}}_t^m) - \mathbf{u})^\top \mathbf{b}_t] + \sum_{t \in \bar{\Lambda}} \left( 8\alpha_t^2 |m_t| + 8\alpha_t \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\ell_\tau - \hat{\ell}_\tau) \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \right] \right). \end{aligned}$$

The work of Van der Hoeven and Cesa-Bianchi (2022) provides a similar result for the multi-armed bandit setting. However, that result does not apply to the more general setting we consider

here as their analysis relies on the fact that in the multi-armed bandit setting the constraint in the Lagrangian of the FTRL objective can be expressed in a simple manner, which is not possible in our setting.

To interpret Lemma 3, consider the multi-armed bandit setting with the standard importance-weighted estimator, no skipping  $\Lambda = \emptyset$ , and regularizer  $R(\mathbf{v}) = \sum_{i=1}^K \frac{1}{\eta} \mathbf{v}(i) \log(\mathbf{v}(i)) - \frac{1}{\gamma} \log(\mathbf{v}(i))$ . The purpose of the log barrier term in the regularizer is to ensure stability of the iterates, as required by the assumptions of the lemma. In this case, if  $\|\ell_t\|_\infty \leq 1$ , then  $\alpha_t$  is  $O(\sqrt{\eta})$ . The quantity  $\beta_t^2$  is a bound on the expectation of the squared local norm of the loss estimate, which is  $O(\eta K)$ . Thus, by choosing  $\mathbf{u} = (1 - \frac{1}{T})\tilde{\mathbf{u}} + \frac{1}{T} \arg \min_{\mathbf{v} \in \mathcal{W}} R(\mathbf{v})$ , we have that the expected regret against  $\tilde{\mathbf{u}}$  is of order

$$\frac{1}{\eta} \log(K) + d_{\max} K \ln(T) + \eta(KT + D) + \sqrt{\eta} \sum_{t=1}^T \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\ell_\tau - \hat{\ell}_\tau) \right\|_{R, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \right], \quad (3)$$

where we used  $\sum_{t=1}^T |m_t| = D$ . The  $d_{\max} K \ln(T)$  term in the above equation comes from the log-barrier part of  $R$ , which—when properly tuned—ensures that the FTRL iterates are close to each other. So far, it seems that we did not manage to separate the cost of delay and bandit feedback because of the final summation in (3). However, due to the delay, if  $\tau, \tau' \in m_t$ , then  $\hat{\ell}_\tau$  and  $\hat{\ell}_{\tau'}$  are independent random variables and  $\ell_\tau$  and  $\ell_{\tau'}$  are their means. Recall that the variance of the sum of independent random variables equals to the sum of their variances. Thus, by applying Jensen's inequality to the square root and using that  $\nabla^2 R(\mathbf{v}) \succeq \text{diag}(\eta \mathbf{v})^{-1}$ , we can see that

$$\begin{aligned} \sqrt{\eta} \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\ell_\tau - \hat{\ell}_\tau) \right\|_{R, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \right] &\leq \sqrt{\eta} \sqrt{\mathbb{E} \left[ \sum_{\tau \in m_t} \left\| (\ell_\tau - \hat{\ell}_\tau) \right\|_{R, \mathbf{w}_t(\hat{\mathbf{L}}_t)}^2 \right]} \\ &\leq 2\sqrt{\eta} \sqrt{\mathbb{E} \left[ \sum_{\tau \in m_t} \left\| (\ell_\tau - \hat{\ell}_\tau) \right\|_{R, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}^2 \right]} \\ &= 2\sqrt{\eta} \sqrt{\mathbb{E} \left[ \sum_{\tau \in m_t} \sum_{i=1}^K \eta \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, i) (\ell_\tau(i) - \hat{\ell}_\tau(i))^2 \right]} \\ &\leq 2\sqrt{\eta^2 |m_t| K}, \end{aligned}$$

where the second inequality is due to Lemma 7, a new result that proves the multi-round stability of FTRL iterates under certain conditions, which can be applied for sufficiently small  $\gamma$ . By using  $\sqrt{\eta |m_t| \eta K} \leq \frac{1}{2}(\eta |m_t| + \eta K)$  we can see that (3) is in fact of order  $\log(K)/\eta + d_{\max} K \ln(T) + \eta(KT + D)$ , which gives a  $O(\sqrt{(KT + D) \log(K)} + d_{\max} K \ln(T))$  bound for an appropriately tuned  $\eta$ .

To conclude, as long as loss estimates  $\hat{\ell}_\tau$  and  $\hat{\ell}_{\tau'}$  are independent for  $\tau, \tau' \in m_t$ , Lemma 3 implies that we have effectively split the cost of delayed feedback and bandit feedback. We formalize the above in Corollary 4, whose proof can be found in Section 3.2.

**Corollary 4** *Under the same assumptions as in Lemma 3, suppose that  $\mathbb{E}[\hat{\ell}_\tau | \mathcal{F}_t] = \ell_\tau$  and that  $\mathbb{E}[(\hat{\ell}_\tau - \ell_\tau)^\top (\nabla^2 R_t(\mathbf{w}_t(\hat{\mathbf{L}}_t)))^{-1} (\hat{\ell}_{\tau'} - \ell_{\tau'}) | \mathcal{F}_t] = 0$  for all  $t \in [T]$  and all  $\tau, \tau' \in m_t$  where*



$\tau' \neq \tau$ . Let  $\Lambda = \emptyset$  and let  $\mathcal{W}_t = \mathcal{W}$ . Then for all  $\mathbf{u} \in \mathcal{W}$

$$\mathbb{E} \left[ \sum_{t=1}^T (\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{u})^\top \boldsymbol{\ell}_t \right] \leq R_T(\mathbf{u}) - \min_{\mathbf{v} \in \mathcal{W}} R_1(\mathbf{v}) + 16 \sum_{t=1}^T \beta_t^2 + 16 \sum_{t=1}^T \alpha_t^2 |m_t|.$$

### 3.2 Analysis Details

In this section we present the proofs of Lemma 3 and Corollary 4. We start by developing the necessary tools in Lemmas 5, 6, and 7. Beginning with the former two, both of which are standard and can be found in various forms in the literature.

**Lemma 5** *Suppose that Assumption 2 holds. Let  $\mathbf{v} \in \mathcal{W}_t$  and  $\mathbf{L}, \mathbf{L}' \in \mathbb{R}^K$  such that  $\mathbf{w}_t(\mathbf{L}'), \mathbf{w}_t(\mathbf{L}) \in \mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2\sqrt{\kappa}})$ , then  $\|\mathbf{w}_t(\mathbf{L}) - \mathbf{w}_t(\mathbf{L}')\|_{R_t, \mathbf{v}}^* \leq 8\|\mathbf{L}' - \mathbf{L}\|_{R_t, \mathbf{v}}$ .*

**Proof** By Taylor's theorem and the optimality of  $\mathbf{w}_t(\mathbf{L}')$  we have that for some  $\zeta$  on the line segment between  $\mathbf{w}_t(\mathbf{L}')$  and  $\mathbf{w}_t(\mathbf{L})$

$$\begin{aligned} & \mathbf{L}'^\top \mathbf{w}_t(\mathbf{L}) + R_t(\mathbf{w}_t(\mathbf{L})) - \mathbf{L}'^\top \mathbf{w}_t(\mathbf{L}') - R_t(\mathbf{w}_t(\mathbf{L}')) \\ & \geq \frac{1}{2} (\mathbf{w}_t(\mathbf{L}') - \mathbf{w}_t(\mathbf{L}))^\top \nabla^2 R_t(\zeta) (\mathbf{w}_t(\mathbf{L}') - \mathbf{w}_t(\mathbf{L})) \\ & \geq \frac{1}{8} (\mathbf{w}_t(\mathbf{L}') - \mathbf{w}_t(\mathbf{L}))^\top \nabla^2 R_t(\mathbf{v}) (\mathbf{w}_t(\mathbf{L}') - \mathbf{w}_t(\mathbf{L})), \end{aligned}$$

where the last inequality is due the assumption on  $\nabla^2 R_t(\mathbf{v})$ , which is applicable because if  $\mathbf{w}_t(\mathbf{L}'), \mathbf{w}_t(\mathbf{L}) \in \mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2\sqrt{\kappa}})$  the line segment between  $\mathbf{w}_t(\mathbf{L}')$  and  $\mathbf{w}_t(\mathbf{L})$  is also in  $\mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2\sqrt{\kappa}})$ . Thus  $\zeta \in \mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2\sqrt{\kappa}})$ .

By adding and subtracting  $\mathbf{L}^\top (\mathbf{w}_t(\mathbf{L}) - \mathbf{w}_t(\mathbf{L}'))$  we have that

$$\begin{aligned} & \mathbf{L}'^\top \mathbf{w}_t(\mathbf{L}) + R_t(\mathbf{w}_t(\mathbf{L})) - \mathbf{L}'^\top \mathbf{w}_t(\mathbf{L}') - R_t(\mathbf{w}_t(\mathbf{L}')) \\ & = (\mathbf{L}' - \mathbf{L})^\top (\mathbf{w}_t(\mathbf{L}) - \mathbf{w}_t(\mathbf{L}')) + \mathbf{L}^\top \mathbf{w}_t(\mathbf{L}) + R_t(\mathbf{w}_t(\mathbf{L})) - \mathbf{L}^\top \mathbf{w}_t(\mathbf{L}') - R_t(\mathbf{w}_t(\mathbf{L}')) \\ & \leq (\mathbf{L}' - \mathbf{L})^\top (\mathbf{w}_t(\mathbf{L}) - \mathbf{w}_t(\mathbf{L}')) \\ & \leq \|\mathbf{L}' - \mathbf{L}\|_{R_t, \mathbf{v}} \|\mathbf{w}_t(\mathbf{L}) - \mathbf{w}_t(\mathbf{L}')\|_{R_t, \mathbf{v}}^*, \end{aligned}$$

where the first inequality is due to the optimality of  $\mathbf{w}_t(\mathbf{L})$  and the second inequality is Hölder's inequality. Thus, we may conclude that

$$\|\mathbf{L}' - \mathbf{L}\|_{R_t, \mathbf{v}} \|\mathbf{w}_t(\mathbf{L}) - \mathbf{w}_t(\mathbf{L}')\|_{R_t, \mathbf{v}}^* \geq \frac{1}{8} \left( \|\mathbf{w}_t(\mathbf{L}) - \mathbf{w}_t(\mathbf{L}')\|_{R_t, \mathbf{v}}^* \right)^2,$$

which concludes the proof after multiplying both sides of the above by  $\frac{8}{\|\mathbf{w}_t(\mathbf{L}) - \mathbf{w}_t(\mathbf{L}')\|_{R_t, \mathbf{v}}^*}$ . ■

**Lemma 6** *Suppose that Assumption 2 holds. Let  $z \subset \mathbb{N}$  be a finite set, and define  $\mathbf{L}' = \mathbf{L} + \sum_{\tau \in z} \mathbf{y}_\tau$ , where  $\mathbf{y}_\tau \in \mathbb{R}^K$ . If  $\sum_{\tau \in z} \sqrt{\kappa} \|\mathbf{y}_\tau\|_{R_{t'}, \mathbf{w}_t(\mathbf{L})} \leq \frac{1}{32}$  and  $\mathcal{W}_t = \mathcal{W}_{t'}$ , then  $\mathbf{w}_{t'}(\mathbf{L}') \in \mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$ .*

**Proof** Because of the strict convexity of all  $R_t$ , to show that  $\mathbf{w}_{t'}(\mathbf{L}') \in \mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$  it suffices to show that for all  $\mathbf{x}$  on the boundary of  $\mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$

$$\mathbf{L}'^\top \mathbf{x} + R_{t'}(\mathbf{x}) \geq \mathbf{L}'^\top \mathbf{w}_t(\mathbf{L}) + R_{t'}(\mathbf{w}_t(\mathbf{L})). \quad (4)$$

To see why the strict convexity of  $R_{t'}$  is sufficient, suppose that all  $\mathbf{x}$  that are on the boundary of  $\mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$  indeed satisfy (4). For the sake of contradiction suppose that  $\mathbf{w}_{t'}(\mathbf{L}')$  is not in  $\mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$ . Let  $\mathbf{z} = (1 - a)\mathbf{w}_t(\mathbf{L}) + a\mathbf{w}_{t'}(\mathbf{L}')$  be the point on the boundary of  $\mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$  on the segment between  $\mathbf{w}_t(\mathbf{L})$  and  $\mathbf{w}_{t'}(\mathbf{L}')$ . Then

$$\begin{aligned} & \mathbf{L}'^\top \mathbf{w}_t(\mathbf{L}) + R_{t'}(\mathbf{w}_t(\mathbf{L})) \\ & \leq \mathbf{L}'^\top \mathbf{z} + R_{t'}(\mathbf{z}) \\ & < (1 - a)(\mathbf{L}'^\top \mathbf{w}_t(\mathbf{L}) + R_{t'}(\mathbf{w}_t(\mathbf{L}))) + a(\mathbf{L}'^\top \mathbf{w}_{t'}(\mathbf{L}') + R_{t'}(\mathbf{w}_{t'}(\mathbf{L}'))) \\ & \leq \mathbf{L}'^\top \mathbf{w}_t(\mathbf{L}) + R_{t'}(\mathbf{w}_t(\mathbf{L})), \end{aligned}$$

where the first inequality holds because we assumed (4) to be true and  $\mathbf{z}$  is on the boundary of  $\mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$  and the last inequality is by definition of  $\mathbf{w}_{t'}(\mathbf{L}')$  and the assumption that  $\mathcal{W}_t = \mathcal{W}_{t'}$ . Thus, we have a contradiction, which implies that if all  $\mathbf{x}$  on the boundary of  $\mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$  satisfy (4), then  $\mathbf{w}_{t'}(\mathbf{L}') \in \mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$ . We proceed by showing that all  $\mathbf{x}$  on the boundary of  $\mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$  satisfy (4). Let  $\mathbf{h} = \mathbf{x} - \mathbf{w}_t(\mathbf{L})$ . Note that

$$\mathbf{w}_t(\mathbf{L}) = \arg \min_{\mathbf{v} \in \mathcal{W}} \{\mathbf{v}^\top \mathbf{L} + R_t(\mathbf{v})\} = \arg \min_{\mathbf{v} \in \mathcal{W}} \{\kappa \mathbf{v}^\top \mathbf{L} + \kappa R_t(\mathbf{v})\}.$$

We have

$$\begin{aligned} \left( \kappa \mathbf{L} + \kappa \nabla R_{t'}(\mathbf{w}_t(\mathbf{L})) \right)^\top \mathbf{h} &= \left( \kappa \mathbf{L} + \kappa \nabla R_t(\mathbf{w}_t(\mathbf{L})) \right)^\top \mathbf{h} + \kappa \left( \nabla R_{t'}(\mathbf{w}_t(\mathbf{L})) - \nabla R_t(\mathbf{w}_t(\mathbf{L})) \right)^\top \mathbf{h} \\ &\geq \kappa \left( \nabla R_{t'}(\mathbf{w}_t(\mathbf{L})) - \nabla R_t(\mathbf{w}_t(\mathbf{L})) \right)^\top \mathbf{h} \\ &\geq -\frac{1}{32} \sqrt{\kappa \mathbf{h}^\top \nabla^2 R_{t'}(\mathbf{w}_t(\mathbf{L})) \mathbf{h}} = -\frac{1}{64}, \end{aligned}$$

where the first inequality is due to the optimality of  $\mathbf{w}_t(\mathbf{L})$ , the second inequality is per Assumption 2(a), implying that  $(\nabla \kappa R_t(\mathbf{w}_t(\mathbf{L})) - \nabla \kappa R_{t'}(\mathbf{w}_t(\mathbf{L})))^\top \mathbf{x} \leq \frac{1}{32} \sqrt{\kappa \mathbf{x}^\top \nabla^2 R_{t'}(\mathbf{w}_t(\mathbf{L})) \mathbf{x}}$ , and the last equality is due to the fact that  $\mathbf{h}$  is a point on the boundary of  $\mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$  and thus  $\|\mathbf{h}\|_{R_{t'}, \mathbf{w}_t(\mathbf{L})}^* = \frac{1}{2\sqrt{\kappa}}$ . Using Taylor's theorem, there exists  $\zeta$  on the segment between  $\mathbf{x}$  and  $\mathbf{w}_t(\mathbf{L})$  such that

$$\begin{aligned} & \kappa \mathbf{L}'^\top \mathbf{x} + \kappa R_{t'}(\mathbf{x}) - \kappa \mathbf{L}'^\top \mathbf{w}_t(\mathbf{L}) - \kappa R_{t'}(\mathbf{w}_t(\mathbf{L})) \\ &= \kappa (\mathbf{L}' - \mathbf{L})^\top \mathbf{h} + \left( \kappa \mathbf{L} + \kappa \nabla R_{t'}(\mathbf{w}_t(\mathbf{L})) \right)^\top \mathbf{h} + \frac{\kappa}{2} \mathbf{h}^\top \nabla^2 R_{t'}(\zeta) \mathbf{h} \\ &\geq \kappa (\mathbf{L}' - \mathbf{L})^\top \mathbf{h} - \frac{1}{64} + \frac{1}{2} \mathbf{h}^\top \nabla^2 R_{t'}(\zeta) \mathbf{h} \\ &\geq \kappa (\mathbf{L}' - \mathbf{L})^\top \mathbf{h} - \frac{1}{64} + \frac{\kappa}{8} \mathbf{h}^\top \nabla^2 R_{t'}(\mathbf{w}_t(\mathbf{L})) \mathbf{h} \\ &= \kappa (\mathbf{L}' - \mathbf{L})^\top \mathbf{h} + \frac{1}{64} \end{aligned} \quad (5)$$

where we also used Assumption 2(a) and that  $\zeta \in \mathcal{D}_{R_{t'}}(\mathbf{w}_t(\mathbf{L}), \frac{1}{2\sqrt{\kappa}})$ . Thus, by applying Hölder's inequality we can see that

$$\begin{aligned} \kappa \mathbf{L}'^\top \mathbf{x} + \kappa R_{t'}(\mathbf{x}) - \kappa \mathbf{L}'^\top \mathbf{w}_t(\mathbf{L}) - \kappa R_{t'}(\mathbf{w}_t(\mathbf{L})) &\geq - \sum_{\tau \in z} \kappa \|\mathbf{y}_\tau\|_{R_{t'}, \mathbf{w}_t(\mathbf{L})} \|\mathbf{h}\|_{R_{t'}, \mathbf{w}_t(\mathbf{L})}^* + \frac{1}{64} \\ &= -\frac{1}{2} \sum_{\tau \in z} \sqrt{\kappa} \|\mathbf{y}_\tau\|_{R_{t'}, \mathbf{w}_t(\mathbf{L})} + \frac{1}{64} \geq 0, \end{aligned}$$

where the equality is due to the fact that  $\|\mathbf{h}\|_{R_{t'}, \mathbf{w}_t(\mathbf{L})}^* = \frac{1}{2\sqrt{\kappa}}$  and the final inequality is due to the assumption that  $\sum_{\tau} \sqrt{\kappa} \|\mathbf{y}_\tau\|_{R_{t'}, \mathbf{w}_t(\mathbf{L})} \leq \frac{1}{32}$ .  $\blacksquare$

The following lemma states that if the local norms of the loss estimates,  $\|\hat{\ell}_t\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)}$ , are small enough, the iterates of FTRL are close across multiple rounds. This is a crucial ingredient in our analysis, as this allows us to use Assumption (a) to control the variance term in Lemma 3. This lemma might be of independent interest.

**Lemma 7** *Suppose that Assumption 1 and Assumption 2 hold. Also suppose that for all  $t, t' \in [T]$ ,  $\sqrt{\kappa} \|\hat{\ell}_t\|_{R_{t'}, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq \frac{1}{128d_{\max}}$ . Then, for all  $t \in [T]$  and all  $\tau \in m_t$  we have that  $\mathbf{w}_t(\hat{\mathbf{L}}_t) \in \mathcal{D}_{R_t}(\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau), \frac{1}{2\sqrt{\kappa}})$ .*

**Proof** We will prove the statement by induction. Assume that there exists a  $t \in [T]$  such that for all  $\tau < t$  and all  $s \in m_\tau$ , it holds that  $\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau) \in \mathcal{D}_{R_\tau}(\mathbf{w}_s(\hat{\mathbf{L}}_s), \frac{1}{2\sqrt{\kappa}})$ . Now pick any  $s \in m_t$ . For the induction step we need to show that  $\mathbf{w}_t(\hat{\mathbf{L}}_t) \in \mathcal{D}_{R_t}(\mathbf{w}_s(\hat{\mathbf{L}}_s), \frac{1}{2\sqrt{\kappa}})$ . The goal is to apply Lemma 6 for which we start by decomposing  $o_t \setminus o_s$  into the losses that were already missing at timestep  $s$  (and were observed later) and the losses that we incurred and observed after the round  $s$ ,

$$\begin{aligned} \sum_{\tau \in o_t \setminus o_s} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_s(\hat{\mathbf{L}}_s)} &= \sum_{\substack{\tau \in o_t \setminus o_s \\ \tau \geq s}} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_s(\hat{\mathbf{L}}_s)} + \sum_{\tau \in m_s \setminus m_t} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_s(\hat{\mathbf{L}}_s)} \\ &\leq 2 \sum_{\substack{\tau \in o_t \setminus o_s \\ \tau \geq s}} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)} + 2 \sum_{\tau \in m_s \setminus m_t} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)} \\ &= 2 \sum_{\tau \in o_t \setminus o_s} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}, \end{aligned}$$

For the inequality, we are applying Lemma 17 using the fact that  $\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau) \in \mathcal{D}_{R_\tau}(\mathbf{w}_s(\hat{\mathbf{L}}_s), \frac{1}{2\sqrt{\kappa}})$  for  $\tau \in m_s$  and  $\mathbf{w}_s(\hat{\mathbf{L}}_s) \in \mathcal{D}_{R_s}(\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau), \frac{1}{2\sqrt{\kappa}})$  for  $\tau \geq s$  (where we follow  $s \in m_\tau$  which follows from  $s \in m_t$  and  $t \geq \tau$ ), both of which hold by the inductive assumption. We continue:

$$2 \sum_{\tau \in o_t \setminus o_s} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)} \leq \frac{2|o_t \setminus o_s|}{128d_{\max}} \leq \frac{1}{32},$$

where the first inequality is per the assumption and the second inequality follows by counting the number of elements in  $o_t \setminus o_s$ , which we do as

$$|o_t \setminus o_s| \leq |\{\hat{\ell}_{t-2d_{\max}}, \dots, \hat{\ell}_{t-1}\}| = 2d_{\max}.$$

Since we have now established that  $\sum_{\tau \in o_t \setminus o_s} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_s(\hat{\mathbf{L}}_s)} \leq \frac{1}{128d_{\max}}$  we can apply Lemma 6 to conclude that  $\mathbf{w}_t(\hat{\mathbf{L}}_t) \in \mathcal{D}_{R_t}(\mathbf{w}_s(\hat{\mathbf{L}}_s), \frac{1}{2\sqrt{\kappa}})$  as  $\mathcal{W}_t = \mathcal{W}_s$ , which holds by Assumption 1. That completes the induction step as we have chosen  $s$  arbitrarily. For the basis of induction it is sufficient to note that  $\mathbf{w}_1(\hat{\mathbf{L}}_1) \in \mathcal{D}_{R_1}(\mathbf{w}_1(\hat{\mathbf{L}}_1), \frac{1}{2\sqrt{\kappa}})$  holds trivially.  $\blacksquare$

Now that we have gathered the necessary tools, we can prove our main Lemma.

**Lemma 3 (RESTATED)** *Suppose that  $\mathbb{E}[\hat{\ell}_t | \mathcal{F}_t] = \ell_t + \mathbf{b}_t$  and suppose that Assumption 1 and Assumption 2 hold. Let  $t \in [T]$  and  $\tau \in m_t \cup \{t\}$ . Suppose that  $\|\ell_\tau\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq \alpha_t$  and  $\mathbb{E}[\|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}^2] \leq \beta_t^2$ . Suppose that for all  $t, t' \in [T]$   $\sqrt{\kappa} \|\hat{\ell}_t\|_{R_{t'}, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq \frac{1}{128d_{\max}}$ . Then for all  $\mathbf{u} \in \mathcal{W}_T$ ,*

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{u})^\top \ell_t \right] &\leq \mathbb{E} \left[ \sum_{t \in \Lambda} (\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{u})^\top \ell_t \right] + R_T(\mathbf{u}) - \min_{v \in \mathcal{W}_1} R_1(v) + \sum_{t \in \bar{\Lambda}} 8\beta_t^2 \\ &\quad - \sum_{t \in \bar{\Lambda}} \mathbb{E}[(\mathbf{w}_t(\hat{\mathbf{L}}_t^m) - \mathbf{u})^\top \mathbf{b}_t] + \sum_{t \in \bar{\Lambda}} \left( 8\alpha_t^2 |m_t| + 8\alpha_t \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\ell_\tau - \hat{\ell}_\tau) \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \right] \right). \end{aligned}$$

**Proof** The first step of the proof is to establish some base facts including that  $\mathbf{w}_t(\hat{\mathbf{L}}_t) \in \mathcal{D}_{R_t}(\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau), \frac{1}{2\sqrt{\kappa}})$  for all  $\tau \in m_t$  and  $\mathbf{w}_t(\hat{\mathbf{L}}_t^m), \mathbf{w}_t(\hat{\mathbf{L}}_t^*) \in \mathcal{D}_{R_t}(\mathbf{w}_t(\hat{\mathbf{L}}_t), \frac{1}{2\sqrt{\kappa}})$ .

Since  $\sqrt{\kappa} \|\hat{\ell}_t\|_{R_{t'}, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq \frac{1}{128d_{\max}}$ , by Lemma 7 we may conclude that  $\mathbf{w}_t(\hat{\mathbf{L}}_t) \in \mathcal{D}_{R_t}(\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau), \frac{1}{2\sqrt{\kappa}})$  for all  $\tau \in m_t$  and all  $t$ , which is also a prerequisite for Lemma 17. We also note that  $\mathbf{w}_t(\hat{\mathbf{L}}_t) \in \mathcal{D}_{R_t}(\mathbf{w}_t(\hat{\mathbf{L}}_t), \frac{1}{2\sqrt{\kappa}})$  holds trivially. Now we can conclude that

$$\sum_{\tau \in m_t} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq \sum_{\tau \in m_t \cup \{t\}} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq 2 \sum_{\tau \in m_t \cup \{t\}} \sqrt{\kappa} \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)} \leq \frac{1}{32},$$

where we used Lemma 17 in the second inequality and the assumption on  $\sqrt{\kappa} \|\hat{\ell}_t\|_{R_{t'}, \mathbf{w}_t(\hat{\mathbf{L}}_t)}$  alongside the fact that  $|m_t \cup \{t\}| \leq d_{\max} + 1 \leq 2d_{\max}$  in the third inequality. By Lemma 6 and Assumption 1 we now know that  $\mathbf{w}_t(\hat{\mathbf{L}}_t^m), \mathbf{w}_t(\hat{\mathbf{L}}_t^*) \in \mathcal{D}_{R_t}(\mathbf{w}_t(\hat{\mathbf{L}}_t), \frac{1}{2})$ .

