
Instance-optimal PAC Algorithms for Contextual Bandits

Zhaoqi Li

Department of Statistics
University of Washington
zli9@uw.edu

Lillian Ratliff

Department of Electrical and Computer Engineering
University of Washington
ratliff@uw.edu

Houssam Nassif

Amazon
houssamn@amazon.com

Kevin Jamieson

Allen School of Computer Science & Engineering
University of Washington
jamieson@cs.washington.edu

Lalit Jain

Foster School of Business
University of Washington
lalitj@uw.edu

Abstract

In the stochastic contextual bandit setting, regret-minimizing algorithms have been extensively researched, but their instance-minimizing best-arm identification counterparts remain seldom studied. In this work, we focus on the stochastic bandit problem in the (ϵ, δ) -PAC setting: given a policy class Π the goal of the learner is to return a policy $\pi \in \Pi$ whose expected reward is within ϵ of the optimal policy with probability greater than $1 - \delta$. We characterize the first *instance-dependent* PAC sample complexity of contextual bandits through a quantity ρ_Π , and provide matching upper and lower bounds in terms of ρ_Π for the agnostic and linear contextual best-arm identification settings. We show that no algorithm can be simultaneously minimax-optimal for regret minimization and instance-dependent PAC for best-arm identification. Our main result is a new instance-optimal and computationally efficient algorithm that relies on a polynomial number of calls to an argmax oracle.

1 Introduction

We consider the stochastic contextual bandit problem in the PAC setting. Fix a distribution ν over a potentially countable¹ set of contexts \mathcal{C} . The action space is \mathcal{A} , and for computational tractability, we assume $|\mathcal{A}|$ is finite. We have a set of policies Π of interest where each policy $\pi \in \Pi$ is a map from contexts to an action space $\pi : \mathcal{C} \rightarrow \mathcal{A}$. The reward function is $r : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}$. At each time $t = 1, 2, \dots$ a context $c_t \sim \nu$ arrives, the learner chooses an action $a_t \in \mathcal{A}$, and receives reward $r_t := r(c_t, a_t) \in \mathbb{R}$ with $\mathbb{E}[r_t | c_t, a_t] = r(c_t, a_t) \in \mathbb{R}$. The value of a policy $V(\pi)$ is the expected reward from playing action $\pi(c)$ in context c : $V(\pi) = \mathbb{E}_{c \sim \nu}[r(c, \pi(c))]$. Given a collection of policies Π , the objective is to identify the optimal policy $\pi_* := \arg \max_{\pi \in \Pi} V(\pi)$, with high probability. Formally, for any $\epsilon > 0$ and $\delta \in (0, 1)$, we seek to characterize the sample complexity of

¹Assuming the set of contexts is countable versus uncountable is for presentation purposes only, since it allows us the notational convenience of letting ν_c denote the probability of context c arriving.

identifying a policy $\pi \in \Pi$ such that $V(\pi) \geq V(\pi_*) - \epsilon$, with probability at least $1 - \delta$. That is, we wish to minimize the total amount of interactions with the environment to learn an ϵ -optimal policy.

We study both the *agnostic* setting, where Π is an arbitrary set of policies with no assumed relationship with the reward function $r(c, a)$; and the *realizable* setting, where the policy class and the reward function follow a linear structure, known as the linear contextual bandit problem. In both cases, we are interested in *instance-dependent* sample complexity bounds. That is, the upper and lower bounds we seek do not simply depend on coarse quantities like $|\Pi|$, $|\mathcal{A}|$, and $1/\epsilon^2$, but more fine-grained relationships between the context distribution ν , geometry of policies Π , and the reward function $r : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}$. Our motivation is that instance-dependent bounds describe the difficulty of a particular problem instance, allowing optimal algorithms to adapt to the true difficulty of the problem, whether easy or hard. We seek algorithms that take advantage of “easy” instances instead of optimizing for the worst-case [23].

1.1 Related work

Minimax regret bounds for general policy classes The vast majority of research in contextual bandits focuses on regret minimization. That is, for a time horizon T , the goal of the player is to minimize $\mathbb{E} \left[\sum_{t=1}^T r(c_t, \pi_*(c_t)) - r(c_t, a_t) \right]$. The landmark algorithm EXP4 for non-stochastic multi-armed bandits [5] achieves a regret bound of $\sqrt{|\mathcal{A}|T \log(|\Pi|)}$. Unfortunately, the running time of EXP4 is linear in $|\Pi|$ which is prohibitive for many problems of interest. The algorithms proposed in [3, 11] achieve the same regret bound with a computational complexity that is only polynomial in T and $\log(|\Pi|)$. Both approaches can be used to obtain an ϵ -optimal policy with probability at least $1 - \delta$ using a sample complexity no more than $\frac{|\mathcal{A}| \log(|\Pi|/\delta)}{\epsilon^2}$. None of these works made any assumption on the connection between the reward function r and the policy class Π (i.e. the agnostic setting).

Instance-dependent regret bounds for general policy classes The epoch-greedy algorithm of [26] achieved the first instance-dependent bounds on regret with a coarse guarantee depending only on the minimum policy gap $\Delta_{\text{pol}} := V(\pi_*) - \max_{\pi \neq \pi_*} V(\pi)$. In the pursuit of more fine-grained regret bounds achievable by computationally efficient algorithms, many authors resort to the *realizability* assumption [14–16, 34]. The learner knows a hypothesis class \mathcal{H} where each $f \in \mathcal{H}$ is a map $f : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}$, and there exists an $f^* \in \mathcal{H}$ such that $r(c, a) = f^*(c, a)$ for all $(c, a) \in \mathcal{C} \times \mathcal{A}$. Under this assumption, [16] proves lower and upper bounds on the instance-dependent regret. Their bounds are in term of the *uniform gap* $\Delta_{\text{uniform}} := \min_{c \in \mathcal{C}} \min_{a \in \mathcal{A}} r(c, \pi_*(c)) - r(c, a)$. In general, for any policy class, they establish matching minimax lower and upper regret bounds of the form $\min \left\{ \sqrt{|\mathcal{A}|T \log(|\mathcal{H}|)}, \frac{|\mathcal{A}| \log(|\mathcal{H}|)}{\Delta_{\text{uniform}}} \mathfrak{C}_{\mathcal{H}}^{\text{pol}} \right\}$, where $\mathfrak{C}_{\mathcal{H}}^{\text{pol}}$ is the *policy disagreement coefficient*, a parameter depending on the geometry of \mathcal{H} and the context distribution ν . That is, these bounds hold with respect to a worst-case family of instances parameterized by Δ_{uniform} and $\mathfrak{C}_{\mathcal{H}}^{\text{pol}}$. Using the standard online-to-batch conversion, this translates to a sample complexity (i.e. the time required to find an ϵ -good policy with constant probability) of roughly $\frac{|\mathcal{A}| \log(|\mathcal{H}|)}{\epsilon \Delta_{\text{uniform}}} \mathfrak{C}_{\mathcal{H}}^{\text{pol}}$. We show in Corollary 2.16 that this sample complexity is at least as large as our bounds. Further, unlike our bounds below, this sample complexity is unbounded as ϵ goes to 0. Recent work refines these kinds of regret bounds further, and provides minimax regret bounds in terms of the *decision-estimation coefficient* [17].

Regret bounds for linear contextual bandits A special case of the realizable case assumes a linear structure for \mathcal{H} . Assume there exists a known feature map $\phi : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}^d$ and an unknown $\theta_* \in \mathbb{R}^d$ such that the true reward function is given as $r(c, a) = \langle \phi(c, a), \theta_* \rangle$. For this setting, popular optimism-based algorithms like LinUCB [27] and Thompson sampling [31, 33] achieve a regret bound of $\min \left\{ d\sqrt{T}, \frac{d^2}{\Delta_{\text{uniform}}} \right\}$ [1]. Appealing to the online-to-batch conversion, this translates to a PAC guarantee of $\frac{d^2}{\epsilon \Delta_{\text{uniform}}}$. More precise instance-dependent upper bounds on regret match instance-dependent lower bounds asymptotically as $T \rightarrow \infty$ [19, 36]. These works are most similar to our setting and have qualitatively similar style algorithms. However, both approaches rely on asymptotics with large problem-dependent terms that may dominate the bounds in finite time. Our work is focused on upper bounds that nearly match lower bounds for all finite times.

Recently, instance-dependent sample complexity results for reinforcement learning in the tabular and linear function approximation settings have appeared [4, 39, 40]. As contextual bandits is a special case of finite-horizon reinforcement learning with a horizon length of 1, their results immediately can

be applied here. However, the cost of this generality is that these algorithms have very large lower order terms (i.e., problem-dependent factors that multiply a $1/\epsilon$ term) making them far from optimal in our setting. Moreover, the leading order term of [39] cannot be related to our lower bounds.

PAC sample complexity for contextual bandits As we will describe, all contextual bandits with an arbitrary policy class can be reduced to PAC learning for linear bandits. Once we made this reduction, our sample complexity analysis draws inspiration from the nearly instance-optimal algorithm for linear best-arm identification [13]. The work in [10] provides a simple regret bound assuming a kernel structure on the reward function, while their bound is minimax and they assume a lower bound on eigenvalues of the covariance matrix of the context distribution. PAC sample complexity of linear contextual bandits was also studied in [41], who shows a minimax guarantee sample complexity that scales with $\frac{d^2}{\epsilon^2} \log(1/\delta)$. Similar to our work, [3] define their action sampling distribution as a convex combination over policies. Our sampling distribution, as well as the optimal sampling distribution, cannot be represented this way and is actually derived from the dual of the optimal experimental design objective.

Contributions. In this work, our contributions include:

1. In the agnostic setting, we introduce a quantity ρ_Π that characterizes the instance-dependent sample complexity of PAC learning for contextual bandits (see Equation 1). We show that ρ_Π appears in information theoretic lower bound on the sample complexity of any PAC algorithm as $\epsilon \rightarrow 0$ in Theorem 2.2. To ground this, we describe it carefully in the setting of the trivial policy class (Section 2.2) and linear policy classes (Section 2.3).
2. We construct an instance on which any regret minimax-optimal algorithm necessarily has a sample complexity that scales quadratically with the optimal sample complexity (Theorem 2.6). This shows that no algorithm can be both regret minimax-optimal and instance-optimal PAC.
3. Finally, we propose Algorithm 3 whose sample complexity nearly matches the lower bound based on ρ_Π . By appealing to an argmax oracle, this algorithm has a runtime polynomial in ρ_Π , $1/\epsilon$, $\log(1/\delta)$, $|\mathcal{A}|$, and $\log(|\Pi|)$, assuming a unit cost of invoking the oracle.

2 Problem statement and main results

More formally, define $\mathcal{F}_t = \sigma(c_1, a_1, r_1, \dots, c_t, a_t, r_t)$ as the natural σ -algebra filtration capturing all observed random variables up to time t . For simplicity, we assume Gaussian noise in some of our analysis. At each time t an *algorithm* defines a *sampling rule* $\mathcal{F}_t \mapsto \mathcal{A}$ which defines a_{t+1} , an $\{\mathcal{F}_t\}_{t \geq 1}$ -adapted stopping time $\tau \in \mathbb{N}$, and a *selection rule* $\mathcal{F}_t \mapsto \Pi$ that is only called once at the stopping time $t = \tau$.

Definition 2.1. Fix $\epsilon \geq 0$ and $\delta \in (0, 1)$. We say an algorithm is (ϵ, δ) -PAC for contextual bandits with policy class Π , if for every instance, at the stopping time $\tau \in \mathbb{N}$ with $\tau < \infty$ almost surely, the algorithm outputs $\hat{\pi} \in \Pi$ satisfying $\mathbb{P}(V(\hat{\pi}) \geq \max_{\pi \in \Pi} V(\pi) - \epsilon) \geq 1 - \delta$.

The *sample complexity* of an (ϵ, δ) -PAC algorithm for contextual bandits is the time at which the algorithm stops and outputs $\hat{\pi}$. As we will discuss, the following quantity governs the sample complexity :

$$\rho_{\Pi, \epsilon} := \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\pi \in \Pi \setminus \pi_*} \frac{\mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{c, \pi(c)}} + \frac{1}{p_{c, \pi_*(c)}} \right) \mathbf{1}\{\pi_*(c) \neq \pi(c)\} \right]}{(\mathbb{E}_{c \sim \nu} [r(c, \pi_*(c)) - r(c, \pi(c))] \vee \epsilon)^2}. \quad (1)$$

Here, for any countable set \mathcal{X} we have that $\Delta_{\mathcal{X}} = \{p \in \mathbb{R}^{|\mathcal{X}|} : \sum_{x \in \mathcal{X}} p_x = 1, p_x \geq 0 \forall x \in \mathcal{X}\}$ so that p_c for every $c \in \mathcal{C}$ defines a probability distribution over actions \mathcal{A} . In addition we use the notation $a \vee b := \max\{a, b\}$. We begin with a necessary condition on the sample complexity for the particular case of exact policy identification ($\epsilon = 0$).

Theorem 2.2 (Lower bound). *Fix $\epsilon = 0$ and $\delta \in (0, 1)$. Moreover, fix a contextual bandit instance $\mu = (\nu, r)$ and a collection of policies Π . Then any $(0, \delta)$ -PAC algorithm for contextual bandits satisfies $\mathbb{E}_\mu[\tau] \geq \rho_{\Pi, 0} \log(1/2.4\delta)$.*

The proof of the lower bound follows from standard information theoretic arguments [24]. The lower bound implicitly applies to learners that know the distribution ν precisely. In practice, such knowledge would never be available however the learner may have a large dataset of offline data.

Assumption 1. Prior to starting the game, the learning algorithm is given a large dataset of contexts $\mathcal{D} = \{c_t\}_{t=1}^T$, where each c_t is drawn IID from ν for all $t \in [T]$, and $T = O(\text{poly}(1/\epsilon, |\mathcal{A}|, \log(1/\delta), \log(|\Pi|)))$.

The above only assumes access to samples from the context distribution, not rewards or the value function. Importantly, since \mathcal{C} could be uncountable, we do not assume \mathcal{D} covers the support of ν . Assumption 1 is satisfied, for example, in an e-commerce setting where the context is the demographic information about visitors to the site for which massive troves of historical data may be available. Other works in PAC learning have made similar assumptions [20]. We would like our algorithm to be computationally efficient in the sense that it makes a polynomial number of calls to what we refer to as argmax oracle. Such an assumption is common in the contextual bandits literature [3, 11, 25].

Definition 2.3 (Argmax oracle (AMO)). The oracle $\text{AMO}(\Pi, \{(c_t, s_t)\}_{t=1}^n)$ is an algorithm that given contexts and cost vectors $(c_1, s_1), \dots, (c_n, s_n) \in \mathcal{C} \times \mathbb{R}^{|\mathcal{A}|}$, returns $\arg \max_{\pi \in \Pi} \sum_{t=1}^n s_t(\pi(c_t))$.

The constrained argmax oracle C-AMO , given an upper bound l on the loss, returns $\arg \max_{\pi \in \Pi} \sum_{t=1}^n s_t(\pi(c_t))$ subject to $\sum_{t=1}^n s_t(\pi(c_t)) \leq l$.

In general we can implement AMO by calling to cost-sensitive classification [6, 11] and C-AMO through a Lagrangian relaxation and a cost-sensitive classification oracle [2, 8]. Our algorithm uses an argmax oracle as a subroutine at most polynomially in ϵ^{-1} , $\log(1/\delta)$, $|\mathcal{A}|$ and $\log(|\Pi|)$. In this sense, it is computationally efficient. The following sufficiency result holds for general $\epsilon \geq 0$.

Theorem 2.4 (Upper bound). Fix $\epsilon \geq 0$ and $\delta \in (0, 1)$. Under Assumption 1, there exists a computationally efficient (ϵ, δ) -PAC algorithm for contextual bandits that satisfies $\tau \leq \rho_{\Pi, \epsilon} \log(|\Pi| \log_2(1/\epsilon)/\delta) \log(1/\Delta_\epsilon)$, where $\Delta_\epsilon = \max\{\epsilon, \min_{\pi \in \Pi \setminus \pi_*} V(\pi_*) - V(\pi)\}$. Furthermore, this sample complexity never exceeds $\frac{|\mathcal{A}|(\log(|\Pi|) + \log(1/\delta)) \log(1/\epsilon)}{\epsilon^2}$.

The second part of the theorem follows from the first, since $\rho_{\Pi, \epsilon} \leq 2|\mathcal{A}|/\epsilon^2$ by taking $p_{c,a} = 1/|\mathcal{A}|$ for all $(c, a) \in \mathcal{C} \times \mathcal{A}$.

2.1 Inefficiency of low-regret algorithms

Computationally efficient algorithms are known to exist, such as ILOVETOCONBANDITS [3], which achieve a minimax-optimal cumulative regret of $\sqrt{T|\mathcal{A}| \log(|\Pi|/\delta)}$. Inspecting the proof in [3], one can extract a sample complexity of $\epsilon^{-2}|\mathcal{A}| \log(|\Pi|/\delta)$ from such results (which is also minimax optimal for PAC). The previous section showed that the sample complexity of our algorithm, Theorem 2.4, nearly matches the instance-dependent lower bound of Theorem 2.2. In other words, our algorithm achieves a nearly optimal instance-dependent PAC sample complexity. However, it is natural to wonder if perhaps with a tighter analysis, the minimax regret optimal algorithm in [3] also obtains the instance-optimal PAC sample complexity. In this section, we show that this is not the case. Indeed, we show that *any* algorithm that is minimax regret optimal must have a sample complexity that is at least quadratic in the optimal PAC sample complexity of some instance.

Definition 2.5 (Hard instance). Fix $m \in \mathbb{N}$, $\Delta \in (0, 1]$ and let $\mathcal{C} = [m]$ with uniform distribution, $\mathcal{A} = \{0, 1\}$. For $i = 1, \dots, m$, let $\pi_i(j) = \mathbf{1}\{i = j\}$ and define $r(i, j) = \Delta \mathbf{1}\{j = \pi_1(i)\}$. Then $V(\pi_1) = \Delta$ and $V(\pi_i) = \Delta(1 - 2/m)$ for all $i \in \mathcal{C} \setminus \{1\}$.

Note that for the hard instance, $m = |\Pi|$. If observations are corrupted by $\mathcal{N}(0, 1)$ additive noise, then a straightforward calculation shows that $\rho_{\Pi, 0} = \frac{4/m}{(2\Delta/m)^2} = m\Delta^{-2}$ for the hard instance.

Theorem 2.6. Fix $\delta \in (0, 1)$ and $\Delta \in (0, 1]$. We say an algorithm is an α -minimax regret algorithm if for some $\alpha > 0$ and all $T \in \mathbb{N}$:

$$\max_{\mu'} \mathbb{E}_{\mu'} \left[\sum_{t=1}^T (r_t(c_t, \pi_*(c_t)) - r_t(c_t, a_t)) \right] = \max_{\mu'} \sum_{c,a} \mathbb{E}_{\mu'} [T_{c,a}(T)] (r(c, \pi_*(c)) - r(c, a)) \leq \sqrt{\alpha |\mathcal{A}| T}$$

where the maximum is taken over all contextual bandit instances $\mu' = (\nu', r')$ and $T_{c,a}(T) = \sum_{t=1}^T \mathbf{1}\{c_t = c, a_t = a\}$. For any α -minimax regret algorithm, it is $(0, \delta)$ -PAC if at a stopping time τ it outputs the optimal policy π_* w. p. at least $1 - \delta$. Any α -minimax regret algorithm that is $(0, \delta)$ -PAC satisfies $\mathbb{E}_{\mu}[\tau] \geq m^2 \Delta^{-2} \log^2(1/2.4\delta)/4\alpha$ for the instance $\mu = (\nu, r)$ defined in 2.5.

