

HIGH PROBABILITY BOUND FOR CROSS-LEARNING CONTEXTUAL BANDITS WITH UNKNOWN CONTEXT DISTRIBUTIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Motivated by applications in online bidding and sleeping bandits, we examine the problem of contextual bandits with cross learning, where the learner observes the loss associated with the action across all possible contexts, not just the current round’s context. Our focus is on a setting where losses are chosen adversarially, and contexts are sampled i.i.d. from a specific distribution. This problem was first studied by Balseiro et al. (2019), who proposed an algorithm that achieves near-optimal regret under the assumption that the context distribution is known in advance. However, this assumption is often unrealistic. To address this issue, Schneider & Zimmert (2023) recently proposed a new algorithm that achieves nearly optimal expected regret. It is well-known that expected regret can be significantly weaker than high-probability bounds. In this paper, we present a novel, in-depth analysis of their algorithm and demonstrate that it actually achieves near-optimal regret with *high probability*. There are steps in the original analysis by Schneider & Zimmert (2023) that lead only to an expected bound by nature. In our analysis, we introduce several new insights. Specifically, we make extensive use of the weak dependency structure between different epochs, which was overlooked in previous analyses. Additionally, standard martingale inequalities are not directly applicable, so we refine martingale inequalities to complete our analysis.

1 INTRODUCTION

In the contextual bandits problem, a learner repeatedly observes a context, chooses an action, and incurs a loss specific to that action. The goal of the learner is to minimize the cumulative loss over the time horizon. The contextual bandits problem is a fundamental problem in online learning having broad applications in fields like online advertising, personalized recommendations, and clinical trials (Li et al., 2010; Kale et al., 2010; Villar et al., 2015).

We consider the cross-learning contextual bandits problem. In this setting, the learner not only observes the loss for the current action under the current context, but also observes the loss for the current action under all other contexts. This problem models many interesting scenarios. One such example is the problem of learning to bid in first-price auctions. In this problem the context is the bidder’s private value for the item, while the action is the bid. The cross-learning structure comes from the fact that the bidder can deduce the utility of the bid under all contexts (i.e., the utility of the bid under different private valuations for the item). Other examples include multi-armed bandits with exogenous costs, dynamic pricing with variable costs, and learning to play in Bayesian games (Balseiro et al., 2019).

Technically, the most interesting setting for the cross-learning contextual bandits problem is when the losses are chosen adversarially but the contexts are i.i.d. samples from an *unknown* distribution ν . Recently, Schneider & Zimmert (2023) gave an algorithm achieving nearly optimal $\tilde{O}(\sqrt{KT})$ expected regret in this scenario.

Schneider & Zimmert (2023) designed a sophisticated algorithm that operates over multiple epochs to achieve near-optimal regret. A key technique in their analysis is to sidestep high-probability bounds and instead focus on bounding the expected summation to improve their results. As a consequence, their analysis only provides a bound that holds in expectation. It is not immediately clear

whether this is due to limitations in the analysis or if the algorithm is inherently suboptimal. In any case, if we aim for a high-probability bound, fundamentally new insights are required.

In this paper, we show that the algorithm indeed achieves nearly optimal $\tilde{O}(\sqrt{KT})$ regret with high-probability. The key contribution of our paper is the following theorem.

Theorem 1 (Informal). *The algorithm in Schneider & Zimmert (2023) yields a regret bound of order $\tilde{O}(\sqrt{KT})$ with high probability for any policy π .*

In this section we only give the informal version of Theorem 1. The formal version can be found in Section 4.

1.1 TECHNICAL OVERVIEW

Our theorem is built on a new and more in-depth analysis of the algorithm in Schneider & Zimmert (2023). This new analysis introduces several new insights. In particular, we exploit the weak dependency structure between different epochs, which was overlooked in previous work. One difficulty of doing so is that standard martingale inequalities are not directly applicable, so we refine martingale inequalities to complete our analysis.

To prepare the readers for our new analysis, we first briefly introduce the algorithm in Schneider & Zimmert (2023). The algorithm in Schneider & Zimmert (2023) is an EXP3-type algorithm. The key novelty in their algorithm is the construction of the loss estimates $\hat{\ell}$ used in the FTRL subroutine. Due to some technical problems we detail later, the algorithm decomposes the time horizon into epochs of equal length. In each epoch e , the algorithm first estimates the probability¹ $f_e(a)$ of observing the reward of each arm a in epoch e by an estimator $\hat{f}_e(a)$, which is constructed exclusively from samples in epoch $e - 1$. Note that thanks to the cross-learning structure, the probability of observing the reward of each arm a is independent of the contexts. The algorithm then constructs the loss estimates as an importance-weighted estimator with $\frac{1}{\hat{f}_e(a)}$ as the importance weight.

Schneider & Zimmert (2023) showed that the performance of the algorithm depends on how well the empirical importance weight $\frac{1}{\hat{f}_e(a)}$ concentrates around the expected importance weight $\frac{1}{f_e(a)}$.

Since the estimator $\hat{f}_e(a)$ is constructed exclusively from samples in a single epoch rather than the entire time horizon, the concentration $|\frac{1}{\hat{f}_e(a)} - \frac{1}{f_e(a)}|$ is not tight enough. To achieve the desired $\tilde{O}(\sqrt{KT})$ regret under a not tight enough concentration, Schneider & Zimmert (2023) bounds only the expected bias of importance estimator $\mathbb{E}[\frac{1}{\hat{f}_e(a)} - \frac{1}{f_e(a)}]$ rather than providing a high-probability bias bound. Bounding only the expected bias gives a small enough bound, however, they can achieve a bound only on the expected regret from a bound on the expected bias.

We overcome this difficulty and show that their algorithm actually achieves a high-probability bound. Our key observation is that different epochs in their algorithm are only weakly dependent on each other. Thus, the bias $\frac{1}{\hat{f}_e(a)} - \frac{1}{f_e(a)}$ for each epoch e is also only weakly dependent on each other. Therefore, although we cannot establish a small enough bound for the bias of a single epoch $\frac{1}{\hat{f}_e(a)} - \frac{1}{f_e(a)}$, we can give a small enough bound for the cumulative bias across all epochs $\sum_e \frac{1}{\hat{f}_e(a)} - \frac{1}{f_e(a)}$. We then use the bound on the cumulative bias to bound the cumulative regret.

In addition to utilizing the weak dependency structure between different epochs, we also address two further technical difficulties to establish our result. The first difficulty is that the existing regret decomposition is too crude to yield a high-probability bound. Schneider & Zimmert (2023) establish an $\tilde{O}(\sqrt{KT})$ expected regret by decomposing the regret into different parts and bounding each part separately. Although their decomposition gives an $\tilde{O}(\sqrt{KT})$ expected regret bound, it is too crude to derive a tight high-probability regret bound, even after utilizing the weak dependency structure. We carefully rearrange the regret decomposition to address this difficulty.

¹For technical reasons, in the actual algorithm, the value $f_e(a)$ actually represents the probability of observing the reward of each arm a in epoch $e + 2$. For ease of understanding, here we instead let it represent the probability of observing the reward of each arm a in each epoch e .

Secondly, we cannot simply apply standard martingale concentration inequalities to $\sum_e \frac{1}{f_e(a)} - \frac{1}{\bar{f}_e(a)}$ to bound its deviation. The main problem is that the random variable $\frac{1}{f_e(a)} - \frac{1}{\bar{f}_e(a)}$ is not almost surely bounded by a constant, which makes standard martingale concentration inequalities inapplicable. We introduce a surrogate sequence of random variables as a bridge to address this problem. We bound the sum over the surrogate sequence, and show that the sum over the real sequence is equal to the surrogate sequence with high probability.

1.2 RELATED WORKS

The cross-learning contextual bandits problem was first proposed in Balseiro et al. (2019). They achieve the nearly optimal $\tilde{O}(\sqrt{KT})$ regret under two scenarios: (1) when both losses and contexts are stochastic, and (2) when losses are adversarial and contexts are stochastic with a known distribution. When losses are adversarial and contexts are stochastic with an unknown distribution, they only achieve the suboptimal $\tilde{O}(K^{1/3}T^{2/3})$ regret. More recently, Schneider & Zimmert (2023) gave a new algorithm that achieves the nearly optimal $\tilde{O}(\sqrt{KT})$ regret in expectation under adversarial losses and stochastic contexts with an unknown distribution.

An important application of the cross-learning contextual bandits problem, which is also the primary motivation for proposing this problem in Balseiro et al. (2019), is to solve the problem of learning to bid in first-price auctions. In this problem the context is the bidder’s private value for the item, while the action is the bid. The cross learning structure comes from the fact that the bidder can deduce the utility of the bid under all contexts (i.e., the utility of the bid under different private valuations for the item).

Balseiro et al. (2019) used the cross-learning contextual bandits problem to model the bidding problem and obtained an $O(T^{3/4})$ regret bound for bidders with an unknown value distribution participating in adversarial first-price auctions, where the only feedback is whether the bidder wins the auction. Later, many works studied different settings of the bidding in first price auctions problem. For example, Han et al. (2020b) considered the problem with censored feedback, where each bidder observes the winning bid. Han et al. (2020a) considered the scenario when the value is also adversarial. Ai et al. (2022); Wang et al. (2023) considered the problem under budget constraints. In all these scenarios, the cross learning structure between different values is an essential component of the analysis.

Another interesting application of the cross-learning contextual bandits problem is the sleeping bandits problem (Kleinberg et al., 2010; Neu & Valko, 2014; Kale et al., 2016; Saha et al., 2020). In this problem, a certain set of arms is unavailable in each round. The sleeping bandits problem is motivated by instances like some items might go out of stock in retail stores or on a certain day some websites could be down. When losses are adversarial and availabilities are stochastic, previous work either requires exponential computing time (Kleinberg et al., 2010; Neu & Valko, 2014) or results in suboptimal regret (Kale et al., 2016; Saha et al., 2020). The first computationally efficient algorithm with optimal regret $\tilde{O}(\sqrt{KT})$ is proposed in Schneider & Zimmert (2023) by modeling the problem as a cross-learning contextual bandit.

We also note that handling unknown context distributions is a common and challenging problem across various contextual bandit problems. For example, in the adversarial linear contextual bandits problem (Neu & Olkhovskaya, 2020), the linear MDP problem (Dai et al., 2023), and the oracle-based adversarial contextual bandits problem (Syrkanis et al., 2016), existing algorithms often rely on knowledge of the context distribution. Removing the reliance on knowledge of the context distribution is typically non-trivial (Liu et al., 2023; Dai et al., 2023).

Recently, Hanna et al. (2023) proposed a method for stochastic linear contextual bandits that maps a multi-context problem to a single-context problem. Unfortunately, their approach cannot be directly applied to our problem for two reasons. First, their method is designed for stochastic bandits, whereas we deal with adversarial bandits. Second, their approach is limited to linear contextual bandits. Whether it can be adapted, with certain modifications, to address our problem remains an intriguing question.

2 PROBLEM STATEMENT

We study a contextual K -armed bandit problem over T rounds, with contexts belonging to the set $[C]$. At the beginning of the problem, an oblivious adversary selects a sequence of losses $\ell_{t,c}(k) \in [0, 1]$ for every round $t \in [T]$, every context $c \in [C]$, and every arm $k \in [K]$. In each round t , we begin by sampling a context $c_t \sim \nu$ i.i.d. from an unknown distribution ν over $[C]$, and we reveal this context to the learner. Based on this context, the learner selects an arm $a_t \in [K]$ to play. The adversary then reveals the function $\ell_{t,c}(a_t)$, and the learner suffers loss $\ell_{t,c_t}(a_t)$. Notably, the learner observes the loss for every context $c \in [C]$, but only for the arm a_t they actually played.