We decompose the regret as follows

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[(\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{u})^\top \ell_t] &= \underbrace{\mathbb{E} \left[ \sum_{t \in \Lambda} (\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{u})^\top \ell_t \right]}_{\text{skipped rounds}} + \underbrace{\sum_{t \in \bar{\Lambda}} -\mathbb{E}[(\mathbf{w}_t(\hat{\mathbf{L}}_t^m) - \mathbf{u})^\top \mathbf{b}_t]}_{\text{bias}} \quad (6) \\ &\quad + \underbrace{\sum_{t \in \bar{\Lambda}} \mathbb{E}[(\mathbf{w}_t(\hat{\mathbf{L}}_t^*) - \mathbf{u})^\top \hat{\ell}_t]}_{\text{cheating regret}} + \sum_{t \in \bar{\Lambda}} \left( \underbrace{\mathbb{E}[(\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{w}_t(\hat{\mathbf{L}}_t^m))^\top \ell_t]}_{H_1} + \underbrace{\mathbb{E}[(\mathbf{w}_t(\hat{\mathbf{L}}_t^m) - \mathbf{w}_t(\hat{\mathbf{L}}_t^*))^\top \hat{\ell}_t]}_{H_2} \right). \end{aligned}$$

By Hölder's inequality and Lemma 5

$$\begin{aligned}
 H_1 &= \mathbb{E} \left[ (\mathbf{w}_t(\widehat{\mathbf{L}}_t) - \mathbf{w}_t(\widehat{\mathbf{L}}_t^m))^\top \boldsymbol{\ell}_t \right] \\
 &\leq \mathbb{E} \left[ \|\mathbf{w}_t(\widehat{\mathbf{L}}_t) - \mathbf{w}_t(\widehat{\mathbf{L}}_t^m)\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)}^* \|\boldsymbol{\ell}_t\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \right] \\
 &\leq \mathbb{E} \left[ 8 \|\widehat{\mathbf{L}}_t - \widehat{\mathbf{L}}_t^m\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \|\boldsymbol{\ell}_t\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \right] \\
 &= \mathbb{E} \left[ 8 \left\| \sum_{\tau \in m_t} \widehat{\boldsymbol{\ell}}_\tau \right\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \|\boldsymbol{\ell}_t\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \right] \\
 &= \mathbb{E} \left[ 8 \left\| \sum_{\tau \in m_t} (\widehat{\boldsymbol{\ell}}_\tau - \boldsymbol{\ell}_\tau) + \sum_{\tau \in m_t} \boldsymbol{\ell}_\tau \right\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \|\boldsymbol{\ell}_t\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \right] \\
 &\leq \mathbb{E} \left[ 8 \left( \left\| \sum_{\tau \in m_t} (\widehat{\boldsymbol{\ell}}_\tau - \boldsymbol{\ell}_\tau) \right\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} + \left\| \sum_{\tau \in m_t} \boldsymbol{\ell}_\tau \right\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \right) \|\boldsymbol{\ell}_t\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \right] \\
 &\leq 8\alpha_t \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\widehat{\boldsymbol{\ell}}_\tau - \boldsymbol{\ell}_\tau) \right\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \right] + 8\alpha_t^2 |m_t|, \tag{7}
 \end{aligned}$$

where the last inequality is due to the triangle inequality and the assumptions on  $\|\boldsymbol{\ell}_\tau\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)}$ . Similarly we bound

$$H_2 = \mathbb{E} \left[ (\mathbf{w}_t(\widehat{\mathbf{L}}_t^m) - \mathbf{w}_t(\widehat{\mathbf{L}}_t^*))^\top \widehat{\boldsymbol{\ell}}_t \right] \leq 8\beta_t^2. \tag{8}$$

By Lemma 18 we have that

$$\text{cheating regret} = \sum_{t \in \overline{\Lambda}} (\mathbf{w}_t(\widehat{\mathbf{L}}_t^*) - \mathbf{u})^\top \widehat{\boldsymbol{\ell}}_t \leq R_T(\mathbf{u}) - \min_{v \in \mathcal{W}_1} R_1(v). \tag{9}$$

By combining equations (7), (8), and (9) with the regret decomposition and leaving the skipped rounds and bias untouched we find

$$\begin{aligned}
 \sum_{t=1}^T \mathbb{E} [(\mathbf{w}_t(\widehat{\mathbf{L}}_t) - \mathbf{u})^\top \boldsymbol{\ell}_t] &= \mathbb{E} \left[ \sum_{t \in \Lambda} (\mathbf{w}_t(\widehat{\mathbf{L}}_t) - \mathbf{u})^\top \boldsymbol{\ell}_t \right] + \sum_{t \in \overline{\Lambda}} -\mathbb{E} [(\mathbf{w}_t(\widehat{\mathbf{L}}_t^m) - \mathbf{u})^\top \mathbf{b}_t] \\
 &+ \sum_{t \in \overline{\Lambda}} \mathbb{E} [(\mathbf{w}_t(\widehat{\mathbf{L}}_t^*) - \mathbf{u})^\top \widehat{\boldsymbol{\ell}}_t] + \sum_{t \in \overline{\Lambda}} \left( \mathbb{E} [(\mathbf{w}_t(\widehat{\mathbf{L}}_t) - \mathbf{w}_t(\widehat{\mathbf{L}}_t^m))^\top \boldsymbol{\ell}_t] + \mathbb{E} [(\mathbf{w}_t(\widehat{\mathbf{L}}_t^m) - \mathbf{w}_t(\widehat{\mathbf{L}}_t^*))^\top \widehat{\boldsymbol{\ell}}_t] \right) \\
 &\leq \mathbb{E} \left[ \sum_{t \in \Lambda} (\mathbf{w}_t(\widehat{\mathbf{L}}_t) - \mathbf{u})^\top \boldsymbol{\ell}_t \right] + R_T(\mathbf{u}) - \min_{v \in \mathcal{W}_1} R_1(v) + \sum_{t \in \overline{\Lambda}} 8\beta_t^2 \\
 &- \sum_{t \in \overline{\Lambda}} \mathbb{E} [(\mathbf{w}_t(\widehat{\mathbf{L}}_t^m) - \mathbf{u})^\top \mathbf{b}_t] + \sum_{t \in \overline{\Lambda}} \left( 8\alpha_t^2 |m_t| + 8\alpha_t \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\boldsymbol{\ell}_\tau - \widehat{\boldsymbol{\ell}}_\tau) \right\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \right] \right).
 \end{aligned}$$

which concludes the proof. ■

We conclude this section with the proof of Corollary 4.

**Corollary 4 (RESTATED)** *Under the same assumptions as in Lemma 3, suppose that  $\mathbb{E}[\hat{\ell}_\tau | \mathcal{F}_t] = \ell_\tau$  and that  $\mathbb{E}[(\hat{\ell}_\tau - \ell_\tau)^\top (\nabla^2 R_t(\mathbf{w}_t(\hat{\mathbf{L}}_t)))^{-1} (\hat{\ell}_{\tau'} - \ell_{\tau'}) | \mathcal{F}_t] = 0$  for all  $t \in [T]$  and all  $\tau, \tau' \in m_t$  where  $\tau' \neq \tau$ . Let  $\Lambda = \emptyset$  and let  $\mathcal{W}_t = \mathcal{W}$ . Then for all  $\mathbf{u} \in \mathcal{W}$*

$$\mathbb{E} \left[ \sum_{t=1}^T (\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{u})^\top \ell_t \right] \leq R_T(\mathbf{u}) - \min_{\mathbf{v} \in \mathcal{W}} R_1(\mathbf{v}) + 16 \sum_{t=1}^T \beta_t^2 + 16 \sum_{t=1}^T \alpha_t^2 |m_t|.$$

**Proof** We are looking to control  $\mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\ell_\tau - \hat{\ell}_\tau) \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \right]$  for a given  $t \in [T]$ . We start by considering

$$\begin{aligned} \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\hat{\ell}_\tau - \ell_\tau) \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)}^2 \right] &= \sum_{\tau \in m_t} \mathbb{E} \left[ \left\| \hat{\ell}_\tau - \ell_\tau \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)}^2 \right] \\ &= \sum_{\tau \in m_t} \left( \mathbb{E} \left[ \left\| \hat{\ell}_\tau \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)}^2 \right] - \mathbb{E} \left[ \left\| \ell_\tau \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)}^2 \right] \right) \\ &\leq \sum_{\tau \in m_t} \mathbb{E} \left[ \left\| \hat{\ell}_\tau \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)}^2 \right], \end{aligned}$$

where we used that  $\mathbb{E}[(\hat{\ell}_\tau - \ell_\tau)^\top (\nabla^2 R_t(\mathbf{w}_t(\hat{\mathbf{L}}_t)))^{-1} (\hat{\ell}_{\tau'} - \ell_{\tau'}) | \mathcal{F}_t] = 0$  for  $\tau \neq \tau'$  in the first equality, and that  $\mathbb{E}[\hat{\ell}_\tau | \mathcal{F}_t] = \ell_\tau$  in the second equality. In turn, the above together with Jensen's inequality implies that

$$\begin{aligned} \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\hat{\ell}_\tau - \ell_\tau) \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \right] &\leq \sqrt{\sum_{\tau \in m_t} \mathbb{E} \left[ \left\| \hat{\ell}_\tau \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)}^2 \right]} \\ &\leq \sqrt{4 \sum_{\tau \in m_t} \mathbb{E} \left[ \left\| \hat{\ell}_\tau \right\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}^2 \right]} \leq \sqrt{4|m_t|\beta_t^2}, \end{aligned} \quad (10)$$

where in the second inequality we used Lemma 7 together with Lemma 17. Finally, the third inequality of (10) is due to the assumptions of Lemma 3. We conclude by substituting this bound into the results of Lemma 3,

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[(\mathbf{w}_t(\hat{\mathbf{L}}_t) - \mathbf{u})^\top \ell_t] &\leq R_T(\mathbf{u}) - \min_{\mathbf{v} \in \mathcal{W}} R_1(\mathbf{v}) + 8 \sum_{t=1}^T \beta_t^2 + \sum_{t=1}^T \left( 8\alpha_t^2 |m_t| + 16\sqrt{|m_t|\alpha_t^2 \beta_t^2} \right) \\ &\leq R_T(\mathbf{u}) - \min_{\mathbf{v} \in \mathcal{W}} R_1(\mathbf{v}) + 16 \sum_{t=1}^T \beta_t^2 + 16\alpha_t^2 \sum_{t=1}^T |m_t|, \end{aligned}$$

where in the last inequality we used that  $\sqrt{ab} \leq \frac{1}{2}(a+b)$  for  $a, b > 0$ . ■

## 4. Combinatorial Semi-Bandits

In this section, we demonstrate how to apply our generic FTRL approach to combinatorial semi-bandits (CMAB) with delayed feedback. As outlined in the introduction, combinatorial semi-bandits extend multi-armed bandits to be able to efficiently deal with combinatorial decision spaces

---

**Algorithm 1:** Delayed FTRL for combinatorial semi-bandits

---

**Input:** Regularizers  $\{R_t\}_{t \geq 1}$  defined in (11), including hyperparams  $\gamma \in (0, 1)$  and  $\{\eta_t\}_{t \geq 1}$ .  
**for**  $t \in [T]$  **do**  
    Observe  $\mathbf{a}_\tau \odot \ell_\tau$  for  $\tau \in o_t \setminus o_{t-1}$ .  
    Find loss estimators  $\hat{\ell}_\tau(i) = \frac{\mathbf{a}_\tau(i)\ell_\tau(i)}{\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, i)}$  for new observations  $\tau \in o_t \setminus o_{t-1}$ .  
    Compute  $\mathbf{w}_t(\hat{\mathbf{L}}_t) = \arg \min_{\mathbf{v} \in \mathcal{W}} \hat{\mathbf{L}}_t^\top \mathbf{v} + R_t(\mathbf{v})$ .  
    Find probability distribution  $\mathbf{p}_t$  such that  $\mathbb{E}_{\mathbf{a} \sim \mathbf{p}_t}[\mathbf{a}] = \mathbf{w}_t(\hat{\mathbf{L}}_t)$ .  
    Draw and play  $\mathbf{a}_t \sim \mathbf{p}_t$ .  
**end for**

---

and have been used in portfolio management (Ni et al., 2023) and recommendation systems (Lou  dec et al., 2015) among others. In combinatorial semi-bandits the learner picks an action  $\mathbf{a}_t \in \mathcal{A}$  at each timestep  $t$ . The actionset  $\mathcal{A} \subseteq \{0, 1\}^K$  is given as part of the problem formulation. The loss of the learner is defined as  $\mathbf{a}_t^\top \ell_t$ , where  $\ell_t \in [-1, 1]^K$ . In the combinatorial semi-bandit setting the feedback function is  $\mathcal{L}(\ell_\tau, \mathbf{a}_\tau) = \mathbf{a}_\tau \odot \ell_\tau$ , where  $\odot$  is the Hadamard (elementwise) vector product. A practical example is a path-finding problem. Consider a directed weighted graph, where the weight on the edges corresponds to some cost associated with traversing an edge. The objective of the learner is to reach a goal state while incurring the least loss. In this setting the dimension of the actions is equal to the number of edges on the graph and the actionset  $\mathcal{A}$  is the set of all valid paths from the starting state to the goal state. The loss is the cost associated with each edge and the feedback is either the individual weights of the edges traversed for semi-bandits or the entire cost of the path taken in full bandits.

We define the pseudo-regret in this setting as

$$\mathcal{R}_T = \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{a}^*)^\top \ell_t \right] \quad \text{with} \quad \mathbf{a}^* = \arg \min_{\mathbf{a} \in \mathcal{A}} \sum_{t=1}^T \mathbf{a}^\top \ell_t.$$

**Algorithm** Algorithm 1 is inspired by the algorithm of Audibert et al. (2014). In any given round  $t$ , Algorithm 1 first computes  $\mathbf{w}_t(\hat{\mathbf{L}}_t)$ , the solution of the FTRL optimization problem of Eq. (1) over the convex hull of the action set, that is with  $\mathcal{W} = \text{Conv}(\mathcal{A})$ . In this setting we are not skipping rounds and the domain is constant.  $\mathbf{w}_t(\hat{\mathbf{L}}_t)$  can be computed efficiently using standard methods from convex optimisation if  $\text{Conv}(\mathcal{A})$  can be described in a polynomial number of linear constraints, see Nemirovski (2004). Then, it constructs a probability distribution  $\mathbf{p}_t$  over  $\mathcal{A}$  such that  $\mathbb{E}_{\mathbf{a} \sim \mathbf{p}_t}[\mathbf{a}] = \mathbf{w}_t(\hat{\mathbf{L}}_t)$ . How to construct  $\mathbf{p}_t$  and if it can be sampled from efficiently depends on the actionset and for many commonly used actionsets, like  $m$ -sets and spanning trees, there exist efficient algorithms. The path finding problem outlined above can also be solved efficiently by relaxing the convex hull of paths in the directed graph to so called unit flows, leading to a runtime of  $O(n^4)$  where  $n$  is the number of nodes in the path finding problem (Koolen et al., 2010). For a more complete discussion on the computational efficiency of FTRL style combinatorial bandit algorithms and for which actionsets  $\mathbf{w}_t(\hat{\mathbf{L}}_t)$  and  $\mathbf{p}_t$  can be obtained efficiently we refer to Koolen et al. (2010), Cesa-Bianchi and Lugosi (2012), and Audibert et al. (2014). The estimator of loss is

given by  $\widehat{\ell}_t(i) = \frac{\mathbf{a}_t(i)\ell_t(i)}{w_t(\widehat{\mathbf{L}}, i)}$ , which is unbiased. We use the regularizer

$$R_t(\mathbf{v}) = \sum_{i=1}^K \left( \frac{1}{\eta_t} \mathbf{v}(i) \log(\mathbf{v}(i)) - \frac{1}{\gamma} \log(\mathbf{v}(i)) \right), \quad (11)$$

where  $\eta_t > 0$  and  $\gamma > 0$  are hyperparameters.

**Main Result and Discussion** The main result of this section is Theorem 8.

**Theorem 8** *Suppose that  $\max_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a}\|_1 \leq B$ . Algorithm 1 with*

$$\eta_t = \min \left\{ \sqrt{\frac{B(1 + \log(\frac{K}{B}))}{16(B \sum_{\tau=1}^t |m_\tau| + Kt)}}, \frac{B^2(1 + \log(\frac{K}{B}))}{128K(Bd_{\max} + K)} \right\}, \quad \gamma = \frac{1}{128\sqrt{B}d_{\max}},$$

*guarantees that*

$$\mathcal{R}_T \leq 12\sqrt{B \left(1 + \log\left(\frac{K}{B}\right)\right) (KT + BD) + 128K^2d_{\max} + 128\sqrt{B}d_{\max}K \log(T)}.$$

The result is based on Corollary 4. After confirming the conditions on the regularizer  $R_t$ , the proof finds  $\alpha_t = \sqrt{\eta_t B}$  and  $\beta_t^2 = \eta_t K$ . The last thing to do is to bound the size of the regularizer  $R_T(\mathbf{u})$  on the comparator  $\mathbf{u}$ , which is a term that also arises from Corollary 4. As  $R_t$  tends to infinity on parts of the boundary of  $\mathcal{W}$  we have to choose a  $\mathbf{u} \neq \mathbf{a}^*$  and we pick  $\mathbf{u}$  as the best point in hindsight in a slightly shrunken actionset. That allows us to bound  $R_T(\mathbf{u})$  in exchange for a small additive bias term. The full proof can be found in Appendix A.

Theorem 9 shows that our results are optimal up to log-factors.

**Theorem 9** *Suppose that  $d_t = d$  for all  $t$  and that  $B \leq K/2$ . Then for any algorithm there exists a sequence of losses such that*

$$\mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{a}^*)^\top \ell_t \right] = \Omega \left( \max \left\{ \sqrt{BKT}, B\sqrt{dT} \right\} \right).$$

The proof for Theorem 9 can be found in Appendix A. When using an action-set constructed of basis vectors, we recover the delayed multi-armed bandit setting, in which we match the optimal upper bound for delayed adversarial bandits due to Zimmert and Seldin (2020) up to constants and log-factors. In the non-delayed setting, we have  $D = 0$  and we recover a bound of  $O(\sqrt{B(1 + \log(\frac{K}{B}))KT})$ , which also matches the optimal upper bound of order by Audibert et al. (2014) up to constants.

## 5. Linear Bandits

In this section, we show how to apply our analysis of FTRL to linear bandits with delayed feedback, which is an instance of our general setting for  $\ell_t \in \mathbb{R}^K$  such that  $\max_t \|\ell_t\|_2 \leq 1$ ,  $\mathcal{A} = \mathcal{W} \subset \mathbb{R}^K$ , and the feedback function is  $\mathcal{L}(\ell, \mathbf{a}) = \ell^\top \mathbf{a}$ . Additionally, we assume that the domain is constant with  $\mathcal{W} \subseteq \mathcal{B}(B)$ , where  $\mathcal{B}(B)$  is an Euclidean ball with radius  $B$ . We are not skipping rounds in this setting.



---

**Algorithm 2:** Delayed FTRL for linear bandits

---

**Input:**  $\nu$ -self concordant barrier  $\Psi$  for  $\mathcal{W}$ , hyperparameters  $\{\eta_t, \gamma_t\}_{t \geq 1}$ .  
**Initialize:**  $\tilde{\Psi}(\cdot) = \Psi(\cdot) - \min_{v \in \mathcal{W}} \Psi(v)$  and  $R_t(\cdot) = \frac{1}{\eta_t} \|\cdot\|_2^2 + \frac{1}{\gamma_t} \tilde{\Psi}(\cdot)$  for  $t \geq 1$ .  
**for**  $t = 1, \dots, T$  **do**  
    Observe  $\mathbf{a}_\tau^\top \boldsymbol{\ell}_\tau$  for  $\tau \in o_t \setminus o_{t-1}$ .  
    Compute  $\hat{\boldsymbol{\ell}}_\tau = K \boldsymbol{\ell}_\tau^\top \mathbf{a}_\tau (\nabla^2 \Psi(\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)))^{1/2} \mathbf{v}_\tau$  for new observations  $\tau \in o_t \setminus o_{t-1}$ .  
    Compute  $\mathbf{w}_t(\hat{\mathbf{L}}_t) = \arg \min_{v \in \mathcal{W}} \hat{\mathbf{L}}_t^\top v + R_t(v)$ .  
    Sample  $\mathbf{v}_t$  uniformly from the unit sphere.  
    Play  $\mathbf{a}_t = \mathbf{w}_t(\hat{\mathbf{L}}_t) + (\nabla^2 \Psi(\mathbf{w}_t(\hat{\mathbf{L}}_t)))^{-1/2} \mathbf{v}_t$ .  
**end for**

---

**Algorithm** Our algorithm for the linear bandit setting is inspired by Abernethy et al. (2008), who provide an algorithm with nearly optimal bounds for the linear bandit setting with an efficient algorithm. For the delayed linear bandit setting we use a regularizer of the form  $R_t(\mathbf{v}) = \frac{1}{\eta_t} \|\mathbf{v}\|_2^2 + \frac{1}{\gamma_t} \tilde{\Psi}(\mathbf{v})$ , where  $\tilde{\Psi}(\mathbf{v}) = \Psi(\mathbf{v}) - \min_{v' \in \mathcal{W}} \Psi(v')$  for a  $\nu$ -self-concordant barrier function  $\Psi$ . For a thorough introduction to self-concordant barriers, we refer the reader to (Nesterov and Nemirovskii, 1994). In Appendix B, we recall the most important properties, which can be found in (Nemirovski and Todd, 2008, Section 2). The main reason for using self-concordant barriers is to adhere to Assumptions 2(a) and 2(b). As detailed in Appendix B, these are standard properties of self-concordant barriers.

Specific examples of self-concordant barriers are  $f(x) = -\log(x)$ , which is 1-self-concordant for the non-negative reals,  $f(\mathbf{x}) = -\log(1 - \|\mathbf{x}\|_2^2)$ , which is 1-self-concordant for the unit ball, the 1-self concordant barrier  $f(\mathbf{x}) = -\log(b - \mathbf{a}^\top \mathbf{x})$  for linear constraints  $\mathbf{a}^\top \mathbf{x} \leq b$ , and the entropic barrier, which is defined as

$$f(\mathbf{x}) = \sup_{\theta \in \mathbb{R}^d} \{\langle \mathbf{x}, \theta \rangle - f^*(\theta)\} \quad \text{where } f^*(\theta) = \ln \left( \int_{\mathcal{W}} \exp(\langle \mathbf{x}, \theta \rangle) d\mathbf{x} \right),$$

which is a  $d$ -self-concordant barrier for any  $\mathcal{W}$ . Unfortunately, even though the entropic barrier is a self-concordant barrier for all domains, it can not always be efficiently computed.

Finally, we turn to the way we choose the action  $\mathbf{a}_t \in \mathcal{A}$  and the construction of the estimator. We use  $\mathbf{a}_t = \mathbf{w}_t(\hat{\mathbf{L}}_t) + (\nabla^2 \Psi(\mathbf{w}_t(\hat{\mathbf{L}}_t)))^{-1/2} \mathbf{v}_t$ , where  $\mathbf{v}_t$  is sampled i.i.d. from the uniform distribution over the unit sphere. To see that  $\mathbf{a}_t \in \mathcal{A}$ , note that  $\mathcal{D}_\Psi(\mathbf{w}, 1) \subseteq \mathcal{W} = \mathcal{A}$  for any  $\mathbf{w} \in \mathcal{W}$  (see Appendix B). Since  $\|\mathbf{a}_t - \mathbf{w}_t(\hat{\mathbf{L}}_t)\|_{\Psi, \mathbf{w}_t(\hat{\mathbf{L}}_t)} = 1$ , we have that  $\mathbf{a}_t \in \mathcal{W}$ . As for the loss estimate, we use  $\hat{\boldsymbol{\ell}}_t = K \boldsymbol{\ell}_t^\top \mathbf{a}_t (\nabla^2 \Psi(\mathbf{w}_t(\hat{\mathbf{L}}_t)))^{1/2} \mathbf{v}_t$ , which can be seen to an unbiased estimator for  $\boldsymbol{\ell}_t$  after observing that  $\mathbb{E}[\mathbf{v}_t \mathbf{v}_t^\top | \mathcal{F}_t] = \frac{1}{K} \mathbf{I}$ .

Given that  $\Psi(\cdot)$ ,  $\nabla \Psi(\cdot)$ , and  $\nabla^2 \Psi(\cdot)$  can be efficiently computed there are two computationally demanding steps in Algorithm 2: the computation of  $\mathbf{w}_t(\hat{\mathbf{L}}_t)$  and the computation of  $(\nabla^2 \Psi(\mathbf{w}_t(\hat{\mathbf{L}}_t)))^{1/2}$  and its inverse.  $(\nabla^2 \Psi(\mathbf{w}_t(\hat{\mathbf{L}}_t)))^{1/2}$  and its inverse can be computed through an eigenvalue decomposition, which can be done in  $O(K^3)$ . Abernethy et al. (2008) show that an approximation of  $\mathbf{w}_t(\hat{\mathbf{L}}_t)$  can be computed in  $O(K^2)$  per round by using the damped Newton method. This approximation maintains the same regret bound up to constants. The implementation as well as an overview of the analysis can be found in Appendix B.1.

**Main Result and Discussion** We arrive at the main result of this section.

**Theorem 10** *Suppose that  $T > 100$  and  $B \geq 1$ . Algorithm 2, run with a  $\nu$ -self-concordant barrier  $\Psi$  and with*

$$\gamma_t = \min \left\{ \frac{1}{256BKd_{\max}}, \sqrt{\frac{\nu \log(1 + \sqrt{T})}{16B^2K^2t}} \right\}$$

$$\eta_t = \min \left\{ \frac{B}{256d_{\max}}, \sqrt{\frac{B^2}{16 \sum_{\tau=1}^t |m_t|}} \right\},$$

*guarantees that, for any  $\mathbf{u} \in \mathcal{W}$ ,*

$$\mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{u})^\top \ell_t \right] \leq 12BK \sqrt{\nu T \log(1 + \sqrt{T})} + 12B\sqrt{D} + 2B\sqrt{T}$$

$$+ 512BKd_{\max}\nu \log(1 + \sqrt{T}).$$

The proof of Theorem 10 can be found in Appendix B. It follows from an application of Corollary 4 and carefully tuning the learning rates.

Let us put Theorem 10 in perspective. For arbitrary  $\mathcal{W}$  we can use the entropic barrier as the regularizer, which means  $\nu = d$  and thus algorithm 2 has a  $\tilde{O}(K^{2/3}\sqrt{T} + \sqrt{D})$  regret bound. For constant delay, i.e.  $d_t = d$ , Ito et al. (2020a) show that continuous exponential weights obtains a  $\tilde{O}(K\sqrt{T} + \sqrt{dT})$  regret bound. Even though this algorithm can be computed in  $\text{poly}(K, T, B)$  time, the algorithm is far from practical. In contrast, (an approximation of) algorithm 2 can be computed in  $O(K^3)$  time, with only a slightly worse regret bound. Huang et al. (2023) provide an algorithm with similar computational complexity as algorithm 2, but their regret bound is  $\tilde{O}(K^{2/3}\sqrt{T} + K^2\sqrt{D})$ , which contains an unnecessary dependence on the dimension  $K$  in the delay term of the regret bound. However, it seems that the regret bound of Huang et al. (2023) can be improved to  $\tilde{O}(K\sqrt{\nu T} + K\sqrt{\nu D})$ . In their terminology: Banker-BOLO is  $(O(\nu \log(T)), K^2)$ -stable, which together with Theorem 4.6 of Huang et al. (2023) leads to a  $\tilde{O}(K\sqrt{\nu T} + K\sqrt{\nu D})$  regret bound. Still, the unnecessary dependence on the dimension  $K$  in the delay term of the regret bound remains.

## 6. Adversarial Markov Decision Processes (MDPs)

In this section, we apply our FTRL approach to adversarial Markov Decision Processes (MDPs) where the transition function is known to the learner in advance. We start with a presentation of the model and regret minimization framework.

A finite-horizon episodic adversarial MDP is defined by  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, p, \{\ell_t\}_{t=1}^T, s_{\text{init}})$ , where  $\mathcal{S}$  and  $\mathcal{A}$  are finite state and action spaces of sizes  $S$  and  $A$ , respectively,  $H$  is the horizon,  $T$  is the number of episodes, and  $s_{\text{init}} \in \mathcal{S}$  is the initial state. The transition function is  $p : [H] \times \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$ , where  $\Delta_{\mathcal{S}}$  is the simplex over the states and  $p(s' \mid h, s, a)$  is the probability of moving to  $s'$  when taking action  $a$  in state  $s$  at time  $h$ . The learner interacts with the environment over  $T$  episodes of length  $H$  each. At the beginning of episode  $t$ , the learner picks a policy  $\pi_t = [H] \times \mathcal{S} \rightarrow \Delta_{\mathcal{A}}$  and starts in the initial state  $s_{t,1} = s_{\text{init}}$ . For each  $h \in [H]$ , the learner observes the current state  $s_{t,h} \in \mathcal{S}$ , draws an action from the policy  $a_{t,h} \sim \pi_t(\cdot \mid h, s_{t,h})$ , and transitions

to the next state  $s_{t,h+1} \sim p(\cdot \mid h, s_{t,h}, a_{t,h})$ . The cost functions  $\ell_t \in [0, 1]^{HS^2A}$  are chosen by an oblivious adversary, and the feedback of episode  $t$  contains the elements of the cost function corresponding to the agent's trajectory  $\{\ell_t(h, s_{t,h}, a_{t,h})\}_{h=1}^H$  (i.e., bandit feedback) and is observed only at the end of episode  $t + d_t$ . The learner's goal is to minimize the value of its policies, where  $V_t^\pi(h, s) = \mathbb{E}[\sum_{h'=h}^H \ell_t(h', s_{h'}, a_{h'}) \mid s_h = s, \pi, p]$  is the value function of policy  $\pi$  with respect to the cost  $\ell_t$ . The performance is measured by the regret, defined as the difference between the cumulative expected cost of the learner and the best fixed policy in hindsight

$$\mathcal{R}_T = \sum_{t=1}^T V_t^{\pi_t}(1, s_{\text{init}}) - \min_{\pi \in \Pi} \sum_{t=1}^T V_t^\pi(1, s_{\text{init}}),$$

where  $\Pi$  is the set of all policies admitted by  $\mathcal{M}$ .