We point out that the minimax regret optimal rate takes $\alpha = \log(m) = \log(|\Pi|)$. Thus, taking $\Delta = 1$ and $\delta = 0.1$, the minimax regret optimal algorithm has a PAC sample complexity of $m^2/\log(m)$; whereas the PAC sample complexity of our algorithm, Theorem 2.4, is just $m \log(m)$. That is, algorithms with optimal minimax regret have a sample complexity that is at least nearly the optimal PAC sample complexity *squared*. This demonstrates that no algorithm can simultaneously be minimax regret optimal and obtain the optimal PAC sample complexity.

2.2 Trivial policy class

As a warm-up to discussing linear policy classes, let us consider the simplest policy class.

Definition 2.7 (Trivial policy class). Assume $|\mathcal{C}| < \infty$ and let $\Pi = \{\pi(c) = a : (c, a) \in \mathcal{C} \times \mathcal{A}\}$ so that $|\Pi| = |\mathcal{A}|^{|\mathcal{C}|}$.

The trivial policy class has the flexibility to predict any action $a \in \mathcal{A}$ individually for each $c \in \mathcal{C}$. This allows us to show that $\rho_{\Pi,0} \leq \max_c \frac{2}{\nu_c} \sum_{a'} \Delta_{c,a'}^{-2}$ (see Appendix A.3). An immediate corollary of Theorem 2.4 is obtained by simply noting that $|\Pi| = |\mathcal{A}|^{|\mathcal{C}|}$.

Corollary 2.8 (Trivial class, upper). Fix $\epsilon > 0$ and $\delta \in (0, 1)$. Let Π be the trivial policy class applied to some fixed \mathcal{C}, \mathcal{A} spaces. Then under Assumption 1 there exists a computationally efficient (ϵ, δ) -PAC algorithm for contextual bandits that satisfies $\tau \leq \min\{A\epsilon^{-2}, \max_c \frac{1}{\nu_c} \sum_{a'} \Delta_{c,a'}^{-2}\} (|\mathcal{C}| \log(|\mathcal{A}|) + \log(1/\delta)) \log(1/\Delta_\epsilon)$, where $\Delta_\epsilon = \max\{\epsilon, \min_{\pi \in \Pi \setminus \pi_*} V(\pi_*) - V(\pi)\}$. Furthermore, this sample complexity never exceeds $\frac{|\mathcal{A}|(|\mathcal{C}| \log(|\mathcal{A}|) + \log(1/\delta))}{\epsilon^2} \log(1/\epsilon)$.

Ignoring log factors, the minimax sample complexity of the trivial class is just $\epsilon^{-2} |\mathcal{A}| (|\mathcal{C}| + \log(1/\delta))$. This is actually a somewhat surprising result, because it says $\lim_{\delta \rightarrow 0} \frac{\mathbb{E}[\tau]}{\log(1/\delta)} \rightarrow \epsilon^{-2} |\mathcal{A}|$ which is *independent* of $|\mathcal{C}|$. To see why this result is somewhat remarkable, if we played a best-arm identification algorithm for each of the $|\mathcal{C}|$ contexts, then this would lead to a sample complexity of $\epsilon^{-2} |\mathcal{C}| \cdot |\mathcal{A}| \log(1/\delta)$. It is somewhat of a surprise that such a natural strategy is not optimal. For intuition for why we can avoid the multiplicative $|\mathcal{C}|$, note that to identify an ϵ -good policy among just two policies (π, π_*) using uniform exploration requires just $\epsilon^{-2} |\mathcal{A}| \log(1/\delta)$ samples. When we have more than two policies, a union bound achieves the claimed result.

The minimax sample complexity of Corollary 2.8 (i.e., the second statement) is nearly tight:

Theorem 2.9 (Trivial class, lower). Fix $\epsilon > 0$ and $\delta \in (0, 1/6)$. Let Π be the trivial policy class applied to some fixed \mathcal{C}, \mathcal{A} spaces. Moreover, fix a contextual bandit instance $\mu = (\nu, \tau)$ and a collection of policies Π . Then any $(0, \delta)$ -PAC algorithm for contextual bandits satisfies $\mathbb{E}_\mu[\tau] \geq \max_c \frac{1}{\nu_c} \sum_a \Delta_{c,a}^{-2} \log(1/2.4\delta)$. Furthermore, $\sup_\mu \mathbb{E}_\mu[\tau] \geq \epsilon^{-2} |\mathcal{A}| (|\mathcal{C}| + \log(1/\delta))$.

2.3 Linear policy class

A particularly compelling model-class of policies is the set of linear policies.

Definition 2.10 (Linear policy class). Fix a feature map $\phi : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}^d$ and assume it is known to the learner. Let $\Pi = \{\pi(c) = \arg \max_{a \in \mathcal{A}} \langle \phi(c, a), \theta \rangle, \forall \theta \in \mathbb{R}^d\}$.

We can consider two settings: the agnostic setting and the realizable setting. In the agnostic setting, there is no assumed relationship between the true reward function $r(c, a)$ and $\phi : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}^d$. In this case, Theorem 2.4 applies directly by taking a cover of Π .

Corollary 2.11 (Agnostic, upper bound). Fix $\epsilon \geq 0$ and $\delta \in (0, 1)$. Let Π be the linear policy class in \mathbb{R}^d . Under Assumption 1 there exists a computationally efficient (ϵ, δ) -PAC algorithm for contextual bandits that satisfies $\tau \leq \rho_{\Pi, \epsilon} \cdot (d \log(1/\epsilon) + \log(1/\delta)) \log(1/\Delta_\epsilon)$ where $\Delta_\epsilon = \max\{\epsilon, \min_{\pi \in \Pi \setminus \pi_*} V(\pi_*) - V(\pi)\}$. Furthermore, this sample complexity never exceeds $\frac{|\mathcal{A}|(d \log(1/\epsilon) + \log(1/\delta))}{\epsilon^2} \log(1/\epsilon)$.

Comparing to the lower bound of Theorem 2.2, the instance dependent upper bound of Corollary 2.11 matches up to a factor of the dimension and negligible log factors. In contrast to the “model-free” feel of the agnostic case, we can also consider a “model-based” type setting, i.e. the realizable setting.

Definition 2.12 (Realizable). We say the linear policy class is *realizable* if there exists a $\theta_* \in \mathbb{R}^d$ such that $r(c, a) = \langle \phi(c, a), \theta_* \rangle$ for all $c \in \mathcal{C}$ and $a \in \mathcal{A}$. Thus, for any $\pi \in \Pi$ we have $V(\pi) = \mathbb{E}_{c \sim \nu}[r(c, \pi(c))] = \mathbb{E}_{c \sim \nu}[\langle \phi(c, \pi(c)), \theta_* \rangle] = \langle \phi_\pi, \theta_* \rangle$ with $\phi_\pi := \mathbb{E}_{c \sim \nu}[\phi(c, \pi(c))]$. Finally, at the start of the game the learner knows this model.

The setting in Definition 2.12 is commonly referred to as the linear contextual bandit problem [1]. Clearly, we have that $\pi_*(c) = \arg \max_{a \in \mathcal{A}} \langle \phi(c, a), \theta_* \rangle$. We begin by defining a quantity fundamental to our sample complexity results:

$$\rho_{\text{lin}, \epsilon} := \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\pi \in \Pi \setminus \pi_*} \frac{\|\phi_\pi - \phi_{\pi_*}\|_{\mathbb{E}_{c \sim \nu}[\sum_{a \in \mathcal{A}} p_{c,a} \phi(c, a) \phi(c, a)^\top]^{-1}}^2}{\langle \phi_{\pi_*} - \phi_\pi, \theta_* \rangle^2 \vee \epsilon^2}.$$

Theorem 2.13 (Realizable, lower bound). Fix $\epsilon = 0$ and $\delta \in (0, 1)$. Let Π be the linear policy class in \mathbb{R}^d and assume it is realizable (see Definitions 2.10 and 2.12). Any $(0, \delta)$ -PAC algorithm in this setting satisfies $\mathbb{E}[\tau] \geq \rho_{\text{lin}, 0} \cdot \log(1/2.4\delta)$.

We now state our nearly matching upper bound. However, in this case we note that the algorithm is not computationally efficient.

Theorem 2.14 (Realizable, upper bound). Fix $\epsilon \geq 0$ and $\delta \in (0, 1)$. Let Π be the linear policy class in \mathbb{R}^d and assume it is realizable (see Definitions 2.10 and 2.12). Under Assumption 1 there exists an (ϵ, δ) -PAC algorithm (see Algorithm 1) for this setting satisfying

$$\tau \leq \rho_{\text{lin}, \epsilon} \cdot (\min\{d \log(1/\epsilon), \log(|\Pi|)\} + \log(1/\delta)) \log(1/\Delta_\epsilon)$$

where $\Delta_\epsilon = \max\{\epsilon, \min_{\pi \in \Pi \setminus \pi_*} \langle \phi_{\pi_*} - \phi_\pi, \theta_* \rangle\} = \max\{\epsilon, \min_{(c,a) \in \mathcal{C} \times \mathcal{A}: \pi_*(c) \neq a} \langle \phi(c, \pi_*(c)) - \phi(c, a), \theta_* \rangle\}$.

Furthermore, this sample complexity never exceeds $\frac{d(d \log(1/\epsilon) + \log(1/\delta)) \log(1/\epsilon)}{\epsilon^2}$.

We remark that the algorithm that achieves this upper bound is very different than popular optimism-based algorithms for linear contextual bandits e.g., UCB or Thompson sampling [1]. Indeed, our algorithm computes an experimental design and is related to instance-dependent linear bandit algorithms developed for best-arm identification [9, 12, 13, 35] and regret minimization [19, 36]. To our knowledge, Theorem 2.14 provides the first instance-dependent sample complexity for the PAC setting of linear contextual bandits. The most relevant work to Theorem 2.14 is the work of [41] which demonstrated a minimax sample complexity of $d^2/\epsilon^2 \log(1/\delta)$. Also, we remark that the lower and upper bounds in this section require an additive Gaussian noise.

Remark 2.15 (Agnostic vs. Realizable). Contrasting the above results, we note that the sample complexity of the agnostic case is always bounded by $|\mathcal{A}|d/\epsilon^2$. whereas it never exceeds d^2/ϵ^2 for the realizable case. This matches the intuition that when the number of actions is much larger than the dimension, assuming realizability can significantly reduce the sample complexity.

2.4 Comparison to the Disagreement Coefficient

The work of [16] provides regret bounds in terms of instance-dependent quantities inspired by the *disagreement coefficient*, a notion of complexity common in the active learning literature [18]. The following corollary relates our sample complexity to these notions of disagreement coefficients.

Define the *policy disagreement coefficient* as

$$\mathfrak{C}_{\Pi}^{\text{pol}}(\epsilon_0) = \sup_{\epsilon \geq \epsilon_0} \frac{\mathbb{E}_{c \sim \nu}[\mathbf{1}\{\exists \pi \in \Pi_\epsilon : \pi(c) \neq \pi_*(c)\}]}{\epsilon}$$

where $\Pi_\epsilon := \{\pi \in \Pi : \mathbb{P}_\nu(\pi(c) \neq \pi_*(c)) \leq \epsilon\}$ and the *cost-sensitive disagreement coefficient* as

$$\mathfrak{C}_{\Pi}^{\text{csc}}(\epsilon_0) = \sup_{\epsilon \geq \epsilon_0} \frac{\mathbb{E}_{c \sim \nu}[\mathbf{1}\{\exists \pi \in \Pi : \pi(c) \neq \pi_*(c), \mathbb{E}_{c \sim \nu}[r(c, \pi_*(c)) - r(c, \pi(c))] \leq \epsilon\}]}{\epsilon}.$$

The AdaCB algorithm of [16] achieves a regret of roughly $R_T = O\left(\min_\delta \left\{ \delta \Delta_{\text{uniform}} T, \frac{|\mathcal{A}| \log(|\Pi|) \mathfrak{C}_{\Pi}^{\text{pol}}(\delta)}{\Delta_{\text{uniform}}} \right\}\right)$ or $R_T = O(\min_\delta \{\delta T, |\mathcal{A}| \log(|\Pi|) \mathfrak{C}_{\Pi}^{\text{csc}}(\delta)\})$. Observe that at time T , given the outputs $\pi_1, \pi_2, \dots, \pi_T$ from AdaCB algorithm, one could return a (randomized) policy $\tilde{\pi}$ which on observing a context, samples from the empirical distribution over

the outputs. By Markov's inequality we have $\tilde{\pi}$, $V(\pi_*) - V(\tilde{\pi}) \leq O(\epsilon)$ with constant probability for $\epsilon = \frac{R_T}{T}$. Therefore, an upper bound on the regret translates to a PAC sample complexity of $\frac{|\mathcal{A}| \log(|\Pi|)}{\epsilon \Delta_{\text{uniform}}} \mathfrak{C}_{\Pi}^{\text{pol}}(\epsilon/\Delta_{\text{uniform}})$ or $\frac{|\mathcal{A}| \log(|\Pi|)}{\epsilon} \mathfrak{C}_{\Pi}^{\text{csc}}(\epsilon)$.

Finally, Corollary 2.16 shows that this sample complexity bound is at least as large as our upper bound, see Appendix A.5 for the proof.

Corollary 2.16. *Recall that $\Delta_{\text{uniform}} := \min_{c \in \mathcal{C}} \min_{a \in \mathcal{A}} r(c, \pi_*(c)) - r(c, a)$. For any $\epsilon_0 > 0$ we have that*

1. $\rho_{\Pi, \epsilon_0} \leq \frac{2|\mathcal{A}|}{\epsilon_0 \Delta_{\text{uniform}}} \mathfrak{C}_{\Pi}^{\text{pol}}(\epsilon_0/\Delta_{\text{uniform}})$;
2. $\rho_{\Pi, \epsilon_0} \leq \frac{2|\mathcal{A}|}{\epsilon_0} \mathfrak{C}_{\Pi}^{\text{csc}}(\epsilon_0)$.

Moreover, for all $\epsilon_0 \geq 0$ we have that $\rho_{\Pi, \epsilon_0} < \infty$ whenever $\Delta_{\text{pol}} := V(\pi_*) - \max_{\pi \neq \pi_*} V(\pi) > 0$.

3 Optimal Algorithms for Contextual Bandits

3.1 Reduction to linear realizability and a simple elimination scheme

The astute reader may have noticed that if we ignore computation, Theorem 2.4 is actually an immediate corollary of Theorem 2.14 by taking $\phi(c, a) = \text{vec}(\mathbf{e}_c \mathbf{e}_a^\top) \in \mathbb{R}^{|\mathcal{C}| \cdot |\mathcal{A}|}$ where \mathbf{e}_i is a one-hot encoded vector so that $r(c, a) = \langle \phi(c, a), \theta_* \rangle$ with $\theta_* \in \mathbb{R}^{|\mathcal{C}| \cdot |\mathcal{A}|}$. This observation is key to our sample complexity results. Recalling $\phi_\pi := \mathbb{E}_{c \sim \nu}[\phi(c, \pi(c))]$ (from Definition 2.12), we have that $V(\pi) = \mathbb{E}_{c \sim \nu}[r(c, \pi(c))] = \mathbb{E}_{c \sim \nu}[\langle \phi(c, \pi(c)), \theta_* \rangle] = \langle \phi_\pi, \theta_* \rangle$. We stress that \mathcal{C} can be uncountable, and thus we would never actually instantiate any of the vectors $\phi(c, a)$.

For notational convenience, define the feasible set of (context, action) probability distributions as $\Omega = \left\{ w \in \Delta_{\mathcal{C} \times \mathcal{A}} : \nu_c = \sum_{a \in \mathcal{A}} w_{a,c} \right\}$. Note that for each context, $p_c := \{w_{c,a}/\nu_c\}_{a \in \mathcal{A}} \in \Delta_{\mathcal{A}}$ defines a probability distribution over actions. Also define $A(w) := \sum_{c,a} w_{c,a} \phi(c, a) \phi(c, a)^\top$ for any $w \in \Omega$. Under this notation, recalling the right hand side from Theorems 2.13 and 2.14 we have

$$\min_{w \in \Omega} \max_{\pi \in \Pi \setminus \pi_*} \frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{\langle \phi_{\pi_*} - \phi_\pi, \theta_* \rangle^2 \vee \epsilon^2} = \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\pi \in \Pi \setminus \pi_*} \frac{\|\phi_\pi - \phi_{\pi_*}\|_{\mathbb{E}_{c \sim \nu}[\sum_{a \in \mathcal{A}} p_{c,a} \phi(c, a) \phi(c, a)^\top]^{-1}}^2}{\langle \phi_{\pi_*} - \phi_\pi, \theta_* \rangle^2 \vee \epsilon^2}$$

To show that the sample complexity of Theorem 2.4 is a corollary of Theorem 2.14, it suffices to show that equation (1) and the above display are equal. To see this, observe

$$\begin{aligned} \|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2 &= \|\mathbb{E}_{c \sim \nu}[\text{vec}(\mathbf{e}_c \mathbf{e}_{\pi(c)}^\top) - \text{vec}(\mathbf{e}_c \mathbf{e}_{\pi_*(c)}^\top)]\|_{A(w)^{-1}}^2 \\ &= \sum_{c,a} \frac{\nu_c^2}{w_{c,a}} (\mathbf{1}\{\pi(c) = a\} + \mathbf{1}\{\pi_*(c) = a\} - 2\mathbf{1}\{\pi(c) = \pi_*(c)\}) \\ &= \mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{c, \pi(c)}} + \frac{1}{p_{c, \pi_*(c)}} \right) \mathbf{1}\{\pi_*(c) \neq \pi(c)\} \right]. \end{aligned}$$

Due to this equivalence, the lower bound of Theorem 2.2 is also a corollary of Theorem 2.13. The lower bound of Theorem 2.13 follows almost immediately from the lower bound argument in [13].

The conclusion of this section is that from a sample complexity analysis alone, all that is left is to prove Theorem 2.14. In the next section we propose an algorithm that achieves this sample complexity but assumes precise knowledge of the context distribution ν (this is relaxed in following sections). While the algorithm is highly impractical for a number of reasons, its analysis provides a great deal of intuition and motivation for our final algorithm.

3.2 A simple, impractical, elimination-style algorithm

Algorithm 1 provides an initial elimination based method for the PAC-contextual bandit problem. The algorithm runs in stages. Before the start of each stage $\ell \in \mathbb{N}$, the algorithm defines a distribution $p_c^{(\ell)} \in \Delta_{\mathcal{A}}$ for each $c \in \mathcal{C}$. At each successive time $t \in [n_\ell]$, it plays random action $a_t \sim p_{c_t}^{(\ell)}$ in response to context $c_t \sim \nu$, and receives random reward r_t with $\mathbb{E}[r_t | c_t, a_t] = \langle \phi(c_t, a_t), \theta_* \rangle$. Observe that

$$\mathbb{E}[\phi(c_t, a_t) r_t] = \mathbb{E}[\phi(c_t, a_t) \phi(c_t, a_t)^\top \theta_*] = \sum_{c \in \mathcal{C}, a \in \mathcal{A}} w_{c,a}^{(\ell)} \phi(c, a) \phi(c, a)^\top \theta_* = A(w^{(\ell)}) \theta_*$$

using the identity $w_{c,a}^{(\ell)} := \nu_c p_{c,a}^{(\ell)}$. Thus, if we set $O_t = A(w^{(\ell)})^{-1} \phi(c_t, a_t) r_t$ then $\mathbb{E}[O_t] = \theta_*$. A straightforward calculation also shows that $\text{Cov}(O_t) = A(w^{(\ell)})^{-1}$ if r_t is perturbed with additive unit variance noise. Thus, an unbiased estimator of $\Delta(\pi, \pi_*) := V(\pi_*) - V(\pi) = \langle \phi_{\pi_*} - \phi_\pi, \theta_* \rangle$ is simply $\langle \phi_{\pi_*} - \phi_\pi, \frac{1}{n_\ell} \sum_t O_t \rangle$ which has variance $\frac{1}{n_\ell} \|\phi_{\pi_*} - \phi_\pi\|_{A(w^{(\ell)})^{-1}}^2$. Intuitively, $\langle \phi_{\pi_*} - \phi_\pi, \frac{1}{n_\ell} \sum_t O_t \rangle = \langle \phi_{\pi_*} - \phi_\pi, \theta_* \rangle \pm \sqrt{\frac{1}{n_\ell} \|\phi_{\pi_*} - \phi_\pi\|_{A(w^{(\ell)})^{-1}}^2}$ so we can safely conclude that a policy π is sub-optimal (i.e., $\pi \neq \pi_*$) if there exists any policy π' such that $\langle \phi_{\pi'} - \phi_\pi, \frac{1}{n_\ell} \sum_t O_t \rangle \gg \sqrt{\frac{1}{n_\ell} \|\phi_{\pi'} - \phi_\pi\|_{A(w^{(\ell)})^{-1}}^2}$. This is the intuition behind Contextual RAGE (Algorithm 1), which inherits its name from the best-arm identification algorithm of [13] that inspired its strategy.