We aim to design learning algorithms that minimize regret. Fix a policy $\pi : [C] \rightarrow [K]$. With a slight abuse of notation, we also denote $\pi_c = e_k \in \Delta([K])$ for each $c \in [C]$. The (unexpected) regret with respect to policy π is

$$\text{Reg}(\pi) = \sum_{t=1}^T \ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}).$$

We aim to upper bound this quantity (for an arbitrary policy π).

Schneider & Zimmert (2023) designed an algorithm that achieves an expected regret bound of $\mathbb{E}[\text{Reg}(\pi)] \leq \tilde{O}(\sqrt{KT})$ for any policy π . We will show that the algorithm in Schneider & Zimmert (2023) actually provides a high-probability regret bound.

3 THE ALGORITHM IN SCHNEIDER AND ZIMMERT (2023)

In this section, we briefly recap the intuition behind the algorithm proposed in Schneider & Zimmert (2023) and redescribe the algorithm formally to prepare the readers for our new analysis.

3.1 INTUITION BEHIND SCHNEIDER AND ZIMMERT (2023)

The algorithm proposed in Schneider & Zimmert (2023) is an EXP3-type algorithm. Similar to the well-known EXP3 algorithm, at each round t , the algorithm generates a distribution using an FTRL subroutine

$$p_{t,c} = \arg \min_{p \in \Delta([K])} \left\langle p, \sum_{s=1}^{t-1} \hat{\ell}_{s,c} \right\rangle - \frac{1}{\eta} F(p)$$

for each context c , where $F(p) = \sum_{i=1}^K p_i \log(p_i)$ is the unnormalized negative entropy, η is a learning rate, and $\hat{\ell}$ are loss estimates to be defined later. The algorithm then essentially samples the action a_t to be played in round t from distribution p_{t,c_t} .

The key novelty in Schneider & Zimmert (2023) lies in the construction of the loss estimates $\hat{\ell}$. An intuitive construction is defined as follows:

$$\tilde{\ell}_{t,c}(a) = \frac{\ell_{t,c}(a)}{\mathbb{E}_{c \sim \nu}[p_{t,c}(a)]} \mathbb{1}(a_t = a).$$

That is, it uses the classic importance-weighted estimator with $\mathbb{E}_{c \sim \nu}[p_{t,c}(a)]$ as the importance². A straightforward analysis shows that this estimator yields a regret bound of $\tilde{O}(\sqrt{KT})$. However, the denominator term $\mathbb{E}_{c \sim \nu}[p_{t,c}(a)]$ is uncomputable because we do not know the distribution of contexts ν . One may attempt to circumvent this issue by replacing the expected importance $\mathbb{E}_{c \sim \nu}[p_{t,c}(a)]$ with the empirical importance $\frac{1}{t} \sum_{s=1}^t p_{t,c_s}(a)$. It is not hard to see that whether we achieve the desired $\tilde{O}(\sqrt{KT})$ regret depends on how well the empirical importance weight $\frac{1}{t} \sum_{s=1}^t p_{t,c_s}(a)$ concentrates around the expected importance weight $\mathbb{E}_{c \sim \nu}[p_{t,c}(a)]$. However, the empirical importance weight $\frac{1}{t} \sum_{s=1}^t p_{t,c_s}(a)$ may not concentrate well around the expected importance

²In this paper we call terms like $\frac{1}{\mathbb{E}_{c \sim \nu}[p_{t,c}(a)]}$ as the *importance weight* and call terms like $\mathbb{E}_{c \sim \nu}[p_{t,c}(a)]$ as the *importance*.

weight $\frac{1}{\mathbb{E}_{c \sim \nu}[p_{t,c}(a)]}$. This is because the probability vector $p_{t,c}$ is not independent of the previous contexts c_s , which makes standard concentration inequalities inapplicable.

To address this difficulty, Schneider & Zimmert (2023) divides the time horizon into epochs of equal length L . At the end of each epoch e , the algorithm stores the FTRL distribution at the current time $t = eL$ in a new distribution s_e ; that is, it takes $s_{e,c}(a) = p_{t,c}(a)$ for each context c and each arm a . The algorithm further decouples the distribution played by the algorithm and the distribution used to estimate the loss vector. For each time t in epoch $e + 1$, the algorithm observes the loss $\ell_{t,c}(a)$ for each arm a and context c with probability $f_e(a) \triangleq \mathbb{E}_{c \sim \nu}[s_{e,c}(a)/2]$. The algorithm then estimates the expected importance $f_e(a)$ using an empirical importance $\hat{f}_e(a)$ constructing solely from contexts in epoch $e + 1$. Finally, the algorithm constructs $\hat{\ell}_{t,c}(a)$ as an importance-weighted estimator with $\hat{f}_e(a)$ serving as the importance.

The advantage of their construction is that the empirical importance weight $\frac{1}{\hat{f}_e(a)}$ concentrates around the expected importance weight $\frac{1}{f_e(a)}$ now. This concentration ensures that the loss estimates $\hat{\ell}_{t,c}(a)$ are good estimates of the true losses $\ell_{t,c}(a)$. And this concentration is achieved because the algorithm constructs the estimator using only samples from epoch $e + 1$, which are independent of the estimand.

3.2 A FORMAL DESCRIPTION OF THE ALGORITHM IN SCHNEIDER AND ZIMMERT (2023)

In this subsection we describe the algorithm in Schneider & Zimmert (2023) formally for the sake of completeness. Readers familiar with Schneider & Zimmert (2023) can skip this subsection safely.

In each round t , the algorithm generates a distribution from an FTRL subroutine:

$$p_{t,c} = \arg \min_{p \in \Delta([K])} \left\langle p, \sum_{s=1}^{t-1} \hat{\ell}_{s,c} \right\rangle - \frac{1}{\eta} F(p)$$

for each context c , where $F(p) = \sum_{i=1}^K p_i \log(p_i)$ is the unnormalized negative entropy, η is the learning rate, and $\hat{\ell}$ are loss estimates to be defined later. The algorithm will not sample the action a_t played in round t directly from p_t but from a distribution q_t to be defined later.

To construct loss estimates $\hat{\ell}$, the algorithm divides the time horizon into epochs of equal length L . We let \mathcal{T}_e to denote the set of rounds in the e -th epoch. At the end of each epoch, the algorithm takes a single snapshot of the underlying FTRL distribution p_t for each context and arm. That is, the algorithm takes

$$s_{e+2,c}(a) = p_{eL,c}(a), \text{ where } s_{1,c}(a) = s_{2,c}(a) = \begin{cases} \frac{1}{|\mathcal{A}_c|} & \text{if } a \in \mathcal{A}_c \\ 0 & \text{otherwise.} \end{cases}$$

For each round $t \in \mathcal{T}_e$, the algorithm observes the loss function of arm a with probability $f_e(a) = \mathbb{E}_{c \sim \nu}[s_{e,c}(a)/2]$. This is guaranteed by the following rejection sampling procedure: we first play an arm according to the distribution

$$q_{t,c_t} = \begin{cases} p_{t,c_t} & \text{if } \forall a \in [K] : p_{t,c_t}(a) \geq s_{e,c_t}(a)/2 \\ s_{e,c_t} & \text{otherwise.} \end{cases}$$

After playing arm a according to q_{t,c_t} , the learner samples a Bernoulli random variable S_t with probability $\frac{s_{e,c_t}(a)}{2q_{t,c_t}(a)}$. If $S_t = 0$, the learner ignores the feedback from this round; otherwise, they use this loss.

The only remaining unspecified part is how to construct the loss estimates. We group all timesteps into consecutive pairs of two. In each pair of consecutive timesteps, we sample from the same distribution and randomly use one to calculate a loss estimate and the other to estimate the sampling frequency. To be precise, let \mathcal{T}_e^f denote the timesteps selected for estimating the sampling frequency and \mathcal{T}_e^ℓ denote the timesteps used to estimate the losses. Then we define

$$\hat{f}_e(a) = \frac{1}{|\mathcal{T}_{e-1}^f|} \sum_{t \in \mathcal{T}_{e-1}^f} \frac{s_{e,c_t}(a)}{2}$$

which is an unbiased estimator of $f_e(a)$. The loss estimators are defined as follows:

$$\widehat{\ell}_{t,c}(a) = \frac{2\ell_{t,c}(a)}{\widehat{f}_e(a) + \frac{3}{2}\gamma} \mathbb{1}(A_t = a \wedge S_t \wedge t \in \mathcal{T}_e^\ell)$$

where γ is a confidence parameter to be specified later.

The algorithm is summarized in Algorithm 1. Furthermore, Schneider & Zimmert (2023) showed that the algorithm achieves an expected regret bound of $\widetilde{O}(\sqrt{KT})$.

Algorithm 1 The algorithm for the cross-learning problem in Schneider & Zimmert (2023)

Input: Parameters $\eta, \gamma > 0$ and $L < T$.

```

282  $\widehat{f}_2 \leftarrow 0$ 
283 for  $t = 1, \dots, L$  do
284   Observe  $c_t$ 
285   Play  $A_t \sim s_{1,c_t}$ 
286    $\widehat{f}_2 \leftarrow \widehat{f}_2 + \frac{s_{2,c_t}}{2L}$ 
287 for  $e = 2, \dots, T/L$  do
288    $\widehat{f}_{e+1} \leftarrow 0$ 
289   for  $t = (e-1)L + 1, t = (e-1)L + 3, \dots, eL - 1$  do
290     Set  $p_{t,c} = \arg \min_{x \in \Delta([K])} \left( \left\langle x, \sum_{s=1}^{t-1} \widehat{\ell}_s(c) \right\rangle - \eta^{-1} F(x) \right)$ 
291     for  $t' = t, t+1$  do
292       Observe  $c_{t'}$ 
293       if  $p_{t,c_{t'}}(a) \geq s_{e,c_{t'}}(a)/2$  for all  $a \in [K]$  then
294         Set  $q_{t',c_{t'}} = p_{t,c_{t'}}$ 
295       else
296         Set  $q_{t',c_{t'}} = s_{e,c_{t'}}$ 
297       Play  $A_{t'} \sim q_{t',c_{t'}}$ 
298       Observe  $\ell_{t',A_{t'}}$ 
299      $t_f, t_\ell \leftarrow \text{RandPerm}(t, t+1)$ 
300      $\widehat{f}_{e+1} \leftarrow \widehat{f}_{e+1} + \frac{s_{e+1,c_{t_f}}}{2(L/2)}$ 
301     Sample  $S_t \sim \mathcal{B} \left( \frac{s_{e,c_{t_\ell}}(A_{t_\ell})}{2q_{t,c_{t_\ell}}(A_{t_\ell})} \right)$ 
302     Set  $\widehat{\ell}_{t_\ell,c}(a) = \frac{2\ell_{t_\ell,c}(a)}{\widehat{f}_e(a) + \frac{3}{2}\gamma} \mathbb{I}(A_t = a, S_t = 1)$ 
303    $s_{e+2} \leftarrow p_t$ 

```

4 MAIN RESULT AND ANALYSIS

The main result of our paper is the following theorem.