Given a policy  $\pi$  and a transition function  $p'$ , the occupancy measure  $\mathbf{q}^{\pi, p'} \in [0, 1]^{HS^2A}$  is a vector, where  $\mathbf{q}^{\pi, p'}(h, s, a, s')$  is the probability to visit state  $s$  at time  $h$ , take action  $a$  and transition to state  $s'$ . We also denote

$$\mathbf{q}^{\pi, p'}(h, s, a) = \sum_{s'} \mathbf{q}^{\pi, p'}(h, s, a, s') \quad \text{and} \quad \mathbf{q}^{\pi, p'}(h, s) = \sum_a \mathbf{q}^{\pi, p'}(h, s, a).$$

By Rosenberg and Mansour (2019b)—see also Zimin and Neu (2013); Dick et al. (2014)—the occupancy measure encodes the policy and the transition function through the relations

$$\pi(a \mid h, s) = \frac{\mathbf{q}^{\pi, p'}(h, s, a)}{\mathbf{q}^{\pi, p'}(h, s)} \quad \text{and} \quad p'(s' \mid h, s, a) = \frac{\mathbf{q}^{\pi, p'}(h, s, a, s')}{\mathbf{q}^{\pi, p'}(h, s, a)}.$$

The set of all occupancy measures with respect to an MDP  $\mathcal{M}$  is denoted by  $\Delta(\mathcal{M})$ , and the set of all policies by  $\Pi = \{\pi : [H] \times \mathcal{S} \rightarrow \Delta_{\mathcal{A}}\}$ . Importantly, the value of a policy from the initial state (i.e., the expected loss in an episode) can be written as the dot product between its occupancy measure and the cost function, i.e.,  $\langle \mathbf{q}^{\pi, p'}, \ell \rangle = \sum_{h,s,a} \mathbf{q}^{\pi, p'}(h, s, a) \ell(h, s, a)$ . Thus, the regret becomes

$$\mathcal{R}_T = \sum_{t=1}^T \langle \mathbf{q}^{\pi_t, p}, \ell_t \rangle - \min_{\mathbf{q} \in \Delta(\mathcal{M})} \sum_{t=1}^T \langle \mathbf{q}, \ell_t \rangle.$$

Whenever  $p'$  is omitted from the notation  $\mathbf{q}^{\pi, p'}$ , it is understood to be the true transition function  $p$ .

With that in hand, the adversarial MDP setting is an instance of the online learning framework where  $\ell_t \in [0, 1]^{HS^2A}$ ,  $\mathcal{A} = \Delta(\mathcal{M})$  as the set of all occupancy measures and the feedback  $\mathcal{L}(\mathbf{w}^{\pi_\tau}, \ell_\tau)$  is the loss over the trajectory  $\{\ell_\tau(h, s_{\tau,h}, a_{\tau,h})\}_{h=1}^H$ .  $\mathcal{W}$  is a (slightly modified) set of occupancy measures which we will define later. Note that in this context,  $\mathbf{w}_t(\mathbf{L})$  is a vector of dimension  $HS^2A$ —we will denote by  $\mathbf{w}_t(\mathbf{L}, h, s, a, s')$  the  $(h, s, a, s')$  element of it and also define  $\mathbf{w}_t(\mathbf{L}, h, s, a) = \sum_{s'} \mathbf{w}_t(\mathbf{L}, h, s, a, s')$ .

**Algorithm** Algorithm 3 is based on the general framework presented in Section 3. To satisfy the stability conditions required for Lemma 3, we employ a hybrid regularization of negative entropy and log-barrier just like in the combinatorial bandit case:

$$R_t(\mathbf{v}) = \frac{1}{\eta_t} \sum_{h,s,a,s'} \mathbf{v}(h, s, a, s') \log \mathbf{v}(h, s, a, s') - \frac{1}{\gamma} \sum_{h,s,a,s'} \log \mathbf{v}(h, s, a, s'). \quad (12)$$

---

**Algorithm 3:** Delayed FTRL for adversarial MDPs
 

---

**Input:** Regularizers  $\{R_t\}_{t \geq 1}$  defined in (12).  
**for**  $t = 1, \dots, T$  **do**  
     Observe feedback  $\ell_t(h, s_{\tau,h}, a_{\tau,h})$  for  $h \in [H], \tau \in o_t \setminus o_{t-1}$ .  
     Compute upper occupancy bounds  $\mathbf{q}_\tau^{\max}(h, s, a) = \max_{\hat{p} \in \mathcal{P}} \mathbf{q}^{\pi_\tau, \hat{p}}(h, s, a)$ .  
     Compute  $\hat{\ell}_\tau(h, s, a) = \frac{\mathbb{I}\{s_{\tau,h}=s, a_{\tau,h}=a\} \ell_\tau(h, s, a)}{\mathbf{q}_\tau^{\max}(h, s, a)}$  for  $\tau \in o_t \setminus o_{t-1}$ .  
     Compute  $\mathbf{w}_t(\hat{\mathbf{L}}_t) = \arg \min_{\mathbf{v} \in \mathcal{W}} \hat{\mathbf{L}}_t^\top \mathbf{v} + R_t(\mathbf{v})$  and policy  $\pi_t(a \mid h, s) = \frac{\mathbf{w}_t(\hat{\mathbf{L}}_t, h, s, a)}{\mathbf{w}_t(\hat{\mathbf{L}}_t, h, s)}$ .  
     Play episode  $t$  with policy  $\pi_t$   
**end for**

---

The main difference is that some of the elements of the occupancy measures may be 0 regardless of the chosen policy (if  $p(s' \mid h, s, a) = 0$  then  $\mathbf{q}^\pi(h, s, a, s') = 0$ ), in which case the regularization is not well-defined. To avoid that, we augment the set of occupancy measures to include occupancy measures for which the associated transition probability differs a little bit from the true transition probabilities

$$\Delta(\mathcal{P}) = \{\mathbf{q}^{\pi, \hat{p}} : \pi \in \Pi, \hat{p} \in \mathcal{P}\} \quad \text{where} \quad \mathcal{P} = \left\{ \hat{p} : \|\hat{p} - p\|_\infty \leq \frac{1}{THSA} \right\}.$$

To complete the presentation of adversarial MDPs as an instance of our online learning framework, we define the constant domain as  $\mathcal{W} = \Delta(\mathcal{P})$ . Also, we are not skipping any rounds. This construction allows us to establish the following properties of  $\mathcal{W}$ :

**Lemma 11**  *$\mathcal{W}$  satisfies the following:*

1. For any  $\mathbf{q} \in \Delta(\mathcal{M})$ , there exists  $\tilde{\mathbf{q}} \in \mathcal{W}$  such that  $\min_{h,s,a,s'} \tilde{\mathbf{q}}(h, s, a, s') \geq \frac{1}{T^3 H^2 S^4 A^2}$  and  $\|\mathbf{q} - \tilde{\mathbf{q}}\|_1 \leq \frac{2H}{T}$ .
2. Given  $\mathbf{v} \in \mathcal{W}$ , let  $\pi$  be defined by  $\pi(a \mid h, s) = \frac{\mathbf{v}(h,s,a)}{\mathbf{v}(h,s)}$  and  $\mathbf{q}^{\max}(h, s, a) = \max_{\hat{p} \in \mathcal{P}} \mathbf{q}^{\pi, \hat{p}}(h, s, a)$ .  
 Then,  $\|\mathbf{q}^\pi - \mathbf{v}\|_1 \leq \frac{2H}{T}$  and  $\|\mathbf{q}^{\max} - \mathbf{v}\|_1 \leq \frac{4H^2 S}{T}$ .

The proof can be found in Appendix C. The importance-weighted loss estimator for Algorithm 3 is inspired by Jin et al. (2020),

$$\hat{\ell}_\tau(h, s, a) = \frac{\mathbb{I}\{s_{\tau,h}=s, a_{\tau,h}=a\} \ell_\tau(h, s, a)}{\mathbf{q}_\tau^{\max}(h, s, a)},$$

where  $\mathbf{q}_\tau^{\max}(h, s, a) = \max_{\hat{p} \in \mathcal{P}} \mathbf{q}^{\pi_\tau, \hat{p}}(h, s, a)$  is an upper bound on the occupancy measure for each state  $h, s, a$  when following policy  $\pi_\tau$ . That means that  $\hat{\ell}_\tau$  is underestimating the actual losses and is a slightly biased estimator.

Note that  $\mathcal{W}$  is a convex set defined by  $O(HS^2A)$  linear equality and inequality constraints. In practice, we can eliminate the equality constraints through a simple re-parameterization, ensuring the variables lie within the linear subspace that satisfies the constraints (Boyd and Vandenberghe, 2004), thereby making the interior of the decision set non-empty. Using that, we can

apply the interior-point method to approximate the solution to the FTRL step with running time  $O(\text{poly}(H, S, A) \log T)$ —Nemirovski (2004); see also Abernethy et al. (2012)—with an error up to  $1/T$  (which affects the regret only by a constant). In addition,  $\mathbf{q}_t^{\max}$  can be computed efficiently as well using dynamic programming (Jin et al., 2020). We note that, while this approach is technically efficient, it becomes impractical when the number of states is significantly large.

**Main Result and Discussion** The main result of this section is Theorem 12.

**Theorem 12** *Suppose that  $T \geq H$ . Algorithm 3 with*

$$\gamma = \frac{1}{128Hd_{\max}} \quad \eta_t = \min \left\{ \frac{\log(SA)}{96HSA\sqrt{SA}d_{\max} + d_{\max}^2}, \frac{\sqrt{\log(SA)}}{\sqrt{SA}t + \sum_{\tau=1}^t |m_t|} \right\}$$

*guarantees*

$$\mathbb{E}[\mathcal{R}_T] \leq 72H\sqrt{\log(SA)(TSA + D)} + 1338d_{\max}H^2S^2A^2\log(HSAT) .$$

The proof relies on yet another regret decomposition given by

$$\mathcal{R}_T = \sum_{t=1}^T \langle \mathbf{q}^{\pi_t} - \mathbf{u}, \ell_t \rangle = \underbrace{\sum_{t=1}^T \langle \mathbf{q}^{\pi_t} - \mathbf{w}_t(\hat{\mathbf{L}}_t), \ell_t \rangle}_{\text{ERROR}} + \underbrace{\sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle}_{\text{REG}} + \underbrace{\sum_{t=1}^T \langle \tilde{\mathbf{u}} - \mathbf{u}, \ell_t \rangle}_{\text{SHIFT-PENALTY}} .$$

$\Delta(\mathcal{P})$  is only slightly larger than  $\Delta(\mathcal{M})$ , and we can easily bound ERROR using the first property in Lemma 11. Since  $R_T(\mathbf{u})$  can be arbitrarily large near the boundary of the domain, we slightly shift  $\mathbf{u}$  to  $\tilde{\mathbf{u}}$  using the first property in Lemma 11 to ensure that (i)  $R_T(\tilde{\mathbf{u}}) \leq \tilde{O}(\frac{HS^2A}{\gamma})$ , and (ii) SHIFT-PENALTY is bounded by  $2H$ . We can not apply Corollary 4 to bound REG because of the bias in our estimator. We apply Lemma 3 instead. By Lemma 25 we can show that  $R_t$  satisfies Assumption 2 and that  $R_T(\tilde{\mathbf{u}}) - \min_{\mathbf{v} \in \mathcal{W}_1} R_1(\mathbf{v})$  is bounded by  $\tilde{O}(\frac{H}{\eta_T} + \frac{HS^2A}{\gamma})$ .

The fact that  $\mathbf{q}_t^{\max}(h, s, a)$  upper bounds both  $\mathbf{w}_t(\hat{\mathbf{L}}_t, h, s, a)$  and  $\mathbf{q}^{\pi_t}(h, s, a)$  allows us to keep local norms related to  $\alpha_t$  and  $\beta_t$  small. In addition, using the second property in Lemma 11, we can also show that the estimator’s bias is only of order  $1/T$  (ignoring  $S, H$  factors). The main part of the remaining of the proof deals with the term  $\|\sum_{\tau \in m_t} (\ell_\tau - \hat{\ell}_\tau)\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)}^2$ . This term, which arises because we are using biased estimators, is not present when applying Corollary 4. The full proof can be found in Appendix C.

The algorithm is optimal, matching the lower bound of Lancewicki et al. (2022a) up to log-factors and improves on previous state-of-the-art regret bounds  $\tilde{O}(H^2S\sqrt{AT} + H(HSA)^{1/4}\sqrt{D})$  by Jin et al. (2022) and  $\tilde{O}(H^2\sqrt{SAT} + H^3\sqrt{D})$  by Lancewicki et al. (2023).

## 7. Adversarial MDPs with Unknown Transitions

In this section, we apply our FTRL approach to adversarial Markov Decision Processes (MDPs) setting detailed in Section 6, for the case that the transition function is unknown to the learner in advance. We show that it yields the first algorithm that handles delay asymptotically optimal in this setting, up to sub-optimality gaps that already exist in the non-delayed setting.

---

**Algorithm 4:** Delayed FTRL for adversarial MDPs with unknown transitions
 

---

**Initialize**  $j = 1$ ,  $\widehat{\mathcal{P}}_1$  as the set of all transition functions,  $\mathcal{W}_0 = \Delta(\widehat{\mathcal{P}}_1)$ .  
 For all  $(h, s, a, s') \in [H] \times \mathcal{S} \times \mathfrak{A} \times \mathcal{S}$  set  $N_0(s'|h, s, a) = N_1(s'|h, s, a) = 0$ .  
**for**  $t = 1, \dots, T$  **do**  
     /\* Transition estimation and epochs \*/  
     Observe trajectories  $(s_{\tau,h}, a_{\tau,h})$  for  $h \in [H], \tau \in \tilde{o}_t \setminus \tilde{o}_{t-1}$ .  
     Update counters:  $N_j(s_{\tau,h+1}|h, s_{\tau,h}, a_{\tau,h}) += 1$  for  $h \in [H], \tau \in \tilde{o}_t \setminus \tilde{o}_{t-1}$ .  
     **if**  $\exists h$  such that  $N_j(h, s_{\tau,h}, a_{\tau,h}) \geq \max\{1, 2N_{j-1}(h, s_{\tau,h}, a_{\tau,h})\}$  **then**  
          $j += 1$   
         For all  $(h, s, a, s') \in \mathcal{S} \times \mathfrak{A} \times \mathcal{S}$ , set  $N_j(s'|h, s, a) = N_{j-1}(s'|h, s, a)$ .  
         Update set  $\widehat{\mathcal{P}}_j$  as in equation (13).  
         Set  $\mathcal{W}_t = \bigcap_{j'=1}^j \Delta(\widehat{\mathcal{P}}_{j'})$ . If  $\mathcal{W}_t = \emptyset$  then set  $\mathcal{W}_t = \Delta(\widehat{\mathcal{P}}_j)$ .  
         Skip all rounds that are missing by adding all elements in  $\tilde{m}_t$  to  $\Lambda$ .  
     **end if**  
     /\* Loss estimation and episode execution \*/  
     If  $\mathcal{W}_t$  is not defined by an epoch change, set  $\mathcal{W}_t = \mathcal{W}_{t-1}$ .  
     Observe feedback  $\ell_t(h, s_{\tau,h}, a_{\tau,h})$  for  $h \in [H], \tau \in o_t \setminus o_{t-1}$ .  
     Compute upper occupancy bounds  $\mathbf{q}_{\tau}^{\max}(h, s, a) = \max_{\hat{p} \in \widehat{\mathcal{P}}_j} \mathbf{q}^{\pi_{\tau}, \hat{p}}(h, s, a)$ .  
     Compute  $\widehat{\ell}_{\tau}(h, s, a) = \frac{\mathbb{I}\{s_{\tau,h}=s, a_{\tau,h}=a\} \ell_{\tau}(h, s, a)}{\mathbf{q}_{\tau}^{\max}(h, s, a) + \xi}$  for new observations  $\tau \in o_t \setminus o_{t-1}$ .  
     Compute  $\mathbf{w}_t(\widehat{\mathbf{L}}_t) = \arg \min_{\mathbf{v} \in \mathcal{W}_t} \widehat{\mathbf{L}}_t^{\top} \mathbf{v} + R_t(\mathbf{v})$  and policy  $\pi_t(a | h, s) = \frac{\mathbf{w}_t(\widehat{\mathbf{L}}_t, h, s, a)}{\mathbf{w}_t(\widehat{\mathbf{L}}_t, h, s)}$ .  
     Play episode  $t$  with policy  $\pi_t$   
**end for**

---

**Algorithm** Algorithm 4 is very similar to the one presented in Section 6 for the known transitions case, with one main differences: In order to estimate the transition function we use a delayed version of the confidence set for the transition function of Jin et al. (2020). The confidence sets are updated in epochs. Specifically, the algorithm maintains counters  $N_j(s' | h, s, a)$  to track the number of visits to state-action pair  $(s, a)$  and transitioning to state  $s'$  at time  $h$ . Mirroring the notation used for occupancy measures we also define  $N_j(h, s, a) = \sum_{s'} N_j(s' | h, s, a)$  as the number of visits to state-action pair  $(s, a)$  at time  $h$ . In each epoch  $j$ , if the counter  $N_j(h, s, a)$  doubles compared to  $N_{j-1}(h, s, a)$  for some triplet  $(h, s, a)$ , a new epoch (epoch  $j + 1$ ) is initiated. The confidence set in epoch  $j$  is defined as

$$\widehat{\mathcal{P}}_j = \left\{ \widehat{p} : \left| \widehat{p}(s'|h, s, a) - \bar{p}_j(s'|h, s, a) \right| \leq \epsilon_j(s'|h, s, a), \forall (h, s', a, s) \in [H] \times \mathcal{S} \times \mathfrak{A} \times \mathcal{S} \right\}, \quad (13)$$

where

$$\epsilon_j(s'|h, s, a) = 2 \sqrt{\frac{\bar{p}_j(s'|h, s, a) \log(HSAT^3)}{\max\{1, N_j(h, s, a) - 1\}}} + \frac{14 \log(HSAT^3)}{3 \max\{1, N_j(h, s, a) - 1\}},$$

for  $\bar{p}_j(s'|h, s, a) = \frac{N_j(s'|h, s, a)}{N_j(h, s, a)}$  being the empirical transition, calculated using the visit counts  $N_j(s'|h, s, a)$  at the beginning of the epoch. The domain is constant in each epoch and is computed as the intersection over all previous  $\Delta(\widehat{\mathcal{P}})$ . That is, if round  $t$  is in epoch  $j$ , then  $\mathcal{W}_t =$

$\bigcap_{j'=1}^j \Delta(\widehat{\mathcal{P}}_{j'})$ . Lemma 13 below shows that the true transition function is contained in our confidence set with high probability.

**Lemma 13** *With probability at least  $1 - 4/T^2$ , we have  $p \in \widehat{\mathcal{P}}_j$  for all  $j$ .*

**Proof** The proof is a straightforward modification of the proof of Lemma 2 of Jin et al. (2020). ■

As a consequence of Lemma 13, the occupancy measure of the benchmark policy is contained in each domain. The reason that we define  $\mathcal{W}_t$  as the intersection of  $\Delta(\mathcal{P}_{j'})$  up to the current epoch is to ensure that  $\mathcal{W}_t \subseteq \mathcal{W}_{t+1}$ . This will later be crucial in the analysis to apply the be-the-leader lemma (Lemma 18 in Appendix F). In order to ensure that Assumption 1, specifically the fact that  $\mathcal{W}_t = \mathcal{W}_\tau$  for all outstanding observations  $\tau \in m_t$ , is met, we skip all outstanding rounds at the beginning of a new epoch. The loss estimator is an importance-weighted estimator with  $\mathbf{q}_\tau^{\max}(h, s, a)$  being an upper confidence estimate for  $\mathbf{q}^{\pi, p}(h, s, a)$ . In addition we add a small bias  $\xi$ , so that the estimator is also bounded under the bad event.

In terms of implementation, we can take the same approach as in the known transition case (Section 6), with the main difference being that the decision set is defined by  $O(H^2 S^3 A^2 \log T)$  linear constraints due to the number of epochs being at most  $HSA \log T$ . Thus, the FTRL solution can be  $1/T$ -approximated with a running time of  $O(\text{poly}(H, S, A, \log T))$ . As noted before, while this approach is technically efficient, it becomes impractical when the number of states is large.

**Main Result and Discussion** The main result of this section is Theorem 14. To slightly simplify the analysis, we choose  $\eta_t = \eta$ . However, a decreasing learning rate is also possible, as shown for MDPs with known transitions in section 6.

**Theorem 14** *Algorithm 4 with  $\gamma = \frac{1}{128\sqrt{H}d_{\max}}$ ,  $\eta = \frac{\sqrt{\log(SA)}}{\sqrt{SAT+D}}$ ,  $\xi = \frac{1}{T}$  and  $T \geq 4$  guarantees,*

$$\begin{aligned} \mathbb{E}[\mathcal{R}_T] &\lesssim H^2 S \sqrt{AT \log(HSAT)} + H \sqrt{D \log(SA)} \\ &\quad + H^3 S^2 A \log(HSAT) d_{\max} + H^3 S^3 A \log^2(HSAT). \end{aligned}$$

The proof relies on the same regret decomposition as in Section 6.

$$\mathcal{R}_T = \mathbb{E} \left[ \underbrace{\sum_{t=1}^T \langle \mathbf{q}^{\pi_t} - \mathbf{w}_t(\widehat{\mathbf{L}}_t), \ell_t \rangle}_{\text{ERROR}} + \underbrace{\sum_{t=1}^T \langle \mathbf{w}_t(\widehat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle}_{\text{REG}} + \underbrace{\sum_{t=1}^T \langle \tilde{\mathbf{u}} - \mathbf{u}, \ell_t \rangle}_{\text{SHIFT-PENALTY}} \right].$$

The ERROR and SHIFT-PENALTY terms are bounded using standard arguments (see Lemma 29 in the appendix). To bound the REG term we will apply Lemma 3 just as in the previous sections. Since we now have a changing domain we need to ensure that all loss estimators that we observe in round  $t$ , that is where  $\tau \in o_t \setminus o_{t-1}$ , use the same domain as round  $t$  does. Since our domains are constant within any epoch, we will simply skip outstanding observations at the start of each epoch, skipping at most  $d_{\max}$  rounds whenever we change epochs. From here applying Lemma 3, bounding the bias of the estimator and bounding the total regret of skipped rounds by  $H \cdot d_{\max} \cdot HSA \log(T)$  yields the desired result. The detailed proof can be found in Appendix D.

The first term in the regret matches the best known regret for adversarial MDPs even without delayed feedback (Jin et al., 2020). The second term matches the lower bound of Lancewicki et al. (2022a) up to logarithmic terms. This improves over the previous state-of-the-art regret bounds  $\tilde{O}(H\sqrt{SAT} + H(HSA)^{1/4}\sqrt{D})$  by Jin et al. (2022) and  $\tilde{O}(H^3 S \sqrt{AT} + H^3 \sqrt{D})$  by Lancewicki et al. (2023).

## 8. Experiments

In this section we are evaluating the performance of Algorithm 1 and Algorithm 2 on synthetic experiments. The full code for the experiments can be found here<sup>1</sup>

### 8.1 Experiments for combinatorial bandits

For the combinatorial bandit setting we split the time horizon of  $T = 10000$  rounds into  $b$  blocks of length  $J$  and the algorithm only receives the feedback for all rounds in a block at the end of that block. As actions we use  $m$ -sets with  $m = 3$  and  $K = 10$ , the losses of in dimension  $i$  are either fixed or oscillating. The fixed arms are always 0, the oscillating arms have a constant loss of  $-1$  in block  $j$  if  $j$  is even and  $0.9$  otherwise. In other words, the oscillating arms are good arms and the constants are the bad arms. We use  $m = 3$  oscillating and  $7$  fixed arms. As mentioned earlier, Algorithm 1 is the first algorithm for delayed adversarial combinatorial bandits. We will therefore compare Algorithm 1 to a standard algorithm not adapted to delay. Namely, an FTRL based version of the algorithm presented in Figure 3 of Audibert et al. (2014), which is the same as running Algorithm 1 with regularizer  $R_t(\mathbf{v}) = \sum_{i=1}^K \frac{1}{\eta_t} \mathbf{v}(i) \log(\mathbf{v}(i))$  and learning rate  $\eta_t = \sqrt{\frac{B(1+\log(\frac{K}{B}))}{16Kt}}$ .

The results of experiments for varying block sizes can be found in Figure 1. Dropping all other dependencies, and assuming a constant delay of  $d_t = d$ , our analysis finds that  $\mathcal{R}_T \lesssim \frac{1}{\eta} + \eta T + \eta d T$ . The delay unaware algorithms tunes  $\eta \approx \frac{1}{\sqrt{T}}$ , leading to a regret bound of  $\mathcal{R}_T \lesssim \sqrt{T} + d\sqrt{T}$ , which matches the roughly linear dependency on delay which we observe for the delay unaware algorithm. When the block size is  $b = 1$ , there is no delay present and the delay unaware method outperforms our algorithm slightly as we are over-regularising. But even for small delays, the delay aware tuning outperforms the non-delayed tuning significantly.

### 8.2 Experiments for linear bandits

In this section we present synthetic experiments for the linear bandit setting. The losses are generated using the same block structure as for the experiments for combinatorial semi-bandits, where the algorithm only observes feedback at the end of the block. There are  $T = 10000$  rounds split into blocks, where the block size is  $J \in \{30, 60, 90, 120, 150\}$ . In each block the losses are either  $(1/\sqrt{K}, \dots, 1/\sqrt{K})$  or  $(-1/\sqrt{K}, \dots, -1/\sqrt{K})$ . As in the combinatorial semi-bandit setting, the sign of the losses alternates between blocks. The dimension is varied between experiments, with  $K \in [10, 20, 40]$ . We implemented Algorithm 2, Algorithm 5, and Banker-BOLO Huang et al. (2023) with  $-\log(1 - \|\mathbf{x}\|_2^2)$  as the 1-self-concordant barrier for the unit ball. We also implemented a version of Banker-BOLO with what we believe to be improved tuning as described in Section 5. This version of Banker-BOLO is denoted by Banker-BOLO-V2. A fifth possible algorithm to compare with is the algorithm of Ito et al. (2020a). This is an instance of continuous exponential weights, which means its computational complexity is  $O(\text{poly}(K, T))$ . However, the degree of this polynomial is high, which means that running this algorithm is infeasible for us.

The results can be found in Figure 2 in the main body, and Figures 3 and 4 in Appendix G. As predicted by theory, the regret grows with the square root of the block size for all algorithms.

1. <https://github.com/LukasZierahn/A-Unified-Analysis-of-Nonstochastic-Delayed-Feedback.git>



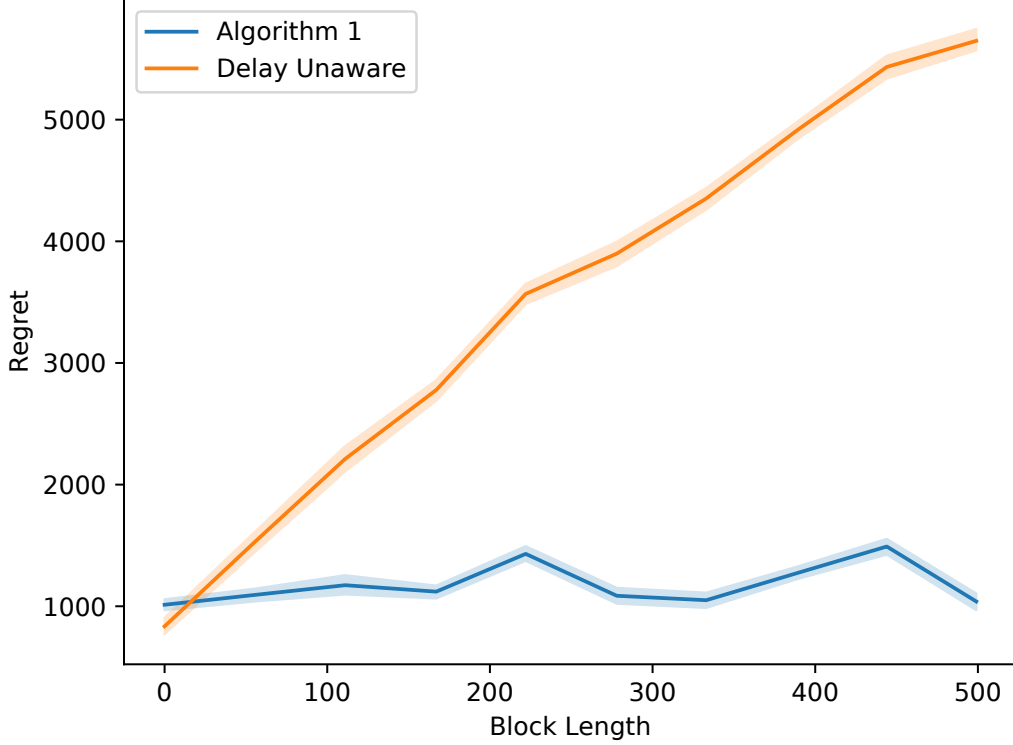


Figure 1: The results for our combinatorial semi-bandits experiments. The solid line is the mean regret in the 20 repetitions over  $T = 10000$  rounds. The shaded areas are 95% confidence intervals.