However, while $\langle \phi_{\pi_*} - \phi_\pi, \frac{1}{n_\ell} \sum_t O_t \rangle$ is unbiased and has controlled variance, it is potentially heavy-tailed because $w_{c,a}^{(\ell)}$ can be arbitrarily small. Instead of trying to control $w_{c,a}^{(\ell)}$ and appealing to Bernstein's inequality, in line 7 we use the robust mean estimator of Catoni [28]. We can then show:

Lemma 3.1. $\pi_* \in \Pi_\ell$ and $\max_{\pi \in \Pi_\ell} \langle \phi_{\pi_*} - \phi_\pi, \theta_* \rangle \leq 4\epsilon_\ell$ for all $\ell > 1$ w.p. at least $1 - \delta$.

The lemma states that if Π_ℓ is the active set of policies still under consideration, the optimal policy π_* is never discarded from Π_ℓ , and moreover, the quality of all policies remaining in Π_ℓ is getting better and better. We are now ready to state the main sample complexity result, with proof in Appendix B.

Theorem 3.2. Fix any policy class $\Pi = \{\pi : \mathcal{C} \rightarrow \mathcal{A}\}_\pi$, distribution over contexts ν , $\delta \in (0, 1)$, $\epsilon \geq 0$, and feature map $\phi : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}^d$ such that $r(c, a) = \langle \phi(c, a), \theta_* \rangle$ (w.l.o.g. one can always take $\phi(c, a) = \text{vec}(\mathbf{e}_c \mathbf{e}_a^\top)$). With probability at least $1 - \delta$, if $\phi_\pi = \mathbb{E}_{c \sim \nu}[\phi(c, \pi(c))]$ and $\pi_* = \arg \max_{\pi} \langle \phi_\pi, \theta_* \rangle$ then Contextual-RAGE returns a policy $\hat{\pi} \in \Pi$ such that $V(\hat{\pi}) \geq V(\pi_*) - \epsilon$ after taking at most

$$c \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)}^2}{(\langle \phi_{\pi_*} - \phi_\pi, \theta_* \rangle \vee \epsilon)^2} \log(\log((\Delta \vee \epsilon)^{-1}) |\Pi| / \delta) \log((\Delta \vee \epsilon)^{-1})$$

samples, where c is an absolute constant and $\Delta = \min_{\pi \in \Pi \setminus \pi_*} V(\pi_*) - V(\pi)$.

3.3 Towards a more efficient algorithm

One major issue with Algorithm 1 is that it explicitly maintains a set of policies Π_ℓ from round to round. Since Π could be exponential in $|\mathcal{A}|$, this is a non-starter for any implementation. As a motivation for our approach, we consider a non-elimination algorithm, Algorithm 2, as an intermediate step. It does not maintain Π_ℓ and instead just solves the optimization problem (2) over Π . The design computed in (2) is chosen to ensure that for all $\pi \in \Pi$, $|\hat{\Delta}_{\ell-1}(\pi, \hat{\pi}_{\ell-1}) - \Delta(\pi, \pi_*)| \leq 2\epsilon_{\ell-1} + \frac{1}{4}\Delta(\pi, \pi_*)$ with high probability (Lemma C.3). Equivalently, we estimate gaps up to a constant factor for policies with $\Delta(\pi, \pi_*) > \epsilon_\ell$, while our gap estimates are bounded by ϵ_ℓ for those policies satisfying $\Delta(\pi, \pi_*) \leq \epsilon_\ell$. This ensures that our choice of $\hat{\pi}_\ell$ is good enough, i.e. satisfies $V(\pi_*) - V(\hat{\pi}_\ell) \leq \epsilon_\ell$ with high probability. The full proof is in Appendix C.

Unfortunately, Algorithm 2 introduces additional problems. It is not clear whether solving (2) is computationally efficient. Also, we need to find an estimator $\hat{\Delta}_l$ that is computationally efficient even if the policy space Π is infinite. In addition, it requires precise knowledge of ν to even define the domain of distributions Ω optimized over, and store the solution $w \in \mathcal{C} \times \mathcal{A}$ explicitly. But in general, such precise knowledge will not be available and is only estimable using past data (Assumption 1).

3.4 An instance-optimal and computationally efficient algorithm.

In this section we provide Algorithm 3, which witnesses the guarantees of Theorem 2.14 for the general agnostic contextual bandit problem. We now address the caveats of the previous approaches.

Access to Offline Data. By Assumption 1, we have access to a large amount of sampled offline contexts \mathcal{D} , where each $c_t \in \mathcal{D}$ is drawn IID from ν . Having access to \mathcal{D} allows us to approximate $\mathbb{E}_{c \sim \nu}[\cdot]$ with expectations over the empirical distribution $\mathbb{E}_{c \sim \nu_{\mathcal{D}}}[\cdot]$, where $\nu_{\mathcal{D}}$ is the uniform distribution over historical data \mathcal{D} . The number of offline contexts we need only scales logarithmically over the size of the policy set Π , more specifically, $\text{poly}(|\mathcal{A}|, \epsilon^{-1}, \log(|\Pi|), \log(1/\delta))$. We quantify the precise number of samples needed in Appendix D.2.

Algorithm 1 Elimination Contextual RAGE**Input:** $\Pi, \phi : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}^d, \delta \in (0, 1)$

- 1: **Initialize** $\Pi_1 = \Pi$
- 2: **for** $\ell = 1, 2, \dots, \lceil \log_2(1/\epsilon) \rceil$ **do**
- 3: $\epsilon_\ell := 2^{-\ell}, \delta_\ell := \delta/(2\ell^2|\Pi|)$
- 4: **Let** n_ℓ **be the minimum value s.t.:**

$$\min_{w \in \Omega} \max_{\pi, \pi' \in \Pi_\ell} \frac{\|\phi_\pi - \phi_{\pi'}\|_{A(w)^{-1}}^2 \log(1/\delta_\ell)}{n_\ell} \leq \epsilon_\ell^2$$

with solution $w^{(\ell)}$.

- 5: **For each** $t \in [n_\ell]$, **get** $c_t \sim \nu$, **pull** $a_t \sim p_{c_t}^{(\ell)}$, **observe reward** r_t
- 6: **Compute** $O_t = A(w^{(\ell)})^{-1} \phi(c_t, a_t) r_t$.
- 7: **For** $\pi, \pi' \in \Pi_\ell$

$$\widehat{\Delta}_\ell(\pi, \pi') = \text{Cat}(\{\langle \phi_\pi - \phi_{\pi'}, O_i \rangle\}_{i=1}^{n_\ell})$$

- 8: **Update**

$$\Pi_{\ell+1} = \Pi_\ell \setminus \{\pi' \in \Pi_\ell \mid \max_{\pi \in \Pi_\ell} \widehat{\Delta}_\ell(\pi, \pi') > \epsilon_\ell\}$$

- 9: **end for**

Output: $\Pi_{\ell+1}$ **Algorithm 2** Non-elimination Contextual RAGE**Input:** $\Pi, \phi : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}^d, \delta \in (0, 1)$

- 1: **Initialize:** $\widehat{\pi}_0 \in \Pi$ **arbitrarily**
- 2: **for** $\ell = 1, 2, \dots, \lceil \log_2(1/\epsilon) \rceil$ **do**
- 3: $\epsilon_\ell := 2^{-\ell}, \delta_\ell := \delta/(2\ell^2|\Pi|)$
- 4: **Let** n_ℓ **be the minimum value s.t.:**

$$\min_{w \in \Omega} \max_{\pi \in \Pi} -\frac{1}{4} \widehat{\Delta}_{\ell-1}(\pi, \widehat{\pi}_{\ell-1}) + \sqrt{\frac{2\|\phi_\pi - \phi_{\widehat{\pi}_{\ell-1}}\|_{A(w)^{-1}}^2 \log(1/\delta_\ell)}{n_\ell}} \leq \epsilon_\ell. \quad (2)$$

with solution $w^{(\ell)}$

- 5: **For each** $t \in [n_\ell]$, **get** $c_t \sim \nu$, **pull** $a_t \sim p_{c_t}^{(\ell)}$, **observe reward** r_t
- 6: **Compute** $O_t = A(w^{(\ell)})^{-1} \phi(c_t, a_t) r_t$.
- 7: **For each** $\pi \in \Pi$, **let**

$$\widehat{\Delta}_\ell(\pi, \widehat{\pi}_{\ell-1}) = \text{Cat}(\{\langle \phi_\pi - \phi_{\widehat{\pi}_{\ell-1}}, O_i \rangle\}_{i=1}^{n_\ell}).$$

- 8: **Set** $\widehat{\pi}_\ell := \arg \min_{\pi \in \Pi} \widehat{\Delta}_\ell(\pi, \widehat{\pi}_{\ell-1})$

- 9: **end for**

Output: $\widehat{\pi}_\ell$

Computing the design efficiently. As described, the context space \mathcal{C} may be infinite so maintaining a distribution $\omega \in \Omega \subset \Delta_{\mathcal{C} \times \mathcal{A}}$ is not possible. To overcome this issue, we consider the dual problem of equation (2). We can remove the square root by noticing that $2\sqrt{xy} = \min_{\gamma > 0} \gamma x + \frac{y}{\gamma}$, and introducing an additional minimization over the variable $\gamma_\pi, \pi \in \Pi$. Then, the dual problem becomes

$$\max_{\lambda \in \Delta_\Pi} \min_{w \in \Omega} \min_{\gamma_\pi \geq 0} \sum_{\pi \in \Pi} \lambda_\pi \left(-\widehat{\Delta}_{\ell-1}(\pi, \widehat{\pi}_{\ell-1}) + \gamma_\pi \|\phi_\pi - \phi_{\widehat{\pi}_{\ell-1}}\|_{A(w)^{-1}}^2 + \frac{\log(1/\delta_\ell)}{2\gamma_\pi n_\ell} \right). \quad (4)$$

Exchanging the order of the minimums on ω and γ , somewhat surprisingly we have the close-form expression (Lemma E.6)

$$\min_{\omega \in \Omega} \sum_{\pi \in \Pi} \lambda_\pi \gamma_\pi \|\phi_\pi - \phi_{\widehat{\pi}_{\ell-1}}\|_{A(w)^{-1}}^2 = \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top t_a^{(c)}(\widehat{\pi}_{\ell-1})} \right)^2 \right],$$

where for $\pi' \in \Pi$, $t_a^{(c)}(\pi') \in \{0, 1\}^{|\Pi|}$ with $[t_a^{(c)}(\pi')]_\pi := \mathbf{1}\{\pi(c) = a, \pi'(c) \neq a\} + \mathbf{1}\{\pi(c) \neq a, \pi'(c) = a\}$ and $[\lambda \odot \gamma]_\pi := \lambda_\pi \gamma_\pi$. Interestingly, this value is achieved at a sampling distribution ω , which is a *non-linear* function of λ rather than a convex combination over policies (as in [3]). Because we have an expectation over contexts, this expectation can be replaced by an empirical estimate using historical data, thus avoiding any issues with an infinite context space. The final algorithm utilizing these observations found is in Algorithm 3.

The main challenge is finding a solution to the design (5). First, we can reduce it to a saddle point problem over (λ, γ) by considering only a dyadic sequence of $n \in \{2^k : k \in \mathbb{N}\}$. We use an alternating ascent/descent method, with the caveat that λ lives in a simplex, and γ in a box. Both spaces are defined over a potentially infinite set of policies Π (in the worst case exponential in $|\mathcal{C}|$).

To handle this, we use the Frank-Wolfe (FW) method on λ . Referring to the iterates of FW as λ^t , FW guarantees that the size of the support of λ^t in each iterate grows by at most 1. Thus, if initialized as a 1-sparse vector, we only need to maintain a sparse λ^t in each iteration. Each iterate of FW computes

$$\arg \max_{\pi \in \Pi} [\nabla_\lambda h_\ell(\lambda, \gamma, n)]_\pi.$$

To do so, we show that we can appeal to a constrained argmax oracle (AMO) to run the Frank-Wolfe algorithm, a similar approach to [3]. To optimize over γ we use a gradient descent procedure. We show that in each iterate, the support of γ is contained in that of λ , and we can quantify the number of steps of gradient descent needed to find an ϵ -good solution. Though $h_\ell(\lambda, \gamma, n)$ might not be convex in γ , we nevertheless are able to argue that it has a unique minima and that gradient descent converges to this minima. We introduce our subroutine and further discuss the above claims in Appendix D.

Regularized Estimator. While Algorithms 1 and 2 use a robust mean estimator as in equation (3), this estimator is impractical with a very large number of policies Π . Instead, we use a regularized IPS estimator that can be computed using historical data and an argmax oracle.

Algorithm 3 Contextual Oracle-efficient Dualized Algorithm (CODA)

Input: policies $\Pi = \{\pi : \mathcal{C} \rightarrow \mathcal{A}\}_\pi$, feature map $\phi : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}^d$, $\delta \in (0, 1)$, historical data $\mathcal{D} = \{\nu_s\}_s$

- 1: initiate $\hat{\pi}_0 \in \Pi$ arbitrarily, $\lambda_0 = \mathbf{e}_{\hat{\pi}_0}$, $\hat{\Delta}_0(\pi)$, γ_0 , γ_{\min} , γ_{\max} appropriately
- 2: **for** $l = 1, 2, \dots$ **do**
- 3: $\epsilon_l = 2^{-l}$, $\delta_l = \delta/(l^2|\Pi|^2)$
- 4: **Define**

$$h_l(\lambda, \gamma, n) = \sum_{\pi \in \Pi} \lambda_\pi \left(-\hat{\Delta}_{l-1}^{\gamma_{l-1}}(\pi, \hat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_l n} \right) + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top t_a^{(c)}(\hat{\pi}_{l-1})} \right)^2 \right].$$

- 5: **Let** $\lambda^l, \gamma^l, n_l = \text{FW-GD}(\Pi, |\mathcal{A}|, \hat{\pi}_{l-1}, \epsilon_l)$. These are the solutions to

$$n_\ell := \min\{n \in \mathbb{N} : \max_{\lambda \in \Delta_\Pi} \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^{|\Pi|}} h_\ell(\lambda, \gamma, n) \leq \epsilon_\ell\} \quad (5)$$

- 6: **For** $i \in [n_\ell]$ **get** $c_i \sim \nu$, pull $a_i \sim p_{c_i}^{(\ell)}$ where $p_{c_s, a_s}^{(\ell)} \propto \sqrt{(\lambda_l \odot \gamma_l)^\top t_{a_s}^{(c_s)}(\hat{\pi}_{l-1})}$, observe rewards r_s
- 7: **For each** $\pi \in \Pi$, **define** the IPS estimator

$$\hat{\Delta}_l^{\gamma_l}(\pi, \hat{\pi}_{l-1}) = \sum_{s=1}^{n_l} \frac{r_s}{p_{c_s, a_s}^{(\ell)} + [\gamma_l]_\pi} (\mathbf{1}\{\hat{\pi}_{l-1}(c_s) = a_s\} - \mathbf{1}\{\pi(c_s) = a_s\})$$

- 8: **set**

$$\hat{\pi}_l = \arg \min_{\pi \in \Pi} \hat{\Delta}_l^{\gamma_l}(\pi, \hat{\pi}_{l-1}) + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\frac{[\gamma_l]_\pi}{p_{c, a}^{(\ell)}} + \frac{[\gamma_l]_\pi}{p_{c, a'}} \right) \mathbf{1}\{\hat{\pi}_{l-1}(c) \neq \pi(c)\} \right] + \frac{\log(1/\delta_l)}{[\gamma_l]_\pi n_l} \quad (6)$$

9: **end for**

Output: $\hat{\pi}_l$

Algorithm 3 puts all the pieces together and Theorem 3.3 shows our main result. Note that for exposition purposes, we have omitted some additional regularization terms in the optimization problems that have no effect on the sample complexity, but ensure finite-time convergence. Appendix E shows the full algorithm and the proof. In what follows, $\text{poly}_1(|\mathcal{A}|, \epsilon^{-1}, \log(1/\delta)) \cdot \log(|\Pi|)$ and $\text{poly}_2(|\mathcal{A}|, \epsilon^{-1}, \log(1/\delta), \log(|\Pi|))$ are polynomials in their arguments.

Theorem 3.3. *Fix set of policies Π , context distribution ν and reward function $r(c, a) \in [0, 1]$. With probability at least $1 - \delta$, provided a history \mathcal{D} whose size exceeds $\text{poly}_1(|\mathcal{A}|, \epsilon^{-1}, \log(1/\delta)) \cdot \log(|\Pi|)$, Algorithm 3 returns a policy $\hat{\pi}$ satisfying $V(\pi_*) - V(\hat{\pi}) \leq \epsilon$ in a number of samples not exceeding $O(\rho_{*, \epsilon} \log(|\Pi| \log_2(1/\Delta_\epsilon)/\delta) \log_2(1/\Delta_\epsilon))$ where $\Delta_\epsilon := \max\{\epsilon, \min_{\pi \in \Pi} V(\pi_*) - V(\pi)\}$.*

In addition, Algorithm 3 is computationally efficient and requires at most $\text{poly}_2(|\mathcal{A}|, \epsilon^{-1}, \log(1/\delta), \log(|\Pi|))$ calls to a constrained argmax oracle.

Conclusion. This work provides the first instance-dependent lower bounds for the (ϵ, δ) -PAC contextual bandit problem. One limitation of this work is that our analysis of Algorithm 3 does not immediately extend to the realizable linear setting. That is, a computationally efficient algorithm that achieves the same bound is not known to exist. In the general agnostic settings discussed in this work, we proposed a computationally efficient algorithm. A second limitation is the assumption that we have access to a large pool of offline data. Because it seems necessary to plan with some information about the context distribution, it is not clear how one would completely remove such an assumption and achieve the same sample complexity bounds. As with any recommender system, there is the potential for unintended consequences from optimizing just a single metric. Moreover, other potential pitfalls can arise, such as negative feedback loops, if our assumptions fail to hold in real-world environments. Such consequences can be mitigated by tracking a diverse set of metrics.

Acknowledgement and Disclosure of Funding This work was supported, in part, by NSF award 1907907.