Theorem 1 (Formal). *For any $\delta \in (0, 1)$, Algorithm 1 with parameters choice $\iota = 2 \log(8KT \frac{1}{\delta})$, $L = \sqrt{\frac{\iota KT}{\log(K)}} = \widetilde{\Theta}(\sqrt{KT \log \frac{1}{\delta}})$, $\gamma = \frac{16\iota}{L} = \widetilde{\Theta}(\sqrt{\frac{\log(1/\delta)}{KT}})$, and $\eta = \frac{\gamma}{2(2L\gamma + \iota)} = \widetilde{\Theta}(1/\sqrt{KT \log(1/\delta)})$ yields a regret bound of*

$$\text{Reg}(\pi) = \widetilde{O} \left(\sqrt{KT \log \frac{1}{\delta}} \right)$$

with probability at least $1 - \delta$ for any policy π .

In what follows, we briefly overview our proof of Theorem 1. The full proof can be found in the appendix.

4.1 REGRET DECOMPOSITION

Denote the set of all timesteps used to estimate the frequency as \mathcal{T}^f and denote the set of all timesteps used to estimate the losses as \mathcal{T}^ℓ . For each $t \in \mathcal{T}_e$, we define $\tilde{\ell}_{t,c}(a) = \frac{2\ell_{t,c}(a)}{\widehat{f}_e(a)+\gamma} \mathbb{1}(A_t = a \wedge S_t \wedge t \in \mathcal{T}_e^\ell)$. We decompose regret $\text{Reg}(\pi) = \sum_{t=1}^T \ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t})$ as follows:

$$\begin{aligned}
\text{Reg}(\pi) &= \underbrace{\sum_{t=1}^T (\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t})) - 2 \sum_{t \in \mathcal{T}^\ell} (\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}))}_{\text{bias}_1} \\
&+ 2 \underbrace{\sum_{t \in \mathcal{T}^\ell} \left(\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - \sum_c \Pr(c) \langle p_{t,c} - \pi_c, \ell_{t,c} \rangle \right)}_{\text{bias}_2} \\
&+ 2 \underbrace{\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c} - \pi_c, \widehat{\ell}_{t,c} \rangle}_{\text{ftrl}} + 2 \underbrace{\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \ell_{t,c} - \tilde{\ell}_{t,c} \rangle}_{\text{bias}_3} \\
&+ 2 \underbrace{\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \tilde{\ell}_{t,c} - \widehat{\ell}_{t,c} \rangle}_{\text{bias}_4} + 2 \underbrace{\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle \pi_c, \widehat{\ell}_{t,c} - \ell_{t,c} \rangle}_{\text{bias}_5}.
\end{aligned}$$

In our decomposition, the bias_1 term refers to the bias introduced by replacing the regret over the entire time horizon with that over \mathcal{T}^ℓ , and the bias_2 term refers to the bias introduced by replacing regret with its linearization. Both of these terms are not hard to bound using standard concentration inequalities. Furthermore, the **ftrl** and bias_3 terms are standard in the analysis of high-probability bounds for bandit algorithms. These two terms are not hard to bound using techniques from EXP3-IX (Neu, 2015; Schneider & Zimmert, 2023). The bias_4 and bias_5 terms correspond to the bias introduced by constructing the importance estimator $\widehat{f}_e(a)$. These two terms are the terms of interest to bound.

Our decomposition is different from the decomposition in Schneider & Zimmert (2023). This difference is essential for deriving a high-probability bound. The key difference lies in the bias_5 term here. This term saves a $\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle \pi_c, \tilde{\ell}_{t,c} - \ell_{t,c} \rangle$ term from the bias_2 term in the decomposition given by Schneider & Zimmert (2023), which is crucial for deriving a high-probability bound.

4.2 IDENTIFYING A PROTOTYPICAL TERM

The terms of interest to bound are bias_4 and bias_5 . These two terms can be bounded using similar methods. Here we take the bias_5 term as a prototypical term and give a sketch of its analysis. Details can be found in the appendix.

To bound bias_5 , we define a filtration $\{\mathcal{H}_t\}_t$ such that the σ -algebra \mathcal{H}_t for each time step t is generated by all randomness before time t . Next, we decompose bias_5 as

$$\begin{aligned}
&\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle \pi_c, \widehat{\ell}_{t,c} - \ell_{t,c} \rangle \\
&= \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \ell_{t,c}(\pi_c) \right) \\
&= \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right) \\
&+ \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] \right).
\end{aligned}$$

In this decomposition, the two terms correspond to different components of the bias of the estimator $\widehat{\ell}_{t,c}$. The first term $\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right)$ corresponds to the bias introduced by constructing the importance estimator $\widehat{f}_e(a)$. The second term $\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] \right)$ corresponds to the bias introduced from the randomness in sampling a_t from q_t . Once again, the analyses of these two terms follow the same principle. We take the first term

$$\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right)$$

as a prototypical term and give a sketch of its analysis for the sake of simplicity. Details can be found in the appendix.

4.3 ANALYSIS OF THE PROTOTYPICAL TERM

To bound the term $\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right)$, we will use the key observation mentioned at Section 1.1: different epochs in Algorithm 1 are only weakly dependent on each other. To use this observation rigorously, we introduce an important technical tool. With a slight abuse of notation, we define a filtration $\{\mathcal{H}_e\}_e$, in which for each epoch e , the σ -algebra \mathcal{H}_e is generated by all randomness in epochs $1, \dots, e-1$ and the randomness in \mathcal{T}_e^ℓ . That is, the σ -algebra \mathcal{H}_e is generated precisely by the context c_t , the random seed used in sampling $a_t \sim q_{t,c_t}$, and the random seed used in sampling $S_t \sim \mathcal{B} \left(\frac{s_{e,c_t}(a_t)}{2q_{t,c_t}(a_t)} \right)$ for $t \leq (e-1)L$ and $t \in \mathcal{T}_e^\ell$. Note that for each epoch e , the σ -algebra \mathcal{H}_e excludes the randomness in \mathcal{T}_e^f . This exclusion is crucial for characterizing the weak dependence structure between epochs.

Given this filtration, we consider the cumulative bias in each epoch. For each epoch e , we define a random variable

$$\text{Bias5}_e \triangleq \sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right).$$

Then the prototypical term $\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right)$ is exactly $\sum_e \text{Bias5}_e$. Our key observation is that, not only

$$\mathbb{E} \left[\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right) \right] \leq 0$$

as shown in Schneider & Zimmert (2023), but also

$$\sum_e \mathbb{E} [\text{Bias5}_e \mid \mathcal{H}_e] \sim - \sum_e \frac{\gamma}{f_e(\pi_c) + \gamma}.$$

This key observation improves the inequality in Schneider & Zimmert (2023) in two ways. Firstly, our bound holds for conditional expectations across epochs, which opens the door to applying martingale concentration inequalities across epochs. Secondly, our new decomposition improves the upper bound from 0 to $-\sum_e \frac{\gamma}{f_e(\pi_c) + \gamma}$. This improvement is essential for deriving a high probabilistic bound.

Given the new bound $\sum_e \mathbb{E} [\text{Bias5}_e \mid \mathcal{H}_e] \sim - \sum_e \frac{\gamma}{f_e(\pi_c) + \gamma}$, we only need to bound the deviation $\sum_e \text{Bias5}_e - \mathbb{E} [\text{Bias5}_e \mid \mathcal{H}_e]$ to get an upper bound on $\sum_e \text{Bias5}_e$. However, we cannot directly apply standard martingale concentration inequalities to $\sum_e \text{Bias5}_e - \mathbb{E} [\text{Bias5}_e \mid \mathcal{H}_e]$. This is because we need to assume that the random variable $|\text{Bias5}_e| \leq 2L$ almost surely to get a tight enough concentration bound when applying standard martingale concentration inequalities. However, this is not the case. The random variable Bias5_e exceeds the constant $2L$ with a small but positive probability. This unboundness prevents us from getting a tight enough concentration bound when applying standard martingale concentration inequalities.

To overcome this problem, we consider the indicator function

$$F_e \triangleq \mathbb{1} \left(\forall a, \left| \widehat{f}_e(a) - f_e(a) \right| \leq 2 \max \left\{ \sqrt{\frac{f_e(a)\iota}{L}}, \frac{\iota}{L} \right\} \right)$$

defined in Schneider & Zimmert (2023). We show that we also have $\sum_e \mathbb{E}[\text{Bias}_{5_e} F_e | \mathcal{H}_e] \sim -\sum_e \frac{\gamma}{f_e(\pi_e) + \gamma}$ and that the random variable $|\text{Bias}_{5_e} F_e| \leq 2L$ almost surely. Thus, we can use standard martingale concentration inequalities to get a tight enough concentration bound on $\sum_e \text{Bias}_{5_e} F_e - \mathbb{E}[\text{Bias}_{5_e} F_e | \mathcal{H}_e]$ and further bound $\sum_e \text{Bias}_{5_e} F_e$. Finally, we have that $\sum_e \text{Bias}_{5_e} F_e = \sum_e \text{Bias}_{5_e}$ with high probability. Thus, a high probability bound on $\sum_e \text{Bias}_{5_e} F_e$ transfers to a high probability bound on $\sum_e \text{Bias}_{5_e}$.

5 CONCLUSIONS

We reanalyze the algorithm proposed by Schneider & Zimmert (2023) and show that it actually achieves near-optimal regret with *high probability* for the cross-learning contextual bandits problem when the losses are chosen adversarially but the contexts are i.i.d. sampled from an *unknown* distribution. Our key technique is utilizing the weak dependency structure between different epochs for an algorithm executing over multiple epochs. It is of interest to investigate that whether this techniques is applicable for deriving high probability bounds for algorithms executing over multiple epochs in other problems.

REFERENCES

- Rui Ai, Chang Wang, Chenchen Li, Jinshan Zhang, Wenhan Huang, and Xiaotie Deng. No-regret learning in repeated first-price auctions with budget constraints, 2022.
- Santiago Balseiro, Negin Golrezaei, Mohammad Mahdian, Vahab Mirrokni, and Jon Schneider. Contextual bandits with cross-learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yan Dai, Haipeng Luo, Chen-Yu Wei, and Julian Zimmert. Refined regret for adversarial mdps with linear function approximation. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 6726–6759. PMLR, 2023.
- Yanjun Han, Zhengyuan Zhou, Aaron Flores, Erik Ordentlich, and Tsachy Weissman. Learning to bid optimally and efficiently in adversarial first-price auctions. *ArXiv*, abs/2007.04568, 2020a.
- Yanjun Han, Zhengyuan Zhou, and Tsachy Weissman. Optimal no-regret learning in repeated first-price auctions. *ArXiv*, abs/2003.09795, 2020b.
- Osama A Hanna, Lin Yang, and Christina Fragouli. Contexts can be cheap: Solving stochastic contextual bandits with linear bandit algorithms. In Gergely Neu and Lorenzo Rosasco (eds.), *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pp. 1791–1821. PMLR, 12–15 Jul 2023.
- Satyen Kale, Lev Reyzin, and Robert E Schapire. Non-stochastic bandit slate problems. In J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta (eds.), *Advances in Neural Information Processing Systems*, volume 23. Curran Associates, Inc., 2010.
- Satyen Kale, Chansoo Lee, and Dávid Pál. Hardness of online sleeping combinatorial optimization problems. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pp. 2181–2189, 2016.
- Robert D. Kleinberg, Alexandru Niculescu-Mizil, and Yogeshwer Sharma. Regret bounds for sleeping experts and bandits. *Machine Learning*, 80:245–272, 2010.