However, it seems that the  $K\sqrt{\nu}\sqrt{D}$  and  $K^{3/2}\sqrt{D}$  terms in the regret bounds of Banker-BOLO-V2 and Banker-BOLO could possibly improved, as we do not see the difference in the regret between our algorithms and the Banker-BOLO algorithm increase as the dimension increases. This is to be expected, as one can probably derive a similar decomposition of the regret for OMD, upon which Banker-BOLO is based, as we did for FTRL. As with FTRL, this would most likely lead to  $\sqrt{D}$  term in the regret bound for the cost of delay, given that the algorithm is appropriately tuned. We see that our algorithms consistently outperform both versions of Banker-BOLO. However, we believe that with the right tuning a version of OMD will perform similarly to our algorithms. It can be seen that Banker-BOLO-V2 has better performance than Banker-BOLO, which is predicted by theory. We also see that the performance of Algorithms 2 and 5 hardly differs, which is also predicted by theory.

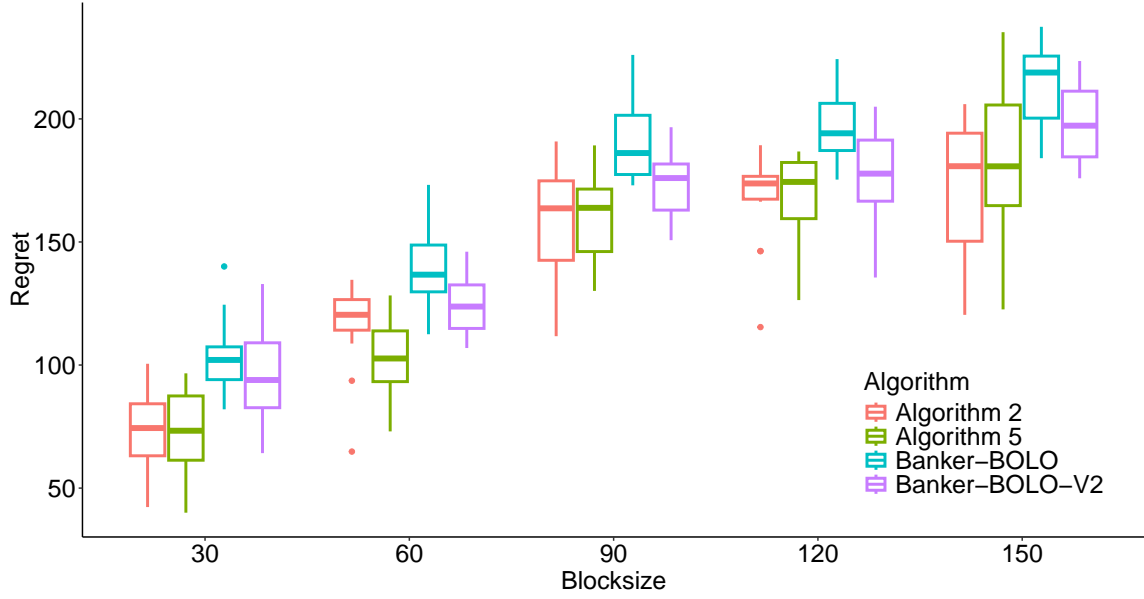


Figure 2: Boxplot of the regret in the linear bandit experiments with 20 repetitions over  $T = 10000$  rounds with  $K = 20$ .

## 9. Conclusion

In sections 4, 6, and 7 we have shown that FTRL leads to optimal regret bounds under delayed feedback in combinatorial semi-bandits, MDPs with known transitions. For MDPs with unknown transitions we provide state-of-the-art results. Furthermore, in section 5 we have provided an efficient algorithm with nearly optimal regret for linear bandits. In section 8 we have shown that Algorithm 1 and Algorithm 2 outperform delay-unaware and previous algorithms respectively on our synthetic datasets.

## Acknowledgments

The authors would like to thank András György for an insightful discussion about the regret decomposition and the suggestion to simplify the decomposition compared to the conference version of this paper. This work was done while LZ was at the Università degli Studi di Milano, Italy and partially done while DvdH was at the University of Milan supported by the MIUR PRIN grant Algorithms, Games, and Digital Markets (ALGADIMAR) and partially done while DvdH was at the University of Amsterdam supported by Netherlands Organization for Scientific Research (NWO), grant number VI.Vidi.192.095. NCB and LZ acknowledge the financial support from: the MUR PRIN grant 2022EKNE5K (Learning in Markets and Society), the FAIR project, funded by the NextGenerationEU program within the PNRR-PE-AI scheme, the EU Horizon CL4-2022-HUMAN-02 research and innovation action under grant agreement 101120237, project ELIAS. TL is supported by the European Research Council (ERC) under the European Union’s Horizon 2020 research and

innovation program (grant agreement No. 882396), by the Israel Science Foundation (grant number 993/17), the Yandex Initiative for Machine Learning at Tel Aviv University and a grant from the Tel Aviv University Center for AI and Data Science (TAD).

## Appendix A. Combinatorial Bandits

In this appendix we proof the main results of Section 4.

**Theorem 8 (RESTATED)** *Suppose that  $\max_{\mathbf{a} \in \mathcal{A}} \|\mathbf{a}\|_1 \leq B$ . Algorithm 1 with*

$$\eta_t = \min \left\{ \sqrt{\frac{B(1 + \log(\frac{K}{B}))}{16(B \sum_{\tau=1}^t |m_t| + Kt)}}, \frac{B^2(1 + \log(\frac{K}{B}))}{128K(Bd_{\max} + K)} \right\}, \quad \gamma = \frac{1}{128\sqrt{B}d_{\max}},$$

*guarantees that*

$$\mathcal{R}_T \leq 12\sqrt{B \left(1 + \log\left(\frac{K}{B}\right)\right) (KT + BD) + 128K^2d_{\max} + 128\sqrt{B}d_{\max}K \log(T)}.$$

**Proof** We start by verifying the conditions of Corollary 4 for  $R_t$ . Because we are not skipping rounds and have a constant actionset of  $\mathcal{W} = \text{Conv}(\mathcal{A})$ , we have that Assumption 1 holds. Next, note that  $R_t$  as specified in (11) does not satisfy Assumption 2(c) because  $R_t(\mathbf{v}) \leq R_{t-1}(\mathbf{v})$ . However, by using regularizer  $\tilde{R}_t(\mathbf{v}) = R_t(\mathbf{v}) - \min_{\mathbf{v} \in \mathcal{W}} R_t(\mathbf{v})$ , and running the algorithm with this regularizer instead, we can see that Assumption 2(c) is satisfied and, crucially, the algorithm produces the same iterates. Note also that the gradients and Hessians of  $R_t$  and  $\tilde{R}_t$  are equivalent. We continue the analysis as if the algorithm is run with regularizer  $\tilde{R}_t$ . By Lemma 26

$$\sqrt{\eta_{t+\delta}} \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \leq \frac{\sqrt{B}d_{\max}}{8\sqrt{K}}.$$

Thus, using Lemma 25 and plugging in  $\gamma$  gives,

$$(\nabla R_t(\mathbf{v}) - \nabla R_t + \delta(\mathbf{v}))^\top \mathbf{y} \leq \sqrt{\gamma \frac{\sqrt{B}d_{\max}}{8}} \sqrt{\mathbf{y}^\top \nabla^2 R_t + \delta(\mathbf{v}) \mathbf{y}} \leq \frac{1}{32} \sqrt{\mathbf{y}^\top \nabla^2 R_t + \delta(\mathbf{v}) \mathbf{y}},$$

for all  $t, \delta \in [d_{\max}]$ ,  $\mathbf{v} \in \mathcal{W}$  and all  $\mathbf{y} \in \mathbb{R}^K$ , which verifies Assumption 2(b) for  $\kappa = \gamma$ . The fact that  $4\nabla^2 R_t(\mathbf{v}) \succeq \nabla^2 R_t(\mathbf{v}') \succeq \frac{1}{4}\nabla^2 R_t(\mathbf{v})$  for all  $\mathbf{v} \in \mathcal{W}$ ,  $\mathbf{v}' \in \mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2\sqrt{\kappa}})$ , and all  $t$  is also shown in Lemma 25, showing that Assumptions 2(a) holds for  $\kappa = \gamma$ .

As the next step, we bound  $\|\ell_\tau\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)}$ . We use  $\sum_{i=1}^K \mathbf{w}_t(\hat{\mathbf{L}}_t, i) \leq B$ ,  $\|\ell\|_\infty \leq 1$ , and  $\eta_t + \gamma \mathbf{w}_t(\hat{\mathbf{L}}_t, i) \geq \gamma \mathbf{w}_t(\hat{\mathbf{L}}_t, i)$  to show that

$$\|\ell_\tau\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} = \sqrt{\sum_{i=1}^K \ell_\tau(i)^2 \frac{\eta_t \gamma \mathbf{w}_t(\hat{\mathbf{L}}_t, i)^2}{\eta_t + \gamma \mathbf{w}_t(\hat{\mathbf{L}}_t, i)}} \leq \underbrace{\sqrt{\eta_t B}}_{\alpha_t}. \quad (14)$$

Next we bound  $\mathbb{E} \left[ \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}^2 \right]$ . By the tower rule we have

$$\mathbb{E} \left[ \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}^2 \right] = \mathbb{E}_{\mathcal{F}_\tau} \left[ \mathbb{E}_{\mathbf{a}_\tau} \left[ \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}^2 | \mathcal{F}_\tau \right] \right]$$

where  $\mathcal{F}_\tau$  is a filtration over all random events observed by the learner as defined in Section 2. Let us consider  $\mathbb{E}_{\mathbf{a}_\tau} \left[ \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}^2 | \mathcal{F}_\tau \right]$  in isolation:

$$\begin{aligned} \mathbb{E}_{\mathbf{a}_\tau} \left[ \|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}^2 | \mathcal{F}_\tau \right] &= \mathbb{E}_{\mathbf{a}_\tau} \left[ \sum_{i=1}^K \left( \frac{\mathbf{a}_\tau(i) \ell_\tau(i)}{\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, i)} \right)^2 \left( \nabla^2 R_t(\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)) \right)^{-1}(i, i) | \mathcal{F}_\tau \right] \\ &= \sum_{i=1}^K \frac{\ell_\tau(i)^2}{\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, i)} \frac{\eta_t \gamma \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, i)^2}{\eta_t + \gamma \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, i)} \leq \underbrace{\eta_t K}_{\beta_t^2}, \end{aligned} \quad (15)$$

where we used that  $\mathbf{a}_\tau^2 = \mathbf{a}_\tau$ ,  $\mathbb{E}_{\mathbf{a}_\tau}[\mathbf{a}_\tau] = \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)$ , and  $\eta_t + \gamma \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, i) \geq \gamma \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, i)$ . We now bound  $\kappa \|\hat{\ell}_t\|_{R_{t'}, \mathbf{w}_t(\hat{\mathbf{L}}_t)}$ :

$$\kappa \|\hat{\ell}_t\|_{R_{t'}, \mathbf{w}_t(\hat{\mathbf{L}}_t)} = \kappa \sqrt{\sum_{i=1}^K \left( \frac{\mathbf{a}_t(i) \ell_t(i)}{\mathbf{w}_t(\hat{\mathbf{L}}_t, i)} \right)^2 \frac{\eta_{t'} \gamma \mathbf{w}_t(\hat{\mathbf{L}}_t, i)^2}{\eta_{t'} + \gamma \mathbf{w}_t(\hat{\mathbf{L}}_t, i)}} \leq \sqrt{\kappa \gamma B} \leq \frac{1}{128 d_{\max}},$$

where we used that  $\kappa = \gamma = \frac{1}{128 \sqrt{B} d_{\max}}$ .

The last requirement is to show that  $\hat{\ell}_\tau$  and  $\hat{\ell}_{\tau'}$  are independent for all  $\tau, \tau' \in m_t$  where  $\tau' \neq \tau$ . Recall that  $\hat{\ell}_\tau(i) = \frac{\mathbf{a}_\tau(i) \ell_\tau(i)}{\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, i)}$ , for all  $i$ . Conditioned on the observed history  $\mathcal{F}_t$ , the only random element of  $\hat{\ell}_\tau$  is  $\mathbf{a}_\tau \sim \mathbf{p}_\tau$ . Since  $\hat{\ell}_{\tau'}$  can not have been used in round  $\tau$  to compute  $\mathbf{p}_\tau$  (and vice versa) we have that  $\hat{\ell}_{\tau'}$  and  $\hat{\ell}_\tau$  are independent. We conclude that

$$\mathbb{E} \left[ (\hat{\ell}_\tau - \ell_\tau)^\top (\nabla^2 R_t(\mathbf{w}_t(\hat{\mathbf{L}}_t)))^{-1} (\hat{\ell}_{\tau'} - \ell_{\tau'}) \middle| \mathcal{F}_t \right] = 0,$$

where we used that  $\hat{\ell}_{\tau'}$  is an unbiased estimator of  $\ell_{\tau'}$ .

We are now in a position to apply Corollary 4. By using  $\alpha_t$  from equation (14) and  $\beta_t^2$  from equation (15) it follows that for any  $\mathbf{u} \in \mathcal{W}$

$$\mathbb{E} \left[ \sum_{t=1}^T (\mathbf{w}_t - \mathbf{u})^\top \ell_t \right] \leq R_T(\mathbf{u}) + 16 \sum_{t=1}^T \beta_t^2 + 16 \sum_{t=1}^T \alpha_t^2 |m_t|. \quad (16)$$

Next we want to bound  $R_T(\mathbf{u})$ , however the negative logarithm component is unbounded and tends to infinity when any element of  $\mathbf{u}$  tends to 0. Thus we cannot compare to  $\mathbf{a}^*$  directly, which might lie on the boundary of  $\mathcal{W}$  and instead we will compare to  $\mathbf{u} = \arg \min_{\mathbf{v} \in \widetilde{\mathcal{W}}} \sum_{t=1}^T \ell_t^\top \mathbf{v}$  where  $\widetilde{\mathcal{W}} = \mathcal{W} \cap \{\mathbf{x} \in \mathbb{R}_+ : \forall i \in [K] \ x(i) \geq \theta\}$  is a shrunken actionset.  $\theta$  now acts as a trade-off between an upper bound on the regularizer and an additional bias-like term that stems from comparing  $\mathbf{a}^*$  to  $\mathbf{u}$  in terms of pseudo-regret. More specifically we can write  $\mathbf{a}^* = \mathbf{u} + \theta \xi$ , for some  $\xi$  with  $\|\xi\|_\infty \leq 1$ .

By Lemma 25, we have

$$R_T(\mathbf{u}) \leq \frac{B \left( 1 + \log \left( \frac{K}{B} \right) \right)}{\eta T} + \frac{K \log \left( \frac{1}{\theta} \right)}{\gamma}. \quad (17)$$

To finish the proof, we start from the regret

$$\begin{aligned}\mathcal{R}_T &= \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{a}^*)^\top \boldsymbol{\ell}_t \right] = \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{u})^\top \boldsymbol{\ell}_t \right] + \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{u} - \mathbf{a}^*)^\top \boldsymbol{\ell}_t \right] \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{u})^\top \boldsymbol{\ell}_t \right] + \mathbb{E} \left[ \sum_{t=1}^T \theta \xi^\top \boldsymbol{\ell}_t \right] \leq \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{u})^\top \boldsymbol{\ell}_t \right] + \theta KT,\end{aligned}$$

where we bound  $\xi^\top \boldsymbol{\ell}_t \leq K$  in the inequality. We continue by using (16)

$$\begin{aligned}\mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{u})^\top \boldsymbol{\ell}_t \right] &\leq R_T(\mathbf{u}) + 16 \sum_{t=1}^T \kappa \beta_t^2 + 16 \sum_{t=1}^T \kappa \alpha_t^2 |m_t| \\ &\leq \frac{B(1 + \log(\frac{K}{B}))}{\eta_T} + \frac{K \log(T)}{\gamma} + 16K \sum_{t=1}^T \eta_t + 16B \sum_{t=1}^T \eta_t |m_t| \\ &\leq \frac{B(1 + \log(\frac{K}{B}))}{\eta_T} + \frac{K \log(T)}{\gamma} + \\ &\quad + 8 \sqrt{B \left( 1 + \log \left( \frac{K}{B} \right) \right) (KT + BD)},\end{aligned}$$

where we first substituted in (17), with  $\theta = \frac{1}{T}$ ,  $\alpha_t^2 = \eta_t B$  and  $\beta_t^2 = \eta_t K$ , and applied Lemma 24. Using the last two inequalities and substituting

$$\begin{aligned}\eta_t &= \min \left\{ \sqrt{\frac{B(1 + \log(\frac{K}{B}))}{16(B \sum_{\tau=1}^t |m_\tau| + Kt)}}, \frac{B^2(1 + \log(\frac{K}{B}))}{128K(Bd_{\max} + K)} \right\} \\ \gamma &= \frac{1}{128\sqrt{B}d_{\max}}\end{aligned}$$

and doing some simplifications yields

$$\mathcal{R}_T \leq 12 \sqrt{B \left( 1 + \log \left( \frac{K}{B} \right) \right) (KT + BD)} + 128K^2d_{\max} + 128\sqrt{B}d_{\max}K \log(T)$$

concluding the proof. ■

We now state a lower bound for the delayed combinatorial semi-bandit setting. This implies that, ignoring terms that are logarithmic in  $T$ , the result of Theorem 8 is optimal. The proof of our lower bound follows from standard arguments in the delayed bandit feedback literature.

**Theorem 9 (RESTATED)** *Suppose that  $d_t = d$  for all  $t$  and that  $B \leq K/2$ . Then for any algorithm there exists a sequence of losses such that*

$$\mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{a}^*)^\top \boldsymbol{\ell}_t \right] = \Omega \left( \max \left\{ \sqrt{BKT}, B\sqrt{dT} \right\} \right).$$

**Proof** By Audibert et al. (2014), we have that any algorithm without delay must suffer at least  $\Omega(\sqrt{BKT})$  regret in the combinatorial semi-bandit setting. Next, we assume full information feedback, which is easier from the point of view of the algorithm. We take inspiration from Langford et al. (2009, Lemma 3). For simplicity we will assume that  $T/d$  is an integer. We divide the  $T$  rounds into  $T/d$  blocks of  $d$  rounds. We take the losses of the lower bound for  $B$ -sets in (Koolen et al., 2010, Section 4), which states that any algorithm in the full information setting must suffer at least  $\Omega(B\sqrt{T'})$  regret after  $T'$  rounds. We take the loss of the first round of the lower bound (Koolen et al., 2010) and copy it  $d$  times, which we use as the losses for the first block. We repeat this process for the remaining blocks. Since the algorithm can not respond to the copied losses, we must have that any algorithm must suffer at least  $\Omega(dB\sqrt{T/d}) = \Omega(B\sqrt{dT})$  regret, which completes the proof.  $\blacksquare$

## Appendix B. Linear Bandits

Recall that a thrice-differentiable function  $\Psi$  is called self-concordant if it is convex and satisfies  $|\nabla^3\Psi(\mathbf{v})[\mathbf{h}, \mathbf{h}, \mathbf{h}]| \leq 2(\nabla^2\Psi(\mathbf{v})[\mathbf{h}, \mathbf{h}])^{3/2}$ , where  $\nabla^3\Psi(\mathbf{v})[\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3] = \frac{\partial^3}{\partial t_1 \partial t_2 \partial t_3} \big|_{t_1=t_2=t_3=0} \Psi(\mathbf{v} + t_1\mathbf{h}_1 + t_2\mathbf{h}_2 + t_3\mathbf{h}_3)$ . A self-concordant function  $\Psi$  is a  $\nu$ -self-concordant barrier if  $|\nabla\Psi(\mathbf{v})[\mathbf{h}]| \leq \sqrt{\nu\nabla^2\Psi(\mathbf{v})[\mathbf{h}, \mathbf{h}]}$ . The following property allows us to satisfy the stability condition of the Hessian in Assumption 2(a): for  $\mathbf{v}, \mathbf{v}' \in \mathcal{W}$ , if  $\|\mathbf{v} - \mathbf{v}'\|_{\Psi, \mathbf{v}}^* < 1$ , then

$$(1 - \|\mathbf{v} - \mathbf{v}'\|_{\Psi, \mathbf{v}}^*)^2 \nabla^2\Psi(\mathbf{v}) \preceq \Psi(\mathbf{v}') \preceq (1 + \|\mathbf{v} - \mathbf{v}'\|_{\Psi, \mathbf{v}}^*)^2 \nabla^2\Psi(\mathbf{v}). \quad (18)$$

Next, given  $\mathbf{y} \in \mathcal{W}$  denote by  $\pi_{\mathbf{y}}(\mathbf{x}) = \inf\{z \geq 0 : \mathbf{y} + z^{-1}(\mathbf{x} - \mathbf{y}) \in \mathcal{W}\}$  the Minkowsky function. We denote by  $\mathcal{W}_\delta = \{\mathbf{v} : \pi_{\mathbf{v}^+}(\mathbf{v}) \leq (1 + \delta)^{-1}\}$ , where  $\mathbf{v}^+ = \arg \min_{\mathbf{v} \in \mathcal{W}} \Psi(\mathbf{v})$  and  $\delta > 0$ . If  $\Psi$  is a  $\nu$ -self-concordant barrier, then for any  $\mathbf{v} \in \mathcal{W}_\delta$

$$\Psi(\mathbf{v}) - \min_{\mathbf{v} \in \mathcal{W}} \Psi(\mathbf{v}) \leq \nu \ln((1 + \delta)\delta^{-1}). \quad (19)$$

This property allows us to show that for any benchmark point  $\tilde{\mathbf{u}} \in \mathcal{W}_\delta$ ,  $R_T(\mathbf{u})$  and is nicely bounded.

**Theorem 10 (RESTATED)** *Suppose that  $T > 100$  and  $B \geq 1$ . Algorithm 2, run with a  $\nu$ -self-concordant barrier  $\Psi$  and with*

$$\begin{aligned} \gamma_t &= \min \left\{ \frac{1}{256BKd_{\max}}, \sqrt{\frac{\nu \log(1 + \sqrt{T})}{16B^2K^2t}} \right\} \\ \eta_t &= \min \left\{ \frac{B}{256d_{\max}}, \sqrt{\frac{B^2}{16 \sum_{\tau=1}^t |m_t|}} \right\}, \end{aligned}$$

*guarantees that, for any  $\mathbf{u} \in \mathcal{W}$ ,*

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{u})^\top \ell_t \right] &\leq 12BK \sqrt{\nu T \log(1 + \sqrt{T})} + 12B\sqrt{D} + 2B\sqrt{T} \\ &\quad + 512BKd_{\max} \nu \log(1 + \sqrt{T}). \end{aligned}$$

**Proof** We start by verifying the assumptions of Corollary 4. Because we are not skipping rounds and have a constant actionset of  $\mathcal{W} = \text{Conv}(\mathcal{A})$ , we have that Assumption 1 holds. Using  $\mathbb{E}[\mathbf{v}_t] = \mathbf{0}$  and  $\mathbb{E}[\mathbf{v}_t \mathbf{v}_t^\top] = \frac{1}{K} \mathbf{I}$  we see that  $\mathbb{E}[\hat{\ell}_t] = \ell_t$ . For  $\tau \leq t$ , observe that the distribution of  $\hat{\ell}_{\tau'}$  is fully determined given  $\mathcal{F}_t$  because  $\mathcal{F}_{\tau'} \subseteq \mathcal{F}_t$ . Furthermore, since  $\hat{\ell}_\tau$  can not be used in round  $\tau'$  because  $\tau$  is not available in round  $t$  due to the delay, we must have that  $\hat{\ell}_{\tau'}$  is independent of  $\hat{\ell}_\tau$ . Thus,

$$\mathbb{E} \left[ (\hat{\ell}_\tau - \ell_\tau)^\top (\nabla^2 R_t(\mathbf{w}_t(\hat{\mathbf{L}}_t)))^{-1} (\hat{\ell}_{\tau'} - \ell_{\tau'}) \middle| \mathcal{F}_t \right] = 0,$$

where we used that  $\mathbb{E}[\hat{\ell}_{\tau'} | \mathcal{F}_t] = \mathbb{E}[\hat{\ell}_\tau | \mathcal{F}_t, \hat{\ell}_\tau] = \ell_t$ .

We now turn to verifying that Assumption 2 holds. Assumption 2(c) holds by definition of  $\eta_t$  and  $\gamma_t$ . Because  $\tilde{\Psi}$  is a self-concordant, if we choose  $\kappa = \frac{1}{256BKd_{\max}}$ ,  $\frac{\kappa}{\gamma_t} \tilde{\Psi}$  is also self-concordant as self-concordance is preserved by scaling of factors exceeding one. Since  $c\|\mathbf{v}\|_2^2$  is self-concordant on  $\mathbb{R}^d$  for any  $c > 0$  and adding two self-concordant barriers yields a self-concordant barrier,  $\kappa R_t$  is also a self-concordant barrier. If  $\mathbf{v}' \in \mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2\sqrt{\kappa}})$  it implies that  $\mathbf{v}' \in \mathcal{D}_{\kappa R_t}(\mathbf{v}, \frac{1}{2})$ . By equation (18), for  $\mathbf{v}' \in \mathcal{D}_{\kappa R_t}(\mathbf{v}, \frac{1}{2})$ , we have  $4\nabla^2 \kappa R_t(\mathbf{v}) \succeq \nabla^2 \kappa R_t(\mathbf{v}') \succeq \frac{1}{4} \nabla^2 \kappa R_t(\mathbf{v})$  or equivalently, for all  $\mathbf{v} \in \mathcal{W}$  and  $\mathbf{v}' \in \mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2\sqrt{\kappa}})$  we have  $4\nabla^2 R_t(\mathbf{v}) \succeq \nabla^2 R_t(\mathbf{v}') \succeq \frac{1}{4} \nabla^2 R_t(\mathbf{v})$ , which verifies Assumption 2(a).

The final condition to check is that  $\kappa (\nabla R_t(\mathbf{v}) - \nabla R_{t+\delta}(\mathbf{v}))^\top \mathbf{y} \leq \frac{1}{32} \sqrt{\kappa \mathbf{y}^\top \nabla^2 R_{t+\delta}(\mathbf{v}) \mathbf{y}}$ . Let  $\mathbf{v} \in \mathcal{W}$ ,  $\mathbf{y} \in \mathbb{R}^K$ , and  $\delta \in [d_{\max}]$ , then

$$\begin{aligned} & \left( \nabla R_t(\mathbf{v}) - \nabla R_{t+\delta}(\mathbf{v}) \right)^\top \mathbf{y} \\ &= \kappa \left( \frac{2}{\eta_t} \mathbf{v} + \frac{1}{\gamma_t} \nabla \Psi(\mathbf{v}) - \frac{2}{\eta_{t+\delta}} \mathbf{v} - \frac{1}{\gamma_{t+\delta}} \nabla \Psi(\mathbf{v}) \right)^\top \mathbf{y} \\ &= 2 \sum_{i=1}^K \kappa \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \mathbf{v}(i) \mathbf{y}(i) + \kappa \left( \frac{1}{\gamma_{t+\delta}} - \frac{1}{\gamma_t} \right) (\nabla \Psi(\mathbf{v}))^\top \mathbf{y}. \end{aligned}$$

By using the Cauchy-Schwarz inequality and the fact that  $\mathcal{W} \subseteq \mathcal{B}(B)$  we can see that

$$\begin{aligned} \sum_{i=1}^K \kappa \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \mathbf{v}(i) \mathbf{y}(i) &\leq \kappa \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) B \sqrt{\sum_{i=1}^K \mathbf{y}(i)^2} \\ &\leq \kappa \sqrt{16 \sum_{\tau=t}^{t+\delta} |m_\tau|} \sqrt{\sum_{i=1}^K \mathbf{y}(i)^2} \\ &\leq 4\kappa d_{\max} \sqrt{\sum_{i=1}^K \mathbf{y}(i)^2} \\ &= \sqrt{\kappa \eta_{t+\delta}} 4d_{\max} \sqrt{\frac{\kappa}{\eta_{t+\delta}} \sum_{i=1}^K \mathbf{y}(i)^2} \\ &\leq \frac{1}{64} \sqrt{\frac{\kappa}{\eta_{t+\delta}} \sum_{i=1}^K \mathbf{y}(i)^2}. \end{aligned}$$

Similarly, since  $\Psi$  is a  $\nu$ -self-concordant barrier and using that  $\log(1 + \sqrt{T}) > 1$  by assumption on  $T$ ,

$$\begin{aligned} \kappa \left( \frac{1}{\gamma_{t+\delta}} - \frac{1}{\gamma_t} \right) (\nabla \Psi(\mathbf{v}))^\top \mathbf{y} &\leq \kappa \sqrt{\nu} \left( \frac{1}{\gamma_{t+\delta}} - \frac{1}{\gamma_t} \right) \sqrt{\mathbf{y}^\top \nabla^2 \Psi(\mathbf{v}) \mathbf{y}} \\ &\leq \kappa \sqrt{16B^2 K^2 d_{\max}} \sqrt{\mathbf{y}^\top \nabla^2 \Psi(\mathbf{v}) \mathbf{y}} \\ &= \sqrt{\kappa \gamma_{t+\delta} 16B^2 K^2 d_{\max}} \sqrt{\frac{\kappa}{\gamma_{t+\delta}} \mathbf{y}^\top \nabla^2 \Psi(\mathbf{v}) \mathbf{y}} \\ &\leq \frac{1}{32\sqrt{2}} \sqrt{\frac{\kappa}{\gamma_{t+\delta}} \mathbf{y}^\top \nabla^2 \Psi(\mathbf{v}) \mathbf{y}}. \end{aligned}$$

By using the above two inequalities and  $\sqrt{a} + \sqrt{b} \leq \sqrt{2}\sqrt{a+b}$  we can see that

$$\begin{aligned} \left( \nabla R_t(\mathbf{v}) - \nabla R_{t+\delta}(\mathbf{v}) \right)^\top \mathbf{y} &\leq \frac{1}{32} \sqrt{\frac{\kappa}{\eta_{t+\delta}} 2 \sum_{i=1}^K \mathbf{y}(i)^2 + \frac{\kappa}{\gamma_{t+\delta}} \mathbf{y}^\top \nabla^2 \Psi(\mathbf{v}) \mathbf{y}} \\ &= \frac{1}{32} \sqrt{\kappa \mathbf{y}^\top \nabla^2 R_{t+\delta}(\mathbf{v}) \mathbf{y}} \end{aligned}$$

Next, pick any  $t' \in [T]$  and observe that because  $\nabla^2 R_{t'}(\mathbf{v}) \succeq \frac{1}{\gamma_{t'}} \nabla^2 \Psi(\mathbf{v})$  we have that

$$\kappa \|\hat{\ell}_t\|_{R_{t'}, \mathbf{w}_t(\hat{\mathbf{L}}_t)}^2 \leq \kappa \gamma_{t'} K^2 (\ell_t^\top \mathbf{a}_t)^2 \mathbf{v}_t^\top \mathbf{v}_t = \kappa \gamma_{t'} K^2 (\ell_t^\top \mathbf{a}_t)^2.$$

Since  $\|\ell_t\|_2 \leq 1$  and  $\mathcal{W} \subseteq \mathcal{B}(B)$ , we have that  $(\ell_t^\top \mathbf{a}_t)^2 \leq B^2$  and thus

$$\sqrt{\kappa} \|\hat{\ell}_t\|_{R_{t'}, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq \sqrt{\kappa} \underbrace{\sqrt{\gamma_{t'} B^2 K^2}}_{\beta_{t'}} \leq \frac{1}{128 d_{\max}},$$

where the last inequality is because  $\gamma_{t'} \leq \frac{1}{128 B K d_{\max}}$ . Let  $\tau \in m_t \cup \{t\}$ , because  $\nabla^2 R_t(\mathbf{v}) \succeq \frac{\kappa}{\eta_t} \mathbf{I}$  and  $\|\ell_\tau\|_2 \leq 1$  we have that

$$\|\ell_\tau\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq \underbrace{\sqrt{\eta_t}}_{\alpha_t}.$$

We have fulfilled all requirements for Corollary 4.