References

- [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. *Advances in neural information processing systems*, 24, 2011.
- [2] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. A reductions approach to fair classification. In *International Conference on Machine Learning*, pages 60–69. PMLR, 2018.
- [3] Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. In *International Conference on Machine Learning*, pages 1638–1646. PMLR, 2014.
- [4] Aymen Al Marjani and Alexandre Proutiere. Adaptive sampling for best policy identification in markov decision processes. In *International Conference on Machine Learning*, pages 7459–7468. PMLR, 2021.
- [5] Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, 32(1):48–77, 2002.
- [6] Alina Beygelzimer, Varsha Dani, Tom Hayes, John Langford, and Bianca Zadrozny. Error limiting reductions between classification tasks. In *Proceedings of the 22nd international conference on Machine learning*, pages 49–56, 2005.
- [7] Sébastien Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [8] Andrew Cotter, Maya Gupta, Heinrich Jiang, Nathan Srebro, Karthik Sridharan, Serena Wang, Blake Woodworth, and Seungil You. Training well-generalizing classifiers for fairness metrics and other data-dependent constraints. In *International Conference on Machine Learning*, pages 1397–1405. PMLR, 2019.
- [9] Rémy Degenne, Pierre Ménard, Xuedong Shang, and Michal Valko. Gamification of pure exploration for linear bandits. In *International Conference on Machine Learning*, pages 2432–2442. PMLR, 2020.
- [10] Aniket Anand Deshmukh, Srinagesh Sharma, James W Cutler, Mark Moldwin, and Clayton Scott. Simple regret minimization for contextual bandits. *arXiv preprint arXiv:1810.07371*, 2018.
- [11] Miroslav Dudík, Daniel Hsu, Satyen Kale, Nikos Karampatziakis, John Langford, Lev Reyzin, and Tong Zhang. Efficient optimal learning for contextual bandits. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, pages 169–178, 2011.
- [12] Tanner Fiez, Sergio Gamez, Arick Chen, Houssam Nassif, and Lalit Jain. Adaptive experimental design and counterfactual inference. In *Workshops of Conference on Recommender Systems (RecSys)*, 2022.
- [13] Tanner Fiez, Lalit Jain, Kevin Jamieson, and Lillian Ratliff. Sequential experimental design for transductive linear bandits. In *Advances in Neural Information Processing Systems*, 2019.
- [14] Dylan Foster, Alekh Agarwal, Miroslav Dudík, Haipeng Luo, and Robert Schapire. Practical contextual bandits with regression oracles. In *International Conference on Machine Learning*, pages 1539–1548. PMLR, 2018.
- [15] Dylan Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual bandits with regression oracles. In *International Conference on Machine Learning*, pages 3199–3210. PMLR, 2020.
- [16] Dylan Foster, Alexander Rakhlin, David Simchi-Levi, and Yunzong Xu. Instance-dependent complexity of contextual bandits and reinforcement learning: A disagreement-based perspective. In *Conference on Learning Theory*, pages 2059–2059. PMLR, 2021.
- [17] Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.

- [18] Steve Hanneke et al. Theory of disagreement-based active learning. *Foundations and Trends® in Machine Learning*, 7(2-3):131–309, 2014.
- [19] Botao Hao, Tor Lattimore, and Csaba Szepesvari. Adaptive exploration in linear contextual bandit. In *International Conference on Artificial Intelligence and Statistics*, pages 3536–3545. PMLR, 2020.
- [20] Tzu-Kuo Huang, Alekh Agarwal, Daniel J Hsu, John Langford, and Robert E Schapire. Efficient and parsimonious agnostic active learning. *Advances in Neural Information Processing Systems*, 28, 2015.
- [21] Martin Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *International Conference on Machine Learning*, pages 427–435. PMLR, 2013.
- [22] Chi Jin, Praneeth Netrapalli, Rong Ge, Sham M Kakade, and Michael I Jordan. On nonconvex optimization for machine learning: Gradients, stochasticity, and saddle points. *Journal of the ACM (JACM)*, 68(2):1–29, 2021.
- [23] Kwang-Sung Jun, Lalit Jain, Blake Mason, and Houssam Nassif. Improved confidence bounds for the linear logistic model and applications to bandits. In *International Conference on Machine Learning (ICML)*, pages 5148–5157, 2021.
- [24] Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm identification in multi-armed bandit models. *The Journal of Machine Learning Research*, 17(1):1–42, 2016.
- [25] Akshay Krishnamurthy, Alekh Agarwal, Tzu-Kuo Huang, Hal Daumé III, and John Langford. Active learning for cost-sensitive classification. In *International Conference on Machine Learning*, pages 1915–1924. PMLR, 2017.
- [26] John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side information. *Advances in neural information processing systems*, 20, 2007.
- [27] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670, 2010.
- [28] Gábor Lugosi and Shahar Mendelson. Mean estimation and regression under heavy-tailed distributions: A survey. *Foundations of Computational Mathematics*, 19(5):1145–1190, 2019.
- [29] Shie Mannor and John N Tsitsiklis. Lower bounds on the sample complexity of exploration in the multi-armed bandit problem. In *Learning Theory and Kernel Machines*, pages 418–432. Springer, 2003.
- [30] John Milnor and David W Weaver. *Topology from the differentiable viewpoint*, volume 21. Princeton university press, 1997.
- [31] Sareh Nabi, Houssam Nassif, Joseph Hong, Hamed Mamani, and Guido Imbens. Bayesian meta-prior learning using Empirical Bayes. *Management Science*, 68(3):1737–1755, 2022.
- [32] Fabian Pedregosa, Geoffrey Negiar, Armin Askari, and Martin Jaggi. Linearly convergent frank-wolfe with backtracking line-search. In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, 2020.
- [33] Daniel Russo. Simple bayesian algorithms for best arm identification. In *Conference on Learning Theory*, pages 1417–1418, 2016.
- [34] David Simchi-Levi and Yunzong Xu. Bypassing the monster: A faster and simpler optimal algorithm for contextual bandits under realizability. *Mathematics of Operations Research*, 2021.
- [35] Marta Soare, Alessandro Lazaric, and Rémi Munos. Best-arm identification in linear bandits. *Advances in Neural Information Processing Systems*, 27, 2014.

- [36] Andrea Tirinzoni, Matteo Pirotta, Marcello Restelli, and Alessandro Lazaric. An asymptotically optimal primal-dual incremental algorithm for contextual linear bandits. *Advances in Neural Information Processing Systems*, 33:1417–1427, 2020.
- [37] Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- [38] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- [39] Andrew Wagenmaker and Kevin Jamieson. Instance-dependent near-optimal policy identification in linear mdps via online experiment design. *arXiv preprint arXiv:2207.02575*, 2022.
- [40] Andrew J Wagenmaker, Max Simchowitz, and Kevin Jamieson. Beyond no regret: Instance-dependent pac reinforcement learning. In *Conference on Learning Theory*, pages 358–418. PMLR, 2022.
- [41] Andrea Zanette, Kefan Dong, Jonathan Lee, and Emma Brunskill. Design of experiments for stochastic contextual linear bandits. *Advances in Neural Information Processing Systems*, 34, 2021.

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? **[Yes]** See Section.
- Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]**
 - (b) Did you describe the limitations of your work? **[Yes]** Please see our conclusion.
 - (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** Please see our conclusion.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - (b) Did you include complete proofs of all theoretical results? **[Yes]**
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[N/A]**
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[N/A]**
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[N/A]**
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[N/A]**

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [N/A]
 - (b) Did you mention the license of the assets? [N/A]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]

 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 1.1 | Related work | 2 |
| 2 | Problem statement and main results | 3 |
| 2.1 | Inefficiency of low-regret algorithms | 4 |
| 2.2 | Trivial policy class | 5 |
| 2.3 | Linear policy class | 5 |
| 2.4 | Comparison to the Disagreement Coefficient | 6 |
| 3 | Optimal Algorithms for Contextual Bandits | 7 |
| 3.1 | Reduction to linear realizability and a simple elimination scheme | 7 |
| 3.2 | A simple, impractical, elimination-style algorithm | 7 |
| 3.3 | Towards a more efficient algorithm | 8 |
| 3.4 | An instance-optimal and computationally efficient algorithm. | 8 |
| A | Lower Bound Results | 16 |
| A.1 | Proof of Theorem 2.2 | 16 |
| A.2 | Proof of Theorem 2.6 | 16 |
| A.3 | Trivial Class: Proof of Theorem 2.9 | 17 |
| A.4 | Proofs of Linear Policy Class | 18 |
| A.5 | Proof for Corollary 2.16 | 20 |
| B | Contextual Regret Proofs Section 3.2 | 21 |
| C | Proof for sample complexity of Algorithm 2 | 22 |
| D | The FW-GD subroutine | 26 |
| D.1 | Proof of computational efficiency | 26 |
| D.2 | Quantify the offline data | 29 |
| E | Proof of Theorem 3.3 | 32 |
| F | Intuition for convergence of duality gap | 39 |
| G | Convergence analysis of FW-GD | 40 |
| G.1 | Statement of the convergence results | 40 |
| G.2 | Technical proofs | 43 |
| G.3 | Convergence of gradient descent | 46 |
| G.4 | Guarantees for strong concavity and local strong convexity | 48 |
| G.5 | Proof of strong duality | 53 |
| H | Useful lemmas | 55 |

Appendix

In the appendix we present algorithms and proofs not included in the main text. Broadly speaking,

- Section A presents proofs for lower bounds;
- Section B-C presents proofs for the proposed computationally inefficient algorithms 1 and 2;
- Section D presents results to justify the computational efficiency of Algorithm 3;
- Section E presents arguments for Algorithm 3 hitting the sample complexity lower bound;
- Section F-H provides technical proofs to argue about convergence of our subroutines.

The table below summarises the notations we used in the proof.

| | |
|--|--|
| $t_a^{(c)}(\pi')$ | $\{\mathbf{1}\{\pi(c) = a, \pi'(c) \neq a\} + \mathbf{1}\{\pi(c) \neq a, \pi'(c) = a\}\}_{\pi \in \Pi} \in \mathbb{R}^\Pi$ |
| S_ℓ | $\{\pi \in \Pi : \langle \phi_{\pi_*} - \phi_\pi, \theta^* \rangle = V(\pi_*) - V(\pi) = \Delta(\pi, \pi_*) \leq \epsilon_\ell\}$ |
| $w(\lambda, \gamma)$ | $[w(\lambda, \gamma)]_{a,c} = \nu_c \cdot p_{c,a} = \nu_c \cdot \frac{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta)}}{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta)}}$ |
| $\widehat{\Delta}_l^\gamma(\pi, \pi')$ | $\sum_{s=1}^{n_l} \frac{r_s}{p_{c_s, a_s} + \gamma_\pi} (\mathbf{1}\{\pi'(c_s) = a_s\} - \mathbf{1}\{\pi(c_s) = a_s\})$ |
| $h_l(\lambda, \gamma, n)$ | $\sum_{\pi \in \Pi} \lambda_\pi \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_\pi n} \right) + \gamma_\pi \mathbb{E}_{c \sim \nu_D} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right)^2 \right]$ |
| $\mathcal{P}_l(w, \gamma)$ | $\max_{\pi \in \Pi} \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \gamma \ \phi_\pi - \phi_{\widehat{\pi}_{l-1}}\ _{A(w)^{-1}}^2 + \frac{\log(1/\delta_l)}{\gamma n_l} \right)$ |

Table 1: Glossary

A Lower Bound Results

A.1 Proof of Theorem 2.2

We quickly point out that the proof of Theorem 2.2 is identical to the proof of the linear policy class case proof of Theorem 2.13. Please see that argument below.

A.2 Proof of Theorem 2.6

Proof of Theorem 2.6. To relate the random stopping time to the regret bound, note that

$$\sum_{c,a} \mathbb{E}_\mu [T_{c,a}(\tau)] (r(c, \pi_*(c)) - r(c, a)) \leq \mathbb{E}_\mu \left[\sqrt{\alpha |\mathcal{A}| \tau} \right] \leq \sqrt{\alpha |\mathcal{A}| \mathbb{E}_\mu [\tau]}$$

where the last inequality follows by Jensen's inequality. Since $\pi_1 := \pi_*$ for our particular instance, if $\bar{c} = \arg \min_{c \in [m]} \mathbb{E}_\mu [T_{c, \pi_c(c)}(\tau)]$ then

$$\begin{aligned} \sum_{c,a} \mathbb{E}_\mu [T_{c,a}(\tau)] (r(c, \pi_1(c)) - r(c, a)) &= \sum_{c,a} \mathbb{E}_\mu [T_{c,a}(\tau)] \Delta \mathbf{1}\{a \neq \pi_1(c)\} \\ &\geq \sum_c \max_a \mathbb{E}_\mu [T_{c,a}(\tau)] \Delta \mathbf{1}\{a \neq \pi_1(c)\} \\ &\geq m \min_c \max_a \mathbb{E}_\mu [T_{c,a}(\tau)] \Delta \mathbf{1}\{a \neq \pi_1(c)\} \\ &= m \mathbb{E}_\mu [T_{\bar{c}, \pi_{\bar{c}}(\bar{c})}(\tau)] \Delta. \end{aligned}$$

Combining the two equations above, and rearranging, we observe that

$$\mathbb{E}_\mu [T_{\bar{c}, \pi_{\bar{c}}(\bar{c})}(\tau)] \leq \frac{1}{m \Delta} \sqrt{\alpha |\mathcal{A}| \mathbb{E}_\mu [\tau]}.$$

Define an instance $\mu' = (\nu, r')$ such that $r'(c, a) = r(c, a)$ for all $(c, a) \in [m] \times \{0, 1\} \setminus (\bar{c}, 1)$, and set $r'(\bar{c}, 1) = r'(\bar{c}, \pi_{\bar{c}}(\bar{c})) = 2\Delta$ under μ' (instead of $r(\bar{c}, \pi_{\bar{c}}(\bar{c})) = 0$ under μ). Note that under μ' ,

we now have that $\pi_{\bar{\epsilon}}$ is the unique optimal policy. If the algorithm is $(0, \delta)$ -PAC then by [24, Lemma 1] we have that

$$\begin{aligned} \log(1/2.4\delta) &\leq \sum_{c,a} KL(\mathcal{N}(r(c,a), 1)|\mathcal{N}(r'(c,a), 1)) \cdot \mathbb{E}_\mu[T_{c,a}(\tau)] \\ &= KL(\mathcal{N}(0, 1)|\mathcal{N}(2\Delta, 1)) \cdot \mathbb{E}_\mu[T_{\bar{\epsilon}, \pi_{\bar{\epsilon}}(\bar{\epsilon})}(\tau)] = 2\Delta^2 \cdot \mathbb{E}_\mu[T_{\bar{\epsilon}, \pi_{\bar{\epsilon}}(\bar{\epsilon})}(\tau)] \\ &\leq 2\Delta^2 \cdot \frac{1}{m\Delta} \sqrt{\alpha|\mathcal{A}|\mathbb{E}_\mu[\tau]} = \sqrt{\frac{4\alpha\mathbb{E}_\mu[\tau]}{m^2\Delta^{-2}}}. \end{aligned}$$

The result follows by rearranging. \square

A.3 Trivial Class: Proof of Theorem 2.9

Firstly note that

$$\begin{aligned} \rho_{\Pi,0}(\Pi, v) &= \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\pi \in \Pi \setminus \pi_*} \frac{\mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{c, \pi(c)}} + \frac{1}{p_{c, \pi_*(c)}} \right) \mathbf{1}\{\pi_*(c) \neq \pi(c)\} \right]}{(\mathbb{E}_{c \sim \nu} [r(c, \pi_*(c)) - r(c, \pi(c))])^2} \\ &= \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\pi \in \Pi \setminus \pi_*} \frac{\sum_{c \in \mathcal{C}} \nu_c \left(\frac{1}{p_{c, \pi(c)}} + \frac{1}{p_{c, \pi_*(c)}} \right) \mathbf{1}\{\pi_*(c) \neq \pi(c)\}}{(\sum_{c \in \mathcal{C}} \nu_c \Delta_{c, \pi(c)} \mathbf{1}\{\pi_*(c) \neq \pi(c)\})^2} \\ &= \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\substack{\alpha \in \{0,1\}^{|\mathcal{C}| \times |\mathcal{A}| \setminus \mathbf{0}} \\ \sum_a \alpha_{c,a} \in \{0,1\}}} \frac{\sum_{c,a} \alpha_{c,a} \nu_c \left(\frac{1}{p_{c, \pi(c)}} + \frac{1}{p_{c, \pi_*(c)}} \right) \mathbf{1}\{\pi_*(c) \neq a\}}{(\sum_{c,a} \alpha_{c,a} \nu_c \Delta_{c, \pi(c)} \mathbf{1}\{\pi_*(c) \neq \pi(c)\})^2} \\ &= \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{c,a: \pi_*(c) \neq a} \frac{\nu_c \left(\frac{1}{p_{c,a}} + \frac{1}{p_{c, \pi_*(c)}} \right)}{(\nu_c \Delta_{c,a})^2} \\ &\leq \max_c \frac{2}{\nu_c} \sum_{a'} \Delta_{c,a'}^{-2}, \end{aligned}$$

where the last equality follows from repeated application of the inequality $\frac{a_1+a_2}{(b_1+b_2)^2} \leq \frac{a_1}{b_1^2} \vee \frac{a_2}{b_2^2}$.

Proof of Theorem 2.9. The proof of the instance-dependent lower bound for $\epsilon = 0$ follows directly from Theorem 2.2. The second minimax statement is, to our best knowledge, novel.

First, note that $\sup_\mu \mathbb{E}_\mu[\tau] \geq \epsilon^{-2} |\mathcal{A}| \log(1/\delta)$ by a reduction to multi-armed bandits by just setting $\nu_1 = 1$ and $\nu_c = 0$ for all $c \neq 1$ [24, 29]. If U denotes the set of instances that achieves this supremum, and V is another set of instances, we note that $\sup_\mu \mathbb{E}_\mu[\tau] = \sup_P \mathbb{E}_{\mu \sim P} \mathbb{E}_\mu[\tau] \geq \frac{1}{2} \sup_{\mu \in U} \mathbb{E}_\mu[\tau] + \frac{1}{2} \sup_{\mu \in V} \mathbb{E}_\mu[\tau]$ for some other set of instances V . Thus, it remains to show that $\sup_\mu \mathbb{E}_\mu[\tau] \geq \epsilon^{-2} |\mathcal{A}| \cdot |\mathcal{C}|$.

Consider the following construction of $|\Pi| = |\mathcal{A}|^{|\mathcal{C}|}$ instances. For each context $c \in \mathcal{C}$ let $\nu_c = 1/|\mathcal{C}|$, and for each $\pi \in \Pi$ let $r_\pi(c, a) = \alpha \epsilon \mathbf{1}\{\pi(c) = a\}$ for some $\alpha > 0$ to be determined later. Clearly, policy π is the unique optimal policy under the reward function $r_\pi(s, a)$. Assume that observations are perturbed by Gaussian $\mathcal{N}(0, 1)$ noise.

Fix $p \in (1/2, 1)$ to be determined later. Let $S := \{c \in \mathcal{C} : \mathbb{P}_{\mu_\pi}(\pi(c) = \hat{\pi}(c)) > p\}$ and suppose $|S| \leq |\mathcal{C}|/8$. Then

$$\begin{aligned}
\mathbb{P}_{\mu_\pi}(V(\pi) - V(\hat{\pi}) \leq \epsilon) &= \mathbb{P}_{\mu_\pi}\left(\frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \alpha \mathbf{1}\{\hat{\pi}(c) \neq \pi(c)\} \leq \epsilon\right) \\
&= \mathbb{P}_{\mu_\pi}\left(\sum_{c \in \mathcal{C}} \mathbf{1}\{\hat{\pi}(c) \neq \pi(c)\} \leq |\mathcal{C}|/\alpha\right) \\
&= \mathbb{P}_{\mu_\pi}\left(\sum_{c \in \mathcal{C}} \mathbf{1}\{\hat{\pi}(c) = \pi(c)\} \geq |\mathcal{C}|(1 - 1/\alpha)\right) \\
&\leq \mathbb{P}_{\mu_\pi}\left(\sum_{c \in \mathcal{C} \setminus S} \mathbf{1}\{\hat{\pi}(c) = \pi(c)\} \geq |\mathcal{C}|(1 - 1/\alpha - 1/8)\right) \\
&\leq \frac{\sum_{c \in \mathcal{C} \setminus S} \mathbb{P}_{\mu_\pi}(\hat{\pi}(c) = \pi(c))}{|\mathcal{C}|(1 - 1/\alpha - 1/8)} \leq \frac{p}{1 - 1/\alpha - 1/8} \leq 5/6
\end{aligned}$$

with $p = 5/8$ and $\alpha = 8$. This implies that for $\delta \in (0, 1/8)$, any (ϵ, δ) -PAC algorithm must satisfy $\min_\pi |\{c \in \mathcal{C} : \mathbb{P}_{\mu_\pi}(\pi(c) = \hat{\pi}(c)) > p\}| \geq |\mathcal{C}|/8$.