- 486 Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. A contextual-bandit approach to per-
487 sonalized news article recommendation. In *Proceedings of the 19th International Conference on*
488 *World Wide Web, WWW '10*, pp. 661–670, New York, NY, USA, 2010. Association for Comput-
489 ing Machinery. ISBN 9781605587998. doi: 10.1145/1772690.1772758.
- 490
491 Haolin Liu, Chen-Yu Wei, and Julian Zimmert. Bypassing the simulator: Near-optimal adversarial
492 linear contextual bandits. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and
493 S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 52086–
494 52131. Curran Associates, Inc., 2023.
- 495 Gergely Neu. Explore no more: Improved high-probability regret bounds for non-stochastic bandits.
496 In Corinna Cortes, Neil D. Lawrence, Daniel D. Lee, Masashi Sugiyama, and Roman Garnett
497 (eds.), *Advances in Neural Information Processing Systems 28: Annual Conference on Neural*
498 *Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pp.
499 3168–3176, 2015.
- 500 Gergely Neu and Julia Olkhovskaya. Efficient and robust algorithms for adversarial linear contextual
501 bandits. In Jacob Abernethy and Shivani Agarwal (eds.), *Proceedings of Thirty Third Conference*
502 *on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pp. 3049–3068.
503 PMLR, 09–12 Jul 2020.
- 504
505 Gergely Neu and Michal Valko. Online combinatorial optimization with stochastic decision sets and
506 adversarial losses. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger
507 (eds.), *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc.,
508 2014.
- 509 Aadirupa Saha, Pierre Gaillard, and Michal Valko. Improved sleeping bandits with stochastic actions
510 sets and adversarial rewards. In *Proceedings of the 37th International Conference on Machine*
511 *Learning, ICML'20*. JMLR.org, 2020.
- 512
513 Jon Schneider and Julian Zimmert. Optimal cross-learning for contextual bandits with unknown
514 context distributions. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine
515 (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 51862–51880. Curran
516 Associates, Inc., 2023.
- 517 Vasilis Syrgkanis, Haipeng Luo, Akshay Krishnamurthy, and Robert E. Schapire. Improved regret
518 bounds for oracle-based adversarial contextual bandits. In Daniel D. Lee, Masashi Sugiyama,
519 Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett (eds.), *Advances in Neural Information*
520 *Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016,*
521 *December 5-10, 2016, Barcelona, Spain*, pp. 3135–3143, 2016.
- 522
523 Sofía S. Villar, Jack Bowden, and James Wason. Multi-armed Bandit Models for the Optimal Design
524 of Clinical Trials: Benefits and Challenges. *Statistical Science*, 30(2):199 – 215, 2015. doi:
525 10.1214/14-STSS04.
- 526 Qian Wang, Zongjun Yang, Xiaotie Deng, and Yuqing Kong. Learning to bid in repeated first-
527 price auctions with budgets. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara
528 Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International*
529 *Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*,
530 pp. 36494–36513. PMLR, 23–29 Jul 2023.

532 A USEFUL LEMMAS

533
534
535 **Lemma 1** (Freedman’s Inequality). *Fix any $\lambda > 0$ and $\delta \in (0, 1)$. Let X_t be a random process*
536 *with respect to a filtration \mathcal{F}_t such that $\mu_t = \mathbb{E}[X_t | \mathcal{F}_{t-1}]$ and $V_t = \mathbb{E}[X_t^2 | \mathcal{F}_{t-1}]$, and satisfying*
537 *$\lambda X_t \leq 1$. Then, with probability at least $1 - \delta$, we have for all t ,*

$$538 \sum_{s=1}^t X_s - \mu_s \leq \lambda \sum_{s=1}^t V_s + \frac{\log(1/\delta)}{\lambda}. \quad 539$$

The next lemma is about the following family of indicator functions.

Definition 1. For each epoch e , we define the following two indicator functions:

$$F_e \triangleq \mathbb{1} \left(\forall a, \left| \widehat{f}_e(a) - f_e(a) \right| \leq 2 \max \left\{ \sqrt{\frac{f_e(a)\iota}{L}}, \frac{\iota}{L} \right\} \right)$$

and

$$L_e \triangleq \mathbb{1} \left(\max_{c \in [C], a \in [K]} \sum_{t \in \mathcal{T}_e} \widetilde{\ell}_{t,c}(a) \leq L + \frac{\iota}{\gamma} \right).$$

We further define the following indicator function:

$$G = \prod_{e=1}^{T/L} F_e L_e.$$

Lemma 2 (Lemma 6 and Lemma 7, Schneider & Zimmert (2023)). For any epoch e , the event F_e holds with probability at least $1 - 2K \exp(-\iota)$, and the event L_e holds with probability at least $1 - K \exp(-\iota)$. Furthermore, the event G holds with probability at least $1 - 3K(T/L) \exp(-\iota)$.

Lemma 3 (Lemma 8, Schneider & Zimmert (2023)). Let $\gamma \geq \frac{4\iota}{L}$, then under event F_e , we have that

$$\frac{1}{2} \leq \frac{f_e(a) + \gamma}{\widehat{f}_e(a) + \frac{3}{2}\gamma} \leq 2.$$

The next lemma is about the following auxiliary probability vector.

Definition 2. For each epoch e and each time step $t \in \mathcal{T}_e$, we define

$$\widetilde{p}_{t,c} \triangleq \arg \min_{p \in \Delta([K])} \left\langle p, \sum_{e'=1}^{e-1} \sum_{s \in \mathcal{T}_{e'}} \widehat{\ell}_{s,c} + \sum_{t' \in \mathcal{T}_e, t' < t} \widetilde{\ell}_{t',c} \right\rangle - \eta^{-1} F(p)$$

where $F(p) = \sum_{i=1}^K p_i \log(p_i)$ is the unnormalized negative entropy.

It is easy to see that $\widetilde{p}_{t,c} \propto s_{e+1,c} \circ \exp \left(-\eta \sum_{t' \in \mathcal{T}_e, t' < t} \widetilde{\ell}_{t',c} \right)$.

Lemma 4 (Lemma 9, Schneider & Zimmert (2023)). If $\gamma \geq \frac{4\iota}{L}$ and $\eta \leq \frac{\log(2)}{5L}$, then under event G , we have for all $t \in \mathcal{T}_e, a \in [K], c \in [C]$ simultaneously

$$2s_{e,c}(a) \geq p_{t,c}(a) \geq s_{e,c}(a)/2 \quad \text{and} \quad 2s_{e,c}(a) \geq \widetilde{p}_{t,c}(a) \geq s_{e,c}(a)/2.$$

This implies that

$$\mathbb{E}_{c \sim \nu} [p_{t,c}(a)] \leq 4f_e(a) \quad \text{and} \quad \mathbb{E}_{c \sim \nu} [\widetilde{p}_{t,c}(a)] \leq 4f_e(a).$$

In addition, this implies that $q_t = p_t$ for all $t \in \mathcal{T}_e$.

Definition 3. We define $p_t(a) \triangleq \mathbb{E}_{c \sim \nu} [p_{t,c}(a)]$ and $\widetilde{p}_t(a) \triangleq \mathbb{E}_{c \sim \nu} [\widetilde{p}_{t,c}(a)]$ for each time step t and each arm a .

Lemma 5 (Lemma 10, Schneider & Zimmert (2023)). If $\gamma \geq \frac{16\iota}{L}$ and $\exp(-\iota) \leq \frac{\gamma}{8K}$, then

$$-\frac{\gamma}{f_e(a)} \leq \mathbb{E} \left[\frac{f_e(a) - \widehat{f}_e(a) - \frac{1}{2}\gamma}{\widehat{f}_e(a) + \frac{3}{2}\gamma} F_e \middle| \mathcal{H}_{e-1} \right] \leq 0.$$

Lemma 6. For any $\eta \leq \frac{\gamma}{2(2L\gamma + \iota)}$, $\gamma \geq \frac{16\iota}{L}$, $\iota \geq \log(8K/\gamma)$, we have

$$\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left\langle p_{t,c} - \widetilde{p}_{t,c}, \widetilde{\ell}_{t,c} - \widehat{\ell}_{t,c} \right\rangle G \leq \frac{98KT\iota}{L} + \frac{\gamma^2 LKT}{\iota}.$$

Proof. The proof of Lemma 6 is contained in the analysis of the bias_3 term in Schneider & Zimmert (2023). \square

Lemma 7. *Decomposing all time steps into consecutive pairs, specifically, decomposing $\{1, 2, \dots, T\}$ into $\{(1, 2), (3, 4), (5, 6), \dots, (t-1, t), \dots, (T-1, T)\}$. Constructing a surrogate loss sequence $\{\tilde{\ell}_s\}_{s=1}^{\frac{T}{2}}$ such that for each surrogate time step s the loss vector $\tilde{\ell}_s$ is uniformly sampled from the pair of true loss vector (ℓ_{2s-1}, ℓ_{2s}) . Denote the time step sampled from the pair $(2s-1, 2s)$ as s_ℓ . For any constant $\delta \in (0, 1)$ and any bandit algorithm such that in each pair of time steps $(t-1, t)$, the algorithm takes actions a_{t-1} and a_t from the same distribution $p_{t-1} = p_t$, we have*

$$\sum_{t=1}^T (\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t})) \leq 2 \sum_{s=1}^{\frac{T}{2}} (\tilde{\ell}_{s,c_{s_\ell}}(a_{s_\ell}) - \tilde{\ell}_{s,c_{s_\ell}}(\pi_{c_{s_\ell}})) + 2\sqrt{T \log(\frac{1}{\delta})}$$

with probability at least $1 - \delta$.

Proof of lemma 7. Consider the sequence of random variable $\{Y_s\}_{s=1}^{\frac{T}{2}}$ such that

$$\begin{aligned} Y_s &= \ell_{2s-1,c_{2s-1}}(a_{2s-1}) - \ell_{2s-1,c_{2s-1}}(\pi_{c_{2s-1}}) \\ &\quad + \ell_{2s,c_{2s}}(a_{2s}) - \ell_{2s,c_{2s}}(\pi_{c_{2s}}) \\ &\quad - 2 \left(\tilde{\ell}_{s,c_{s_\ell}}(a_{s_\ell}) - \tilde{\ell}_{s,c_{s_\ell}}(\pi_{c_{s_\ell}}) \right). \end{aligned}$$

Consider the filtration $\{\tilde{H}_s\}_{s=1}^{\frac{T}{2}}$ such that for each s the σ -field \tilde{H}_s is generated by the randomness within c_t and a_t for $t \leq 2s$ and the randomness within sampling from pair $(2\tau-1, 2\tau)$ for each $\tau \leq s$. It is easy to see that the sequence $\{Y_s\}_{s=1}^{\frac{T}{2}}$ forms a martingale difference sequence adapted to the filtration $\{\tilde{H}_s\}_{s=1}^{\frac{T}{2}}$. Moreover, it is also easy to see that $|Y_s| \leq 2$. Using Azuma-Hoeffding's inequality, for any constant $\delta \in (0, 1)$, we have

$$\sum_{s=1}^{\frac{T}{2}} Y_s \leq 2\sqrt{T \log(\frac{1}{\delta})}$$

with probability at least $1 - \delta$. This completes the proof of the lemma. \square

B DETAILED PROOF OF THEOREM 1

B.1 DECOMPOSITION

As we mentioned in Section 4, we decompose the regret as

$$\begin{aligned} \text{Reg}(\pi) &= \underbrace{\sum_{t=1}^T \ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - 2 \sum_{t \in \mathcal{T}^\ell} \ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t})}_{\text{bias}_1} \\ &\quad + 2 \underbrace{\sum_{t \in \mathcal{T}^\ell} \left(\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - \sum_c \Pr(c) \langle p_{t,c} - \pi_c, \ell_{t,c} \rangle \right)}_{\text{bias}_2} \\ &\quad + 2 \underbrace{\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c} - \pi_c, \hat{\ell}_{t,c} \rangle}_{\text{ftrl}} + 2 \underbrace{\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \ell_{t,c} - \tilde{\ell}_{t,c} \rangle}_{\text{bias}_3} \\ &\quad + 2 \underbrace{\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \rangle}_{\text{bias}_4} + 2 \underbrace{\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle \pi_c, \hat{\ell}_{t,c} - \ell_{t,c} \rangle}_{\text{bias}_5}. \end{aligned}$$

We bound these terms one by one.