Let  $\tilde{\mathbf{u}} = \frac{\mathbf{u} - \mathbf{v}^+}{1 + \frac{1}{\sqrt{T}}} + \mathbf{v}^+ \in \mathcal{W}_{\frac{1}{\sqrt{T}}}$ , with  $\mathbf{v}^+ = \arg \min_{\mathbf{v} \in \mathcal{W}} R_1(\mathbf{v})$ . By equation 19 and  $\mathcal{W} \subseteq \mathcal{B}(B)$  we have that

$$\begin{aligned} R_T(\tilde{\mathbf{u}}) - R_1(\mathbf{v}^+) &\leq \frac{B^2}{\eta_T} + \frac{1}{\gamma_T} \tilde{\Psi}(\tilde{\mathbf{u}}) - \frac{1}{\gamma_1} \tilde{\Psi}(\mathbf{v}^+) \\ &\leq \frac{B^2}{\eta_T} + \frac{\nu}{\gamma_T} \log \left( \frac{1 + \xi}{\xi} \right), \end{aligned} \tag{20}$$

where we used that  $\tilde{\Psi}(\mathbf{v}^+) \geq 0$ . Furthermore, by using that  $\mathcal{W} \in \mathcal{B}(B)$  and  $\|\ell_t\|_2 \leq 1$ , we have that

$$\sum_{t=1}^T (\tilde{\mathbf{u}} - \mathbf{u})^\top \ell_t = \sum_{t=1}^T \left( 1 - \frac{1}{1 + \frac{1}{\sqrt{T}}} \right) (\tilde{\mathbf{u}} - \mathbf{v}^+)^\top \ell_t \leq 2TB \left( \frac{\frac{1}{\sqrt{T}}}{1 + \frac{1}{\sqrt{T}}} \right) \leq 2B\sqrt{T}.$$



---

**Algorithm 5:** Efficient implementation of delayed FTRL for linear bandits
 

---

**Input:**  $\nu$ -self concordant barrier  $\Psi$  for  $\mathcal{W}$ , hyperparameters  $\eta, \gamma$ .  
 Set  $\mathbf{z}_1 = \arg \min_{\mathbf{v}} \Psi_1(\mathbf{v})$   
**for**  $t = 1, \dots, T$  **do**  
     Observe  $\mathbf{a}_\tau^\top \boldsymbol{\ell}_\tau$  for  $\tau \in o_t \setminus o_{t-1}$ .  
     Find loss estimators  $\hat{\boldsymbol{\ell}}_\tau = K \boldsymbol{\ell}_\tau^\top \mathbf{a}_\tau (\nabla^2 \Psi(\mathbf{z}_\tau))^{1/2} \mathbf{v}_\tau$  for new observations  $\tau \in o_t \setminus o_{t-1}$ .  
     Compute  $\mathbf{z}_t = DN(\Psi_{t-1}, \mathbf{z}_{t-1})$   
     Play  $\mathbf{a}_t = \mathbf{z}_t + (\nabla^2 \Psi(\mathbf{z}_t))^{-1/2} \mathbf{v}_t$ , where  $\mathbf{v}_t$  is uniformly sampled from the unit sphere.  
**end for**

---

Thus, by Corollary 4 with  $\alpha_t^2 = \eta_t$  and  $\beta_t^2 = \gamma_t B^2 K^2$  we obtain

$$\begin{aligned}
 \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \mathbf{u})^\top \boldsymbol{\ell}_t \right] &\leq \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{a}_t - \tilde{\mathbf{u}})^\top \boldsymbol{\ell}_t \right] + 2B\sqrt{T} \\
 &\leq R_T(\mathbf{u}) - R_1(\mathbf{v}^+) + 16B^2 K^2 \sum_{t=1}^T \gamma_t + 16 \sum_{t=1}^T \eta_t |m_t| + 2B\sqrt{T} \\
 &\leq \frac{B^2}{\eta_T} + \frac{\nu}{\gamma_T} \log(1 + \sqrt{T}) + 8BK \sqrt{\nu T \log(1 + \sqrt{T})} + 8B\sqrt{D} + 2B\sqrt{T} \\
 &\leq 512BK d_{\max} \nu \log(1 + \sqrt{T}) + 12BK \sqrt{\nu T \log(1 + \sqrt{T})} + 12B\sqrt{D} + 2B\sqrt{T},
 \end{aligned}$$

where in the second inequality we used  $\sum_{t=1}^T \gamma_t \leq \frac{\sqrt{\nu T \log(1 + \sqrt{T})}}{2BK}$  and  $\sum_{t=1}^T \eta_t |m_t| \leq \frac{1}{2} B\sqrt{D}$ , both of which follow from Lemma 24, and the last inequality follows from simplifications.  $\blacksquare$

### B.1 Efficient Implementation

In this section we will use fixed learning rates  $\eta_t = \eta$  and  $\gamma_t = \gamma$  for simplicity. Define

$$\begin{aligned}
 \Phi_t(\mathbf{v}) &= \gamma \hat{\mathbf{L}}_t^\top \mathbf{v} + \gamma R(\mathbf{v}) = \gamma \hat{\mathbf{L}}_t^\top \mathbf{v} + \frac{\gamma}{\eta} \|\mathbf{v}\|_2^2 + \tilde{\Phi}(\mathbf{v}) \\
 e(\Phi_t, \mathbf{v}) &= -(\nabla^2 \Phi_t(\mathbf{v}))^{-1} \nabla \Phi_t(\mathbf{v}) \\
 \lambda(\Phi_t, \mathbf{v}) &= \sqrt{\nabla \Phi_t(\mathbf{v})^\top (\nabla^2 \Phi_t(\mathbf{v}))^{-1} \nabla \Phi_t(\mathbf{v})} \\
 DN(\Phi_t, \mathbf{v}) &= \mathbf{v} - \frac{1}{1 + \lambda(\Phi_t, \mathbf{v})} e(\Phi_t, \mathbf{v}) \\
 \mathbf{z}_t^* &= \arg \min_{\mathbf{v}} \Phi_{t-1}(\mathbf{v}).
 \end{aligned}$$

The following facts can be found in Nemirovski and Todd (2008)

$$\lambda(\Phi_t, DN(\Phi_t, \mathbf{v})) \leq 2\lambda(\Phi_t, \mathbf{v})^2 \quad (21)$$

$$\|\mathbf{v} - \mathbf{z}_t^*\|_{\Phi_{t-1}, \mathbf{z}_t^*} \leq \frac{\lambda(\Phi_{t-1}, \mathbf{v})}{1 - 2\lambda(\Phi_{t-1}, \mathbf{v})} \quad \text{if } \lambda(\Phi_{t-1}, \mathbf{v}) < \frac{1}{2} \quad (22)$$

Algorithm 5 is a simple modification of Algorithm 2 in section 9 of Abernethy et al. (2008). Abernethy et al. (2008) show that in the non-delayed setting, given the previous iterate, it takes essentially one iteration of the damped Newton method to compute  $w_t(\hat{L}_t)$ . If an easily computed self-concordant barrier is available, the computational complexity of the damped Newton method is  $O(K^2)$ . Since we can compute  $(\nabla^2 R_t(z_t))^{1/2}$  and its inverse in  $O(K^3)$  time by means of an eigenvalue decomposition, the total runtime is  $O(K^3)$ . In what follows we provide a modification of Lemma 7 by Abernethy et al. (2008) to the delayed setting. In what follows we will show that  $z_t^*$  is close to  $z_t$  as measured in local distance. While this may seem arbitrary, we have that  $z_t^* = \arg \min_v \Phi_t(v) = \arg \min_v \frac{1}{\gamma} \Phi_t(v)$ , which in turn is the FTRL objective we have been working with throughout this paper. Thus, showing that  $z_t^*$  is close to  $z_t$  implies that  $z_t$  will have a similar regret bound as we would obtain from  $z_t^*$ , as argued by Abernethy et al. (2008). With Lemma 15 in hand, one can follow the steps provided by Abernethy et al. (2008) to see that the regret of Algorithm 5 is of the same order as that of Algorithm 2.

**Lemma 15** *Suppose that  $\eta_t = \eta > 0$  and  $\gamma_t = \gamma \leq \frac{1}{162K^2d_{\max}}$ . It holds that for all  $t$*

$$\lambda(\Phi_t, z_t)^2 \leq 9\gamma^2 K^2 d_{\max} \quad \text{and} \quad \|z_t - z_t^*\|_{\Phi_{t-1}, z_t^*} \leq 648\gamma^2 K^2 d_{\max}.$$

**Proof** The proof is by induction on  $t$ . The base case holds by definition. Suppose the statement holds for  $t-1$ . Using  $(x+y)^\top A(x+y) \leq 2x^\top Ax + 2y^\top Ay$  we get that

$$\begin{aligned} \lambda(\Phi_t, z_t)^2 &= \nabla \Phi_t(z_t)^\top (\nabla^2 \Phi_t(z_t))^{-1} \nabla \Phi_t(z_t) \\ &= \left( \nabla \Phi_{t-1}(z_t) + \gamma \sum_{\tau \in m_t} \hat{\ell}_\tau \right)^\top \left( \frac{\gamma}{\eta} \mathbf{I} + \nabla^2 \Psi(z_t) \right)^{-1} \left( \nabla \Phi_{t-1}(z_t) + \gamma \sum_{\tau \in m_t} \hat{\ell}_\tau \right) \\ &\leq 2\gamma^2 \sum_{\tau \in m_t} \hat{\ell}_\tau^\top \left( \frac{\gamma}{\eta} \mathbf{I} + \nabla^2 \Psi(z_t) \right)^{-1} \hat{\ell}_\tau \\ &\quad + 2 \underbrace{\nabla \Phi_{t-1}(z_t)^\top (\nabla^2 \Psi_{t-1}(z_t))^{-1} \nabla \Phi_{t-1}(z_t)}_{\lambda(\Phi_{t-1}, z_t)^2}. \end{aligned}$$

Now, with a minor modification of Lemma 7 we can see that  $z_t^* \in \mathcal{D}_{R_t}(z_{t-\delta}^*, \frac{1}{2})$  and  $z_t^* \in \mathcal{D}_{R_t}(z_{t-\delta}, \frac{1}{2})$  for all  $\delta \in [\min\{d_{\max}, T-t\}]$ . In turn, this implies that

$$\begin{aligned} &2\gamma^2 \sum_{\tau \in m_t} \hat{\ell}_\tau^\top \left( \frac{\gamma}{\eta} \mathbf{I} + \nabla^2 \Psi(z_t) \right)^{-1} \hat{\ell}_\tau \\ &\leq 2\gamma^2 \sum_{\tau \in m_t} \hat{\ell}_\tau^\top (\nabla^2 \Psi(z_t))^{-1} \hat{\ell}_\tau \\ &\leq 8\gamma^2 \sum_{\tau \in m_t} \hat{\ell}_\tau^\top (\nabla^2 \Psi(z_\tau))^{-1} \hat{\ell}_\tau \quad \text{(equation (18))} \\ &\leq 8\gamma^2 K^2 |m_t| \leq 8\gamma^2 K^2 d_{\max}. \quad \text{(by def. of } \hat{\ell}_\tau) \end{aligned}$$

By equation (21) and the induction assumption we have that

$$\lambda(\Phi_{t-1}, z_t)^2 \leq 2\lambda(\Phi_{t-1}, z_{t-1})^4 \leq 162\gamma^4 K^4 d_{\max}^2. \quad (23)$$

Thus, we can apply the assumption  $\gamma^2 \leq \frac{1}{162K^2d_{\max}}$  to find that

$$\lambda(\Phi_t, \mathbf{z}_t)^2 \leq 8\gamma^2 K^2 d_{\max} + 162\gamma^4 K^4 d_{\max}^2 \leq 9\gamma^2 K d_{\max},$$

after which we have proven the induction step for the first claim. For the second claim, we start with equation (22) and then the fact that  $\lambda(\Phi_{t-1}, \mathbf{z}_t)^2 \leq \frac{1}{16}$  which follows by equation (23) and the assumption that  $\gamma^2 \leq \frac{1}{162K^2d_{\max}}$ , then using equation (21) and finally applying the first claim yields

$$\|\mathbf{z}_t - \mathbf{z}_t^*\|_{\Phi_{t-1}, \mathbf{z}_t^*} \leq \frac{\lambda(\Phi_{t-1}, \mathbf{z}_t)}{1 - 2\lambda(\Phi_{t-1}, \mathbf{z}_t)} \leq 2\lambda(\Phi_{t-1}, \mathbf{z}_t) \leq 4\lambda(\Phi_{t-1}, \mathbf{z}_{t-1})^2 \leq 648\gamma^2 K^2 d_{\max}.$$

■

## Appendix C. Adversarial Markov Decision Processes (MDPs)

**Lemma 11 (RESTATEd)**  $\mathcal{W}$  satisfies the following:

1. For any  $\mathbf{q} \in \Delta(\mathcal{M})$ , there exists  $\tilde{\mathbf{q}} \in \mathcal{W}$  such that  $\min_{h,s,a,s'} \tilde{\mathbf{q}}(h, s, a, s') \geq \frac{1}{T^3 H^2 S^4 A^2}$  and  $\|\mathbf{q} - \tilde{\mathbf{q}}\|_1 \leq \frac{2H}{T}$ .
2. Given  $\mathbf{v} \in \mathcal{W}$ , let  $\pi$  be defined by  $\pi(a \mid h, s) = \frac{v(h,s,a)}{v(h,s)}$  and  $\mathbf{q}^{\max}(h, s, a) = \max_{\hat{p} \in \mathcal{P}} \mathbf{q}^{\pi, \hat{p}}(h, s, a)$ .

Then,  $\|\mathbf{q}^\pi - \mathbf{v}\|_1 \leq \frac{2H}{T}$  and  $\|\mathbf{q}^{\max} - \mathbf{v}\|_1 \leq \frac{4H^2 S}{T}$ .

**Proof** We start by proving the first claim. Define  $\tilde{p} : [H] \times \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$  by  $\tilde{p}(s' \mid h, s, a) = (1 - \frac{1}{THSA})p(s' \mid h, s, a) + \frac{1}{THS^2A}$  and notice that  $\tilde{p} \in \mathcal{P}$  since  $|p(s' \mid h, s, a) - \tilde{p}(s' \mid h, s, a)| \leq \frac{1}{THSA}$ . Next, let  $\pi_u$  be the uniformly random policy, and define  $\tilde{\mathbf{q}} = (1 - \frac{1}{T})\mathbf{q} + \frac{1}{T}\mathbf{q}^{\pi_u, \tilde{p}}$ . It holds that  $\tilde{\mathbf{q}} \in \mathcal{W}$  because  $\mathcal{W}$  is a convex set. Moreover, notice that  $\mathbf{q}^{\pi_u, \tilde{p}}(h, s, a, s') \geq \frac{1}{(THS^2A)^2A}$  which implies that  $\tilde{\mathbf{q}}(h, s, a, s') \geq \frac{1}{T^3 H^2 S^4 A^2}$ . Finally,

$$\begin{aligned} \|\mathbf{q} - \tilde{\mathbf{q}}\|_1 &= \sum_{h,s,a,s'} |\mathbf{q}(h, s, a, s') - \tilde{\mathbf{q}}(h, s, a, s')| \\ &= \sum_{h,s,a,s'} \left| \frac{1}{T}\mathbf{q}(h, s, a, s') - \frac{1}{T}\mathbf{q}^{\pi_u, \tilde{p}}(h, s, a, s') \right| \\ &\leq \frac{1}{T} \sum_{h,s,a,s'} \mathbf{q}(h, s, a, s') + \frac{1}{T} \sum_{h,s,a,s'} \mathbf{q}^{\pi_u, \tilde{p}}(h, s, a, s') = \frac{2H}{T}. \end{aligned}$$

Now we prove the second claim. Define loss function  $\tilde{\ell}(h, s, a) = \text{sign}(\mathbf{q}^\pi(h, s, a) - \mathbf{v}(h, s, a))$  and note that  $\|\mathbf{q}^\pi - \mathbf{v}\|_1 = V^{\pi, p, \tilde{\ell}}(1, s_{\text{init}}) - V^{\pi, \hat{p}, \tilde{\ell}}(1, s_{\text{init}})$  for some  $\hat{p} \in \mathcal{P}$ . Combining the value difference lemma (see, e.g., Shani et al. (2020)) with  $\|p - \hat{p}\|_\infty \leq \frac{1}{THSA}$  proves that  $\|\mathbf{q}^\pi - \mathbf{v}\|_1 \leq \frac{2H}{T}$ . Now, let  $\hat{p}^{h,s}$  be the transition function that corresponds to  $\mathbf{u}(h, s)$ . We have that,  $\|\hat{p}^{h,s} - \hat{p}\|_\infty \leq \|\hat{p}^{h,s} - p\|_\infty + \|p - \hat{p}\|_\infty \leq \frac{2}{THSA}$ . Thus, using the same argument as the above,  $\|\mathbf{u} - \mathbf{v}\|_1 \leq \sum_{h,s} \|\mathbf{q}^{\pi, \hat{p}^{h,s}} - \mathbf{v}\|_1 \leq \frac{4H^2 S}{T}$ . ■

**Theorem 12 (RESTATED)** Suppose that  $T \geq H$ . Algorithm 3 with

$$\gamma = \frac{1}{128Hd_{\max}} \quad \eta_t = \min \left\{ \frac{\log(SA)}{96HSA\sqrt{SAd_{\max} + d_{\max}^2}}, \frac{\sqrt{\log(SA)}}{\sqrt{SA t + \sum_{\tau=1}^t |m_t|}} \right\}$$

guarantees

$$\mathbb{E}[\mathcal{R}_T] \leq 72H\sqrt{\log(SA)(TSA + D)} + 1338d_{\max}H^2S^2A^2\log(HSAT) .$$

**Proof** As in the proof of Theorem 8, the regularizer  $R_t$  as specified in (12) does not satisfy Assumption 2(c) because we can have  $R_t(\mathbf{v}) \leq R_{t-1}(\mathbf{v})$ , but, as argued in the proof of Theorem 8, we can overcome this issue in a relative straightforward manner via the regularizer  $\tilde{R}_t(\mathbf{v}) = R_t(\mathbf{v}) - \min_{\mathbf{v}' \in \mathcal{W}} R_t(\mathbf{v}')$ , which has no impact on the iterates. We continue by decomposing the regret as

$$\mathcal{R}_T = \sum_{t=1}^T \langle \mathbf{q}^{\pi_t} - \mathbf{u}, \ell_t \rangle = \underbrace{\sum_{t=1}^T \langle \mathbf{q}^{\pi_t} - \mathbf{w}_t(\hat{\mathbf{L}}_t), \ell_t \rangle}_{\text{ERROR}} + \underbrace{\sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle}_{\text{REG}} + \underbrace{\sum_{t=1}^T \langle \tilde{\mathbf{u}} - \mathbf{u}, \ell_t \rangle}_{\text{SHIFT-PENALTY}},$$

where, by the second property in Lemma 11, ERROR is bounded by  $2H$ ; and by the first property  $\tilde{\mathbf{u}} \in \mathcal{W}$  exists such that both SHIFT-PENALTY is bounded by  $2H$  and  $\min_{h,s,a,s'} \tilde{\mathbf{u}}(h,s,a,s') \geq \frac{1}{T^3H^2S^4A^2}$ , which we will choose as our comparator to ensure that the regularisation is always bounded. For REG we use Lemma 3 with  $\kappa = \gamma$ . The fact that  $4\nabla^2 R_t(\mathbf{v}) \succeq \nabla^2 R_t(\mathbf{v}') \succeq \frac{1}{4}\nabla^2 R_t(\mathbf{v})$  for all  $\mathbf{v} \in \mathcal{W}$  and  $\mathbf{v}' \in \mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2\sqrt{\gamma}})$  and all  $t$  comes directly from Lemma 25.

For  $\delta \in d_{\max}$  we have that

$$\sqrt{\eta_{t+\delta}} \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \leq \sqrt{\eta_{t+\delta}} \frac{\sqrt{SAd_{\max} + d_{\max}^2}}{\sqrt{\log(SA)}} \leq \frac{1}{\sqrt{32HSA}} (SAd_{\max} + d_{\max}^2)^{1/4}.$$

Thus, by definition of  $\gamma$ , Lemma 25 also implies that that

$$\gamma (\nabla R_t(\mathbf{v}) - \nabla R_{t+\delta}(\mathbf{v}))^\top y \leq \frac{1}{32} \sqrt{\gamma y^\top \nabla^2 R_{t+\delta}(\mathbf{v}) y}.$$

Next pick any  $t \in [T]$ ,  $\tau \in m_t$ , then

$$\|\ell_\tau\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \leq \sqrt{\eta_t \sum_{h,s,a} \mathbf{w}_t(\hat{\mathbf{L}}_t, h, s, a) \ell(h, s, a)^2} \leq \sqrt{\eta_t \sum_{h,s,a} \mathbf{w}_t(\hat{\mathbf{L}}_t, h, s, a)} = \underbrace{\sqrt{\eta_t H}}_{\alpha_t}.$$

Since  $\mathbf{q}_\tau^{\max}(h, s, a) = \max_{\hat{p} \in \mathcal{P}} \mathbf{q}^{\pi_\tau, \hat{p}}(h, s, a) \geq \max\{\mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, h, s, a), \mathbf{q}^{\pi_\tau}(h, s, a)\}$  we have that

$$\begin{aligned} \mathbb{E}[\|\hat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau)}^2] &= \eta_t \mathbb{E} \left[ \sum_{h,s,a} \mathbf{w}_\tau(\hat{\mathbf{L}}_\tau, h, s, a) \hat{\ell}_\tau(h, s, a)^2 \right] \\ &\leq \eta_t \mathbb{E} \left[ \sum_{h,s,a} \frac{\mathbb{E}[\mathbb{I}\{s_{t,h} = s, a_{t,h} = a\} \mid \mathcal{F}_\tau]}{\mathbf{q}_\tau^{\max}(h, s, a)} \right] \\ &= \eta_t \mathbb{E} \left[ \sum_{h,s,a} \frac{\mathbf{q}^{\pi_\tau}(h, s, a)}{\mathbf{q}_\tau^{\max}(h, s, a)} \right] \leq \underbrace{\eta_t HSA}_{\beta_t^2}. \end{aligned}$$

Finally, pick any  $t, t'$ ,

$$\begin{aligned}\sqrt{\gamma}\|\widehat{\ell}_t\|_{R_{t'}, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} &\leq \gamma \sqrt{\sum_{h,s,a} \mathbf{w}_t(\widehat{\mathbf{L}}_t, h, s, a)^2 \widehat{\ell}_t(h, s, a)^2} \leq \gamma \sqrt{\sum_{h,s,a} \mathbb{I}\{s_{t,h} = s, a_{t,h} = a\}} \\ &= \gamma \sqrt{H} \leq \frac{1}{128d_{\max}},\end{aligned}$$

where the second inequality is due to the fact that  $\mathbf{q}_t^{\max}(h, s, a) \geq \mathbf{q}^{\pi_t}(h, s, a)$  and that  $\gamma = \kappa$  and the last inequality is by definition of  $\gamma$ . Thus, applying Lemma 3 with  $\mathbf{b}_t = \mathbb{E}[\widehat{\ell}_t - \ell_t \mid \mathcal{F}_t]$ , we get

$$\begin{aligned}\text{REG} &\leq \underbrace{R_T(\tilde{\mathbf{u}}) - \min_{\mathbf{v} \in \mathcal{W}} R_1(\mathbf{v})}_{\text{PENALTY}} + 8HSA \sum_{t=1}^T \eta_t + 8H \sum_{t=1}^T \eta_t |m_t| + \underbrace{\sum_{t=1}^T \mathbb{E}[\mathbf{w}_t(\widehat{\mathbf{L}}_t^m)^\top (\ell_t - \widehat{\ell}_t)]}_{\text{BIAS}_1} \\ &\quad + \underbrace{\sum_{t=1}^T \mathbb{E}[\tilde{\mathbf{u}}^\top (\widehat{\ell}_t - \ell_t)]}_{\text{BIAS}_2} + \underbrace{8\sqrt{H} \sum_{t=1}^T \sqrt{\eta_t} \mathbb{E}[\|\sum_{\tau \in m_t} (\ell_\tau - \widehat{\ell}_\tau)\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)}]}_{\text{DRIFT}}.\end{aligned}$$

Recall that  $\min_{h,s,a,s'} \tilde{\mathbf{u}}(h, s, a, s') \geq \frac{1}{T^3 H^2 S^4 A^2}$ . Thus, using the third fact of Lemma 25 with  $b = \frac{1}{T^3 H^2 S^4 A^2}$ ,  $K = HS^2A$  and  $B = H$ , we conclude

$$\begin{aligned}\text{PENALTY} &\leq \frac{H(1 + \log(S^2A))}{\eta_T} + \frac{HS^2A \log(T^3 H^2 S^4 A^2)}{\gamma} \\ &\leq \frac{4H \log(SA)}{\eta_T} + \frac{4HS^2A \log(HSAT)}{\gamma}.\end{aligned}$$

Now we deal with the primary term that makes up DRIFT, for each  $t$

$$\begin{aligned}\|\sum_{\tau \in m_t} (\ell_\tau - \widehat{\ell}_\tau)\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)}^2 &= \left(\sum_{\tau \in m_t} \ell_\tau - \widehat{\ell}_\tau\right)^\top \nabla^{-2} R_t(\mathbf{w}_t(\widehat{\mathbf{L}}_t)) \left(\sum_{\tau \in m_t} \ell_\tau - \widehat{\ell}_\tau\right) \\ &= \sum_{\tau \in m_t} \|(\ell_\tau - \widehat{\ell}_\tau)\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)}^2 + \sum_{\tau \in m_t} \sum_{\tau' \in m_t \setminus \{\tau\}} (\ell_\tau - \widehat{\ell}_\tau)^\top \nabla^{-2} R_t(\mathbf{w}_t(\widehat{\mathbf{L}}_t)) (\ell_{\tau'} - \widehat{\ell}_{\tau'}). \quad (24)\end{aligned}$$