Assume the algorithm is permutation invariant (note that any reasonable algorithm satisfies this, including UCB, Thompson Sampling, elimination, etc.). Let $\mu_\pi^{(i)} = (\nu, r_0)$ where $r_\pi^{(i)}(c, i) = r_\pi^{(i)}(c, \pi(c)) = \alpha\epsilon$, and $r_\pi^{(i)}(c, j) = 0$ for $j \notin \{i, \pi(c)\}$. Note that $\mathbb{P}_{\mu_\pi}(\pi(c) = \hat{\pi}(c)) \geq p = 5/6$ and also by the symmetric algorithm assumption that $\mathbb{P}_{\mu_\pi^{(i)}}(\pi(c) = \hat{\pi}(c)) \leq 1/2$ because there are two identical best-arms. Note that $\sum_{j \in \mathcal{A}} \mathbb{E}_{\mu_\pi}[T_{c,j}] KL(\mu_\pi(j), \mu_\pi^{(i)}(j)) = \mathbb{E}_{\mu_\pi}[T_{c,i}] \alpha^2 \epsilon^2 / 2$ for $i \neq \pi(c)$. Putting these two pieces together and applying Lemma 1 of [24], we have:

$$\begin{aligned}
\mathbb{E}_{\mu_\pi}[T_{c,i}] \alpha^2 \epsilon^2 / 2 &= \sum_{j \in \mathcal{A}} \mathbb{E}_{\mu_\pi}[T_{c,j}] KL(\mu_\pi(j), \mu_\pi^{(i)}(j)) \\
&\geq d(\mathbb{P}_{\mu_\pi}(\pi(c) = \hat{\pi}(c)), \mathbb{P}_{\mu_\pi^{(i)}}(\pi(c) = \hat{\pi}(c))) \\
&\geq d(5/6, 1/2) = \frac{1}{6} \log(5^5/3^6) \geq 1/10.
\end{aligned}$$

Thus, $\mathbb{E}_{\mu_\pi}[\sum_{i \neq \pi_*(c)} T_{c,i}] \geq \frac{1}{5} \alpha^{-2} \epsilon^{-2} (|\mathcal{A}| - 1)$ and this must occur on at least $|\mathcal{C}|/8$ contexts. Pick one context c of these arbitrarily. Then

$$\frac{1}{5} \alpha^{-2} \epsilon^{-2} (|\mathcal{A}| - 1) \leq \mathbb{E}_{\mu_\pi} \left[\sum_{i \neq \pi_*(c)} T_{c,i} \right] = \mathbb{E}_{\mu_\pi} \left[\sum_{t=1}^{\tau} \mathbf{1}\{c_t = c\} \right] = \mathbb{E}_{\mu_\pi}[\tau] \nu_c = \mathbb{E}_{\mu_\pi}[\tau] / |\mathcal{C}|.$$

Consequently, $\mathbb{E}[\tau] \geq \frac{1}{5} \alpha^{-2} \epsilon^{-2} (|\mathcal{A}| - 1) |\mathcal{C}|$. □

A.4 Proofs of Linear Policy Class

Recall a quantity fundamental to our sample complexity results:

$$\rho_{\text{lin}, \epsilon} := \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\pi \in \Pi \setminus \pi_*} \frac{\|\phi_\pi - \phi_{\pi_*}\|_{\mathbb{E}_{c \sim \nu} [\sum_{a \in \mathcal{A}} p_{c,a} \phi(c,a) \phi(c,a)^\top]^{-1}}^2}{\langle \phi_{\pi_*} - \phi_\pi, \theta_* \rangle^2 \vee \epsilon^2}. \quad (7)$$

Proof of Corollary 2.11. Consider an ϵ -cover of Π and denote it as Π' . Since we are only interested in finding an ϵ -good policy, it is sufficient to find an ϵ -good policy of Π' . Therefore, we can replace Π with Π' in the statement of Theorem 2.4. By inspecting the difference between the statement of the corollary and Theorem 2.4, it is left to show that we can replace $\log(|\Pi'|)$ with $d \log(1/\epsilon)$. In what follows, we will construct an ϵ -cover of Π . Let $\Theta \subset \mathbb{R}^d$ denote the space of θ . Since the reward function $r(c, a) \in [0, 1]$ is bounded for any $c \in \mathcal{C}$ and $a \in \mathcal{A}$, without loss of generality, we assume $\|\theta\|_2 \leq 1$ so $\Theta \subset \mathcal{B}^d$ where \mathcal{B}^d is the unit ball of dimension d . Let Θ' be an ϵ -net of Θ . For any $\theta' \in \Theta'$, define the policy $\pi_{\theta'}$ such that $\pi_{\theta'} := \arg \max_{a \in \mathcal{A}} \langle \phi(c, a), \theta' \rangle$. Then, define $\Pi' = \{\pi_{\theta'} : \theta' \in \Theta'\}$. First, Π' is an ϵ -cover of Π since Θ' is an ϵ -cover of Θ . Also, $|\Pi'| = |\Theta'|$ by construction. By classical argument on covering numbers [38], we have that $|\Theta'| \leq (3/\epsilon)^d$ so $\log(|\Theta'|) \leq d \log(3/\epsilon) = O(d \log(1/\epsilon))$. □

We quickly point out that the proof of Theorem 2.2 is identical to the proof of the linear policy class case proof of Theorem 2.13.

Proof of Theorem 2.13. For any $\theta \in \mathbb{R}^d$ let $\mathbb{P}_\theta(\cdot)$ and $\mathbb{E}_\theta[\cdot]$ denote the probability and expectation laws under θ and ν such that $c_t \sim \nu$ and playing action $a_t \in \mathcal{A}$ results in reward $r_t \sim \mathcal{N}(\langle \phi(c_t, a_t), \theta \rangle, 1)$. If an algorithm is $(0, \delta)$ -PAC then $\sup_{\theta \in \mathbb{R}^d} \mathbb{P}_\theta(V(\hat{\pi}(c)) < V(\pi_*(c))) \leq \delta$. Now, of course, under θ we have that

$$\begin{aligned} V(\hat{\pi}(c)) < V(\pi_*(c)) &\iff \mathbb{E}_{c \sim \nu}[\langle \theta, \phi(c, \hat{\pi}(c)) - \phi(c, \pi_*(c)) \rangle] < 0 \\ &\iff \langle \theta, \phi_{\hat{\pi}} - \phi_{\pi_*} \rangle < 0 \\ &\iff \exists c : \nu_c \langle \theta, \phi(c, \hat{\pi}(c)) - \phi(c, \pi_*(c)) \rangle < 0. \end{aligned}$$

Fix $\theta_* \in \mathbb{R}^d$ and recall that under θ we have that $\pi_*(c) = \arg \max_{a \in \mathcal{A}} \langle \phi(c, a), \theta \rangle$. Fix any $\theta \in \mathbb{R}^d$ and $\max_{c, a} \nu_c \langle \theta, \phi(c, a) - \phi(c, \pi_*(c)) \rangle > 0$. Then by [24, Lemma 1] we have that

$$\begin{aligned} &d(\mathbb{P}_{\theta_*}(\hat{\pi} = \pi_*), \mathbb{P}_\theta(\hat{\pi} = \pi_*)) \\ &\leq \sum_{c', a'} \mathbb{E}_{\theta_*}[T_{c', a'}(\tau)] KL(\mathcal{N}(\langle \theta_*, \phi(c', a') \rangle, 1) | \mathcal{N}(\langle \theta, \phi(c', a') \rangle, 1)) \\ &= \sum_{c', a'} \mathbb{E}_{\theta_*}[T_{c', a'}(\tau)] \|\theta_* - \theta\|_{\phi(c', a')\phi(c', a')^\top}^2 / 2 \\ &= \mathbb{E}_{\theta_*}[\tau] \sum_{c', a'} \frac{\mathbb{E}_{\theta_*}[T_{c', a'}(\tau)]}{\mathbb{E}_{\theta_*}[\tau]} \|\theta_* - \theta\|_{\phi(c', a')\phi(c', a')^\top}^2 / 2 \\ &\leq \max_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \mathbb{E}_{\theta_*}[\tau] \sum_{c', a'} \nu_{c'} p_{c', a'} \|\theta_* - \theta\|_{\phi(c', a')\phi(c', a')^\top}^2 / 2 \\ &= \max_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \mathbb{E}_{\theta_*}[\tau] \|\theta_* - \theta\|_{\mathbb{E}_{c \sim \nu}[\sum_a p_{c, a} \phi(c, a)\phi(c, a)^\top]}^2 / 2 \end{aligned}$$

where the last inequality follows from Wald's identity:

$$\sum_{a' \in \mathcal{A}} \mathbb{E}_{\theta_*}[T_{c', a'}(\tau)] = \sum_{a' \in \mathcal{A}} \mathbb{E}_{\theta_*} \left[\sum_{t=1}^{\tau} \mathbf{1}\{a_t = a', c_t = c'\} \right] = \mathbb{E}_{\theta_*} \left[\sum_{t=1}^{\tau} \mathbf{1}\{c_t = c'\} \right] = \mathbb{E}_{\theta_*}[\tau] \nu_{c'}.$$

Noting that $d(\mathbb{P}_{\theta_*}(\hat{\pi} = \pi_*), \mathbb{P}_\theta(\hat{\pi} = \pi_*)) \geq \log(1/2.4\delta)$ and we can minimize over θ , given the conditions, we have that

$$\begin{aligned} \log(1/2.4\delta) &\leq \max_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \min_{\theta : \exists c : \nu_c \langle \theta, \phi(c, a) - \phi(c, \pi_*(c)) \rangle > 0} \mathbb{E}_{\theta_*}[\tau] \|\theta_* - \theta\|_{\mathbb{E}_{c \sim \nu}[\sum_a p_{c, a} \phi(c, a)\phi(c, a)^\top]}^2 / 2 \\ &= \mathbb{E}_{\theta_*}[\tau] \max_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \min_{\substack{c, a \in \mathcal{C} \times \mathcal{A} \\ \pi_*(c) \neq a}} \frac{\langle \phi(c, \pi_*(c)) - \phi(c, a), \theta_* \rangle^2}{2 \|\phi(c, a) - \phi(c, \pi_*(c))\|_{\mathbb{E}_{c \sim \nu}[\sum_a p_{c, a} \phi(c, a)\phi(c, a)^\top]}^{-1}}. \end{aligned}$$

After rearranging we conclude that

$$\mathbb{E}_{\theta_*}[\tau] \geq \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\substack{c, a \in \mathcal{C} \times \mathcal{A} \\ \pi_*(c) \neq a}} \frac{2 \|\phi(c, a) - \phi(c, \pi_*(c))\|_{\mathbb{E}_{c \sim \nu}[\sum_a p_{c, a} \phi(c, a)\phi(c, a)^\top]}^{-1}}{\langle \phi(c, \pi_*(c)) - \phi(c, a), \theta_* \rangle^2} \log(1/2.4\delta).$$

To see that equation (7) is a lower bound, follow the exact same sequence of steps but taking any $\theta \in \mathbb{R}^d$ and $\max_{\pi \in \Pi} \mathbb{E}_{c \sim \nu}[\langle \theta, \phi(c, \pi(c)) - \phi(c, \pi_*(c)) \rangle] > 0$. \square

Proof of Theorem 2.14 To see the second part of the theorem statement, observe that

$$\begin{aligned}
& \max_{\pi \in \Pi \setminus \pi_*} \|\phi_\pi - \phi_{\pi_*}\|_{\mathbb{E}_{c \sim \nu} [\sum_{a \in \mathcal{A}} p_{c,a} \phi(c,a) \phi(c,a)^\top]^{-1}}^2 \\
&= \max_{\pi \in \Pi \setminus \pi_*} \|\mathbb{E}_{c \sim \nu} [\phi(c, \pi(c)) - \phi(c, \pi_*(c))]\|_{\mathbb{E}_{c \sim \nu} [\sum_{a \in \mathcal{A}} p_{c,a} \phi(c,a) \phi(c,a)^\top]^{-1}}^2 \\
&\leq \max_{\pi \in \Pi \setminus \pi_*} \mathbb{E}_{c \sim \nu} \left[\|\phi(c, \pi(c)) - \phi(c, \pi_*(c))\|_{\mathbb{E}_{c \sim \nu} [\sum_{a \in \mathcal{A}} p_{c,a} \phi(c,a) \phi(c,a)^\top]^{-1}}^2 \right] \\
&\leq \max_{\pi \in \Pi} 4 \mathbb{E}_{c \sim \nu} \left[\|\phi(c, \pi(c))\|_{\mathbb{E}_{c \sim \nu} [\sum_{a \in \mathcal{A}} p_{c,a} \phi(c,a) \phi(c,a)^\top]^{-1}}^2 \right] \\
&= \max_{q \in \Delta_\Pi} 4 \mathbb{E}_{c \sim \nu} \left[\sum_{\pi \in \Pi} q_\pi \|\phi(c, \pi(c))\|_{\mathbb{E}_{c \sim \nu} [\sum_{a \in \mathcal{A}} p_{c,a} \phi(c,a) \phi(c,a)^\top]^{-1}}^2 \right] \\
&= \max_{q \in \Delta_\Pi} 4 \operatorname{Tr} \left(\mathbb{E}_{c \sim \nu} \left[\sum_{\pi \in \Pi} q_\pi \phi(c, \pi(c)) \phi(c, \pi(c))^\top \right] \mathbb{E}_{c \sim \nu} \left[\sum_{a \in \mathcal{A}} p_{c,a} \phi(c,a) \phi(c,a)^\top \right]^{-1} \right) \\
&\leq 4d
\end{aligned}$$

where the last line takes $p_{c,a} = \sum_{\pi \in \Pi} \mathbf{1}\{\pi(c) = a\} q_\pi$, which is at least as good as the minimizing choice in the theorem.

A.5 Proof for Corollary 2.16

Proof. Observe that

$$\begin{aligned}
\rho_{\Pi, \epsilon_0} &:= \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\pi \in \Pi \setminus \pi_*} \frac{\mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{c, \pi(c)}} + \frac{1}{p_{c, \pi_*(c)}} \right) \mathbf{1}\{\pi_*(c) \neq \pi(c)\} \right]}{(\mathbb{E}_{c \sim \nu} [r(c, \pi_*(c)) - r(c, \pi(c))] \vee \epsilon_0)^2} \\
&= \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\epsilon \geq \epsilon_0} \max_{\pi \in \Pi \setminus \pi_* : \Delta(\pi) \leq \epsilon} \frac{\mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{c, \pi(c)}} + \frac{1}{p_{c, \pi_*(c)}} \right) \mathbf{1}\{\pi_*(c) \neq \pi(c)\} \right]}{\epsilon^2} \\
&= \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\epsilon \geq \epsilon_0} \max_{\pi \in \Pi \setminus \pi_* : \Delta(\pi) \leq \epsilon} \frac{\mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{c, \pi(c)}} + \frac{1}{p_{c, \pi_*(c)}} \right) \mathbf{1}\{\pi_*(c) \neq \pi(c), \Delta(\pi) \leq \epsilon\} \right]}{\epsilon^2} \\
&\leq \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\epsilon \geq \epsilon_0} \max_{\pi \in \Pi \setminus \pi_* : \Delta(\pi) \leq \epsilon} \frac{\mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{c, \pi(c)}} + \frac{1}{p_{c, \pi_*(c)}} \right) \mathbf{1}\{\exists \pi \in \Pi : \pi_*(c) \neq \pi(c), \Delta(\pi) \leq \epsilon\} \right]}{\epsilon^2} \\
&\stackrel{(i)}{\leq} \max_{\epsilon \geq \epsilon_0} \max_{\pi \in \Pi \setminus \pi_* : \Delta(\pi) \leq \epsilon} \frac{\mathbb{E}_{c \sim \nu} [(|\mathcal{A}| + |\mathcal{A}|) \mathbf{1}\{\exists \pi \in \Pi : \pi_*(c) \neq \pi(c), \Delta(\pi) \leq \epsilon\}]}{\epsilon^2} \\
&= \max_{\epsilon \geq \epsilon_0} \frac{2|\mathcal{A}| \mathbb{E}_{c \sim \nu} [\mathbf{1}\{\exists \pi \in \Pi : \pi_*(c) \neq \pi(c), \Delta(\pi) \leq \epsilon\}]}{\epsilon^2} \leq \frac{2|\mathcal{A}|}{\epsilon_0} \mathfrak{C}_{\Pi}^{\text{csc}}(\epsilon_0),
\end{aligned}$$

where (i) follows from taking $p_c \in \Delta_{\mathcal{A}}$ to be the uniform distribution over all actions for each $c \in \mathcal{C}$. To relate this to the policy disagreement coefficient, note that

$$\begin{aligned}
\Delta(\pi) &= \mathbb{E}_{c \sim \nu} [r(c, \pi_*(c)) - r(c, \pi(c))] \geq \mathbb{E}_{c \sim \nu} [\mathbf{1}\{\pi(c) \neq \pi_*(c)\} (\min_{c \in \mathcal{C}} \min_{a \in \mathcal{A}} r(c, \pi_*(c)) - r(c, a))] \\
&= \mathbb{P}_\nu(\pi(c) \neq \pi_*(c)) \Delta_{\text{uniform}}.
\end{aligned}$$

Therefore,

$$\begin{aligned}
& \max_{\epsilon \geq \epsilon_0} \frac{2|\mathcal{A}| \mathbb{E}_{c \sim \nu} [\mathbf{1}\{\exists \pi \in \Pi : \pi_*(c) \neq \pi(c), \Delta(\pi) \leq \epsilon\}]}{\epsilon^2} \\
&\leq \max_{\epsilon \geq \epsilon_0} \frac{2|\mathcal{A}| \mathbb{E}_{c \sim \nu} \left[\mathbf{1}\{\exists \pi \in \Pi : \pi_*(c) \neq \pi(c), \mathbb{P}_\nu(\pi(c) \neq \pi_*(c)) \leq \frac{\epsilon}{\Delta_{\text{uniform}}}\} \right]}{\epsilon^2} \\
&\leq \frac{2|\mathcal{A}|}{\epsilon_0 \Delta_{\text{uniform}}} \mathfrak{C}_{\Pi}^{\text{pol}}(\epsilon_0 / \Delta_{\text{uniform}}).
\end{aligned}$$