The **bias**₁, **bias**₂, **ftrl**, and **bias**₃ terms are not hard to bound. The terms of interest to bound are **bias**₄ and **bias**₅. We first bound these two terms. In these two terms, the **bias**₅ term is the one easier to bound. We first bound **bias**₅ to provide some intuition for our readers.

B.2 UPPER BOUNDING bias_5

We first bound the fifth term $\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle \pi_c, \widehat{\ell}_{t,c} - \ell_{t,c} \rangle$. We decompose the fifth term into two components:

$$\begin{aligned} & \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle \pi_c, \widehat{\ell}_{t,c} - \ell_{t,c} \rangle \\ &= \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \ell_{t,c}(\pi_c) \right) \\ &= \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right) \\ & \quad + \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] \right). \end{aligned}$$

We bound these two components separately.

For each epoch e , we define a random variable

$$\text{Bias5}_e \triangleq \sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right).$$

We rewrite the term $\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right)$ as $\sum_{e=1}^{T/L} \text{Bias5}_e$. Recall our key observation: different epochs are only weakly dependent on each other. We bound the summation over epochs $\sum_{e=1}^{T/L} \text{Bias5}_e$ by leveraging the weak dependence structure between $\{\text{Bias5}_e\}_e$.

The sequence of random variables $\{\text{Bias5}_e\}_e$ has the following properties:

1. For each epoch e , the random variable Bias5_e is measurable under σ -algebra \mathcal{H}_e .
2. For each epoch e , we have³

$$\begin{aligned} & \mathbb{E} \left[\text{Bias5}_e \cdot F_e \mid \mathcal{H}_{e-1} \right] \\ &= \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^\ell} \ell_{t,c}(\pi_c) \mathbb{E} \left[\left(\frac{f_e(\pi_c)}{\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma} - 1 \right) F_e \mid \mathcal{H}_{e-1} \right]. \end{aligned}$$

We further have

$$\begin{aligned} & \mathbb{E} \left[\left(\frac{f_e(\pi_c)}{\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma} - 1 \right) F_e \mid \mathcal{H}_{e-1} \right] \\ &= \mathbb{E} \left[\left(\frac{f_e(\pi_c)}{\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma} - \frac{f_e(\pi_c)}{f_e(\pi_c) + \gamma} + \frac{f_e(\pi_c)}{f_e(\pi_c) + \gamma} - 1 \right) F_e \mid \mathcal{H}_{e-1} \right] \\ &= \mathbb{E} \left[\frac{f_e(\pi_c)}{f_e(\pi_c) + \gamma} \frac{\left(f_e(\pi_c) - \widehat{f}_e(\pi_c) - \frac{1}{2}\gamma \right)}{\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma} F_e \mid \mathcal{H}_{e-1} \right] - \frac{\gamma}{f_e(\pi_c) + \gamma} \mathbb{E} \left[F_e \mid \mathcal{H}_{e-1} \right] \\ &\leq - \frac{\gamma}{f_e(\pi_c) + \gamma} \mathbb{E} \left[F_e \mid \mathcal{H}_{e-1} \right] \tag{Lemma 5} \\ &\leq - \frac{\gamma}{f_e(\pi_c) + \gamma} (1 - 2K \exp(-\iota)). \tag{Lemma 2} \end{aligned}$$

³Readers familiar with Schneider & Zimmert (2023) may wonder why we do not directly consider $\text{Bias5}_e G$ but consider $\text{Bias5}_e F_e$ instead. This is because there is a small flaw in the argument of Schneider & Zimmert (2023). Schneider & Zimmert (2023) essentially argues that $\mathbb{E}[\text{Bias5}_e G \mid \mathcal{H}_{e-1}] = \mathbb{E}[\text{Bias5}_e \mid \mathcal{H}_{e-1}] \mathbb{E}[G \mid \mathcal{H}_{e-1}]$. However, this equality may not hold since the indicator G depends on Bias5_e and these two terms are not conditionally independent given \mathcal{H}_{e-1} . This is why we consider $\text{Bias5}_e F_e$ here instead.

Thus we have

$$\begin{aligned} & \mathbb{E} [\text{Bias5}_e \cdot F_e \mid \mathcal{H}_{e-1}] \\ & \leq - \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^t} \ell_{t,c}(\pi_c) \frac{\gamma}{f_e(\pi_c) + \gamma} (1 - 2K \exp(-t)). \end{aligned}$$

3. For each epoch e , we have

$$\begin{aligned} \text{Bias5}_e F_e & \leq \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^t} \mathbb{E} [\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1}] F_e \\ & = \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^t} \ell_{t,c}(\pi_c) \frac{f_e(\pi_c)}{\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma} F_e \\ & \leq \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^t} \ell_{t,c}(\pi_c) \frac{2f_e(\pi_c)}{f_e(\pi_c) + \gamma} \quad (\text{Lemma 3}) \\ & \leq 2L = \frac{32t}{\gamma}. \end{aligned}$$

4. For each epoch e , we have

$$\begin{aligned} & \mathbb{E} [(\text{Bias5}_e F_e)^2 \mid \mathcal{H}_{e-1}] \\ & = \mathbb{E} \left[\left(\sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^t} \ell_{t,c}(\pi_c) \frac{f_e(\pi_c) - \widehat{f}_e(\pi_c) - \frac{3}{2}\gamma}{\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma} F_e \right)^2 \mid \mathcal{H}_{e-1} \right] \\ & \leq \sum_c \Pr(c) \mathbb{E} \left[\left(\sum_{t \in \mathcal{T}_e^t} \ell_{t,c}(\pi_c) \frac{f_e(\pi_c) - \widehat{f}_e(\pi_c) - \frac{3}{2}\gamma}{\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma} F_e \right)^2 \mid \mathcal{H}_{e-1} \right] \\ & = \sum_c \Pr(c) \left(\sum_{t \in \mathcal{T}_e^t} \ell_{t,c} \right)^2 \mathbb{E} \left[\left(\frac{f_e(\pi_c) - \widehat{f}_e(\pi_c) - \frac{3}{2}\gamma}{\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma} F_e \right)^2 \mid \mathcal{H}_{e-1} \right] \\ & \leq L \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^t} \ell_{t,c} \mathbb{E} \left[\left(\frac{f_e(\pi_c) - \widehat{f}_e(\pi_c) - \frac{3}{2}\gamma}{\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma} F_e \right)^2 \mid \mathcal{H}_{e-1} \right] \\ & \leq 4L \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^t} \ell_{t,c} \mathbb{E} \left[\left(\frac{f_e(\pi_c) - \widehat{f}_e(\pi_c) - \frac{3}{2}\gamma}{f_e(\pi_c) + \gamma} F_e \right)^2 \mid \mathcal{H}_{e-1} \right] \quad (\text{Lemma 3}) \\ & \leq \sum_c \Pr(c) \frac{4L}{(f_e(\pi_c) + \gamma)^2} \sum_{t \in \mathcal{T}_e^t} \ell_{t,c}(\pi_c) \mathbb{E} \left[\left(f_e(\pi_c) - \widehat{f}_e(\pi_c) - \frac{3}{2}\gamma \right)^2 \mid \mathcal{H}_{e-1} \right] \\ & = \sum_c \Pr(c) \frac{4L}{(f_e(\pi_c) + \gamma)^2} \sum_{t \in \mathcal{T}_e^t} \ell_{t,c}(\pi_c) \left(\mathbb{E} \left[\left(f_e(\pi_c) - \widehat{f}_e(\pi_c) \right)^2 \mid \mathcal{H}_{e-1} \right] + \frac{9}{4}\gamma^2 \right) \\ & \leq \sum_c \Pr(c) \frac{4L}{(f_e(\pi_c) + \gamma)^2} \sum_{t \in \mathcal{T}_e^t} \sum_{t \in \mathcal{T}_e^t} \ell_{t,c}(\pi_c) \left(\frac{f_e(\pi_c)}{L} + \frac{9}{4}\gamma^2 \right) \\ & \leq 4 \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^t} \ell_{t,c}(\pi_c) \left(\frac{1}{f_e(\pi_c) + \gamma} + \frac{9L\gamma}{4(f_e(\pi_c) + \frac{3}{2}\gamma)} \right) \\ & \leq 4 \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^t} \ell_{t,c}(\pi_c) \frac{36t}{f_e(\pi_c) + \gamma}. \end{aligned}$$

Given these properties, we use Freedman's inequality to get that for any $0 < \lambda < \frac{\gamma}{32\iota}$, with probability at least $1 - \delta$, we have

$$\sum_e \text{Bias5}_e F_e - \mathbb{E}[\text{Bias5}_e F_e | \mathcal{H}_{e-1}] \leq 4\lambda \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^\ell} \ell_{t,c}(\pi_c) \frac{36\iota}{f_e(\pi_c) + \frac{3}{2}\gamma} + \frac{\log(1/\delta)}{\lambda}.$$

We further have that event $\{\sum_e \text{Bias5}_e F_e = \sum_e \text{Bias5}_e\}$ holds if event G holds. Combining these two facts, we get that the inequality

$$\sum_e \text{Bias5}_e - \mathbb{E}[\text{Bias5}_e F_e | \mathcal{H}_{e-1}] \leq 4\lambda \sum_c \Pr(c) \sum_{t \in \mathcal{T}_e^\ell} \ell_{t,c}(\pi_c) \frac{36\iota}{f_e(\pi_c) + \gamma} + \frac{\log(1/\delta)}{\lambda}$$

holds with probability at least $\Pr(G) - \delta$.