We continue by bounding the first term on the right hands side in expectation:

$$\begin{aligned}\sum_{\tau \in m_t} \mathbb{E}[\|(\ell_\tau - \widehat{\ell}_\tau)\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)}^2] &\leq 4 \sum_{\tau \in m_t} \mathbb{E}[\|(\ell_\tau - \widehat{\ell}_\tau)\|_{R_t, \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)}^2] \\ &\leq 4 \sum_{\tau \in m_t} \mathbb{E}[\|\ell_\tau\|_{R_t, \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)}^2] + 4 \sum_{\tau \in m_t} \mathbb{E}[\|\widehat{\ell}_\tau\|_{R_t, \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)}^2] \\ &\leq 4|m_t|(\alpha_t^2 + \beta_t^2) \leq 8\eta_t HSA|m_t|.\end{aligned}$$

We bound the second term on the right hand side of (24) by first applying the law of total expectation,

$$\begin{aligned}
 & \mathbb{E} \left[ \sum_{\tau \in m_t} \sum_{\tau' \in m_t \setminus \{\tau\}} (\ell_\tau - \widehat{\ell}_\tau) \nabla^{-2} R_t(\mathbf{w}_t(\widehat{\mathbf{L}}_t)) (\ell_{\tau'} - \widehat{\ell}_{\tau'}) \right] \\
 &= \mathbb{E} \left[ \eta_t \sum_{\tau \in m_t} \sum_{\tau' \in m_t \setminus \{\tau\}} \sum_{h,s,a} \mathbf{w}_t(\widehat{\mathbf{L}}_t, h, s, a) (\ell_\tau(h, s, a) - \mathbb{E}[\widehat{\ell}_\tau(h, s, a) \mid \mathcal{F}_t]) (\ell_{\tau'}(h, s, a) \right. \\
 &\quad \left. - \mathbb{E}[\widehat{\ell}_{\tau'}(h, s, a) \mid \mathcal{F}_t]) \right] \\
 &+ \mathbb{E} \left[ \gamma \sum_{\tau \in m_t} \sum_{\tau' \in m_t \setminus \{\tau\}} \sum_{h,s,a} \mathbf{w}_t(\widehat{\mathbf{L}}_t, h, s, a)^2 (\ell_\tau(h, s, a) - \mathbb{E}[\widehat{\ell}_\tau(h, s, a) \mid \mathcal{F}_t]) (\ell_{\tau'}(h, s, a) \right. \\
 &\quad \left. - \mathbb{E}[\widehat{\ell}_{\tau'}(h, s, a) \mid \mathcal{F}_t]) \right].
 \end{aligned}$$

Then, since  $\ell_{\tau'}(h, s, a) - \mathbb{E}[\widehat{\ell}_{\tau'}(h, s, a) \mid \mathcal{F}_t] \in [0, 1]$ , we can bound the first term on the right-hand-side above by

$$\begin{aligned}
 & \mathbb{E} \left[ \eta_t \sum_{\tau \in m_t} \sum_{\tau' \in m_t \setminus \{\tau\}} \sum_{h,s,a} \mathbf{w}_t(\widehat{\mathbf{L}}_t, h, s, a) (\ell_\tau(h, s, a) - \mathbb{E}[\widehat{\ell}_\tau(h, s, a) \mid \mathcal{F}_t]) \right] \\
 &\leq \mathbb{E} \left[ \eta_t |m_t| \sum_{\tau \in m_t} \sum_{h,s,a} \mathbf{w}_t(\widehat{\mathbf{L}}_t, h, s, a) \ell_\tau(h, s, a) \frac{\mathbf{q}_\tau^{\max}(h, s, a) - \mathbf{q}^{\pi_\tau}(h, s, a)}{\mathbf{q}_\tau^{\max}(h, s, a)} \right] \\
 &\leq 2 \mathbb{E} \left[ \eta_t |m_t| \sum_{\tau \in m_t} \sum_{h,s,a} \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau, h, s, a) \frac{\mathbf{q}_\tau^{\max}(h, s, a) - \mathbf{q}^{\pi_\tau}(h, s, a)}{\mathbf{q}_\tau^{\max}(h, s, a)} \right] \\
 &\leq 2 \mathbb{E} \left[ \eta_t |m_t| \sum_{\tau \in m_t} \sum_{h,s,a} \mathbf{q}_\tau^{\max}(h, s, a) - \mathbf{q}^{\pi_\tau}(h, s, a) \right] \\
 &\leq 2 \mathbb{E} \left[ \eta_t |m_t| \sum_{\tau \in m_t} \|\mathbf{q}_\tau^{\max} - \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)\|_1 + \|\mathbf{q}^{\pi_\tau} - \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)\|_1 \right] \leq 12 \eta_t |m_t|^2 \frac{H^2 S}{T},
 \end{aligned}$$

where the third inequality is by equation (37), the fourth inequality is since  $\mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau, h, s, a) \leq \mathbf{q}_\tau^{\max}(h, s, a)$ , and the fifth inequality is by Lemma 11. Likewise we can see that

$$\begin{aligned}
 & \mathbb{E} \left[ \sum_{\tau \in m_t} \sum_{\tau' \in m_t \setminus \{\tau\}} \sum_{h,s,a} \mathbf{w}_t(\widehat{\mathbf{L}}_t, h, s, a)^2 (\ell_\tau(h, s, a) - \mathbb{E}[\widehat{\ell}_\tau(h, s, a) \mid \mathcal{F}_t]) (\ell_{\tau'}(h, s, a) \right. \\
 &\quad \left. - \mathbb{E}[\widehat{\ell}_{\tau'}(h, s, a) \mid \mathcal{F}_t]) \right] \leq 12 |m_t|^2 \frac{H^2 S}{T}.
 \end{aligned}$$

These inequalities combined with Jensen's inequality gives

$$\begin{aligned}
 8\sqrt{H} \cdot \text{DRIFT} &\leq 32H \sum_{t=1}^T \eta_t \sqrt{SA|m_t|} + \sum_{t=1}^T 96H(\eta_t + \sqrt{\eta_t \gamma})|m_t| \sqrt{\frac{H^2 S}{T}} \\
 &\leq 32H \sum_{t=1}^T \eta_t (SA + |m_t|) + H\sqrt{ST} \\
 &\leq 65H \sqrt{\log(SA)(TSA + D)},
 \end{aligned}$$

where we used that  $\eta, \gamma \leq \frac{1}{96Hd_{\max}}$ . Next, we bound  $\text{BIAS}_1$ . Let  $\mathcal{G}_t$  be the history of all episodes in  $[t-1]$ , and note that  $\mathbf{w}_t(\hat{\mathbf{L}}_t)$ ,  $\mathbf{q}_t^{\max}$  and  $\mathbf{w}_t(\hat{\mathbf{L}}_t^m)$  are all determined by  $\mathcal{G}_t$ . Therefore,

$$\begin{aligned}
 \text{BIAS}_1 &= \mathbb{E} \left[ \sum_{t,h,s,a} \mathbf{w}_t(\hat{\mathbf{L}}_t^m, h, s, a) (\ell_t(h, s, a) - \mathbb{E}[\hat{\ell}_t(h, s, a) \mid \mathcal{G}_t]) \right] \\
 &= \mathbb{E} \left[ \sum_{t,h,s,a} \mathbf{w}_t(\hat{\mathbf{L}}_t^m, h, s, a) \ell_t(h, s, a) \left( 1 - \frac{\mathbf{q}^{\pi_t}(h, s, a)}{\mathbf{q}_{t,h}^{\max}(s, a)} \right) \right] \\
 &\leq \mathbb{E} \left[ \sum_{t,h,s,a} \mathbf{w}_t(\hat{\mathbf{L}}_t^m, h, s, a) \frac{|\mathbf{q}_{t,h}^{\max}(s, a) - \mathbf{q}^{\pi_t}(h, s, a)|}{\mathbf{q}_t^{\max}(h, s, a)} \right].
 \end{aligned}$$

Now, as in the proof of Lemma 3,  $\mathbf{w}_t(\hat{\mathbf{L}}_t^m) \in \mathcal{D}_{R_t}(\mathbf{w}_t(\hat{\mathbf{L}}_t), \frac{1}{2\sqrt{\gamma}})$ . Thus, by equation (37),  $\mathbf{w}_t(\hat{\mathbf{L}}_t^m, h, s, a) \leq 2\mathbf{w}_t(\hat{\mathbf{L}}_t, h, s, a) \leq 2\mathbf{q}_t^{\max}(h, s, a)$ . Therefore,

$$\text{BIAS}_1 \leq 2 \sum_{t=1}^T \mathbb{E}[\|\mathbf{q}_t^{\max} - \mathbf{w}_t(\hat{\mathbf{L}}_t)\|_1] \leq 2 \sum_{t=1}^T \mathbb{E}[\|\mathbf{q}_t^{\max} - \mathbf{w}_t(\hat{\mathbf{L}}_t)\|_1 + \|\mathbf{q}^{\pi_t} - \mathbf{w}_t(\hat{\mathbf{L}}_t)\|_1] \leq 12H^2 S.$$

where the last is by article 2 in Lemma 11.

Recall that by definition  $\mathbf{q}_t^{\max}(h, s, a) \geq \mathbf{w}^{\pi_t}(h, s, a)$ . Thus,  $\mathbb{E}[\hat{\ell}_t(h, s, a) \mid \mathcal{F}_t] \leq \ell_t$  and  $\text{BIAS}_2 \leq 0$ . Putting everything together gives

$$\begin{aligned}
 \mathbb{E}[\mathcal{R}_T] &\leq 16H^2 S + 71H \sqrt{\log(SA)(TSA + D)} + \frac{4H \log(SA)}{\eta_T} + \frac{4HS^2 A \log(HSAT)}{\gamma} \\
 &\leq 16H^2 S + 72H \sqrt{\log(SA)(TSA + D)} + 512d_{\max} H^2 S^2 A \log(HSAT) \\
 &\quad + 800d_{\max} H^2 S^2 A^2 \\
 &\leq 72H \sqrt{\log(SA)(TSA + D)} + 1338d_{\max} H^2 S^2 A^2 \log(HSAT).
 \end{aligned}$$

■

## Appendix D. Adversarial MDPs with Unknown Transitions

**Theorem 14 (RESTATED)** *Algorithm 4 with  $\gamma = \frac{1}{128\sqrt{H}d_{\max}}$ ,  $\eta = \frac{\sqrt{\log(SA)}}{\sqrt{SAT+D}}$ ,  $\xi = \frac{1}{T}$  and  $T \geq 4$  guarantees,*

$$\begin{aligned} \mathbb{E}[\mathcal{R}_T] &\lesssim H^2 S \sqrt{AT \log(HSAT)} + H \sqrt{D \log(SA)} \\ &\quad + H^3 S^2 A \log(HSAT) d_{\max} + H^3 S^3 A \log^2(HSAT). \end{aligned}$$

**Proof** We introduce  $\varsigma$ , the good event, where  $p \in \mathcal{P}_{j'}$  for all  $j' \leq j_t$  and the compliment of the good event  $\varsigma^c$ . By Lemma 13 we have that the good event holds with probability of at least  $1 - 4/T^2$ . We then start by decomposing the regret in the same way as in the known transition setting, let  $\mathbf{u}$  be any comparator

$$\begin{aligned} \mathcal{R}_T &= \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{q}^{\pi_t} - \mathbf{u}, \ell_t \rangle \right] \\ &= \mathbb{E} \left[ \underbrace{\sum_{t=1}^T \langle \mathbf{q}^{\pi_t} - \mathbf{w}_t(\hat{\mathbf{L}}_t), \ell_t \rangle}_{\text{ERROR}} + \underbrace{\sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle}_{\text{REG}} + \underbrace{\sum_{t=1}^T \langle \tilde{\mathbf{u}} - \mathbf{u}, \ell_t \rangle}_{\text{SHIFT-PENALTY}} \right]. \end{aligned} \quad (25)$$

Under the good event and by Lemma 11 there exists an  $\tilde{\mathbf{u}} \in \mathcal{W}_T$  such that SHIFT-PENALTY is bounded by  $2H$ . That allows us to bound

$$\text{SHIFT-PENALTY} = \mathbb{E} \left[ \mathbb{I}\{\varsigma\} \sum_{t=1}^T \langle \tilde{\mathbf{u}} - \mathbf{u}, \ell_t \rangle + \mathbb{I}\{\varsigma^c\} \sum_{t=1}^T \langle \tilde{\mathbf{u}} - \mathbf{u}, \ell_t \rangle \right] \leq 2H + 4 \frac{HSA}{T}, \quad (26)$$

where we also used that  $\langle \tilde{\mathbf{u}} - \mathbf{u}, \ell_t \rangle \leq 1$  as well as the probability of the bad event being greater or equal to  $4/T^2$ . The ERROR can be bound using standard tools in the analysis of MDPs with unknown transitions (Lemma 29) by

$$\text{ERROR} \lesssim \sqrt{H^4 S^2 AT \log(HSAT^3)} + H^3 S^3 A \log^2(HSAT^3) + H^3 S^2 A d_{\max}. \quad (27)$$

Bounding the REG will be the main challenge of this proof as we now estimate the transition function and consequently have a changing domain  $\mathcal{W}_t$ . We are looking to apply Lemma 3, our analysis is structured around epochs and to make sure that  $\mathcal{W}_t = \mathcal{W}_\tau$  whenever  $\tau \in m_t$ , as is required by Assumption 1, we will only change our  $\mathcal{W}_t$  once each epoch and not use any delayed information from any previous epoch. We define  $\mathcal{E} = \{t : j_t \neq j_{t-1}\}$  be the set of rounds in which a new epoch starts and if we are changing epoch in round  $t \in \mathcal{E}$ , then we skip all outstanding observations,  $m_t \subseteq \Lambda$ .

$\mathcal{W}_t \subseteq \mathcal{W}_{t-1}$  holds by the construction of  $\mathcal{W}_t$  and  $\mathcal{W}_T$  is non-empty under the good event, fulfilling the rest of Assumption 1. Our regularizer fulfils Assumptions 2 under the same caveates as in the previous section. We split REG into the good and bad event

$$\text{REG} = \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle \right] = \mathbb{E} \left[ \mathbb{I}\{\varsigma\} \sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle \right] + \mathbb{E} \left[ \mathbb{I}\{\varsigma^c\} \sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle \right]$$



and bound the bad event first. By Lemma 13 the good event  $\varsigma$  happens with a probability of at least  $1 - \frac{4}{T^2}$  and we have that

$$\mathbb{E} \left[ \mathbb{I}\{\varsigma^c\} \sum_{t=1}^T (\mathbf{w}_t(\widehat{\mathbf{L}}_t) - \tilde{\mathbf{u}})^\top \widehat{\boldsymbol{\ell}}_t \right] \leq \mathbb{E} \left[ \mathbb{I}\{\varsigma^c\} \sum_{t=1}^T \|\widehat{\boldsymbol{\ell}}_t\|_1 \right] \leq \mathbb{E} [\mathbb{I}\{\varsigma^c\} 4HT^2] \leq 4H, \quad (28)$$

where we used Hölders inequality on  $(\mathbf{w}_t(\widehat{\mathbf{L}}_t) - \tilde{\mathbf{u}})^\top \widehat{\boldsymbol{\ell}}_t$ , upper bounded  $\|\mathbf{w}_t(\widehat{\mathbf{L}}_t) - \tilde{\mathbf{u}}\|_\infty \leq 1$  and finally upper bound all  $H$  non-zero elements of  $\widehat{\boldsymbol{\ell}}_t$  with  $\widehat{\ell}_t(h, s, a) \leq \frac{1}{\xi} = T$ . Under the good event we already showed that Assumption 1 and Assuming 2 are satisfied, both of which are required for Lemma 3. Next we bound  $\sqrt{\kappa} \|\widehat{\boldsymbol{\ell}}_t\|_{R, \mathbf{w}_t(\widehat{\mathbf{L}}_t)}$  for  $\kappa = \gamma$ . We have that

$$\begin{aligned} \sqrt{\kappa} \|\widehat{\boldsymbol{\ell}}_t\|_{R, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} &\leq \sqrt{\kappa} \sqrt{\gamma \sum_{h,s,a} \mathbf{w}_t(\widehat{\mathbf{L}}_t, h, s, a)^2 \widehat{\ell}_t(h, s, a)^2} \leq \sqrt{\kappa \gamma} \sqrt{\sum_{h,s,a} \mathbb{I}\{s_{t,h} = s, a_{t,h} = a\}} \\ &= \gamma \sqrt{H} = \frac{1}{128d_{\max}}, \end{aligned}$$

where the second inequality is due to the fact that  $\mathbf{q}_t^{\max}(h, s, a) \geq \mathbf{w}_t(\widehat{\mathbf{L}}_t, h, s, a)$  and the last inequality is by definition of  $\gamma$ . For all  $\tau \in m_t \cup \{t\}$ , we have that

$$\|\boldsymbol{\ell}_\tau\|_{R, \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)}^2 \leq \eta \sum_{h,s,a} \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau, h, s, a) \ell_\tau(h, s, a)^2 \leq \underbrace{\eta H}_{\alpha^2}. \quad (29)$$

We re-define the filtration over all past events observed by the learner to include state information  $\mathcal{F}_t = \{(\tau, s_{\tau,h}, a_{\tau,h}, h, \ell_\tau(h, s_{\tau,h}, a_{\tau,h})) : \tau + d_\tau < t, h \in [H]\}$ . Then using that on the good event  $\varsigma$ ,  $\mathbf{q}_\tau^{\max}(h, s, a) \geq \mathbf{q}^{\pi_\tau}(h, s, a)$ , we can bound

$$\begin{aligned} \mathbb{E}[\|\widehat{\boldsymbol{\ell}}_\tau\|_{R, \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)}^2 | \mathcal{F}_t, \varsigma] &\leq \mathbb{E} \left[ \eta \sum_{h,s,a} \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau, h, s, a) \widehat{\ell}_\tau(h, s, a)^2 \mid \mathcal{F}_t, \varsigma \right] \\ &\leq \mathbb{E} \left[ \eta \sum_{h,s,a} \frac{\mathbb{I}\{s_{\tau,h} = s, a_{\tau,h} = a\}}{\mathbf{q}^{\pi_\tau}(h, s, a)} \mid \mathcal{F}_\tau, \varsigma \right] = \underbrace{\eta H S A}_{\beta^2}. \end{aligned} \quad (30)$$

As in the proof of Theorem 8, the regularizer  $R_t$  as specified in (12) does not satisfy Assumption 2(c) because we can have  $R_t(\mathbf{v}) \leq R_{t-1}(\mathbf{v})$ , but, as argued in the proof of Theorem 8, we can overcome this issue in a relative straightforward manner via the regularizer  $\tilde{R}_t(\mathbf{v}) = R_t(\mathbf{v}) - \min_{\mathbf{v}' \in \mathcal{W}} R_t(\mathbf{v}')$ , which has no impact on the iterates. This already puts us in a position to apply Lemma 3 to find that

$$\begin{aligned} \mathbb{E} \left[ \mathbb{I}\{\varsigma\} \sum_{t=1}^T \langle \mathbf{w}_t(\widehat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \boldsymbol{\ell}_t \rangle \right] &\leq \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{w}_t(\widehat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \boldsymbol{\ell}_t \rangle \mid \varsigma \right] \\ &\leq \underbrace{\mathbb{E} \left[ \sum_{t \in \Lambda} (\mathbf{w}_t(\widehat{\mathbf{L}}_t) - \tilde{\mathbf{u}})^\top \boldsymbol{\ell}_t \mid \varsigma \right]}_{\text{SKIPPED ROUNDS}} + \underbrace{R_T(\tilde{\mathbf{u}}) - \min_{\mathbf{v} \in \mathcal{W}_1} R_1(\mathbf{v})}_{\text{PENALTY}} + \sum_{t \in \bar{\Lambda}} 8\beta_t^2 \\ &\quad - \underbrace{\sum_{t \in \bar{\Lambda}} \mathbb{E} [(\mathbf{w}_t(\widehat{\mathbf{L}}_t^m) - \tilde{\mathbf{u}})^\top \mathbf{b}_t \mid \varsigma]}_{\text{BIAS}} + \sum_{t \in \bar{\Lambda}} \underbrace{\left( 8\alpha_t^2 |m_t| + 8\alpha_t \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\boldsymbol{\ell}_\tau - \widehat{\boldsymbol{\ell}}_\tau) \right\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \mid \varsigma \right] \right)}_{\text{MISSING ESTIMATES}} \end{aligned} \quad (31)$$

Since we only start a new episode when any counter  $N_j$  doubles, we only start an new episode logarithmically often. This implies that  $|\Lambda| \leq d_{\max} HSA \log(T)$ , which means that the cost for the SKIPPED ROUNDS is upper bounded by:

$$\sum_{t \in \Lambda} (\mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}})^\top \ell_t \leq \sum_{t \in \Lambda} \mathbf{w}_t(\hat{\mathbf{L}}_t)^\top \ell_t - \tilde{\mathbf{u}}^\top \ell_t \leq d_{\max} H^2 SA \log(T). \quad (32)$$

where we used that  $\ell_t \in [0, 1]$  per assumption.

We now bound the BIAS term. We know that  $\mathbf{w}_t(\hat{\mathbf{L}}_t^m) \in \mathcal{D}_{R_t}(\mathbf{w}_t(\hat{\mathbf{L}}_t), \frac{1}{2\sqrt{\gamma}})$  when  $t \in \bar{\Lambda}$ , which is also shown in the proof of Lemma 3. By Lemma 19 we have that  $\mathbf{w}_t(\hat{\mathbf{L}}_t^m, h, s, a) \leq 2\mathbf{w}_t(\hat{\mathbf{L}}_t, h, s, a)$ . By definition  $\mathbf{q}_t^{\max}(h, s, a) \geq \mathbf{w}_t(\hat{\mathbf{L}}_t, h, s, a)$  and thus,

$$\begin{aligned} & -\mathbb{E} \left[ \sum_{t \in \bar{\Lambda}} \mathbf{w}_t(\hat{\mathbf{L}}_t^m)^\top \mathbf{b}_t \right] \\ &= \mathbb{E} \left[ \sum_{t \in \bar{\Lambda}} \sum_{h,s,a} \mathbf{w}_t(\hat{\mathbf{L}}_t^m, h, s, a) \ell_t(h, s, a) \left( 1 - \frac{\mathbf{q}^{\pi_t}(h, s, a)}{\mathbf{q}_t^{\max}(h, s, a) + \xi} \right) \right] \\ &= \mathbb{E} \left[ \sum_{t \in \bar{\Lambda}} \sum_{h,s,a} \frac{\mathbf{w}_t(\hat{\mathbf{L}}_t^m, h, s, a) \ell_t(h, s, a)}{\mathbf{q}_t^{\max}(h, s, a) + \xi} (\mathbf{q}_t^{\max}(h, s, a) - \mathbf{q}^{\pi_t}(h, s, a) + \xi) \right] \\ &\leq 2 \mathbb{E} \left[ \sum_{t \in \bar{\Lambda}} \sum_{h,s,a} \frac{\mathbf{w}_t(\hat{\mathbf{L}}_t, h, s, a) \ell_t(h, s, a)}{\mathbf{q}_t^{\max}(h, s, a) + \xi} (|\mathbf{q}_t^{\max}(h, s, a) - \mathbf{q}^{\pi_t}(h, s, a)| + \xi) \right] \\ &\leq 2 \mathbb{E} \left[ \sum_{t \in \bar{\Lambda}} \left( \sum_{h,s,a} |\mathbf{q}_t^{\max}(h, s, a) - \mathbf{q}^{\pi_t}(h, s, a)| + \xi HSA \right) \right] \\ &\lesssim \sqrt{H^4 S^2 AT \log(HSAT^3)} + H^3 S^3 A \log^2(HSAT^3) + H^3 S^2 A d_{\max} + HSA, \quad (33) \end{aligned}$$

where the last inequality is due to  $\xi = 1/T$  and Lemma 29, where we take an expectation over the event that the inequality in Lemma 29 holds similar to how we have been treating the good event and its complementary in other equations. For the second term in the BIAS, note that  $\mathbb{E}[\tilde{\mathbf{u}}^\top \mathbf{b}_t \mid \varsigma] \leq 0$ , since under the good event we have  $\mathbf{q}_\tau^{\max}(h, s, a) \geq \mathbf{q}^{\pi_\tau}(h, s, a)$  and thus,  $\mathbf{b}_t(h, s, a) \leq 0$ .

Next is the MISSING ESTIMATES term. We can not simply use the same argument as in Corollary 4 out of the box because  $\hat{\ell}_t$  is a biased estimator, so we add and subtract the bias  $\mathbf{b}_t$  and use the triangle inequality to find

$$\mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\ell_\tau - \hat{\ell}_\tau) \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \mid \varsigma \right] \leq \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\ell_\tau + \mathbf{b}_\tau - \hat{\ell}_\tau) \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} + \left\| \sum_{\tau \in m_t} \mathbf{b}_\tau \right\|_{R_t, \mathbf{w}_t(\hat{\mathbf{L}}_t)} \mid \varsigma \right]$$

Now we recognise that we are not using any information from rounds that we have not seen yet and thus  $\hat{\ell}_\tau$  and  $\hat{\ell}_{\tau'}$  are independent if  $\tau, \tau' \in m_t$ . Furthermore,  $\hat{\ell}_t$  is an unbiased estimator of  $\ell_t + \mathbf{b}_t$ ,

which allows us to use the exact same arguments as in Corollary 4 for the first term to find

$$\begin{aligned} \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\ell_\tau + \mathbf{b}_\tau - \widehat{\ell}_\tau) \right\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \middle| \varsigma \right] &\leq \sqrt{\sum_{\tau \in m_t} \left( \mathbb{E} [\|\widehat{\ell}_\tau\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)}^2] - \mathbb{E} [\|\ell_\tau + \mathbf{b}_\tau\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)}^2] \right)} \\ &\leq \sqrt{4|m_t|\beta_t^2} \end{aligned}$$

For the second term we start with the triangle inequality and Lemma 19,

$$\begin{aligned} \mathbb{E} \left[ \left\| \sum_{\tau \in m_t} \mathbf{b}_\tau \right\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \middle| \varsigma \right] &\leq \mathbb{E} \left[ \sum_{\tau \in m_t} \|\mathbf{b}_\tau\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \middle| \varsigma \right] \\ &\leq 2 \mathbb{E} \left[ \sum_{\tau \in m_t} \|\mathbf{b}_\tau\|_{R_t, \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)} \middle| \varsigma \right] \\ &\stackrel{(a)}{\leq} 2 \mathbb{E} \left[ \sum_{\tau \in m_t} \left( \|\ell_\tau + \mathbf{b}_\tau\|_{R_t, \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)} + \|\ell_\tau\|_{R_t, \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)} \right) \middle| \varsigma \right] \\ &\stackrel{(b)}{\leq} 2 \mathbb{E} \left[ \sum_{\tau \in m_t} \left( 2\|\ell_\tau\|_{R_t, \mathbf{w}_\tau(\widehat{\mathbf{L}}_\tau)} \right) \middle| \varsigma \right] \leq 4|m_t|\alpha_t, \end{aligned}$$

where we added and subtracted  $\ell_\tau$  and used the triangle inequality again in inequality (a). Inequality (b) holds as  $\mathbf{b}_\tau \leq 0$ , which holds as  $\mathbf{q}_\tau^{\max}(h, s, a) \geq \mathbf{q}^{\pi_\tau}(h, s, a)$  under the good event for all  $(h, s, a)$ . At the same time we have that  $\widehat{\ell}_\tau \geq 0$  by the construction of  $\widehat{\ell}_\tau$ , showing that  $\ell_\tau + \mathbf{b}_\tau = \mathbb{E}[\widehat{\ell}_\tau]$  is non-negative. Together both of those facts allow us to conclude that  $|\ell_\tau + \mathbf{b}_\tau| = \ell_\tau + \mathbf{b}_\tau \leq \ell_\tau$ , which we use in inequality (b).