□

B Contextual Rage Proofs Section 3.2

Proof of Lemma 3.1. For any $\mathcal{V} \subseteq \Pi$ and $\pi \in \mathcal{V}$, define the Catoni estimator as $\widehat{\phi}_{\pi_*, \pi, \ell}(\mathcal{V})$ and define the event

$$\mathcal{E}_{\pi, \ell}(\mathcal{V}) = \{|\widehat{\phi}_{\pi_*, \pi, \ell}(\mathcal{V}) - \langle \phi_{\pi_*} - \phi_{\pi}, \theta_* \rangle| \leq \epsilon_{\ell}\}$$

where it is implicit that $\widehat{\phi}_{\pi_*, \pi, \ell} := \widehat{\phi}_{\pi_*, \pi, \ell}(\mathcal{V})$ is the resulting estimate after round ℓ if Π_{ℓ} had been equal to \mathcal{V} . Define $w_{\ell}(\mathcal{V})$ such that $[w_{\ell}(\mathcal{V})]_{c, a} = \nu_c p_{c, a}^{(\ell)}$ and $\tau_{\ell}(\mathcal{V})$ to be the number of samples in ℓ th round analogously. By the properties of the Catoni estimator, we have for any $\mathcal{V} \subseteq \Pi$ with probability at least $1 - \frac{\delta}{2\ell^2|\Pi|}$ that

$$\begin{aligned} |\widehat{\phi}_{\pi_*, \pi, \ell}(\mathcal{V}) - \langle \phi_{\pi_*} - \phi_{\pi}, \theta_* \rangle| &\leq \|\phi_{\pi_*} - \phi_{\pi}\|_{A(w_{\ell}(\mathcal{V}))^{-1}} \sqrt{\frac{2 \log(2\ell^2|\Pi|/\delta)}{\tau_{\ell}(\mathcal{V}) - \log(2\ell^2|\Pi|/\delta)}} \\ &\leq \sqrt{\frac{\|\phi_{\pi_*} - \phi_{\pi}\|_{A(w_{\ell}(\mathcal{V}))^{-1}}^2}{2\epsilon_{\ell}^{-2} \rho(w_{\ell}(\mathcal{V}), \mathcal{V}) \log(2\ell^2|\Pi|/\delta)}} \sqrt{2 \log(2\ell^2|\Pi|/\delta)} = \epsilon_{\ell}. \end{aligned}$$

Consequently,

$$\begin{aligned} \mathbb{P}\left(\bigcup_{\ell=1}^{\infty} \bigcup_{\pi \in \Pi_{\ell}} \{\mathcal{E}_{\pi, \ell}^c(\Pi_{\ell})\}\right) &\leq \sum_{\ell=1}^{\infty} \mathbb{P}\left(\bigcup_{\pi \in \Pi_{\ell}} \{\mathcal{E}_{\pi, \ell}^c(\Pi_{\ell})\}\right) \\ &= \sum_{\ell=1}^{\infty} \sum_{\mathcal{V} \subseteq \Pi} \mathbb{P}\left(\bigcup_{\pi \in \mathcal{V}} \{\mathcal{E}_{\pi, \ell}^c(\mathcal{V})\}, \Pi_{\ell} = \mathcal{V}\right) \\ &= \sum_{\ell=1}^{\infty} \sum_{\mathcal{V} \subseteq \Pi} \mathbb{P}\left(\bigcup_{\pi \in \mathcal{V}} \{\mathcal{E}_{\pi, \ell}^c(\mathcal{V})\}\right) \mathbb{P}(\Pi_{\ell} = \mathcal{V}) \\ &\leq \sum_{\ell=1}^{\infty} \sum_{\mathcal{V} \subseteq \Pi} \frac{\delta|\mathcal{V}|}{2\ell^2|\Pi|} \mathbb{P}(\Pi_{\ell} = \mathcal{V}) \leq \delta. \end{aligned}$$

Thus, assume $\bigcap_{\ell=1}^{\infty} \bigcap_{\pi \in \Pi_{\ell}} \{\mathcal{E}_{\pi, \ell}(\Pi_{\ell})\}$ holds. For any $\pi \in \Pi_{\ell}$ we have

$$\begin{aligned} \widehat{\phi}_{\pi, \pi_*, \ell} &= \widehat{\phi}_{\pi, \pi_*, \ell} - \langle \phi_{\pi} - \phi_{\pi_*}, \theta_* \rangle + \langle \phi_{\pi_*}, \theta_* \rangle \\ &\leq \epsilon_{\ell} + \langle \phi_{\pi} - \phi_{\pi_*}, \theta_* \rangle \leq \epsilon_{\ell} \end{aligned}$$

which implies that π_* would survive to round $\ell + 1$. And for any $\pi' \in \Pi_{\ell}$ such that $\langle \phi_{\pi_*} - \phi_{\pi'}, \theta_* \rangle > 2\epsilon_{\ell}$ we have

$$\begin{aligned} \max_{\pi \in \Pi_{\ell}} \widehat{\phi}_{\pi, \pi', \ell} &\geq \widehat{\phi}_{\pi_*, \pi', \ell} \\ &= \langle \phi_{\pi'} - \phi_{\pi_*}, \theta_* \rangle - \widehat{\phi}_{\pi', \pi_*, \ell} + \langle \phi_{\pi_*} - \phi_{\pi'}, \theta_* \rangle \\ &> -\epsilon_{\ell} + 2\epsilon_{\ell} = \epsilon_{\ell} \end{aligned}$$

which implies this π' would be kicked out. Note that this implies that $\max_{\pi \in \Pi_{\ell+1}} \langle \phi_{\pi_*} - \phi_{\pi}, \theta_* \rangle \leq 2\epsilon_{\ell} = 4\epsilon_{\ell+1}$. \square

Proof of Theorem 3.2. Define $S_{\ell} = \{\pi \in \Pi : \langle \phi_{\pi_*} - \phi_{\pi}, \theta_* \rangle \leq 4\epsilon_{\ell}\}$. The above lemma implies that with probability at least $1 - \delta$ we have $\bigcap_{\ell=1}^{\infty} \{\Pi_{\ell} \subseteq S_{\ell}\}$. Observe that if for any $\mathcal{V} \subseteq \Pi$ we define $\rho(w, \mathcal{V}) := \min_{w \in \Omega} \max_{\pi, \pi' \in \mathcal{V}} \|\phi_{\pi} - \phi_{\pi'}\|_{A(w)^{-1}}^2$ then

$$\rho(w^{(\ell)}, \Pi_{\ell}) = \min_{w \in \Omega} \max_{\pi, \pi' \in \Pi_{\ell}} \|\phi_{\pi} - \phi_{\pi'}\|_{A(w)^{-1}}^2 \leq \min_{w \in \Omega} \max_{\pi, \pi' \in S_{\ell}} \|\phi_{\pi} - \phi_{\pi'}\|_{A(w)^{-1}}^2 = \rho(S_{\ell}).$$

For $\ell \geq \lceil \log_2(4\Delta^{-1}) \rceil$ we have that $S_\ell = \{\pi_*\}$, thus the sample complexity to identify π_* is

$$\begin{aligned} \sum_{\ell=1}^{\lceil \log_2(4\Delta^{-1}) \rceil} \tau_\ell &= \sum_{\ell=1}^{\lceil \log_2(4\Delta^{-1}) \rceil} \lceil 4\epsilon_\ell^{-2} \rho(w^{(\ell)}, \Pi_\ell) \log(2\ell^2 |\Pi|/\delta) \rceil \\ &\leq \sum_{\ell=1}^{\lceil \log_2(4\Delta^{-1}) \rceil} 4\epsilon_\ell^{-2} \rho(S_\ell) \log(2\ell^2 |\Pi|/\delta) + 1 \\ &\leq c \log(\log(\Delta^{-1}) |\Pi|/\delta) \sum_{\ell=1}^{\lceil \log_2(4\Delta^{-1}) \rceil} \epsilon_\ell^{-2} \rho(S_\ell) \end{aligned}$$

for some absolute constant $c > 0$. We now note that

$$\begin{aligned} \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{(\langle \phi_{\pi_*} - \phi_\pi, \theta^* \rangle)^2} &= \min_{w \in \Omega} \max_{\ell \leq \lceil \log_2(4\Delta^{-1}) \rceil} \max_{\pi \in S_\ell} \frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{(\langle \phi_{\pi_*} - \phi_\pi, \theta^* \rangle)^2} \\ &\geq \frac{1}{\lceil \log_2(4\Delta^{-1}) \rceil} \min_{w \in \Omega} \sum_{\ell=1}^{\lceil \log_2(4\Delta^{-1}) \rceil} \max_{\pi \in S_\ell} \frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{(\langle \phi_{\pi_*} - \phi_\pi, \theta^* \rangle)^2} \\ &\geq \frac{1}{16 \lceil \log_2(4\Delta^{-1}) \rceil} \sum_{\ell=1}^{\lceil \log_2(4\Delta^{-1}) \rceil} \epsilon_\ell^{-2} \min_{w \in \Omega} \max_{\pi \in S_\ell} \|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2 \\ &\geq \frac{1}{64 \lceil \log_2(4\Delta^{-1}) \rceil} \sum_{\ell=1}^{\lceil \log_2(4\Delta^{-1}) \rceil} \epsilon_\ell^{-2} \min_{w \in \Omega} \max_{\pi, \pi' \in S_\ell} \|\phi_\pi - \phi_{\pi'}\|_{A(w)^{-1}}^2 \\ &= \frac{1}{64 \lceil \log_2(4\Delta^{-1}) \rceil} \sum_{\ell=1}^{\lceil \log_2(4\Delta^{-1}) \rceil} \epsilon_\ell^{-2} \rho(S_\ell) \end{aligned}$$

where we have used the fact that $\max_{\pi, \pi' \in S_\ell} \|\phi_\pi - \phi_{\pi'}\|_{A(w)^{-1}}^2 \leq 4 \max_{\pi \in S_\ell} \|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2$ by the triangle inequality. \square

C Proof for sample complexity of Algorithm 2

In this section we provide a proof for the sample complexity of Algorithm 2.

Theorem C.1. *Under \mathcal{E} , for all $\ell \in \mathbb{N}$, the following holds:*

1. $\hat{\pi}_\ell \in S_\ell := \{\pi \in \Pi : V(\pi_*) - V(\pi) \leq \epsilon_\ell\}$;
2. $n_\ell \lesssim \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\pi_*} - \phi_\pi\|_{A(w)^{-1}}^2 \log(1/\delta_\ell)}{\epsilon_\ell^2 + \Delta(\pi)^2}$.

Without loss of generality, we assume that $\forall t$, the reward $r_t \in [0, 1]$. Note that by the result about Cati estimator in [28], we have for all $\ell \in \mathbb{N}$ and $\pi, \pi' \in \Pi$, that

$$|\text{Cat}(\{\langle \phi_\pi - \phi_{\pi'}, O_t \rangle\}_{t=1}^{n_\ell}) - \langle \phi_\pi - \phi_{\pi'}, \theta_* \rangle| \leq \|\phi_\pi - \phi_{\pi'}\|_{A(w^{(\ell)})^{-1}} \sqrt{\frac{2 \log(2\ell^2 |\Pi|/\delta)}{n_\ell - \log(2\ell^2 |\Pi|/\delta)}}.$$

Therefore, in the ℓ th round, we have for any $\pi, \pi' \in \Pi$,

$$\begin{aligned} \left| \widehat{\Delta}_\ell(\pi, \pi') - \Delta(\pi, \pi') \right| &= |\text{Cat}(\{\langle \phi_\pi - \phi_{\pi'}, O_i \rangle\}_{i=1}^{n_\ell}) - \langle \phi_\pi - \phi_{\pi'}, \theta_* \rangle| \\ &\leq \sqrt{\frac{2 \|\phi_\pi - \phi_{\pi'}\|_{A(w^{(\ell)})^{-1}}^2 \log(2\ell^2 |\Pi|/\delta)}{n_\ell}}. \end{aligned} \quad (8)$$

Then, let $\delta_\ell = \frac{\delta}{2\ell^2 |\Pi|}$ we define the event

$$\mathcal{E}_\ell = \bigcap_{\pi, \pi' \in \Pi} \left\{ \left| \widehat{\Delta}_\ell(\pi, \pi') - \Delta(\pi, \pi') \right| \leq \sqrt{\frac{2 \|\phi_\pi - \phi_{\pi'}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_\ell)}{n_\ell}} \right\},$$

and $\mathcal{E} = \bigcap_{l=0}^{\infty} \mathcal{E}_l$. First, by equation 8, we have that \mathcal{E} happens with probability at least $1 - \delta$. In order to show the sample complexity lower bound, we use proof by induction. Note that in a step of Lemma C.4, we can show that $n_l \lesssim \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\hat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w)^{-1}}^2 \log(1/\delta_l)}{\epsilon_l^2 + \Delta(\pi)^2}$, so we induct on this result. Assume in round $l - 1$, $\hat{\pi}_{l-1} \in S_{l-1} = \{\pi \in \Pi : \Delta(\pi, \pi_*) \leq \epsilon_{l-1}\}$ and $n_{l-1} \lesssim \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\hat{\pi}_{l-2}} - \phi_{\pi}\|_{A(w)^{-1}}^2 \log((l-1)^2 |\Pi|^2 / \delta)}{\epsilon_{l-1}^2 + \Delta(\pi)^2}$. Then, the following lemma gives us an upper bound on the UCB.

Lemma C.2. *We have for any $\pi \in \Pi$,*

$$\sqrt{\frac{\|\phi_{\hat{\pi}_l} - \phi_{\pi}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{n_l}} \leq \frac{1}{28} \left(4\epsilon_l + \hat{\Delta}_{l-1}(\pi, \hat{\pi}_{l-1}) \right).$$

Proof. By definition of n_l and $w^{(\ell)}$ and $\pi^{(\ell)}$ being the saddle point, we have

$$\begin{aligned} & -\frac{1}{4} \hat{\Delta}_{l-1}(\pi^{(\ell)}, \hat{\pi}_{l-1}) + 28 \sqrt{\frac{2 \|\phi_{\pi^{(\ell)}} - \phi_{\hat{\pi}_{l-1}}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{n_l}} \\ &= \max_{\pi \in \Pi} -\frac{1}{4} \hat{\Delta}_{l-1}(\pi, \hat{\pi}_{l-1}) + \sqrt{\frac{1568 \|\phi_{\pi} - \phi_{\hat{\pi}_{l-1}}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{n_l}} \leq \epsilon_l. \end{aligned}$$

Solving for n_l gives us

$$n_l \geq \max_{\pi \in \Pi} \frac{1568 \|\phi_{\pi} - \phi_{\hat{\pi}_{l-1}}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{(4\epsilon_l + \hat{\Delta}_{l-1}(\pi, \hat{\pi}_{l-1}))^2}.$$

We have for any $\pi \in \Pi$,

$$\begin{aligned} 2n_l &\geq 3136 \max_{\pi \in \Pi} \frac{\|\phi_{\hat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{(4\epsilon_l + \hat{\Delta}_{l-1}(\pi, \hat{\pi}_{l-1}))^2} \\ &\geq 1568 \frac{\|\phi_{\hat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{(4\epsilon_l + \hat{\Delta}_{l-1}(\pi, \hat{\pi}_{l-1}))^2} \\ &\quad + 1568 \frac{\|\phi_{\hat{\pi}_{l-1}} - \phi_{\hat{\pi}_l}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{(4\epsilon_l + \hat{\Delta}_{l-1}(\hat{\pi}_l, \hat{\pi}_{l-1}))^2} \\ &\stackrel{(i)}{\geq} 1568 \frac{\left(\|\phi_{\hat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w^{(\ell)})^{-1}}^2 + \|\phi_{\hat{\pi}_{l-1}} - \phi_{\hat{\pi}_l}\|_{A(w^{(\ell)})^{-1}}^2 \right) \log(1/\delta_l)}{\max\{(4\epsilon_l + \hat{\Delta}_{l-1}(\hat{\pi}_l, \hat{\pi}_{l-1}))^2, (4\epsilon_l + \hat{\Delta}_{l-1}(\pi, \hat{\pi}_{l-1}))^2\}} \\ &\stackrel{(ii)}{\geq} 1568 \frac{\|\phi_{\hat{\pi}_l} - \phi_{\pi}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{\max\{(4\epsilon_l + \hat{\Delta}_{l-1}(\hat{\pi}_l, \hat{\pi}_{l-1}))^2, (4\epsilon_l + \hat{\Delta}_{l-1}(\pi, \hat{\pi}_{l-1}))^2\}}. \end{aligned}$$

where (i) holds by lower bounding the ratio with a larger denominator, and (ii) holds by triangular inequality. Therefore, using the fact that $\hat{\Delta}(\pi, \hat{\pi}_{l-1}) \geq 0$ for any $\pi \in \Pi$ since $\hat{\pi}_{l-1} = \arg \max_{\pi \in \Pi} \hat{V}_{l-1}(\pi)$, we have $\sqrt{\max\{(4\epsilon_l + \hat{\Delta}_{l-1}(\hat{\pi}_l, \hat{\pi}_{l-1}))^2, (4\epsilon_l + \hat{\Delta}_{l-1}(\pi, \hat{\pi}_{l-1}))^2\}} = \max\{4\epsilon_l + \hat{\Delta}_{l-1}(\hat{\pi}_l, \hat{\pi}_{l-1}), 4\epsilon_l + \hat{\Delta}_{l-1}(\pi, \hat{\pi}_{l-1})\}$, so we have

$$\sqrt{\frac{\|\phi_{\hat{\pi}_l} - \phi_{\pi}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{n_l}} \leq \frac{1}{28} \left(4\epsilon_l + \max\{\hat{\Delta}_{l-1}(\pi, \hat{\pi}_{l-1}), \hat{\Delta}_{l-1}(\hat{\pi}_l, \hat{\pi}_{l-1})\} \right).$$

□

With the above results, the following lemma controls the difference between the empirical gap and the true gap.

Lemma C.3. *With inductive hypotheses, we have for any $\pi \in \Pi$,*

$$|\widehat{\Delta}_{l-1}(\pi, \widehat{\pi}_{l-1}) - \Delta(\pi, \pi_*)| \leq 2\epsilon_{l-1} + \frac{1}{4}\Delta(\pi, \pi_*).$$

Proof. We prove this by induction. First, in round $l = 0$, this holds by choosing a sufficiently large n_0 . Then, in round $l - 1$,

$$\begin{aligned} & |\widehat{\Delta}_{l-1}(\pi, \widehat{\pi}_{l-1}) - \Delta(\pi, \pi_*)| \\ &= |\widehat{\Delta}_{l-1}(\pi, \widehat{\pi}_{l-1}) - \Delta(\pi, \widehat{\pi}_{l-1}) - \Delta(\widehat{\pi}_{l-1}, \pi_*)| \\ &\leq \sqrt{\frac{2 \|\phi_\pi - \phi_{\widehat{\pi}_{l-1}}\|_{A(w^{(l-1))^{-1}}}^2 \log(1/\delta_{l-1})}{n_{l-1}}} + \epsilon_{l-1} \\ &\stackrel{(i)}{\leq} \frac{\sqrt{2}}{28} \left(4\epsilon_{l-1} + \max\{\widehat{\Delta}_{l-2}(\pi, \widehat{\pi}_{l-2}), \widehat{\Delta}_{l-2}(\widehat{\pi}_{l-1}, \widehat{\pi}_{l-2})\} \right) + \epsilon_{l-1} \\ &\stackrel{(ii)}{\leq} \frac{\sqrt{2}}{28} \left(4\epsilon_{l-1} + 2\epsilon_{l-2} + \frac{5}{4}\Delta(\pi, \widehat{\pi}_{l-2}) + 2\epsilon_{l-2} + \frac{5}{4}\Delta(\widehat{\pi}_{l-1}, \widehat{\pi}_{l-2}) \right) + \epsilon_{l-1} \\ &\leq \frac{\sqrt{2}}{28} \left(4\epsilon_{l-1} + 4\epsilon_{l-2} + \frac{5}{4}\Delta(\pi, \pi_*) + \frac{5}{4}\Delta(\widehat{\pi}_{l-1}, \pi_*) \right) + \epsilon_{l-1} \\ &\leq \frac{\sqrt{2}}{28} \left(4\epsilon_{l-1} + 4\epsilon_{l-2} + \frac{5}{4}\Delta(\pi, \pi_*) + \frac{5}{4}\epsilon_{l-1} \right) + \epsilon_{l-1} \\ &\leq 2\epsilon_{l-1} + \frac{1}{4}\Delta(\pi, \pi_*), \end{aligned}$$

where (i) follows from the preceding lemma and (ii) follows from the inductive hypothesis that

$$|\widehat{\Delta}_{l-2}(\pi, \widehat{\pi}_{l-2}) - \Delta(\pi, \pi_*)| \leq 2\epsilon_{l-2} + \frac{1}{4}\Delta(\pi, \pi_*).$$

□

We make use of these two lemmas to state a lower bound on n_l .