We now bound the second component

$$\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] \right).$$

The second term $\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] \right)$ has the following properties:

1. The sequence of random variables $\left\{ \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] \right) \right\}_{t \in \mathcal{T}^\ell}$ forms a martingale difference sequence with respect to the filtration $\{\mathcal{H}_t\}_t$.
2. Each random variable $\sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] \right)$ satisfies $\left| \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] \right) \right| \leq \frac{1}{\gamma}$.
3. Each random variable $\sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] \right)$ satisfies

$$\begin{aligned} \text{Var} \left[\sum_c \Pr(c) \widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] &\leq \sum_c \Pr(c) \text{Var} \left[\widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] \\ &= \sum_c \Pr(c) \frac{f_e(\pi_c) - f_e^2(\pi_c)}{\left(\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma \right)^2} \ell_{t,c}(\pi_c)^2 \\ &\leq \sum_c \Pr(c) \frac{f_e(\pi_c)}{\left(\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma \right)^2} \ell_{t,c}(\pi_c). \end{aligned}$$

Applying Freedman's inequality to the sequence of random variables $\left\{ \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] \right) \right\}_{t \in \mathcal{T}^\ell}$, we get that for each $\delta \in (0, 1)$ and each $0 < \lambda < \gamma$, with probability at least $1 - \delta$, we have

$$\begin{aligned} &\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] \right) \\ &\leq \lambda \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \frac{f_e(\pi_c)}{\left(\widehat{f}_e(\pi_c) + \frac{3}{2}\gamma \right)^2} \ell_{t,c}(\pi_c) + \frac{1}{\lambda} \log\left(\frac{1}{\delta}\right). \end{aligned}$$

By assuming that event G holds, we further get that with probability at least $\Pr(G) - \delta$, the following inequality holds:

$$\begin{aligned} &\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) | \mathcal{H}_{t-1} \right] \right) \\ &\leq 4\lambda \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \frac{f_e(\pi_c)}{\left(f_e(\pi_c) + \gamma \right)^2} \ell_{t,c}(\pi_c) + \frac{1}{\lambda} \log\left(\frac{1}{\delta}\right). \end{aligned} \quad (\text{Lemma 3})$$

Combining all previous inequalities, we get that for any $0 < \lambda_1 < \frac{\gamma}{32\iota}$ and $0 < \lambda_2 < \gamma$, the following inequality holds with probability at least $\Pr(G) - 2\delta$:

$$\begin{aligned}
& \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle \pi_c, \widehat{\ell}_{t,c} - \ell_{t,c} \rangle \\
&= \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] - \ell_{t,c}(\pi_c) \right) \\
&\quad + \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] \right) \\
&= \sum_e \text{Bias5}_e - \mathbb{E}[\text{Bias5}_e F_e \mid \mathcal{H}_{e-1}] \\
&\quad + \sum_e \mathbb{E} [\text{Bias5}_e F_e \mid \mathcal{H}_{e-1}] \\
&\quad + \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left(\widehat{\ell}_{t,c}(\pi_c) - \mathbb{E} \left[\widehat{\ell}_{t,c}(\pi_c) \mid \mathcal{H}_{t-1} \right] \right) \\
&\leq 4\lambda_1 \sum_c \Pr(c) \sum_{t \in \mathcal{T}^\ell} \ell_{t,c}(\pi_c) \frac{36\iota}{f_e(\pi_c) + \gamma} + \frac{\log(1/\delta)}{\lambda_1} \\
&\quad - \sum_c \Pr(c) \sum_{t \in \mathcal{T}^\ell} \ell_{t,c}(\pi_c) \frac{\gamma}{f_e(\pi_c) + \gamma} (1 - 2K \exp(-\iota)) \\
&\quad + 4\lambda_2 \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \frac{f_e(\pi_c)}{(f_e(\pi_c) + \gamma)^2} \ell_{t,c}(\pi_c) + \frac{1}{\lambda_2} \log\left(\frac{1}{\delta}\right).
\end{aligned}$$

Note that in the previous analysis, we combined two good events each happening with probability at least $\Pr(G) - \delta$. The combined good event happens with probability $\Pr(G) - 2\delta$ rather than the vanilla union bound $1 - 2(1 - \Pr(G) + \delta)$. This is because, in both events, the $\Pr(G)$ term comes from assuming event G happens. Thus, in the combined event, we can simply assume event G happens and count the corresponding bad event G^c only once. We will use this small trick repeatedly in the following analysis.

We pick $\lambda_1 = \frac{\gamma}{8 \cdot 36\iota}$ and $\lambda_2 = \frac{\gamma}{8}$ to get that the following inequality holds with probability at least $\Pr(G) - 2\delta$:

$$\begin{aligned}
& \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle \pi_c, \widehat{\ell}_{t,c} - \ell_{t,c} \rangle \\
&\leq 4\lambda_1 \sum_c \Pr(c) \sum_{t \in \mathcal{T}^\ell} \ell_{t,c}(\pi_c) \frac{36\iota}{f_e(\pi_c) + \gamma} + \frac{\log(1/\delta)}{\lambda_1} \\
&\quad - \sum_c \Pr(c) \sum_{t \in \mathcal{T}^\ell} \ell_{t,c}(\pi_c) \frac{\gamma}{f_e(\pi_c) + \gamma} (1 - 2K \exp(-\iota)) \\
&\quad + 4\lambda_2 \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \frac{f_e(\pi_c)}{(f_e(\pi_c) + \gamma)^2} \ell_{t,c}(\pi_c) + \frac{1}{\lambda_2} \log\left(\frac{1}{\delta}\right) \\
&= \left(\frac{8 \cdot 36\iota}{\gamma} + \frac{8}{\gamma} \right) \log(1/\delta) + \sum_c \Pr(c) \sum_{t \in \mathcal{T}^\ell} \ell_{t,c}(\pi_c) \frac{\gamma}{f_e(\pi_c) + \gamma} 2K \exp(-\iota) \\
&\leq \left(\frac{8 \cdot 36\iota}{\gamma} + \frac{8}{\gamma} \right) \log(1/\delta) + KT \exp(-\iota).
\end{aligned}$$

B.3 UPPER BOUNDING bias_4

We then bound the forth term $\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \rangle$. Similar to the previous analysis, we decompose it as follows:

$$\begin{aligned} & \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \rangle \\ &= \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \mathbb{E} [\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} | \mathcal{H}_{t-1}] \rangle \\ & \quad + \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} - \mathbb{E} [\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} | \mathcal{H}_{t-1}] \rangle. \end{aligned}$$

We bound these two components separately.

Similar to the previous analysis, we decompose the first component

$$\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \mathbb{E} [\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} | \mathcal{H}_{t-1}] \rangle$$

as

$$\sum_e \sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \langle p_{t,c}, \mathbb{E} [\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} | \mathcal{H}_{t-1}] \rangle.$$

For each epoch e we define a random variable

$$\text{Bias4}_e \triangleq \sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \langle p_{t,c}, \mathbb{E} [\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} | \mathcal{H}_{t-1}] \rangle.$$

We need to bound $\sum_e \text{Bias4}_e$.

We decompose $\sum_e \text{Bias4}_e$ as $\sum_e \text{Bias4}_e F_e L_e + \sum_e \text{Bias4}_e (1 - F_e L_e)$. As usual we have that $\sum_e \text{Bias4}_e (1 - F_e L_e) = 0$ whenever event G holds. Thus we can focus on bounding $\sum_e \text{Bias4}_e F_e L_e$. Firstly we bound

$$\sum_e \mathbb{E} [\text{Bias4}_e F_e L_e | \mathcal{H}_{e-1}].$$

We have

$$\begin{aligned} & \sum_e \mathbb{E} [\text{Bias4}_e F_e L_e | \mathcal{H}_{e-1}] \\ &= \sum_e \mathbb{E} \left[\sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \langle \tilde{p}_{t,c}, \mathbb{E} [\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} | \mathcal{H}_{t-1}] \rangle F_e L_e | \mathcal{H}_{e-1} \right] \\ & \quad + \sum_e \mathbb{E} \left[\sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \langle p_{t,c} - \tilde{p}_{t,c}, \mathbb{E} [\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} | \mathcal{H}_{t-1}] \rangle F_e L_e | \mathcal{H}_{e-1} \right]. \end{aligned}$$

By Lemma 6, the latter term

$$\sum_e \mathbb{E} \left[\sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \langle p_{t,c} - \tilde{p}_{t,c}, \mathbb{E} [\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} | \mathcal{H}_{t-1}] \rangle F_e L_e | \mathcal{H}_{e-1} \right]$$

is bounded by $\frac{98KT}{L} + \frac{\gamma^2 LKT}{L}$. Furthermore, condition on \mathcal{H}_{e-1} , the indicator function F_e is affected only by randomness within time steps $t \in \mathcal{T}_e^f$, thus the indicator function F_e is conditional

independent with the probability vector $\tilde{p}_{t,c}$. We have

$$\begin{aligned}
& \sum_e \mathbb{E} \left[\sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \left\langle \tilde{p}_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle F_e \mid \mathcal{H}_{e-1} \right] \\
&= \sum_e \sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \left\langle \mathbb{E} [\tilde{p}_{t,c} \mid \mathcal{H}_{e-1}], \mathbb{E} \left[(\tilde{\ell}_{t,c} - \hat{\ell}_{t,c}) F_e \mid \mathcal{H}_{e-1} \right] \right\rangle \\
&\leq \sum_e \sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \sum_a \frac{\tilde{p}_{t,c}(a) \gamma}{f_e(a)}. \tag{Lemma 5}
\end{aligned}$$

By Lemma 4, whenever event G holds, the ratio $\frac{\tilde{p}_{t,c}(a)}{f_e(a)} \leq 4$. Thus we have

$$\sum_e \mathbb{E} \left[\sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \left\langle \tilde{p}_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle F_e \mid \mathcal{H}_{e-1} \right] \leq 4\gamma KT$$

whenever event G holds.

We further have $\left| \left\langle \tilde{p}_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle \right| \leq \frac{1}{\gamma}$. Thus we have

$$\begin{aligned}
& \sum_e \mathbb{E} \left[\sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \left\langle \tilde{p}_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle F_e (L_e - 1) \mid \mathcal{H}_{e-1} \right] \\
&\leq \sum_e \mathbb{E} \left[\sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \frac{1}{\gamma} F_e (L_e - 1) \mid \mathcal{H}_{e-1} \right] \\
&\leq \frac{1}{\gamma} \sum_e \mathbb{E} \left[\sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) |L_e - 1| \mid \mathcal{H}_{e-1} \right] \\
&\leq \frac{K \exp(-\iota) T}{\gamma}. \tag{Lemma 2}
\end{aligned}$$

Thus we have

$$\begin{aligned}
& \sum_e \mathbb{E} \left[\sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \left\langle \tilde{p}_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle F_e L_e \mid \mathcal{H}_{e-1} \right] \\
&\leq 4\gamma KT + \frac{K \exp(-\iota) T}{\gamma}
\end{aligned}$$

whenever event G holds.

Thus we have

$$\sum_e \mathbb{E} [\text{Bias}_{4_e} F_e L_e \mid \mathcal{H}_{e-1}] \leq 4K\gamma T + \frac{K \exp(-\iota) T}{\gamma} + \frac{98KT\iota}{L} + \frac{\gamma^2 LKT}{\iota}$$

whenever event G holds.

We then only need to bound the concentration term

$$\sum_e \text{Bias}_{4_e} F_e L_e - \mathbb{E} [\text{Bias}_{4_e} F_e L_e \mid \mathcal{H}_{e-1}].$$

For each random variable $\text{Bias}_{4_e} F_e L_e$, we have

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

$$\begin{aligned}
& \text{Bias4}_e F_e L_e \\
&= \sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \left\langle p_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle F_e L_e \\
&= \sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \sum_a p_{t,c}(a) f_e(a) \ell_{t,c}(a) \left(\frac{1}{f_e(a) + \gamma} - \frac{1}{\hat{f}_e(a) + \frac{3}{2}\gamma} \right) F_e L_e \\
&= \sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \sum_a \tilde{p}_{t,c}(a) f_e(a) \ell_{t,c}(a) \frac{\hat{f}_e(a) - f_e(a) + \frac{1}{2}\gamma}{(f_e(a) + \gamma)(\hat{f}_e(a) + \frac{3}{2}\gamma)} F_e L_e.
\end{aligned}$$

Thus we have

$$\begin{aligned}
& |\text{Bias4}_e F_e L_e| \\
&\leq \left| \sum_{t \in \mathcal{T}_e^\ell} \sum_c \Pr(c) \sum_a \tilde{p}_{t,c}(a) f_e(a) \ell_{t,c}(a) \frac{\hat{f}_e(a) - f_e(a) + \frac{1}{2}\gamma}{(f_e(a) + \gamma)(\hat{f}_e(a) + \frac{3}{2}\gamma)} \right| F_e L_e \\
&\leq \sum_{t \in \mathcal{T}_e^\ell} \sum_a \tilde{p}_t(a) \left| \frac{\hat{f}_e(a) - f_e(a) + \frac{1}{2}\gamma}{\hat{f}_e(a) + \frac{3}{2}\gamma} \right| F_e L_e \\
&\leq 8 \sum_{t \in \mathcal{T}_e^\ell} \sum_a \max \left\{ \sqrt{\frac{f_e(a)\iota}{L}}, \frac{\iota}{L} \right\} \\
&\leq 8L \left(\sqrt{\frac{K\iota}{L}} + \frac{K\iota}{L} \right) \\
&= 8(\sqrt{KL\iota} + K\iota).
\end{aligned}$$