Putting the last two equations together lets us bound the MISSING ESTIMATES term

$$\mathbb{E} \left[ \left\| \sum_{\tau \in m_t} (\ell_\tau - \widehat{\ell}_\tau) \right\|_{R_t, \mathbf{w}_t(\widehat{\mathbf{L}}_t)} \middle| \varsigma \right] \leq \sqrt{4|m_t|\beta_t^2} + 4|m_t|\alpha_t. \quad (34)$$

The last thing to bound is the PENALTY term. Using the third fact of Lemma 25 with  $b = \frac{1}{T^3 H^2 S^4 A^2}$ ,  $K = HS^2A$  and  $B = H$ , we conclude

$$\begin{aligned} \text{PENALTY} &\leq \frac{H(1 + \log(S^2A))}{\eta} + \frac{HS^2A \log(T^3 H^2 S^4 A^2)}{\gamma} \\ &\leq \frac{4H \log(SA)}{\eta} + \frac{4HS^2A \log(HSAT)}{\gamma}. \end{aligned} \quad (35)$$

Putting things together for the REG term gives

$$\begin{aligned}
 \text{REG} &= \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle \right] \\
 &= \mathbb{E} \left[ \mathbb{I}\{\varsigma\} \sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle \right] + \mathbb{E} \left[ \mathbb{I}\{\varsigma^c\} \sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle \right] \\
 &\leq \underbrace{\frac{4H \log(SA)}{\eta} + \frac{4HS^2A \log(HSAT)}{\gamma}}_{\text{PENALTY}} + \sum_{t \in \bar{\Lambda}} 8\beta_t^2 \\
 &\quad + \sum_{t \in \bar{\Lambda}} \left( 8\alpha_t^2 |m_t| + 8\alpha_t \left( \underbrace{\sqrt{4|m_t|\beta_t^2} + 4|m_t|\alpha_t}_{\text{MISSING ESTIMATES}} \right) \right) + \underbrace{d_{\max} H^2 SA \log(T)}_{\text{SKIPPED ROUNDS}} \\
 &\quad - \underbrace{\sum_{t \in \bar{\Lambda}} \mathbb{E} [(\mathbf{w}_t(\hat{\mathbf{L}}_t^m) - \mathbf{u})^\top \mathbf{b}_t | \varsigma]}_{\text{BIAS}} + 4H \\
 &\leq \frac{4H \log(SA)}{\eta} + \frac{4HS^2A \log(HSAT)}{\gamma} + 136\eta HSAT \\
 &\quad + 148\eta HD + d_{\max} H^2 SA \log(T) \\
 &\quad - \sum_{t \in \bar{\Lambda}} \mathbb{E} [(\mathbf{w}_t(\hat{\mathbf{L}}_t^m) - \mathbf{u})^\top \mathbf{b}_t | \varsigma] + 4H,
 \end{aligned}$$

where we used equations (28), (31), (32), (34), (35) in the first inequality and plugged in the values of  $\alpha$  and  $\beta$  we found in equations (29) and (30) and also used that  $\sqrt{ab} \leq \frac{1}{2}(a+b)$  for  $a, b > 0$ . We plug in the learning rates to find

$$\begin{aligned}
 \text{REG} &\lesssim \frac{H \log(SA)}{\eta} + \frac{HS^2A \log(HSAT)}{\gamma} + \eta HSAT \\
 &\quad + \eta HD + d_{\max} H^2 SA \log(T) \\
 &\quad - \sum_{t \in \bar{\Lambda}} \mathbb{E} [(\mathbf{w}_t(\hat{\mathbf{L}}_t^m) - \mathbf{u})^\top \mathbf{b}_t | \varsigma] + H \\
 &\lesssim H \sqrt{SAT \log(SA)} + H \sqrt{D \log(SA)} \\
 &\quad + d_{\max} \sqrt{H} S^2 A \log(HSAT) + d_{\max} H^2 SA \log(T) \\
 &\quad - \sum_{t \in \bar{\Lambda}} \mathbb{E} [(\mathbf{w}_t(\hat{\mathbf{L}}_t^m) - \mathbf{u})^\top \mathbf{b}_t | \varsigma] + H. \tag{36}
 \end{aligned}$$

We can finally put everything together, starting from the regret again

$$\begin{aligned}
 \mathcal{R}_T &= \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{q}^{\pi_t} - \mathbf{w}_t(\hat{\mathbf{L}}_t), \ell_t \rangle + \sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle + \sum_{t=1}^T \langle \tilde{\mathbf{u}} - \mathbf{u}, \ell_t \rangle \right] \quad (\text{Eqn (25)}) \\
 &\lesssim \mathbb{E} \left[ \sum_{t=1}^T \langle \mathbf{w}_t(\hat{\mathbf{L}}_t) - \tilde{\mathbf{u}}, \ell_t \rangle \right] + \sqrt{H^4 S^2 A T \log(HSAT^3)} \\
 &\quad + H^3 S^3 A \log^2(HSAT^3) + d_{\max} H^3 S^2 A \quad (\text{Eqns (26) and (27)}) \\
 &\lesssim \sqrt{H^4 S^2 A T \log(HSAT^3)} + H \sqrt{D \log(SA)} + d_{\max} H^3 S^2 A \quad (\text{Eqns (36) and (33)}) \\
 &\quad + H^3 S^3 A \log^2(HSAT^3) + d_{\max} \sqrt{H} S^2 A \log(HSAT) + d_{\max} H^2 S A \log(T),
 \end{aligned}$$

which concludes the proof.  $\blacksquare$

## Appendix E. Doubling with Delayed Feedback

In this section we show how to handle unknown problem parameters. For simplicity of presentation we assume that only  $d_{\max}$  is unknown. The case of unknown  $T$  and  $D$  can be done in a similar fashion (e.g., see Bistritz et al. (2019); Lancewicki et al. (2022a)).

---

### Algorithm 6: Doubling procedure

---

**Input:**  $T, D$  and algorithm  $ALG$  (for known  $T, D$  and  $d_{\max}$ ).

Set epoch index  $e = 1$  and initialize  $ALG$  with  $T, D$  and  $2^e$  as  $d_{\max}$ .

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

**if**  $\max_{j \in \mathcal{O}_t} d_j \geq 2^e$  **then**

        Start a new epoch  $e = e + 1$ , and re-initiate  $ALG$  with  $T, D$  and  $2^e$  as  $d_{\max}$ .

**end if**

    Play according to  $ALG$ .

**end for**

---

**Theorem 16** *Let  $ALG$  be an algorithm for known  $T, D$  and  $d_{\max}$  and assume that  $ALG$  guarantees regret of  $R_{T,D}(d_{\max})$  whenever initiated properly. Then, running Algorithm 6 with unknown  $d_{\max}$  guarantees regret,*

$$\mathcal{R}_T \leq 2R_{T,D}(2d_{\max}) \log T + 2Md_{\max} \log T,$$

where  $M = \max_{t \in [T], \mathbf{a}, \tilde{\mathbf{a}} \in \mathcal{A}} (\mathbf{a} - \tilde{\mathbf{a}})^\top \ell_t$  is the maximal regret per round (e.g., in Section 6,  $M \leq H$ ).

**Proof** Let  $\mathcal{T}_e = \{t : 2^{e-1} \leq \max_{j \in \mathcal{O}_t} d_j \leq 2^e\}$  be the set of indices of epoch  $e$ , and let  $\tilde{\mathcal{T}}_e = \{t \in \mathcal{T}_e : d_t \leq 2^e\}$  be the indices of epoch  $e$  with delay  $\leq 2^e$ . The regret in rounds  $t \in \tilde{\mathcal{T}}_e$  is at most  $R_{T,D}(2^e) \leq R_{T,D}(2d_{\max})$  since the maximal delay in these rounds is indeed bounded by  $2^e$ . In addition, the regret in  $\mathcal{T}_e \setminus \tilde{\mathcal{T}}_e$  is at most  $Md_{\max}$  since  $|\mathcal{T}_e \setminus \tilde{\mathcal{T}}_e| \leq d_{\max}$ . Thus, the total regret in epoch  $e$  is at most,

$$\underbrace{R_{T,D}(2d_{\max})}_{\text{Regret in } \tilde{\mathcal{T}}_e} + \underbrace{Md_{\max}}_{\text{Regret in } \mathcal{T}_e \setminus \tilde{\mathcal{T}}_e}.$$

Finally, the total number of epochs is at most  $\log d_{max} + 1 \leq 2 \log T$  and thus, the total regret is bounded by,

$$\mathcal{R}_T \leq 2R_{T,D}(2d_{max}) \log T + 2Md_{max} \log T.$$

■

## Appendix F. Auxiliary Lemmas

**Lemma 17** *Let  $t \in [T]$  and suppose that  $4\nabla^2 R_t(\mathbf{u}) \succeq \nabla^2 R_t(\mathbf{u}') \succeq \frac{1}{4}\nabla^2 R_t(\mathbf{u})$  for all  $\mathbf{u} \in \mathcal{W}_t$  and  $\mathbf{u}' \in \mathcal{D}_{R_t}(\mathbf{u}, \frac{1}{2})$ . Let  $\mathbf{v}' \in \mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2})$  or  $\mathbf{v} \in \mathcal{D}_{R_t}(\mathbf{v}', \frac{1}{2})$ , then*

$$\|x\|_{R_t, \mathbf{v}'} \leq 2\|x\|_{R_t, \mathbf{v}},$$

for all  $x \in \mathbb{R}^K$ .

**Proof** First consider that if  $\mathbf{v}' \in \mathcal{D}_{R_t}(\mathbf{v}, \frac{1}{2})$  then  $\nabla^2 R_t(\mathbf{v}') \succeq \frac{1}{4}\nabla^2 R_t(\mathbf{v})$  and thus

$$(\nabla^2 R_t(\mathbf{v}'))^{-1} \preceq 4(\nabla^2 R_t(\mathbf{v}))^{-1}.$$

We can arrive to the same inequality if  $\mathbf{v} \in \mathcal{D}_{R_t}(\mathbf{v}', \frac{1}{2})$ , by using  $4\nabla^2 R_t(\mathbf{v}') \succeq \nabla^2 R_t(\mathbf{v})$ . We can then follow directly

$$\|x\|_{R_t, \mathbf{v}'} = \sqrt{x^\top (\nabla^2 R_t(\mathbf{v}'))^{-1} x} \leq 2\sqrt{x^\top (\nabla^2 R_t(\mathbf{v}))^{-1} x} = 2\|x\|_{R_t, \mathbf{v}}.$$

■

**Lemma 18 (Be-The-Leader Lemma)** *Let  $\tilde{\mathbf{w}}_t(\hat{\mathbf{L}}_t^*) = \arg \min_{\mathbf{w} \in \mathcal{W}_t} \mathbf{w}^\top \hat{\mathbf{L}}_t^* + R_t(\mathbf{w})$ . Suppose that  $R_t(\mathbf{v}) \leq R_{t+1}(\mathbf{v})$  for all  $\mathbf{v} \in \mathcal{W}_t$  and all  $t \in [T]$  and that  $\mathcal{W}_t \subseteq \mathcal{W}_{t-1}$  is a non-empty compact convex set. Then, for any fixed  $\mathbf{u} \in \mathcal{W}_T$ , we have that*

$$\sum_{t \in \bar{\Lambda}} \hat{\ell}_t^\top (\tilde{\mathbf{w}}_t(\hat{\mathbf{L}}_t^*) - \mathbf{u}) \leq R_T(\mathbf{u}) - \min_{\mathbf{v} \in \mathcal{W}_1} R_1(\mathbf{v})$$

**Proof** We will prove the statement by induction on  $T$ . For the induction step, assume that

$$\sum_{t \in \bar{\Lambda} \cap [T-1]} \hat{\ell}_t^\top \tilde{\mathbf{w}}_t(\hat{\mathbf{L}}_t^*) + R_1(\mathbf{w}_1(\hat{\mathbf{L}}_1)) \leq \sum_{t \in \bar{\Lambda} \cap [T-1]} \hat{\ell}_t^\top \mathbf{v} + R_{T-1}(\mathbf{v})$$

for any  $\mathbf{v} \in \mathcal{W}_T$ . If  $\bar{\Lambda} \cap [T-1] = \bar{\Lambda} \cap [T]$  the induction step holds. Otherwise  $T \in \bar{\Lambda}$  and adding  $\hat{\ell}_T^\top \tilde{\mathbf{w}}_T(\hat{\mathbf{L}}_T^*)$  to both sides of the above inequality and setting  $\mathbf{v} = \tilde{\mathbf{w}}_T(\hat{\mathbf{L}}_T^*)$  on the right-hand side of the above inequality we find

$$\begin{aligned} \sum_{t \in \bar{\Lambda} \cap [T]} \hat{\ell}_t^\top \tilde{\mathbf{w}}_t(\hat{\mathbf{L}}_t^*) + R_1(\mathbf{w}_1(\hat{\mathbf{L}}_1)) &\leq \sum_{t \in \bar{\Lambda} \cap [T]} \hat{\ell}_t^\top \tilde{\mathbf{w}}_t(\hat{\mathbf{L}}_t^*) + R_{T-1}(\tilde{\mathbf{w}}_T(\hat{\mathbf{L}}_T^*)) \\ &\leq \sum_{t \in \bar{\Lambda} \cap [T]} \hat{\ell}_t^\top \tilde{\mathbf{w}}_T(\hat{\mathbf{L}}_T^*) + R_T(\tilde{\mathbf{w}}_T(\hat{\mathbf{L}}_T^*)) \\ &\leq \sum_{t \in \bar{\Lambda} \cap [T]} \hat{\ell}_t^\top \mathbf{u} + R_T(\mathbf{u}), \end{aligned}$$

which proves the induction step after reordering and observing that the base case holds by definition of  $w_1(\widehat{L}_1)$ . The statment is proven after applying  $R_T \geq R_\tau$ , which holds for all  $\tau \in \overline{\Lambda}$ , once. ■

**Lemma 19** *Let  $\mathcal{V} \subseteq \{\mathbf{x} \in \mathbb{R}^n : \forall i \in [n], \mathbf{x}(i) > 0\}$ . Let  $R : \mathcal{V} \rightarrow \mathbb{R}$  be some twice-differentiable convex function, and let  $\phi(\mathbf{v}) = -\frac{1}{\gamma} \sum_{i=1}^n \log \mathbf{v}(i)$  be the log barrier with  $\gamma \in (0, 1)$ . Assume that for any  $\mathbf{v} \in \mathcal{V}$ ,  $\nabla^2 R(\mathbf{v}) \succeq \nabla^2 \phi(\mathbf{v})$ . Then for any  $\mathbf{v}' \in \mathcal{D}_R(\mathbf{v}, \frac{1}{2\sqrt{\gamma}})$  and all  $i \in [n]$ ,*

$$\frac{1}{2}\mathbf{v}(i) \leq \mathbf{v}'(i) \leq 2\mathbf{v}(i) .$$

**Proof** Since  $\nabla^2 R(\mathbf{v}) \succeq \nabla^2 \phi(\mathbf{v})$ , for any  $\mathbf{v}' \in \mathcal{D}_R(\mathbf{v}, \frac{1}{2\sqrt{\gamma}})$ ,

$$(\|\mathbf{v}' - \mathbf{v}\|_{\phi, \mathbf{v}}^*)^2 \leq (\|\mathbf{v}' - \mathbf{v}\|_{R, \mathbf{v}}^*)^2 \leq \frac{1}{4\gamma} .$$

On the other hand,

$$(\|\mathbf{v}' - \mathbf{v}\|_{\phi, \mathbf{v}}^*)^2 = \sum_{j=1}^n \frac{(\mathbf{v}'(j) - \mathbf{v}(j))^2}{\gamma \mathbf{v}(j)^2} \geq \frac{(\mathbf{v}'(i) - \mathbf{v}(i))^2}{\gamma \mathbf{v}(i)^2} .$$

Thus,  $|\mathbf{v}'(i) - \mathbf{v}(i)| \leq \frac{1}{2}\mathbf{v}(i)$  which implies that  $\frac{1}{2}\mathbf{v}(i) \leq \mathbf{v}'(i) \leq 2\mathbf{v}(i)$ . ■

**Lemma 20** *Let  $a, b \in \mathbb{R}_+$  such that  $a \geq b$ , then*

$$\sqrt{a} - \sqrt{b} \leq \sqrt{a - b} .$$

**Proof** We show directly

$$\sqrt{a} - \sqrt{b} = \sqrt{(\sqrt{a} - \sqrt{b})^2} = \sqrt{a + b - 2\sqrt{ab}} \leq \sqrt{a - b} .$$

■

**Lemma 21**  $\log(x)^2 \leq \frac{1}{x}$  for all  $0 < x \leq 1$ .

**Proof** Note that since  $\log(x) \leq 0 \leq 1/x$  for  $0 < x \leq 1$ ,  $\log(x)^2 \leq \frac{1}{x}$  is equivalent to  $-\log(x) \leq \frac{1}{\sqrt{x}}$  which we rearrange to  $-\sqrt{x} \log(x) \leq 1$ . We maximise the function on the lefthandside on  $x \in (0, 1]$ , taking a derivative yields

$$\frac{\partial -\sqrt{x} \log(x)}{\partial x} = -\frac{1}{2\sqrt{x}} \log(x) - \frac{1}{\sqrt{x}} .$$

Setting the derivative to 0 gives  $x = e^{-2}$  as a possible maximum and  $-\sqrt{e^{-2}} \log(e^{-2}) = \frac{2}{e} < 1$ . The supremum of  $-\sqrt{x} \log(x)$  may also lie on the boundary of  $(0, 1]$  but  $-\sqrt{1} \log(1) = 0 < 1$  and

$$\lim_{x \rightarrow 0+} -\sqrt{x} \log(x) = \lim_{x \rightarrow 0+} \frac{-\log(x)}{\frac{1}{\sqrt{x}}} = \lim_{x \rightarrow 0+} \frac{x^{-\frac{1}{2}}}{2x^{-\frac{3}{2}}} = 0 < 1 ,$$

where we also used L'Hôpital's rule. We conclude that  $-\sqrt{x} \log(x) \leq 1$  for all  $0 < x \leq 1$ . ■

**Lemma 22** *Let  $a, b, c \in \mathbb{R}$ . Let  $b \geq c$ , then*

$$\max\{a, b\} - \max\{a, c\} \leq b - c$$

**Proof** If  $a \geq b$ , then  $\max\{a, b\} - \max\{a, c\} = a - a = 0 \leq b - c$ .

If  $b \geq a \geq c$ , then  $\max\{a, b\} - \max\{a, c\} = b - a \leq b - c$ .

If  $b, c \geq a$ , then  $\max\{a, b\} - \max\{a, c\} = b - c$ . ■

**Lemma 23 (Part of Lemma 14 from Gaillard et al. (2014))** *Let  $a_1, \dots, a_m \in \mathbb{R}_+$  and call  $s_i = a_1 + \dots + a_i$ . Let  $f : (0, \infty) \rightarrow [0, \infty]$  be a non-increasing function. Then*

$$\sum_{i=1}^m a_i f(s_i) \leq \int_{a_1}^{s_m} f(x) dx$$

**Proof**

$$\sum_{i=1}^m a_i f(s_i) = \sum_{i=1}^m \int_{s_{i-1}}^{s_i} f(s_i) dx \leq \sum_{i=1}^m \int_{s_{i-1}}^{s_i} f(x) dx = \int_{a_1}^{s_m} f(x) dx,$$

where we used a telescoping sum in the first equality, the fact that  $f$  is non-increasing in the inequality and another telescoping sum in the last equality. ■

**Lemma 24**

$$\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq 2\sqrt{T}, \quad \sum_{t=1}^T \frac{|m_t|}{\sqrt{\sum_{\tau=1}^t |m_\tau|}} \leq 2\sqrt{D}, \quad \sum_{t=1}^T t^{-\frac{1}{4}} \leq \frac{4}{3} T^{\frac{3}{4}}.$$

**Proof** By Lemma 23 with  $a_1, \dots, a_T = 1$  and  $f(x) = \frac{1}{\sqrt{x}}$

$$\sum_{t=1}^T \frac{1}{\sqrt{t}} \leq \int_1^T \frac{1}{\sqrt{x}} dx \leq 2\sqrt{T}.$$

We replace  $m_t$  by  $\tilde{m}_t$ , where we used that  $|m_t| \leq |\tilde{m}_t|$  which holds with probability one, then by Lemma 23 with  $a_i = |\tilde{m}_i|$  and  $f(x) = \frac{1}{\sqrt{x}}$

$$\sum_{t=1}^T \frac{|m_t|}{\sqrt{\sum_{\tau=1}^t |m_\tau|}} \leq \sum_{t=1}^T \frac{|\tilde{m}_t|}{\sqrt{\sum_{\tau=1}^t |\tilde{m}_\tau|}} \leq \int_1^D \frac{1}{\sqrt{x}} dx \leq 2\sqrt{D}.$$



One last time by Lemma 23 with  $a_1, \dots, a_T = 1$  and  $f(x) = x^{-\frac{1}{4}}$

$$\sum_{t=1}^T t^{-\frac{1}{4}} \leq \int_1^T x^{-\frac{1}{4}} dx \leq \frac{4}{3} T^{\frac{3}{4}}.$$

■

**Lemma 25** *Let  $\mathcal{V}(b) \subseteq \{\mathbf{x} : 0 \leq b \leq \mathbf{x}(i) \leq 1\}$  and let  $\Gamma_t(\mathbf{v}) = \sum_{i=1}^K \left( \frac{1}{\eta_t} \mathbf{v}(i) \log(\mathbf{v}(i)) - \frac{1}{\gamma} \log(\mathbf{v}(i)) \right)$  for some  $\gamma, \eta_t > 0$  and  $\eta_t \geq \eta_{t+1}$ , then*

$$4\nabla^2 \Gamma_t(\mathbf{v}) \succeq \nabla^2 \Gamma_t(\mathbf{v}') \succeq \frac{1}{4} \nabla^2 \Gamma_t(\mathbf{v}),$$

*for all  $\mathbf{v}', \mathbf{v} \in \mathcal{V}(b)$ ,  $\mathbf{v}' \in \mathcal{D}_{\Gamma_t}(\mathbf{v}, \frac{1}{2\sqrt{\gamma}})$ , and all  $t$ . Furthermore, if there exists an  $\lambda > 0$  and  $\eta_t$  is such that  $\sqrt{\eta_{t+\delta}} \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \leq \frac{\sqrt{\lambda}}{\sqrt{K}}$  for a given  $t$  and  $\delta \in [d_{\max}]$ , then*

$$(\nabla \Gamma_t(\mathbf{v}) - \nabla \Gamma_{t+\delta}(\mathbf{v}))^\top y \leq \sqrt{\lambda} \sqrt{y^\top \nabla^2 \Gamma_{t+\delta}(\mathbf{v}) y},$$

*for all  $\mathbf{v} \in \mathcal{V}(b)$  and all  $\mathbf{y} \in \mathbb{R}^K$ . Finally, if  $b' \geq 0$ ,  $b > 0$ , and  $\|\mathbf{v}\|_1 \leq B$  for some  $B > 0$  and all  $\mathbf{v}(i) \in \mathcal{V}(b')$ , then for all  $\mathbf{u} \in \mathcal{V}(b)$*

$$\Gamma_T(\mathbf{u}) - \min_{\mathbf{v} \in \mathcal{V}(b')} \Gamma_T(\mathbf{v}) \leq \frac{B(1 + \log(\frac{K}{B}))}{\eta_T} + \frac{K \log(\frac{1}{b})}{\gamma}.$$

**Proof** We start with the first statement and we state the derivatives of  $\Gamma_t$

$$\begin{aligned} \Gamma_t(\mathbf{v}) &= \sum_{i=1}^K \left( \frac{1}{\eta_t} \mathbf{v}(i) \log(\mathbf{v}(i)) - \frac{1}{\gamma} \log(\mathbf{v}(i)) \right) \\ (\nabla \Gamma_t(\mathbf{v}))(i) &= \frac{1}{\eta_t} \log(\mathbf{v}(i)) - \frac{1}{\gamma \mathbf{v}(i)} \\ (\nabla^2 \Gamma_t(\mathbf{v}))(i, i) &= \frac{1}{\eta_t \mathbf{v}(i)} + \frac{1}{\gamma \mathbf{v}^2(i)}, \end{aligned}$$

where  $(\nabla^2 \Gamma_t(\mathbf{v}))(j, i) = 0$  if  $j \neq i$ . Now, we have that

$$\frac{1}{2\sqrt{\gamma}} \geq \|\mathbf{v} - \mathbf{v}'\|_{\Gamma_t, \mathbf{v}} \geq \frac{|\mathbf{v}(i) - \mathbf{v}'(i)|}{\sqrt{\gamma \mathbf{v}(i)}} = \frac{1}{\sqrt{\gamma}} \left| 1 - \frac{\mathbf{v}'(i)}{\mathbf{v}(i)} \right|,$$

or equivalently,  $\left| 1 - \frac{\mathbf{v}'(i)}{\mathbf{v}(i)} \right| \leq \frac{1}{2}$ . If  $\mathbf{v}(i) \geq \mathbf{v}'(i)$  then  $\left| 1 - \frac{\mathbf{v}'(i)}{\mathbf{v}(i)} \right| = 1 - \frac{\mathbf{v}'(i)}{\mathbf{v}(i)}$  and we re-arrange to find  $\mathbf{v}'(i) \geq \frac{1}{2} \mathbf{v}(i)$ . Likewise, if  $\mathbf{v}(i) \leq \mathbf{v}'(i)$  then we can see that  $\mathbf{v}'(i) \leq \frac{3}{2} \mathbf{v}(i)$ . Thus, we can conclude that for  $\mathbf{v}' \in \mathcal{D}_{\Gamma_t}(\mathbf{v}, \frac{1}{2\sqrt{\gamma}})$

$$\frac{1}{2} \mathbf{v}(i) \leq \mathbf{v}'(i) \leq \frac{3}{2} \mathbf{v}(i). \quad (37)$$

Using these properties and the second derivative of  $\Gamma_t$  as written above we can verify the first statement as

$$4\nabla^2\Gamma_t(\mathbf{v}) = 4 \operatorname{diag} \left( \frac{\gamma}{\eta_t \mathbf{v}} + \frac{1}{\mathbf{v}^2} \right) \succeq \operatorname{diag} \left( \frac{\gamma}{\eta_t \mathbf{v}'} + \frac{1}{\mathbf{v}'^2} \right) = \nabla^2\Gamma_t(\mathbf{v}'),$$

where the division is meant elementwise and  $\nabla^2\Gamma_t(\mathbf{v}') \succeq \frac{1}{4}\nabla^2\Gamma_t(\mathbf{v})$  goes through analogously. Next for the second statement, we first pick any  $\mathbf{y} \in \mathbb{R}^K$  and then establish that

$$\sum_{i=1}^K -|y(i)| \log(\mathbf{v}(i)) \leq \sqrt{K} \sqrt{\sum_{i=1}^K y(i)^2 \log(\mathbf{v}(i))^2} \leq \sqrt{K} \sqrt{\sum_{i=1}^K y(i)^2 \frac{1}{\mathbf{v}(i)}}$$

using the AM-QM inequality, which holds as all  $-|y(i)| \log(\mathbf{v}(i))$  are real positive numbers and using the fact that  $\log(x)^2 \leq \frac{1}{x}$  for  $0 < x \leq 1$  as shown by Lemma 21. Using the above equation gives

$$\begin{aligned} (\nabla\Gamma_t(\mathbf{v}) - \nabla\Gamma_{t+\delta}(\mathbf{v}))^\top \mathbf{y} &\leq \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \sum_{i=1}^K -|y(i)| \log(\mathbf{v}(i)) \\ &\leq \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \sqrt{K} \sqrt{\sum_{i=1}^K y(i)^2 \frac{1}{\mathbf{v}(i)}} \\ &= \sqrt{\eta_{t+\delta}} \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \sqrt{K} \sqrt{\sum_{i=1}^K \frac{1}{\eta_{t+\delta}} y(i)^2 \frac{1}{\mathbf{v}(i)}} \\ &\leq \sqrt{\lambda} \sqrt{\sum_{i=1}^K \left( \frac{1}{\eta_{t+\delta}} y(i)^2 \frac{1}{\mathbf{v}(i)} + \frac{1}{\gamma} \frac{y(i)^2}{\mathbf{v}(i)^2} \right)} \\ &= \sqrt{\lambda} \sqrt{\mathbf{y}^\top \nabla^2\Gamma_{t+\delta}(\mathbf{v}) \mathbf{y}}, \end{aligned}$$

where we only used  $|y(i)| \geq y(i)$  in the first inequality, the above equation in the second inequality, and the assumption on  $\eta_t$  and  $\lambda$  and the fact that  $\frac{1}{\gamma} \frac{y(i)^2}{\mathbf{v}(i)^2} \geq 0$  in the last inequality.