Lemma C.4. *Under \mathcal{E} , the choice for n_l in the algorithm satisfies*

$$n_l \lesssim \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\pi_*} - \phi_\pi\|_{A(w)^{-1}}^2 \log(1/\delta_l)}{\epsilon_l^2 + \Delta(\pi)^2}.$$

Proof. By inductive hypothesis on n_{l-1} and under \mathcal{E}_l , we have for any $\pi \in \Pi$,

$$\begin{aligned} \Delta(\pi, \pi_*) &= \Delta(\pi, \widehat{\pi}_{l-1}) + \Delta(\widehat{\pi}_{l-1}, \pi_*) \\ &\stackrel{(i)}{\leq} \widehat{\Delta}_{l-1}(\pi, \widehat{\pi}_{l-1}) + \sqrt{\frac{2 \|\phi_{\widehat{\pi}_{l-1}} - \phi_\pi\|_{A(w^{(l-1))^{-1}}}^2 \log((l-1)^2 |\Pi|^2 / \delta)}{n_{l-1}}} + \epsilon_{l-1} \\ &\stackrel{(ii)}{\leq} \widehat{\Delta}_{l-1}(\pi, \widehat{\pi}_{l-1}) + \frac{\sqrt{2}}{28} \left(4\epsilon_{l-1} + \widehat{\Delta}_{l-2}(\pi, \widehat{\pi}_{l-2}) \right) + \epsilon_{l-1} \\ &\leq \widehat{\Delta}_{l-1}(\pi, \widehat{\pi}_{l-1}) + \frac{\sqrt{2}}{28} \left(4\epsilon_{l-1} + \frac{5}{4}\Delta(\pi, \pi_*) + 2\epsilon_{l-2} \right) + \epsilon_{l-1} \\ &\leq \widehat{\Delta}_{l-1}(\pi, \widehat{\pi}_{l-1}) + \frac{1}{4}\Delta(\pi, \pi_*) + 2\epsilon_{l-1}. \end{aligned}$$

where (i) follows from \mathcal{E}_{l-1} and (ii) follows from Lemma C.2. Therefore,

$$\begin{aligned}
& \min_{w \in \Omega} \max_{\pi \in \Pi} -\frac{1}{4} \widehat{\Delta}_{l-1}(\pi, \widehat{\pi}_{l-1}) + 28 \sqrt{\frac{2 \|\phi_\pi - \phi_{\widehat{\pi}_{l-1}}\|_{A(w)^{-1}}^2 \log(1/\delta_l)}{n_l}} \\
& \leq \min_{w \in \Omega} \max_{\pi \in \Pi} -\frac{3}{16} \Delta(\pi, \pi_*) + \frac{1}{2} \epsilon_l + 28 \sqrt{\frac{2 \|\phi_\pi - \phi_{\widehat{\pi}_{l-1}}\|_{A(w)^{-1}}^2 \log(1/\delta_l)}{n_l}} \\
& \leq \min_{w \in \Omega} \max_{\pi \in \Pi} \left(-\frac{3}{16} \Delta(\pi, \pi_*) + 28 \sqrt{\frac{2 \|\phi_{\pi_*} - \phi_\pi\|_{A(w)^{-1}}^2 \log(1/\delta_l)}{n_l}} \right. \\
& \quad \left. + 28 \sqrt{\frac{2 \|\phi_{\pi_*} - \phi_{\widehat{\pi}_{l-1}}\|_{A(w)^{-1}}^2 \log(1/\delta_l)}{n_l}} \right) + \frac{1}{2} \epsilon_l \\
& \leq \min_{w \in \Omega} \max_{\pi \in \Pi} \left(-\frac{3}{16} \Delta(\pi, \pi_*) + 28 \sqrt{\frac{2 \|\phi_{\pi_*} - \phi_\pi\|_{A(w)^{-1}}^2 \log(1/\delta_l)}{n_l}} \right. \\
& \quad \left. + 28 \sqrt{\frac{\max_{\pi' \in S_{l-1}} 2 \|\phi_{\pi_*} - \phi_{\pi'}\|_{A(w)^{-1}}^2 \log(1/\delta_l)}{n_l}} \right) + \frac{1}{2} \epsilon_l
\end{aligned}$$

which is less than ϵ_l whenever

$$n_l \gtrsim \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\pi_*} - \phi_\pi\|_{A(w)^{-1}}^2 \log(1/\delta_l)}{\epsilon_l^2 + \Delta(\pi, \pi_*)^2}.$$

□

Then we finish our first goal. The next goal is to show that $\widehat{\pi}_l \in S_l$.

Lemma C.5. *Under \mathcal{E}_l , we have $\Delta(\widehat{\pi}_l, \pi_*) \leq \epsilon_l$.*

Proof. On \mathcal{E}_l , we have

$$\begin{aligned}
& \Delta(\widehat{\pi}_l, \widehat{\pi}_{l-1}) \\
& \leq \widehat{\Delta}_l(\widehat{\pi}_l, \widehat{\pi}_{l-1}) + \sqrt{\frac{2 \|\phi_{\widehat{\pi}_l} - \phi_{\widehat{\pi}_{l-1}}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{n_l}} \quad (\text{by event } \mathcal{E}_l) \\
& \leq \widehat{\Delta}_l(\pi_*, \widehat{\pi}_{l-1}) + \sqrt{\frac{2 \|\phi_{\widehat{\pi}_l} - \phi_{\widehat{\pi}_{l-1}}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{n_l}} \quad (\text{by minimality of } \widehat{\pi}_l) \\
& \leq \Delta(\pi_*, \widehat{\pi}_{l-1}) + \sqrt{\frac{2 \|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi_*}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{n_l}} + \sqrt{\frac{2 \|\phi_{\widehat{\pi}_l} - \phi_{\widehat{\pi}_{l-1}}\|_{A(w^{(\ell)})^{-1}}^2 \log(1/\delta_l)}{n_l}} \\
& \quad (\text{by event } \mathcal{E}_l) \\
& \leq \Delta(\pi_*, \widehat{\pi}_{l-1}) + \frac{\sqrt{2}}{28} \left(4\epsilon_l + \widehat{\Delta}_{l-1}(\pi_*, \widehat{\pi}_{l-1}) + 4\epsilon_l + \widehat{\Delta}_{l-1}(\widehat{\pi}_l, \widehat{\pi}_{l-1}) \right) \quad (\text{by Lemma C.2}) \\
& \leq \Delta(\pi_*, \widehat{\pi}_{l-1}) + \frac{\sqrt{2}}{28} \left(4\epsilon_l + 2\epsilon_{l-1} + \frac{5}{4} \Delta(\pi_*, \widehat{\pi}_{l-1}) + 4\epsilon_l + 2\epsilon_{l-1} + \frac{5}{4} \Delta(\widehat{\pi}_l, \widehat{\pi}_{l-1}) \right) \\
& \quad (\text{by Lemma C.3}) \\
& \leq \Delta(\pi_*, \widehat{\pi}_{l-1}) + \frac{3}{56} \left(8\epsilon_{l-1} + \frac{5}{4} \Delta(\widehat{\pi}_l, \pi_*) \right).
\end{aligned}$$

Therefore, $\frac{209}{224} \Delta(\widehat{\pi}_l, \pi_*) \leq \frac{6}{7} \epsilon_l$ and $\Delta(\widehat{\pi}_l, \pi_*) \leq \epsilon_l$, so $\widehat{\pi}_l \in S_l$. □

D The FW-GD subroutine

We now introduce the FW-GD subroutine that solves the optimization problem of equation (5). Note that its objective has three variables. We first reduce it to a max-min problem over (λ, γ) by considering n in a dyadic sequence. This is good enough as we only need to find the optimal n up to a constant factor. Then, we combine the Frank-Wolfe algorithm [21] for minimizing over λ with the gradient descent algorithm [7] which minimizes over γ . Algorithm 4 shows the full subroutine. In line 10, we use the standard gradient descent subroutine combining with a clipping on λ , with details in Algorithm 7.

Algorithm 4 FW-GD

Input: Π policy sets, number of actions $|\mathcal{A}|$, $\hat{\pi}_{l-1} \in \Pi$, $\eta_l > 0$, $K \in \mathbb{N}$, threshold ϵ_l , γ_{\min} , γ_{\max}

- 1: Initialize $n_1 = 1$, $L = |\mathcal{A}|^2 \frac{((1+\eta_l)\gamma_{\max})^{5/2}}{\eta_l^{3/2}\gamma_{\min}^2}$
- 2: **for** $r = 1, 2, \dots$ **do**
- 3: Initialize $\lambda^0 = \mathbf{e}_0 \in \mathbb{R}^\Pi$, $\gamma^0 = \mathbf{1}_{|\Pi|} \cdot \sqrt{\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n_r}} \in \mathbb{R}^{|\Pi|}$ // Never explicitly materialized
- 4: **for** $t = 0, 1, 2, \dots, K$ **do**
- 5: Compute

$$\pi_t = \arg \max_{\pi \in \Pi} [\nabla_\lambda h_l(\lambda^t, \gamma^t, n_r)]_\pi \quad (9)$$
- 6: Set the FW-gap

$$g_t = \langle \nabla_\lambda h_l(\lambda^t, \gamma^t, n_r), \mathbf{e}_{\pi_t} - \lambda^t \rangle = [\nabla_\lambda h_l(\lambda^t, \gamma^t, n_r)]_{\pi_t} - \sum_{\pi \in \text{supp}(\lambda^t)} [\nabla_\lambda h_l(\lambda^t, \gamma^t, n_r)]_\pi$$
- 7: Set $\beta_t = \min \left\{ \frac{g_t}{L \|\lambda^t - \mathbf{e}_{\pi_t}\|_1^2}, 1 \right\}$
- 8: Set $\kappa_t = \frac{\epsilon_l}{(t+1)^2}$
- 9: Set $\lambda^{t+1} = (1 - \beta_t)\lambda^t + \beta_t \mathbf{e}_{\pi_t}$ // Only 1-sparse updates recorded
- 10: Set $\gamma^{t+1} = \text{GD}(\lambda^t, n_r, \kappa_t)$ // Only differences from γ_0 recorded
- 11: **end for**
- 12: **if** $h_l(\lambda^{K+1}, \gamma^{K+1}, n_r) \leq \epsilon_l$ **then**
- 13: **break**
- 14: **else**
- 15: $n_{r+1} = 2 \cdot n_r$
- 16: **end if**
- 17: **end for**

Output: $\lambda^{K+1} \in \Delta_\Pi$, $\gamma^{K+1} \in \mathbb{R}_+^{|\Pi|}$, n_r

In this section, we will mainly focus on showing that the algorithm is computationally efficient with access to an argmax oracle (Definition 2.3), i.e. the second part of Theorem 3.3. Specifically, Section D.1 quantifies the number of oracle calls, and Section D.2 quantifies the number of offline data needed in order to approximate the expectation over the context distribution. We leave the convergence analysis of the algorithm in Section G. The main result for this section is stated below.

Theorem D.1. *Let T_l be the number of iterations for FW-GD in the l th round. Then, Algorithm 3 is computationally efficient and requires at most $O(\sum_{l=1}^{\log_2(1/\epsilon)} T_l^2 |\mathcal{D}|)$ calls to a constrained argmax oracle, with the size of the history \mathcal{D} exceeding $\text{poly}(\epsilon^{-1}, \log |\Pi|, \gamma_{\max}, \gamma_{\min}^{-1}, \eta^{-1}, |\mathcal{A}|, \log(1/\delta))$ with probability at least $1 - \delta$, where poly denotes some polynomial.*

The size of the history follows directly from Lemma D.6 and D.7. We will see that $\eta, \gamma_{\max}, \gamma_{\min}$ all scale at most polynomially on $|\mathcal{A}|$ and ϵ^{-1} , and thus we get the statement in Theorem D.1. The bound on the number of oracle calls follows directly from Lemma D.2 and the fact that $T_{l-1} \leq T_l$. We will see in Theorem G.1 that $T_l = \text{poly}(|\mathcal{A}|, \epsilon_l^{-1})$, which shows that the total number of oracle calls is at most $\text{poly}(|\mathcal{A}|, \epsilon^{-1}, \log(1/\delta), \log(|\Pi|))$.

D.1 Proof of computational efficiency

In this section, we address the technical issues on computational efficiency of our algorithm. Fix an iteration t and let T_l be the number of iterations for FW-GD in the l th round.

Lemma D.2. Equation (9) can be computed with $(t + T_{l-1})|\mathcal{D}|$ calls to a cost-sensitive classification oracle.

Proof. We consider the t th iteration of the l th round for some n_r . In this iteration, we compute

$$\begin{aligned} [\nabla_{\lambda} h_l(\lambda^t, \gamma^t, n_r)]_{\pi} &= \sum_{i=1}^{n_l} \frac{r_i}{p_{c_i, a_i}^{(\ell)} + [\gamma^{l-1}]_{\pi}} (\mathbf{1}\{\pi(c_i) = a_i\} - \mathbf{1}\{\widehat{\pi}_{l-1}(c_i) = a_i\}) + \frac{\log(1/\delta_l)}{[\gamma^t]_{\pi} n} \\ &\quad + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda^t \odot \gamma^t)^{\top} (t_a^{(c)} + \eta_l)} \right) \left(\sum_{a' \in \mathcal{A}} \frac{[\gamma^t]_{\pi} (t_{a'}^{(c)} + \eta_l)_{\pi}}{\sqrt{(\lambda^t \odot \gamma^t)^{\top} (t_{a'}^{(c)} + \eta_l)}} \right) \right]. \end{aligned}$$

Define $\gamma_0 := \sqrt{\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n_r}}$. Initially, each coordinate of γ^t is γ_0 . In round t of the algorithm, at most t coordinates of γ will change, and these coordinates will be in $\text{supp}(\lambda^t)$. Also, for any $j \notin \text{supp}(\lambda^{l-1})$, $\gamma_j^{l-1} = \gamma_0$. Therefore, let $t_a^{(c)}(\cdot, \widehat{\pi}_{l-1}) \in \mathbb{R}^{|\Pi|}$, in round l ,

$$\begin{aligned} &\text{argmax}_{\pi \in \Pi \setminus (\text{supp}(\lambda^t) \cup \text{supp}(\lambda^{l-1}))} [\nabla_{\lambda} h_l(\lambda^t, \gamma^t, n_r)]_{\pi} \\ &= \text{argmax}_{\pi \in \Pi \setminus (\text{supp}(\lambda^t) \cup \text{supp}(\lambda^{l-1}))} \sum_{i=1}^{n_l} \frac{r_i}{p_{c_i, a_i}^{(\ell)} + \gamma_0} \mathbf{1}\{\pi(c_i) = a_i\} + \frac{\log(1/\delta_l)}{\gamma_0 n_r} \\ &\quad + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda^t \odot \gamma^t)^{\top} (t_a^{(c)}(\widehat{\pi}_{l-1}) + \eta_l)} \right) \left(\sum_{a' \in \mathcal{A}} \frac{\gamma_0 (t_{a'}^{(c)}(\widehat{\pi}_{l-1}) + \eta_l)_{\pi}}{\sqrt{(\lambda^t \odot \gamma^t)^{\top} (t_{a'}^{(c)}(\widehat{\pi}_{l-1}) + \eta_l)}} \right) \right] \\ &= \text{argmax}_{\pi \in \Pi \setminus (\text{supp}(\lambda^t) \cup \text{supp}(\lambda^{l-1}))} \sum_{i=1}^{n_l} \frac{r_i}{p_{c_i, a_i}^{(\ell)} + \gamma_0} \mathbf{1}\{\pi(c_i) = a_i\} \\ &\quad + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_{a' \in \mathcal{A}} \frac{\sum_{a \in \mathcal{A}} \sqrt{(\lambda^t \odot \gamma^t)^{\top} (t_a^{(c)}(\widehat{\pi}_{l-1}) + \eta_l)}}{\sqrt{(\lambda^t \odot \gamma^t)^{\top} (t_{a'}^{(c)}(\widehat{\pi}_{l-1}) + \eta_l)}} \gamma_0 t_{a'}^{(c)}(\widehat{\pi}_{l-1})_{\pi} \right] \\ &= \text{argmax}_{\pi \in \Pi \setminus (\text{supp}(\lambda^t) \cup \text{supp}(\lambda^{l-1}))} \sum_{i=1}^{n_l + |\mathcal{D}|} L_i(\pi(c_i)) \end{aligned}$$

which is a cost-sensitive classification problem with cost vector

$$L_i(a) = \begin{cases} \frac{r_i}{p_{c_i, a_i}^{(\ell)} + \gamma_0} \mathbf{1}\{a = a_i\} & \text{for } i = 1, \dots, n_l \\ \left(\frac{\gamma_0}{s_{a, c_i}} + \frac{\gamma_0}{s_{\widehat{\pi}_{l-1}(c_i), c_i}} \right) \mathbf{1}\{a \neq \widehat{\pi}_{l-1}(c_i)\} & \text{for } i = n_l + 1, \dots, n_l + |\mathcal{D}| \end{cases}$$

where $s_{a, c} = \frac{\sqrt{(\lambda^t \odot \gamma^t)^{\top} (t_a^{(c)}(\widehat{\pi}_{l-1}) + \eta_l)}}{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda^t \odot \gamma^t)^{\top} (t_{a'}^{(c)}(\widehat{\pi}_{l-1}) + \eta_l)}}$. Note that $s_{a, c}$ is computable since λ^t has at most t non-zero elements in step t . Then, let $\pi^{\#} := \text{supp}(\lambda^t) \cup \text{supp}(\lambda^{l-1})$, we have

$$\begin{aligned} &\text{argmax}_{\pi \in \Pi} [\nabla_{\lambda} h_l(\lambda^t, \gamma^t, n_r)]_{\pi} \\ &= \text{argmax} \left\{ \text{argmax}_{\pi \in \Pi^{\#}} [\nabla_{\lambda} h_l(\lambda^t, \gamma^t, n_r)]_{\pi}, \text{argmax}_{\pi \in \Pi \setminus \Pi^{\#}} [\nabla_{\lambda} h_l(\lambda^t, \gamma^t, n_r)]_{\pi} \right\}. \end{aligned}$$

The first piece could be found directly since $\text{supp}(\lambda^t) \cup \text{supp}(\lambda^{l-1}) \leq t + T_{l-1}$. The second piece could be computed with $(t + T_{l-1})|\mathcal{D}|$ calls to a constrained cost-sensitive classification oracle, stated in Lemma D.3 below. \square

Lemma D.3. For any set $B_t \subset \Pi$, we can compute $\text{argmax}_{\pi \in \Pi \setminus B_t} [\nabla_{\lambda} h_l(\lambda^t, \gamma^t, n_r)]_{\pi}$ using $|B_t| \cdot |\mathcal{D}|$ calls to a constrained cost-sensitive classification oracle defined in Definition 2.3.