Applying Azuma-Hoeffding's inequality to

$$\sum_e \text{Bias4}_e F_e L_e - \mathbb{E} [\text{Bias4}_e F_e L_e \mid \mathcal{H}_{e-1}],$$

we get that for any $\delta > 0$, with probability at least $1 - \delta$, we have

$$\begin{aligned}
& \sum_e \text{Bias4}_e F_e L_e - \mathbb{E} [\text{Bias4}_e F_e L_e \mid \mathcal{H}_{e-1}] \\
&\leq 8(\sqrt{KL\iota} + K\iota) \sqrt{2 \frac{T}{L} \log\left(\frac{\delta}{2}\right)} \\
&= 8 \left(\sqrt{2KT\iota \log\left(\frac{\delta}{2}\right)} + \sqrt{2 \frac{TK^2}{L} \iota \log\left(\frac{\delta}{2}\right)} \right).
\end{aligned}$$

Thus we have

$$\begin{aligned}
& \sum_e \text{Bias4}_e F_e L_e \\
&\leq 4K\gamma T + \frac{K \exp(-\iota)T}{\gamma} + \frac{98KT\iota}{L} + \frac{\gamma^2 LKT}{\iota} \\
&\quad + 8 \left(\sqrt{2KT\iota \log\left(\frac{\delta}{2}\right)} + \sqrt{2 \frac{TK^2}{L} \iota \log\left(\frac{\delta}{2}\right)} \right).
\end{aligned}$$

with probability at least $\Pr(G) - \delta$.

1026 We then bound the second term

1027

1028

1029

1030

1031

1032

1033

1034 For each time step $t \in \mathcal{T}_e^\ell$, we define an indicator function

1035

1036

1037

1038

1039

1040 By Lemma 4, event G implies J_t .

1041

1042 Similar to previous analysis, we decompose the first term as

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057 Since the auxiliary probability vector $\tilde{p}_{t,c}$ is determined at time $t - 1$, the indicator function J_t is also determined at time $t - 1$. Furthermore, the indicator function F_e is determined at epoch $e - 1$. Thus the product of indicator functions $F_e J_t$ is measurable under filtration \mathcal{H}_{t-1} . We have

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068 Thus the sequence of random variables

1069

1070

1071

1072

1073

1074

1075

1076 forms a martingale difference sequence under the filtration $\{\mathcal{H}_t\}_t$.

1077

1078 We further have that the term

1079

1080

$$\sum_c \Pr(c) \left\langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} - \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle F_e J_t$$

1080 satisfies

$$\begin{aligned}
& \left| \sum_c \Pr(c) \left\langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} - \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle \right| F_e J_t \\
& \leq \left| \sum_c \Pr(c) \sum_a p_{t,c}(a) \ell_{t,c}(a) \left(\frac{\mathbb{1}(a_t = a)}{f_e(a) + \gamma} - \frac{f_e(a)}{f_e(a) + \gamma} \right) \right| F_e J_t \\
& \quad + \left| \sum_c \Pr(c) \sum_a p_{t,c}(a) \ell_{t,c}(a) \left(\frac{\mathbb{1}(a_t = a)}{\hat{f}_e(a) + \frac{3}{2}\gamma} - \frac{f_e(a)}{\hat{f}_e(a) + \frac{3}{2}\gamma} \right) \right| F_e J_t \\
& \leq \sum_c \Pr(c) p_{t,c}(a_t) \frac{\ell_{t,c}(a_t)}{f_e(a_t) + \gamma} F_e J_t \\
& \quad + \sum_c \Pr(c) \sum_a p_{t,c}(a) \ell_{t,c}(a) \frac{f_e(a)}{f_e(a) + \gamma} F_e J_t \\
& \quad + \sum_c \Pr(c) p_{t,c}(a_t) \frac{\ell_{t,c}(a_t)}{\hat{f}_e(a_t) + \frac{3}{2}\gamma} F_e J_t \\
& \quad + \sum_c \Pr(c) \sum_a p_{t,c}(a) \ell_{t,c}(a) \frac{f_e(a)}{\hat{f}_e(a) + \frac{3}{2}\gamma} F_e J_t \\
& \leq \left(\frac{p_t(a_t)}{f_e(a_t) + \gamma} + 1 + \frac{p_t(a_t)}{\hat{f}_e(a_t) + \frac{3}{2}\gamma} + \sum_a p_t(a) \frac{f_e(a)}{\hat{f}_e(a) + \frac{3}{2}\gamma} \right) F_e J_t \\
& \leq 4 + 1 + 8 + 2 = 15. \tag{Lemma 3}
\end{aligned}$$

1107 Applying Azuma-Hoeffding's inequality to the sequence of random variables

$$\left\{ \sum_c \Pr(c) \left\langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} - \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle F_e J_t \right\}_{t \in \mathcal{T}^\ell},$$

1113 we get that for any $\delta > 0$, the inequality

$$\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left\langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} - \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle F_e J_t \leq 15 \sqrt{T \log\left(\frac{1}{\delta}\right)}$$

1120 holds with probability at least $1 - \delta$.

1121 On the other hand, note that event G implies event $F_e J_t$, we get that

$$\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left\langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} - \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle (1 - F_e J_t) = 0$$

1127 whenever event G holds. Thus we get that for any $\delta > 0$, inequality

$$\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left\langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} - \mathbb{E} \left[\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle \leq 15 \sqrt{T \log\left(\frac{1}{\delta}\right)}$$

1133 holds with probability at least $\Pr(G) - 2\delta$.

Combining previous results, we get that the following inequality holds with probability at least $\Pr(G) - 2\delta$:

$$\begin{aligned}
& \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} \rangle \\
&= \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \mathbb{E} [\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} | \mathcal{H}_{t-1}] \rangle \\
& \quad + \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \tilde{\ell}_{t,c} - \hat{\ell}_{t,c} - \mathbb{E} [\tilde{\ell}_{t,c} - \hat{\ell}_{t,c} | \mathcal{H}_{t-1}] \rangle \\
&\leq 4K\gamma T + \frac{K \exp(-\iota)T}{\gamma} + \frac{98KT\iota}{L} + \frac{\gamma^2 LKT}{\iota} \\
& \quad + 8 \left(\sqrt{2KT\iota \log\left(\frac{\delta}{2}\right)} + \sqrt{2\frac{TK^2}{L}\iota \log\left(\frac{\delta}{2}\right)} \right) \\
& \quad + 15\sqrt{T \log\left(\frac{1}{\delta}\right)}.
\end{aligned}$$

B.4 UPPER BOUNDING REMAINING TERMS

The remaining terms are the **bias**₁, **bias**₂, **ftl**, and **bias**₃ terms. These terms are not hard to bound using techniques in standard EXP3-IX analysis (Neu, 2015). We write down these analyses in the sake of completeness.

B.4.1 UPPER BOUNDING **bias**₁

Applying Lemma 7 on the loss sequence used in calculating loss estimates, we get that

$$\sum_{t=1}^T \ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - 2 \sum_{t \in \mathcal{T}^\ell} \ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) \leq 2\sqrt{T \log\left(\frac{1}{\delta}\right)}.$$

B.4.2 UPPER BOUNDING **bias**₂

Here we bound the **bias**₂ term

$$\sum_{t \in \mathcal{T}^\ell} \left(\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - \sum_c \Pr(c) \langle p_{t,c} - \pi_c, \ell_{t,c} \rangle \right).$$

Whenever event G holds, we have $q_t = p_t$. Thus we assume event G holds and replace the **bias**₂ term by

$$\sum_{t \in \mathcal{T}^\ell} \left(\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - \sum_c \Pr(c) \langle q_{t,c} - \pi_c, \ell_{t,c} \rangle \right).$$

The new term have the following properties:

- The sequence of random variables

$$\left\{ \ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - \sum_c \Pr(c) \langle q_{t,c} - \pi_c, \ell_{t,c} \rangle \right\}_{t \in \mathcal{T}^\ell}$$

adapts to the filtration $\{\mathcal{H}_t\}_{t \in \mathcal{T}^\ell}$.

- The sequence of random variables satisfies

$$\mathbb{E} \left[\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - \sum_c \Pr(c) \langle q_{t,c} - \pi_c, \ell_{t,c} \rangle \middle| \mathcal{H}_{t-1} \right] = 0.$$

- Each random variable satisfies $\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - \sum_c \Pr(c) \langle q_{t,c} - \pi_c, \ell_{t,c} \rangle \in [-2, 2]$.

By applying Azuma-Hoeffding inequality, we get that for any $\delta \in (0, 1)$, the following inequality holds with probability at least $1 - \delta$:

$$\sum_{t \in \mathcal{T}^\ell} \left(\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - \sum_c \Pr(c) \langle q_{t,c} - \pi_c, \ell_{t,c} \rangle \right) \leq 2\sqrt{T \log\left(\frac{1}{\delta}\right)}$$

Thus the following inequality holds with probability at least $\Pr(G) - \delta$:

$$\sum_{t \in \mathcal{T}^\ell} \left(\ell_{t,c_t}(a_t) - \ell_{t,c_t}(\pi_{c_t}) - \sum_c \Pr(c) \langle p_{t,c} - \pi_c, \ell_{t,c} \rangle \right) \leq 2\sqrt{T \log\left(\frac{1}{\delta}\right)}.$$

B.4.3 UPPER BOUNDING FTRL

Here we bound the **ftrl** term

$$\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c} - \pi_c, \hat{\ell}_{t,c} \rangle.$$

By the standard analysis of FTRL algorithms, the **ftrl** term satisfies

$$\begin{aligned} & \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c} - \pi_c, \hat{\ell}_{t,c} \rangle \\ & \leq \sum_c \Pr(c) \left(\frac{1}{\eta} \log K + \frac{\eta}{2} \sum_{t \in \mathcal{T}^\ell} \langle p_{t,c}, \hat{\ell}_{t,c}^2 \rangle \right). \end{aligned}$$

Here $\hat{\ell}_{t,c}^2$ denotes the vector formed by squaring each component of $\hat{\ell}_{t,c}$.

By Lemma 3, under event G , we have $\hat{\ell}_{t,c} \leq 2\tilde{\ell}_{t,c}$. Thus assuming event G holds, we can focus on upper bounding

$$\sum_c \Pr(c) \left(\frac{1}{\eta} \log K + 2\eta \sum_{t \in \mathcal{T}^\ell} \langle p_{t,c}, \tilde{\ell}_{t,c}^2 \rangle \right).$$

It suffices to upper bound $\sum_c \Pr(c) \sum_{t \in \mathcal{T}^\ell} \langle p_{t,c}, \tilde{\ell}_{t,c}^2 \rangle$. We have

$$\begin{aligned} & \sum_c \Pr(c) \sum_{t \in \mathcal{T}^\ell} \langle p_{t,c}, \tilde{\ell}_{t,c}^2 \rangle \\ & = \sum_c \Pr(c) \sum_{t \in \mathcal{T}^\ell} p_{t,c}(a_t) \tilde{\ell}_{t,c}^2(a_t) \\ & = \sum_c \Pr(c) \sum_{t \in \mathcal{T}^\ell} p_{t,c}(a_t) \frac{\ell_{t,c}^2(a_t)}{(f_e(a_t) + \gamma)^2} \\ & \leq \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) p_{t,c}(a_t) \frac{1}{(f_e(a_t) + \gamma)^2} \\ & = \sum_{t \in \mathcal{T}^\ell} \frac{p_t(a_t)}{(f_e(a_t) + \gamma)^2}. \end{aligned}$$

By Lemma 4, under event G , we have

$$\sum_{t \in \mathcal{T}^\ell} \frac{p_t(a_t)}{(f_e(a_t) + \gamma)^2} \leq 2 \sum_{t \in \mathcal{T}^\ell} \frac{1}{f_e(a_t) + \gamma}.$$

We then focus on upper bounding $\sum_{t \in \mathcal{T}^\ell} \frac{1}{f_e(a_t) + \gamma}$.