For the last statement we start with the negative entropy component of  $\Gamma_T$ . Without loss of generality we may assume that  $\mathbf{v}(i) > 0$  as we may define  $-\mathbf{v}(i) \log(\mathbf{v}(i)) = 0$ . We can bound the negative entropy component of  $\Gamma_T$  as

$$\begin{aligned} -\sum_{i=1}^K \mathbf{v}(i) \log \mathbf{v}(i) &= \|\mathbf{v}\|_1 \sum_{i=1}^K \frac{\mathbf{v}(i)}{\|\mathbf{v}\|_1} \log \frac{1}{\mathbf{v}(i)} \\ &\leq \|\mathbf{v}\|_1 \log \left( \sum_{i=1}^K \frac{\mathbf{v}(i)}{\|\mathbf{v}\|_1} \frac{1}{\mathbf{v}(i)} \right) \\ &\leq \|\mathbf{v}\|_1 \log \left( \frac{K}{\|\mathbf{v}\|_1} \right) + \|\mathbf{v}\|_1 \leq B \left( 1 + \log \left( \frac{K}{B} \right) \right), \end{aligned}$$

where we used Jensen's inequality in the second step and the fact that  $x \log(\frac{K}{x}) + x$  is increasing on  $x \in [1, K]$  in the last inequality. Set  $\mathbf{v}^+ = \arg \min_{\mathbf{v} \in \mathcal{W}(b)} \Gamma_T(\mathbf{v})$ , then

$$\begin{aligned} \Gamma_T(\mathbf{u}) - \Gamma_T(\mathbf{v}^+) &= \sum_{i=1}^K \left( \frac{\mathbf{u}(i)}{\eta_T} \log(\mathbf{u}(i)) - \frac{1}{\gamma} \log(\mathbf{u}(i)) \right) \\ &\quad - \sum_{i=1}^K \left( \frac{\mathbf{v}^+(i)}{\eta_T} \log(\mathbf{v}^+(i)) - \frac{1}{\gamma} \log(\mathbf{v}^+(i)) \right) \\ &\leq \frac{B(1 + \log(\frac{K}{B}))}{\eta_T} + \frac{K \log(\frac{1}{b})}{\gamma} \end{aligned}$$

where we used the fact that  $b \leq \mathbf{u}(i) \leq 1$  since  $\mathbf{u}(i) \in \mathcal{V}(b)$ ,  $\frac{1}{\eta_1} \leq \frac{1}{\eta_T}$ , and the fact that  $-\log(x)$  is a decreasing function and non-negative for  $x \in (0, 1]$ .  $\blacksquare$

**Lemma 26** Let  $\eta_t = \min \left\{ a, \frac{1}{\sqrt{bt+c \sum_{\tau=1}^t |m_\tau|}} \right\}$  for some  $a, b, c \in \mathbb{R}_+$ . If  $a \leq \frac{d}{bd_{\max} + cd_{\max}^2}$  for some  $d \in \mathbb{R}$ , then

$$\sqrt{\eta_{t+\delta}} \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \leq \sqrt{d}.$$

**Proof** We start by showing that

$$\begin{aligned} \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} &\leq \sqrt{b(t+\delta) + c \sum_{\tau=1}^{t+\delta} |m_\tau|} - \sqrt{bt + c \sum_{\tau=1}^t |m_\tau|} \\ &\leq \sqrt{b\delta + c \sum_{\tau=t+1}^{t+\delta} |m_\tau|} \\ &\leq \sqrt{bd_{\max} + cd_{\max}^2}, \end{aligned}$$

where we used our assumption on  $\eta_t$  together with  $\max\{x, y\} - \max\{x, z\} \leq y - z$  (Lemma 22) in the first inequality,  $\sqrt{a} - \sqrt{b} \leq \sqrt{a-b}$  (Lemma 20) in the second inequality and the fact that  $\delta \leq d_{\max}$  and  $|m_t| \leq d_{\max}$  in the third inequality. And from here we can see that

$$\sqrt{\eta_{t+\delta}} \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \leq \sqrt{\eta_{t+\delta}} \sqrt{bd_{\max} + cd_{\max}^2} \leq \sqrt{d}.$$

$\blacksquare$

**Lemma 27** Let  $\eta_t = \min \left\{ a, \frac{1}{\sqrt{bt+c \sum_{\tau=1}^t |m_\tau|}} \right\}$  for some  $a, b, c \in \mathbb{R}_+$ . If  $a \leq \frac{d}{\sqrt{bd_{\max} + cd_{\max}^2}}$  for some  $d \in \mathbb{R}$ , then

$$\sqrt{\eta_{t+\delta}} \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \leq \sqrt{d} (bd_{\max} + cd_{\max}^2)^{1/4}.$$

**Proof** We start by showing that

$$\begin{aligned} \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} &\leq \sqrt{b(t+\delta) + c \sum_{\tau=1}^{t+\delta} |m_\tau|} - \sqrt{bt + c \sum_{\tau=1}^t |m_\tau|} \\ &\leq \sqrt{b\delta + c \sum_{\tau=t+1}^{t+\delta} |m_\tau|} \\ &\leq \sqrt{bd_{\max} + cd_{\max}^2}, \end{aligned}$$

where we used our assumption on  $\eta_t$  together with  $\max\{x, y\} - \max\{x, z\} \leq y - z$  (Lemma 22) in the first inequality,  $\sqrt{a} - \sqrt{b} \leq \sqrt{a-b}$  (Lemma 20) in the second inequality and the fact that  $\delta \leq d_{\max}$  and  $|m_t| \leq d_{\max}$  in the third inequality. And from here we can see that

$$\sqrt{\eta_{t+\delta}} \left( \frac{1}{\eta_{t+\delta}} - \frac{1}{\eta_t} \right) \leq \sqrt{\eta_{t+\delta}} \sqrt{bd_{\max} + cd_{\max}^2} \leq \sqrt{d}.$$

■

**Lemma 28 (Lemma D.11 of Jin et al. (2022); see also Lemma 4 of Jin et al. (2020))** *With probability  $1 - \delta$ , for any collection of transition functions  $\{p_{t,h}^s\}_{s \in \mathcal{S}}$  such that  $p_t^s \in \hat{P}_j$*

$$\begin{aligned} \sum_{t=1}^T \sum_{h=1}^H \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} |q_t^{p_t^s, \pi_t}(h, s, a) - q_t^{\pi_t}(h, s, a)| &\lesssim H \sum_{t=1}^T \sum_{h=1}^H \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \epsilon_t(h, s, a) q^{\pi_t}(h, s, a) \\ &+ HS \sum_{t=1}^T \sum_{1 \leq h < \tilde{h} \leq H} \sum_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}} \sum_{\tilde{s} \in \mathcal{S}, \tilde{a} \in \mathcal{A}} \epsilon_t(s' | h, s, a) q^{\pi_t}(h, s, a) \\ &\cdot \min \left\{ 2, \sum_{\tilde{s}' \in \mathcal{S}} \epsilon_t(\tilde{s}' | \tilde{h}, \tilde{s}, \tilde{a}) \right\} q^{\pi_t}(\tilde{h}, \tilde{s}, \tilde{a} | s'; h+1) + H^3 S^2 Ad_{\max} \end{aligned} \quad (38)$$

where  $q^{\pi_t}(\tilde{h}, \tilde{s}, \tilde{a} | \tilde{s}'; h)$  be the probability to visit  $(\tilde{s}, \tilde{a})$  in time  $\tilde{h}$  given that we visited  $\tilde{s}'$  in time  $h$ .

**Lemma 29 (Lemma D.12 of Jin et al. (2022) adapted to epochs)** *With probability  $1 - 10/T$ , for any collection of transition functions  $\{p_t^{h,s}\}_{h \in [H], s \in \mathcal{S}}$  such that  $p_t^{h,s} \in \hat{P}_j$*

$$\begin{aligned} \sum_{t=1}^T \sum_{h=1}^H \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} |q_t^{p_t^{h,s}, \pi_t}(h, s, a) - q_t^{\pi_t}(h, s, a)| \\ \lesssim \sqrt{H^4 S^2 AT \log(HSAT^3)} + H^3 S^3 A \log^2(HSAT^3) + H^3 S^2 Ad_{\max}. \end{aligned}$$

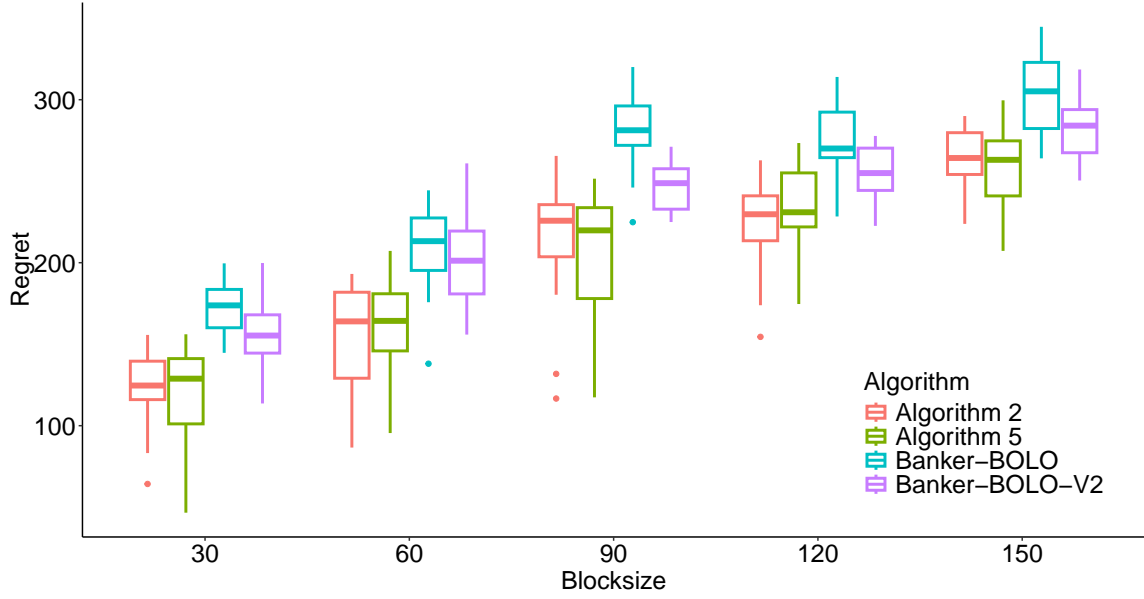


Figure 3: Boxplot of the regret over 20 repetitions over  $T = 10000$  rounds with  $K = 10$ .

**Proof** Following the exact same steps as in the proof of Lemma E.5 in Lancewicki et al. (2022b), with probability of at least  $1 - \delta$ , the first term in (38) can be bounded by,

$$\begin{aligned} \sum_{t=1}^T \sum_{h=1}^H \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \epsilon_t(h, s, a) q^{\pi_t}(h, s, a) &\lesssim \sqrt{S \log(HSAT^3)} \sum_{t,h,s,a} \frac{\mathbb{I}\{s_{t,h} = s, a_{t,h} = a\}}{\sqrt{N_{j(t)}(h, s, a) \vee 1}} \\ &+ S \log(HSAT^3) \sum_{t,h,s,a} \frac{\mathbb{I}\{s_{t,h} = s, a_{t,h} = a\}}{N_{j(t)}(h, s, a) \vee 1} + HS \log^2(HSAT^3) \end{aligned} \quad (39)$$

Similarly, the second summation (38) is bounded by,

$$HS \log^2(HSAT^3) \sum_{t,h,s,a} \frac{\mathbb{I}\{s_{t,h} = s, a_{t,h} = a\}}{N_{j(t)}(h, s, a) \vee 1} \quad (40)$$

Finally, (39) and (40) are bounded by  $O(HS \sqrt{AT \log(HSAT^3)} + HS^2 \log^2(HSAT^3))$  and  $O(H^2 S^2 A \log^2(HSAT^3))$  respectively using standard arguments - see for example the proof of Lemma 10 in Jin et al. (2020).  $\blacksquare$

## Appendix G. Further Results of the Experiments

### References

Jacob Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In *Conference on Learning Theory*, 2008.

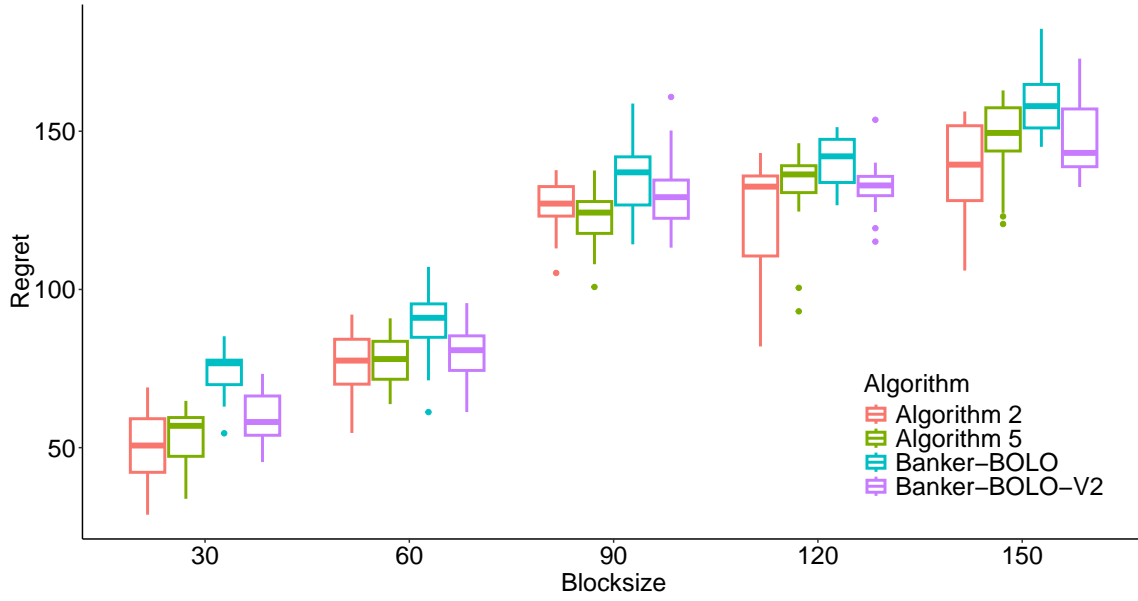


Figure 4: Boxplot of the regret over 20 repetitions over  $T = 10000$  rounds with  $K = 40$ .

Jacob D Abernethy, Elad Hazan, and Alexander Rakhlin. Interior-point methods for full-information and bandit online learning. *IEEE Transactions on Information Theory*, 58(7):4164–4175, 2012.

Alekh Agarwal and John C Duchi. Distributed delayed stochastic optimization. In *IEEE Conference on Decision and Control*, pages 5451–5452, 2012.

Jean-Yves Audibert, Sébastien Bubeck, and Gábor Lugosi. Regret in online combinatorial optimization. *Mathematics of Operations Research*, 39(1):31–45, 2014.

Baruch Awerbuch and Robert D Kleinberg. Adaptive routing with end-to-end feedback: Distributed learning and geometric approaches. In *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*, pages 45–53, 2004.

Itai Bistriz, Zhengyuan Zhou, Xi Chen, Nicholas Bambos, and Jose Blanchet. Online Exp3 learning in adversarial bandits with delayed feedback. In *Advances in Neural Information Processing Systems*, pages 11349–11358, 2019.

Stephen Boyd and Lieven Vandenbergh. *Convex optimization*. Cambridge university press, 2004.

Sébastien Bubeck and Ronen Eldan. The entropic barrier: a simple and optimal universal self-concordant barrier. In *Conference on Learning Theory*, pages 279–279, 2015.

Nicolo Cesa-Bianchi and Gábor Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78(5):1404–1422, 2012.

Nicolò Cesa-Bianchi, Claudio Gentile, Yishay Mansour, and Alberto Minora. Delay and cooperation in nonstochastic bandits. In *Conference on Learning Theory*, pages 605–622, 2016.

- Nicolo Cesa-Bianchi, Claudio Gentile, and Yishay Mansour. Delay and cooperation in nonstochastic bandits. *The Journal of Machine Learning Research*, 20(1):613–650, 2019.
- Wei Chen, Yajun Wang, and Yang Yuan. Combinatorial multi-armed bandit: General framework and applications. In Sanjoy Dasgupta and David McAllester, editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 151–159, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. URL <https://proceedings.mlr.press/v28/chen13a.html>.
- Alon Cohen, Amit Daniely, Yoel Drori, Tomer Koren, and Mariano Schain. Asynchronous stochastic optimization robust to arbitrary delays. *Advances in Neural Information Processing Systems*, 34:9024–9035, 2021.
- Richard Combes, Mohammad Sadegh Talebi Mazraeh Shahi, and Alexandre Proutiere. Combinatorial bandits revisited. In *Advances in Neural Information Processing Systems*, 2015.
- Thomas M Cover. Universal portfolios. *Mathematical finance*, 1(1):1–29, 1991.
- Yan Dai, Haipeng Luo, and Liyu Chen. Follow-the-perturbed-leader for adversarial Markov Decision Processes with bandit feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 11437–11449. Curran Associates, Inc., 2022.
- Varsha Dani and Thomas P Hayes. Robbing the bandit: Less regret in online geometric optimization against an adaptive adversary. In *SODA*, volume 6, pages 937–943, 2006.
- Varsha Dani, Sham M Kakade, and Thomas Hayes. The price of bandit information for online optimization. *Advances in Neural Information Processing Systems*, 20, 2007.
- Travis Dick, Andras Gyorgy, and Csaba Szepesvari. Online learning in Markov Decision Processes with changing cost sequences. In *International Conference on Machine Learning*, pages 512–520. PMLR, 2014.
- Miroslav Dudik, Daniel Hsu, Satyen Kale, Nikos Karampatziakis, John Langford, Lev Reyzin, and Tong Zhang. Efficient optimal learning for contextual bandits. In *Conference on Uncertainty in Artificial Intelligence*, pages 169–178, 2011.
- Stephen G. Eick. The two-armed bandit with delayed responses. *The Annals of Statistics*, 1988.
- Eyal Even-Dar, Sham M Kakade, and Yishay Mansour. Online Markov Decision Processes. *Mathematics of Operations Research*, 34(3):726–736, 2009.
- Genevieve E Flaspohler, Francesco Orabona, Judah Cohen, Soukayna Mouatadid, Miruna Oprescu, Paulo Orenstein, and Lester Mackey. Online learning with optimism and delay. In *International Conference on Machine Learning*, pages 3363–3373, 2021.
- Manegueu Anne Gael, Claire Vernade, Alexandra Carpentier, and Michal Valko. Stochastic bandits with arm-dependent delays. In *International Conference on Machine Learning*, pages 3348–3356, 2020.

- Yi Gai, Bhaskar Krishnamachari, and Rahul Jain. Combinatorial network optimization with unknown variables: Multi-armed bandits with linear rewards and individual observations. *IEEE/ACM Transactions on Networking*, 20(5):1466–1478, 2012. doi: 10.1109/TNET.2011.2181864.
- Pierre Gaillard, Gilles Stoltz, and Tim van Erven. A second-order bound with excess losses. In *Proceedings of The 27th Conference on Learning Theory*, volume 35 of *Proceedings of Machine Learning Research*, pages 176–196. PMLR, 13–15 Jun 2014.
- Andras Gyorgy and Pooria Joulani. Adapting to delays and data in adversarial multi-armed bandits. In *International Conference on Machine Learning*, pages 3988–3997, 2021.
- András György, Tamás Linder, Gábor Lugosi, and György Ottucsák. The on-line shortest path problem under partial monitoring. *Journal of Machine Learning Research*, 8(10), 2007.
- András György, Tamás Linder, Gábor Lugosi, and György Ottucsák. The on-line shortest path problem under partial monitoring. *Journal of Machine Learning Research*, 8(79):2369–2403, 2007. URL <http://jmlr.org/papers/v8/gyoergy07a.html>.
- Elad Hazan and Zohar Karnin. Volumetric spanners: an efficient exploration basis for learning. *Journal of Machine Learning Research*, 2016.
- Dirk van der Hoeven and Nicolo Cesa-Bianchi. Nonstochastic bandits and experts with arm-dependent delays. In *International Conference on Artificial Intelligence and Statistics*, 2022.
- Dirk van der Hoeven, Lukas Zierahn, Tal Lancewicki, Aviv Rosenberg, and Nicolò Cesa-Bianchi. A unified analysis of nonstochastic delayed feedback for combinatorial semi-bandits, linear bandits, and mdps. In *Conference on Learning Theory*, pages 1285–1321, 2023.
- Benjamin Howson, Ciara Pike-Burke, and Sarah Filippi. Optimism and delays in episodic reinforcement learning. In *International Conference on Artificial Intelligence and Statistics*, pages 6061–6094, 2023a.
- Benjamin Howson, Ciara Pike-Burke, and Sarah Filippi. Delayed feedback in generalised linear bandits revisited. In *International Conference on Artificial Intelligence and Statistics*, pages 6095–6119, 2023b.
- Jiatai Huang, Yan Dai, and Longbo Huang. Banker online mirror descent: A universal approach for delayed online bandit learning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 13814–13844. PMLR, 23–29 Jul 2023.
- Shinji Ito, Daisuke Hatano, Hanna Sumita, Kei Takemura, Takuro Fukunaga, Naonori Kakimura, and Ken-Ichi Kawarabayashi. Delay and cooperation in nonstochastic linear bandits. In *Advances in Neural Information Processing Systems*, pages 4872–4883, 2020a.
- Shinji Ito, Shuichi Hirahara, Tasuku Soma, and Yuichi Yoshida. Tight first-and second-order regret bounds for adversarial linear bandits. In *Advances in Neural Information Processing Systems*, pages 2028–2038, 2020b.



- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(4), 2010.
- Chi Jin, Tiancheng Jin, Haipeng Luo, Suvrit Sra, and Tiancheng Yu. Learning adversarial Markov decision processes with bandit feedback and unknown transition. In *International Conference on Machine Learning*, pages 4860–4869, 2020.
- Tiancheng Jin, Tal Lancewicki, Haipeng Luo, Yishay Mansour, and Aviv Rosenberg. Near-optimal regret for adversarial MDP with delayed bandit feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 33469–33481. Curran Associates, Inc., 2022.
- Pooria Joulani, Andras Gyorgy, and Csaba Szepesvári. Online learning under delayed feedback. In *International Conference on Machine Learning*, pages 1453–1461, 2013.
- Pooria Joulani, Andras Gyorgy, and Csaba Szepesvári. Delay-tolerant online convex optimization: Unified analysis and adaptive-gradient algorithms. In *AAAI Conference on Artificial Intelligence*, 2016.
- Pooria Joulani, András György, and Csaba Szepesvári. A modular analysis of adaptive (non-) convex optimization: Optimism, composite objectives, variance reduction, and variational bounds. *Theoretical Computer Science*, 808:108–138, 2020.
- Satyen Kale, Lev Reyzin, and Robert E Schapire. Non-stochastic bandit slate problems. In *Advances in Neural Information Processing Systems*, 2010.
- Konstantinos V Katsikopoulos and Sascha E Engelbrecht. Markov decision processes with delays and asynchronous cost collection. *IEEE transactions on automatic control*, 48(4):568–574, 2003.
- Wouter M Koolen, Manfred K Warmuth, and Jyrki Kivinen. Hedging structured concepts. In *Conference on Learning Theory*, pages 93–105, 2010.
- Branislav Kveton, Zheng Wen, Azin Ashkan, and Csaba Szepesvari. Tight Regret Bounds for Stochastic Combinatorial Semi-Bandits. In Guy Lebanon and S. V. N. Vishwanathan, editors, *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, volume 38 of *Proceedings of Machine Learning Research*, pages 535–543, San Diego, California, USA, 09–12 May 2015. PMLR. URL <https://proceedings.mlr.press/v38/kveton15.html>.
- Tal Lancewicki, Shahar Segal, Tomer Koren, and Yishay Mansour. Stochastic multi-armed bandits with unrestricted delay distributions. In *International Conference on Machine Learning*, pages 5969–5978, 2021.
- Tal Lancewicki, Aviv Rosenberg, and Yishay Mansour. Learning adversarial Markov decision processes with delayed feedback. In *AAAI Conference on Artificial Intelligence*, volume 36, pages 7281–7289, 2022a.
- Tal Lancewicki, Aviv Rosenberg, and Yishay Mansour. Cooperative online learning in stochastic and adversarial MDPs. In *International Conference on Machine Learning*, pages 11918–11968, 2022b.

- Tal Lancelwicki, Aviv Rosenberg, and Dmitry Sotnikov. Delay-adapted policy optimization and improved regret for adversarial MDP with delayed bandit feedback. In *International Conference on Machine Learning*, pages 18482–18534. PMLR, 2023.
- John Langford, Alex Smola, and Martin Zinkevich. Slow learners are fast. In *Advances in neural information processing systems*, 2009.
- Tor Lattimore, Branislav Kveton, Shuai Li, and Csaba Szepesvari. Toprank: A practical algorithm for online stochastic ranking. In *Advances in Neural Information Processing Systems*, 2018.
- Nick Littlestone and Manfred K Warmuth. The weighted majority algorithm. *Information and computation*, 108(2):212–261, 1994.
- Jonathan Lou  dec, Max Chevalier, Josiane Mothe, Aur  lien Garivier, and S  bastien Gerchinovitz. A multiple-play bandit algorithm applied to recommender systems. In *28th International Florida Artificial Intelligence Research Society (FLAIRS 2015)*, pages 67–72, Hollywood, United States, May 2015. URL <https://hal.science/hal-04077707>.
- Haipeng Luo, Chen-Yu Wei, and Chung-Wei Lee. Policy optimization in adversarial MDPs: Improved exploration via dilated bonuses. In *Advances in Neural Information Processing Systems*, 2021.
- Saeed Masoudian, Julian Zimmert, and Yevgeny Seldin. A best-of-both-worlds algorithm for bandits with delayed feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 11752–11762. Curran Associates, Inc., 2022.
- H Brendan McMahan and Avrim Blum. Online geometric optimization in the bandit setting against an adaptive adversary. In *Conference on Learning Theory*, pages 109–123, 2004.
- Arkadi Nemirovski. Interior point polynomial time methods in convex programming. *Lecture notes*, 42(16):3215–3224, 2004.
- Arkadi S Nemirovski and Michael J Todd. Interior-point methods for optimization. *Acta Numerica*, 17:191–234, 2008.
- Yurii Nesterov and Arkadii Nemirovskii. *Interior-point polynomial algorithms in convex programming*. SIAM, 1994.
- Gergely Neu, Andr  s Gy  rgy, Csaba Szepesv  ri, and Andr  s Antos. Online Markov Decision Processes under bandit feedback. In *Advances in Neural Information Processing Systems*, 2010.
- Gergely Neu, Andr  s Gy  rgy, Csaba Szepesv  ri, and Andr  s Antos. Online Markov Decision Processes under bandit feedback. *IEEE Trans. Automat. Contr.*, 59(3):676–691, 2014.
- He Ni, Hao Xu, Dan Ma, and Jun Fan. Contextual combinatorial bandit on portfolio management. *Expert Systems with Applications*, 221:119677, 2023.
- Francesco Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*, 2019.

- Ciara Pike-Burke, Shipra Agrawal, Csaba Szepesvari, and Steffen Grunewalder. Bandits with delayed, aggregated anonymous feedback. In *International Conference on Machine Learning*, pages 4105–4113, 2018.
- Kent Quanrud and Daniel Khashabi. Online learning with adversarial delays. In *Advances in Neural Information Processing Systems*, pages 1270–1278, 2015.
- Alexander Rakhlin and Karthik Sridharan. Online learning with predictable sequences. In *Conference on Learning Theory*, pages 993–1019, 2013.
- Aviv Rosenberg and Yishay Mansour. Online stochastic shortest path with bandit feedback and unknown transition function. In *Advances in Neural Information Processing Systems*, pages 2209–2218, 2019a.
- Aviv Rosenberg and Yishay Mansour. Online convex optimization in adversarial Markov decision processes. In *International Conference on Machine Learning*, pages 5478–5486, 2019b.
- Aviv Rosenberg and Yishay Mansour. Stochastic shortest path with adversarially changing costs. In Zhi-Hua Zhou, editor, *International Joint Conference on Artificial Intelligence*, pages 2936–2942, 2021.
- Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.
- Lior Shani, Yonathan Efroni, Aviv Rosenberg, and Shie Mannor. Optimistic policy optimization with bandit feedback. In *International Conference on Machine Learning*, pages 8604–8613, 2020.
- Tobias Sommer Thune, Nicolò Cesa-Bianchi, and Yevgeny Seldin. Nonstochastic multiarmed bandits with unrestricted delays. In *Advances in Neural Information Processing Systems*, pages 6541–6550, 2019.
- Taishi Uchiya, Atsuyoshi Nakamura, and Mineichi Kudo. Algorithms for adversarial bandit problems with multiple plays. In *International Conference on Algorithmic Learning Theory*, pages 375–389, 2010.
- Claire Vernade, Alexandra Carpentier, Tor Lattimore, Giovanni Zappella, Beyza Ermis, and Michael Brueckner. Linear bandits with stochastic delayed feedback. In *International Conference on Machine Learning*, pages 9712–9721, 2020.
- Volodimir G Vovk. Aggregating strategies. In *Proceedings of the 3rd Annual Workshop on Computational Learning Theory*, 1990.
- Thomas J Walsh, Ali Nouri, Lihong Li, and Michael L Littman. Learning and planning in environments with delayed feedback. *Autonomous Agents and Multi-Agent Systems*, 18(1):83, 2009.
- Marcelo J Weinberger and Erik Ordentlich. On delayed prediction of individual sequences. *IEEE Transactions on Information Theory*, 48(7):1959–1976, 2002.

- Zheng Wen, Branislav Kveton, and Azin Ashkan. Efficient learning in large-scale combinatorial semi-bandits. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1113–1122, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/wen15.html>.
- Zhengyuan Zhou, Renyuan Xu, and Jose Blanchet. Learning in generalized linear contextual bandits with stochastic delays. In *Advances in Neural Information Processing Systems*, pages 5197–5208, 2019.
- Alexander Zimin and Gergely Neu. Online learning in episodic Markovian decision processes by relative entropy policy search. In *Advances in Neural Information Processing Systems*, pages 1583–1591, 2013.
- Julian Zimmert and Tor Lattimore. Return of the bias: Almost minimax optimal high probability bounds for adversarial linear bandits. In *Conference on Learning Theory*, pages 3285–3312, 2022.
- Julian Zimmert and Yevgeny Seldin. An optimal algorithm for adversarial bandits with arbitrary delays. In *International Conference on Artificial Intelligence and Statistics*, pages 3285–3294, 2020.
- Julian Zimmert, Haipeng Luo, and Chen-Yu Wei. Beating stochastic and adversarial semi-bandits optimally and simultaneously. In *International Conference on Machine Learning*, pages 7683–7692. PMLR, 2019.