We have

- The sum of conditional expectations $\sum_{t \in \mathcal{T}^\ell} \mathbb{E} \left[\frac{1}{f_e(a_t) + \gamma} \mid \mathcal{H}_{t-1} \right] \leq KT$.
- Each term $\frac{1}{f_e(a_t) + \gamma} \leq \frac{1}{\gamma}$.
- The sum of conditional quadratic expectations

$$\begin{aligned}
& \sum_{t \in \mathcal{T}^\ell} \mathbb{E} \left[\left(\frac{1}{f_e(a_t) + \gamma} \right)^2 \mid \mathcal{H}_{t-1} \right] \\
&= \sum_{t \in \mathcal{T}^\ell} \sum_a \frac{f_e(a)}{(f_e(a) + \gamma)^2} \\
&\leq \sum_{t \in \mathcal{T}^\ell} \sum_a \frac{1}{f_e(a) + \gamma}.
\end{aligned}$$

By Freedman's inequality, we have that for any $\lambda \in (0, \gamma]$ and any $\delta \in (0, 1)$, the following inequality holds with probability at least $1 - \delta$:

$$\sum_{t \in \mathcal{T}^\ell} \frac{1}{f_e(a_t) + \gamma} \leq \lambda \sum_{t \in \mathcal{T}^\ell} \sum_a \frac{1}{f_e(a) + \gamma} + \frac{1}{\lambda} \log\left(\frac{1}{\delta}\right) + KT.$$

We pick $\lambda = \gamma$ to get that

$$\sum_{t \in \mathcal{T}^\ell} \frac{1}{f_e(a_t) + \gamma} \leq 2KT + \frac{1}{\gamma} \log\left(\frac{1}{\delta}\right)$$

with probability at least $1 - \delta$.

Substituting this inequality back, we get that the following inequality holds with probability at least $\Pr(G) - \delta$:

$$\begin{aligned}
& \sum_c \Pr(c) \left(\frac{1}{\eta} \log K + \frac{\eta}{2} \sum_{t \in \mathcal{T}^\ell} \langle p_{t,c}, \hat{\ell}_{t,c}^2 \rangle \right) \\
&\leq \frac{1}{\eta} \log K + 4\eta \sum_{t \in \mathcal{T}^\ell} \frac{1}{f_e(a_t) + \gamma} \\
&\leq \frac{1}{\eta} \log K + 4\eta \left(2KT + \frac{1}{\gamma} \log\left(\frac{1}{\delta}\right) \right).
\end{aligned}$$

B.4.4 UPPER BOUNDING bias_3

Here we bound the bias_3 term

$$\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \ell_{t,c} - \tilde{\ell}_{t,c} \rangle.$$

We decompose it as

$$\begin{aligned}
& \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \ell_{t,c} - \tilde{\ell}_{t,c} \rangle \\
&= \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \ell_{t,c} - \mathbb{E}[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1}] \rangle \\
&\quad + \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \mathbb{E}[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1}] - \tilde{\ell}_{t,c} \rangle.
\end{aligned}$$

1296 The first component satisfies

$$\begin{aligned}
1297 & \\
1298 & \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left\langle p_{t,c}, \ell_{t,c} - \mathbb{E} \left[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] \right\rangle \\
1299 & \\
1300 & = \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \sum_a p_{t,c}(a) \ell_{t,c}(a) \frac{\gamma}{f_e(a) + \gamma} \\
1301 & \\
1302 & \leq \sum_{t \in \mathcal{T}^\ell} \sum_a p_t(a) \frac{\gamma}{f_e(a) + \gamma}.
\end{aligned}$$

1303 Assuming event G holds, we have

$$\begin{aligned}
1304 & \\
1305 & \sum_{t \in \mathcal{T}^\ell} \sum_a p_t(a) \frac{\gamma}{f_e(a) + \gamma} \\
1306 & \leq 2 \sum_{t \in \mathcal{T}^\ell} \sum_a \gamma \leq 2KT\gamma. \tag{Lemma 4}
\end{aligned}$$

1307 We then bound the second component $\sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left\langle p_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] - \tilde{\ell}_{t,c} \right\rangle$. For each

1308 time step $t \in \mathcal{T}^\ell$, we define an indicator function

$$1309 L_t \triangleq \mathbb{1}(\forall a, p_t(a) \leq 4f_e(a)).$$

1310 By Lemma 4, event G implies event L_t . Thus we have

$$\begin{aligned}
1311 & \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left\langle p_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] - \tilde{\ell}_{t,c} \right\rangle \\
1312 & = \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left\langle p_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] - \tilde{\ell}_{t,c} \right\rangle L_t.
\end{aligned}$$

1313 under event G . We then assume event G holds and focus on upper bounding

$$1314 \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \left\langle p_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] - \tilde{\ell}_{t,c} \right\rangle L_t.$$

1315 Since the probability vector $p_{t,c}$ is determined at time $t - 1$, the indicator function L_t is also determined at time $t - 1$. Thus the summand random variable

$$1316 \sum_c \Pr(c) \left\langle p_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] - \tilde{\ell}_{t,c} \right\rangle L_t$$

1317 satisfies

$$\begin{aligned}
1318 & \mathbb{E} \left[\sum_c \Pr(c) \left\langle p_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] - \tilde{\ell}_{t,c} \right\rangle L_t \mid \mathcal{H}_{t-1} \right] \\
1319 & = \mathbb{E} \left[\sum_c \Pr(c) \left\langle p_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] - \tilde{\ell}_{t,c} \right\rangle \mid \mathcal{H}_{t-1} \right] L_t \\
1320 & = 0.
\end{aligned}$$

1321 Thus the sequence of random variables

$$1322 \left\{ \sum_c \Pr(c) \left\langle p_{t,c}, \mathbb{E} \left[\tilde{\ell}_{t,c} \mid \mathcal{H}_{t-1} \right] - \tilde{\ell}_{t,c} \right\rangle L_t \right\}_{t \in \mathcal{T}^\ell}$$

1323 forms a martingale difference sequence with respect to the filtration $\{\mathcal{H}_t\}_{t \in \mathcal{T}^\ell}$.

1350 The summand random variable further satisfies

$$\begin{aligned}
1351 & \\
1352 & \sum_c \Pr(c) \langle p_{t,c}, \tilde{\ell}_{t,c} \rangle L_t \\
1353 & \\
1354 & = \sum_c \Pr(c) p_{t,c}(a_t) \frac{\ell_{t,c}(a_t)}{f_e(a_t) + \gamma} L_t \\
1355 & \\
1356 & \leq \sum_c \Pr(c) p_{t,c}(a_t) \frac{1}{f_e(a_t) + \gamma} L_t \\
1357 & \\
1358 & = \frac{p_t(a_t)}{f_e(a_t) + \gamma} L_t \\
1359 & \\
1360 & \leq 4. \\
1361 & \\
1362 &
\end{aligned}$$

1363 Then by Azuma-Hoeffding inequality, the following inequality holds with probability at least $1 - \delta$:

$$1364 \\
1365 \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \mathbb{E} [\tilde{\ell}_{t,c} | \mathcal{H}_{t-1}] - \tilde{\ell}_{t,c} \rangle L_t \leq 4\sqrt{T \log(\frac{1}{\delta})}. \\
1366 \\
1367$$

1368 Thus we have

$$1369 \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \mathbb{E} [\tilde{\ell}_{t,c} | \mathcal{H}_{t-1}] - \tilde{\ell}_{t,c} \rangle \leq 4\sqrt{T \log(\frac{1}{\delta})} \\
1370 \\
1371$$

1372 with probability at least $\Pr(G) - \delta$.

1373 Combining the first component and the second component, we get that

$$1374 \\
1375 \sum_{t \in \mathcal{T}^\ell} \sum_c \Pr(c) \langle p_{t,c}, \ell_{t,c} - \tilde{\ell}_{t,c} \rangle \leq 4\sqrt{T \log(\frac{1}{\delta})} + 2KT\gamma \\
1376 \\
1377$$

1378 with probability at least $\Pr(G) - \delta$.

1380 B.5 COMBINING THE PIECES

1381 Combining all previous arguments, we get that

$$\begin{aligned}
1382 & \text{Reg}(\pi) \\
1383 & \leq 2\sqrt{T \log(\frac{1}{\delta})} \\
1384 & + 4\sqrt{T \log(\frac{1}{\delta})} \\
1385 & + 2\frac{1}{\eta} \log K + 8\eta \left(2KT + \frac{1}{\gamma} \log(\frac{1}{\delta}) \right) \\
1386 & + 8\sqrt{T \log(\frac{1}{\delta})} + 4KT\gamma \\
1387 & + 8K\gamma T + 2\frac{K \exp(-\iota)T}{\gamma} + 2\frac{98KT\iota}{L} + 2\frac{\gamma^2 LKT}{\iota} \\
1388 & + 16 \left(\sqrt{2KT\iota \log(\frac{\delta}{2})} + 2\sqrt{2\frac{TK^2}{L}\iota \log(\frac{\delta}{2})} \right) \\
1389 & + 30\sqrt{T \log(\frac{1}{\delta})} \\
1390 & + 2\left(\frac{8 \cdot 36\iota}{\gamma} + 2\frac{8}{\gamma}\right) \log(1/\delta) + 2KT \exp(-\iota) \\
1391 & \\
1392 & \\
1393 & \\
1394 & \\
1395 & \\
1396 & \\
1397 & \\
1398 & \\
1399 & \\
1400 & \\
1401 & \\
1402 & \\
1403 &
\end{aligned}$$

with probability at least $\Pr(G) - 8\delta$.

1404 Taking $\iota = 2 \log(8KT\frac{1}{\delta})$, $L = \sqrt{\frac{\iota KT}{\log(K)}} = \tilde{\Theta}(\sqrt{KT \log \frac{1}{\delta}})$, $\gamma = \frac{16\iota}{L} = \tilde{\Theta}(\sqrt{\frac{\log(1/\delta)}{KT}})$, and
1405
1406 $\eta = \frac{\gamma}{2(2L\gamma + \iota)} = \tilde{\Theta}(1/\sqrt{KT \log(1/\delta)})$, it is easy to see that $\text{Reg}(\pi) = \tilde{O}(\sqrt{KT \log \frac{1}{\delta}})$ with
1407 probability at least $\Pr(G) - 8\delta \geq 1 - 9\delta$ for any policy π and any $\delta \in (0, 1)$.
1408

1409 The final step is rescaling the probability constant by a factor of $1/9$, which gives that $\text{Reg}(\pi) =$
1410 $\tilde{O}(\sqrt{KT \log \frac{1}{\delta}})$ with probability at least $1 - \delta$ and ends the proof of Theorem 1.
1411

1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457