Proof. Algorithm 5 below shows that we could compute this argmax via the C-AMO oracle. First, by construction of the algorithm, we have that $\pi_e \notin B_t$, so $\pi_e \in \Pi \setminus B_t$. It remains to show that π_e achieves the maximum. We prove this via contradiction. Assume that there is some other $\pi' \neq \pi_e$ that satisfies $\pi' \notin B_t$ and $\nabla_\lambda[h_l(\lambda, \gamma, n)]_{\pi'} > \nabla_\lambda[h_l(\lambda, \gamma, n)]_{\pi_e}$. By construction of our algorithm, we know that $\nabla_\lambda[h_l(\lambda, \gamma, n)]_{\pi_k}$ is non-increasing in k . We find the largest $0 \leq j \leq i-1$ such that

$$\nabla_\lambda[h_l(\lambda, \gamma, n)]_{\pi_{j+1}} \leq \nabla_\lambda[h_l(\lambda, \gamma, n)]_{\pi'} \leq \nabla_\lambda[h_l(\lambda, \gamma, n)]_{\pi_j}.$$

First, since j is the largest, we have $\nabla_\lambda[h_l(\lambda, \gamma, n)]_{\pi_{j+1}} < \nabla_\lambda[h_l(\lambda, \gamma, n)]_{\pi'}$, i.e. the first inequality is strict. By assumption that $\pi' \notin B_t$ and $\pi' \neq \pi_e$, we have $\pi' \neq \pi_k, \forall 0 \leq k \leq i$. So $\exists c_0 \in \mathcal{D}$ such that $\pi'(c_0) \neq \pi_j(c_0)$. Then we get a contradiction since in iteration j , at line 6 we should return π'_{c_0} instead of π_{j+1} . Therefore, there does not exist such π' and π_e achieves the maximum. \square

Algorithm 5 Constrained cost-sensitive classification

Input: policy set Π , set of policies to avoid B_t , objective function h_l , context history \mathcal{D} , tolerance ϵ

1: $\pi_0 = \operatorname{argmax}_{\pi \in \Pi} [\nabla_\lambda h_l(\lambda, \gamma, n)]_\pi, i = 0$

2: **while** $\pi_i \in B_t$ **do**

3: **for** $c \in \mathcal{D}$ **do**

4: compute $\pi'_c = \operatorname{argmax}_{\substack{\pi \in \Pi \\ \pi(c) \neq \pi_i(c)}} [\nabla_\lambda h_l(\lambda, \gamma, n)]_\pi$ s.t. $[\nabla_\lambda h_l(\lambda, \gamma, n)]_\pi \leq [\nabla_\lambda h_l(\lambda, \gamma, n)]_{\pi_i}$

5: **end for**

6: $\pi_{i+1} = \operatorname{argmax}_{c \in \mathcal{D}} [\nabla_\lambda h_l(\lambda, \gamma, n)]_{\pi'_c}$

7: $i = i + 1$

8: **end while**

9: $\pi_e = \pi_i$

Output: π_e

Lemma D.4. We can compute equation (6) with $T_l|\mathcal{D}|$ calls to a constrained argmax oracle.

Proof. We follow the proof technique in Lemma D.2 and break the argmin into two pieces with $\pi \in \operatorname{supp}(\lambda^l)$ and $\pi \in \Pi \setminus \operatorname{supp}(\lambda^l)$. We only show how to compute the second piece as the first piece could be compute directly. We know that $\widehat{\Delta}_i^{\gamma^l}(\pi, \widehat{\pi}_{l-1}) = \sum_{i=1}^{n_l} \frac{r_i}{p_{c_i, a_i}^{(\ell)} + [\gamma^l]_\pi} (\mathbf{1}\{\widehat{\pi}_{l-1}(c_i) = a_i\} - \mathbf{1}\{\pi(c_i) = a_i\})$. Then, similar to proof of Lemma D.2, let $\gamma_\pi = \gamma_0$ for all $\pi \in \Pi \setminus \operatorname{supp}(\lambda^l)$, we have

$$\begin{aligned} & \operatorname{argmin}_{\pi \in \Pi \setminus \operatorname{supp}(\lambda^l)} \widehat{\Delta}_i^{\gamma^l}(\pi, \widehat{\pi}_{l-1}) + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\frac{[\gamma^l]_\pi}{p_{c,a}^{(\ell)}} + \frac{[\gamma^l]_\pi}{s_{a',c}} \right) \mathbf{1}\{\widehat{\pi}_{l-1}(c) \neq \pi(c)\} \right] + \frac{\log(1/\delta_l)}{[\gamma^l]_\pi n_l} \\ &= \operatorname{argmin}_{\pi \in \Pi \setminus \operatorname{supp}(\lambda^l)} \sum_{i=1}^{n_l} \frac{r_i}{p_{c_i, a_i}^{(\ell)} + [\gamma^l]_\pi} (\mathbf{1}\{\widehat{\pi}_{l-1}(c_i) = a_i\} - \mathbf{1}\{\pi(c_i) = a_i\}) \\ & \quad + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\frac{[\gamma^l]_\pi}{p_{c,a}^{(\ell)}} + \frac{[\gamma^l]_\pi}{p_{c,a'}^{(\ell)}} \right) \mathbf{1}\{\widehat{\pi}_{l-1}(c) \neq \pi(c)\} \right] \\ &= \operatorname{argmin}_{\pi \in \Pi \setminus \operatorname{supp}(\lambda^l)} \sum_{i=1}^{n_l} \frac{r_i}{p_{c_i, a_i}^{(\ell)} + \gamma_0} \mathbf{1}\{\pi(c_i) = a_i\} \\ & \quad + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\frac{\gamma_0}{p_{c,a}^{(\ell)}} + \frac{\gamma_0}{p_{c,a'}^{(\ell)}} \right) \mathbf{1}\{\widehat{\pi}_{l-1}(c) \neq \pi(c)\} \right] \\ &= \operatorname{argmin}_{\pi \in \Pi \setminus \operatorname{supp}(\lambda^l)} \sum_{i=1}^{n_l} \frac{r_i}{p_{c_i, a_i}^{(\ell)} + \gamma_0} \mathbf{1}\{\pi(c_i) = a_i\} \\ & \quad - \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\frac{\gamma_0}{p_{c,a}^{(\ell)}} + \frac{\gamma_0}{p_{c,a'}^{(\ell)}} \right) \mathbf{1}\{\widehat{\pi}_{l-1}(c) \neq \pi(c)\} \right] \end{aligned}$$

which is a cost-sensitive classification problem with cost vector

$$L_i(a) = \begin{cases} \frac{r_i}{p_{c_i, a_i}^{(\ell)} + \gamma_0} \mathbf{1}\{a = a_i\} & \text{for } i = 1, \dots, n_l \\ - \left(\frac{\gamma_0}{p_{c_i, a}^{(\ell)}} + \frac{\gamma_0}{p_{c_i, \hat{\pi}_{l-1}(c_i)}^{(\ell)}} \right) \mathbf{1}\{a \neq \hat{\pi}_{l-1}(c_i)\} & \text{for } i = n_l + 1, \dots, n_l + |\mathcal{D}|. \end{cases}$$

□

D.2 Quantify the offline data

We first prove a general result for an empirical process bound of the difference of the expectation and the truth in Lemma D.5.

Lemma D.5. *Let $m = |\mathcal{D}|$ and define some set $\mathcal{K} \subset \gamma_{\max} \Delta_{\Pi}$. Consider some function $u : \mathcal{C} \times \mathcal{K} \rightarrow \mathbb{R}$ with $c, \kappa \mapsto u(c, \kappa)$ and define $\mathcal{F} \triangleq \{c \mapsto u(c, \kappa) : \kappa \in \mathcal{K}\}$. If*

1. *u satisfies that for any $c \in \mathcal{C}$ and $\kappa \in \mathcal{K}$, $u(c, \kappa) \in [0, b]$ where $b < \infty$ is a uniform upper bound;*
2. *there exists $L < \infty$ such that $\|u(\cdot, \kappa_1) - u(\cdot, \kappa_2)\|_{\mathcal{F}} \leq L \|\kappa_1 - \kappa_2\|_1$.*

Then, with probability at least $1 - \delta$,

$$\sup_{\kappa \in \mathcal{K}} |\mathbb{E}_{c \sim \nu_{\mathcal{D}}} [u(c, \kappa)] - \mathbb{E}_{c \sim \nu} [u(c, \kappa)]| \leq \sqrt{\frac{b^2}{2m} \log \left(\frac{2}{\delta} \right)} + \frac{16}{\sqrt{m}} L \gamma_{\max} \sqrt{2k \log(3e|\Pi|/k)}.$$

Proof. By the bounded condition on u we have $\{\mathbb{E}_{c \sim \nu_{\mathcal{D}}} [u(c, \kappa)] : \kappa \in \mathcal{K}\}$ satisfies the bounded difference property with parameter b . Then we use McDiarmid's inequality to get with probability at least $1 - \delta$,

$$\begin{aligned} & \sup_{\kappa \in \mathcal{K}} |\mathbb{E}_{c \sim \nu_{\mathcal{D}}} [u(c, \kappa)] - \mathbb{E}_{c \sim \nu} [u(c, \kappa)]| \\ & \leq \sqrt{\frac{b^2}{2m} \log \left(\frac{2}{\delta} \right)} + \mathbb{E} \left[\sup_{\kappa \in \mathcal{K}} |\mathbb{E}_{c \sim \nu_{\mathcal{D}}} [u(c, \kappa)] - \mathbb{E}_{c \sim \nu} [u(c, \kappa)]| \right]. \end{aligned}$$

Also, note that by definition of \mathcal{F} and classical results on entropy integral [37],

$$\mathbb{E} \left[\sup_{\kappa \in \mathcal{K}} |\mathbb{E}_{c \sim \nu_{\mathcal{D}}} [u(c, \kappa)] - \mathbb{E}_{c \sim \nu} [u(c, \kappa)]| \right] \leq \frac{8}{\sqrt{n}} \sup_Q \int_0^{\infty} \sqrt{\log N(\mathcal{F}, L_2(Q), \epsilon) d\epsilon},$$

where $N(\mathcal{F}, L_2(Q), \epsilon)$ is the covering number. By condition 2 and property of covering numbers,

$$\sup_Q N(\mathcal{F}, L_2(Q), \epsilon) \leq N(\mathcal{F}, \|\cdot\|_{\mathcal{F}}, \epsilon) \leq N(\mathcal{K}, \|\cdot\|_1, \epsilon/L).$$

Denote B_1^k as the l_1 ball with dimension k . We know that for $\epsilon \leq 1$, $N(B_1^k, \|\cdot\|_1, \epsilon) \leq \left(\frac{3}{\epsilon}\right)^k$. Since $\mathcal{K} \subset \gamma_{\max} \Delta_{\Pi}^{(k)} \subset \gamma_{\max} B_1^k$, and there are $\binom{\Pi}{k}$ ways to choose such a support $\gamma_{\max} B_1^k$, by union bound over k -dimensional subspaces we have

$$\begin{aligned} N(\mathcal{K}, \|\cdot\|_1, \epsilon/L) & \leq \binom{\Pi}{k} N(\gamma_{\max} B_1^k, \|\cdot\|_1, \epsilon/L) \\ & \leq \binom{\Pi}{k} N(B_1^k, \|\cdot\|_1, \epsilon/(L\gamma_{\max})) \\ & \leq \left(\frac{e|\Pi|}{k} \right)^k \left(\frac{3L\gamma_{\max}}{\epsilon} \right)^k \leq \left(\frac{3L\gamma_{\max} e|\Pi|}{\epsilon k} \right)^k. \end{aligned}$$

Therefore,

$$\begin{aligned}
\sup_Q \int_0^\infty \sqrt{\log N(\mathcal{F}, L_2(Q), \epsilon)} d\epsilon &\leq \int_0^\infty \sqrt{\log N(\mathcal{K}, \|\cdot\|_1, \epsilon/L)} d\epsilon \\
&\leq \int_0^{L\gamma_{\max}} \sqrt{k \log \left(\frac{3L\gamma_{\max}e|\Pi|}{\epsilon k} \right)} d\epsilon \\
&= L\gamma_{\max} \int_0^1 \sqrt{k \log \left(\frac{3e|\Pi|}{\epsilon k} \right)} d\epsilon \\
&\leq L\gamma_{\max} \sqrt{\int_0^1 k \log \left(\frac{3e|\Pi|}{\epsilon k} \right) d\epsilon} \\
&\leq L\gamma_{\max} \sqrt{2k \log(3e|\Pi|/k)}.
\end{aligned}$$

Combining all results yields

$$\begin{aligned}
\mathbb{E} \left[\sup_{\lambda \in \Delta_\Pi^{(k)}} |\mathbb{E}_{c \sim \nu_{\mathcal{D}}} [u(c, \kappa)] - \mathbb{E}_{c \sim \nu} [u(c, \kappa)]| \right] &\leq \frac{16}{\sqrt{m}} \sup_Q \int_0^\infty \sqrt{\log N(\mathcal{F}, L_2(Q), \epsilon)} d\epsilon \\
&\leq \frac{16}{\sqrt{m}} L\gamma_{\max} \sqrt{2k \log(3e|\Pi|/k)}.
\end{aligned}$$

Therefore, our result follows. \square

Then, we take two special kind of $u(c, \kappa)$, and get the bounds for our estimate of the expectation over ν with the offline history \mathcal{D} .

Lemma D.6. *Let $m = |\mathcal{D}|$. Then, with probability at least $1 - \delta$, we have*

$$\begin{aligned}
&\sup_{(\lambda, \gamma) \in \gamma_{\max} \Delta_\Pi^{(k)}} \left| \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right)^2 \right] - \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right)^2 \right] \right| \\
&\leq \sqrt{\frac{|\mathcal{A}|^4 \gamma_{\max}^2 (1 + \eta_l)^2}{2m} \log \left(\frac{2}{\delta} \right)} + \frac{16}{\sqrt{m}} |\mathcal{A}|^2 \gamma_{\max} \sqrt{\frac{2k(1 + \eta_l) \gamma_{\max}}{\eta_l \gamma_{\min}} \log \left(\frac{3e|\Pi|}{k} \right)}.
\end{aligned}$$

Proof. Define $\kappa \in \mathcal{K}$ such that $\kappa_\pi = \lambda_\pi \gamma_\pi$. Then, $\mathcal{K} \subset \gamma_{\max} \Delta_\Pi$ since $\sum_{\pi \in \Pi} \kappa_\pi = \sum_{\pi \in \Pi} \lambda_\pi \gamma_\pi \leq \gamma_{\max}$. Then, let $u(c, \kappa) = \left(\sum_{a \in \mathcal{A}} \sqrt{\kappa^\top (t_a^{(c)} + \eta_l)} \right)^2$. We aim to use the result of Lemma D.5 to get our bound. First, since for any $\kappa \in \mathcal{K}$ and any $c \in \mathcal{D}$, $u(c, \kappa) \in [|\mathcal{A}|^2 \gamma_{\min} \eta_l, |\mathcal{A}|^2 (1 + \eta_l) \gamma_{\max}]$,

so condition 1 is satisfied. Also, note that $u(c, \kappa)$ is Lipschitz in κ , i.e.

$$\begin{aligned}
& \|u(\cdot, \kappa_1) - u(\cdot, \kappa_2)\|_{\mathcal{F}} \\
&= \sup_{c \in \mathcal{C}} |u(c, \kappa_1) - u(c, \kappa_2)| \\
&= \sup_{c \in \mathcal{C}} \left| \left(\sum_{a \in \mathcal{A}} \sqrt{\kappa_1^\top (t_a^{(c)} + \eta_l)} \right)^2 - \left(\sum_{a \in \mathcal{A}} \sqrt{\kappa_2^\top (t_a^{(c)} + \eta_l)} \right)^2 \right| \\
&\leq \sup_{c \in \mathcal{C}} \left| \left(\sum_{a \in \mathcal{A}} \sqrt{\kappa_1^\top (t_a^{(c)} + \eta_l)} + \sqrt{\kappa_2^\top (t_a^{(c)} + \eta_l)} \right) \left(\sum_{a \in \mathcal{A}} \sqrt{\kappa_1^\top (t_a^{(c)} + \eta_l)} - \sqrt{\kappa_2^\top (t_a^{(c)} + \eta_l)} \right) \right| \\
&= \sup_{c \in \mathcal{C}} \left(\sum_{a \in \mathcal{A}} \sqrt{\kappa_1^\top (t_a^{(c)} + \eta_l)} + \sqrt{\kappa_2^\top (t_a^{(c)} + \eta_l)} \right) \left(\sum_{a \in \mathcal{A}} \frac{|(\kappa_1 - \kappa_2)^\top t_a^{(c)}|}{\sqrt{\kappa_1^\top (t_a^{(c)} + \eta_l)} + \sqrt{\kappa_2^\top (t_a^{(c)} + \eta_l)}} \right) \\
&\leq \sup_{c \in \mathcal{C}} \left(\sum_{a \in \mathcal{A}} \sqrt{\kappa_1^\top (t_a^{(c)} + \eta_l)} + \sqrt{\kappa_2^\top (t_a^{(c)} + \eta_l)} \right) \left(\sum_{a \in \mathcal{A}} \frac{\|\kappa_1 - \kappa_2\|_1}{\sqrt{\kappa_1^\top (t_a^{(c)} + \eta_l)} + \sqrt{\kappa_2^\top (t_a^{(c)} + \eta_l)}} \right) \\
&\leq |\mathcal{A}|^2 \sqrt{\frac{(1 + \eta_l) \gamma_{\max}}{\eta_l \gamma_{\min}}} \|\kappa_1 - \kappa_2\|_1.
\end{aligned}$$

Therefore, condition 2 is satisfied with $L = |\mathcal{A}|^2 \sqrt{\frac{(1 + \eta_l) \gamma_{\max}}{\eta_l \gamma_{\min}}}$. Plugging in the result in Lemma D.5, we get

$$\begin{aligned}
& \sup_{\lambda \in \Delta_{\Pi}^{(k)}} |\mathbb{E}_{c \sim \nu_{\mathcal{D}}} [u(c, \kappa)] - \mathbb{E}_{c \sim \nu} [u(c, \kappa)]| \\
&\leq \sqrt{\frac{|\mathcal{A}|^4 \gamma_{\max}^2 (1 + \eta_l)^2 \log\left(\frac{2}{\delta}\right)}{2m}} + \frac{16}{\sqrt{m}} |\mathcal{A}|^2 \gamma_{\max} \sqrt{\frac{2k(1 + \eta_l) \gamma_{\max}}{\eta_l \gamma_{\min}} \log\left(\frac{3e|\Pi|}{k}\right)}.
\end{aligned}$$

□

Lemma D.7. For any $\pi \in \Pi$, with probability at least $1 - \delta$,

$$\begin{aligned}
& \sup_{(\lambda, \gamma) \in \gamma_{\max} \Delta_{\Pi}} \left| \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_{a \in \mathcal{A}} \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)}} (\gamma_{\pi} [t_a^{(c)}]_{\pi}) \right] \right. \\
& \quad \left. - \mathbb{E}_{c \sim \nu} \left[\sum_{a \in \mathcal{A}} \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)}} (\gamma_{\pi} [t_a^{(c)}]_{\pi}) \right] \right| \\
&\leq \gamma_{\max} \left(\sqrt{\frac{|\mathcal{A}|^4 (1 + \eta_l) \gamma_{\max}}{2\eta_l \gamma_{\min} m}} \log\left(\frac{2}{\delta}\right) + \frac{8|\mathcal{A}|^2 \gamma_{\max}}{\sqrt{m} (\eta_l \gamma_{\min})^{3/2}} \sqrt{2k \log(3e|\Pi|/k)} \right).
\end{aligned}$$

Proof. First, note that

$$\frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)}} (\gamma_{\pi} [t_a^{(c)}]_{\pi}) \leq \gamma_{\max} \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)}} [t_a^{(c)}]_{\pi}.$$

Then, we define $u(c, \kappa) = \sum_{a \in \mathcal{A}} \frac{\sum_{a' \in \mathcal{A}} \sqrt{\kappa^\top(t_{a'}^{(c)} + \eta)}}{\sqrt{\kappa^\top(t_a^{(c)} + \eta)}} [t_a^{(c)}]_\pi$. First, note that for any $c \in \mathcal{C}$ and $\kappa \in \mathcal{K}$, $u(c, \kappa) \in \left[0, |\mathcal{A}|^2 \frac{\sqrt{(1+\eta)\gamma_{\max}}}{\sqrt{\eta\gamma_{\min}}}\right]$, so condition 1 in Lemma D.5 is satisfied. Also,

$$\begin{aligned}
& \|u(c, \kappa_1) - u(c, \kappa_2)\|_{\mathcal{F}} = \sup_{c \in \mathcal{C}} |u(c, \kappa_1) - u(c, \kappa_2)| \tag{10} \\
&= \sup_{c \in \mathcal{C}} \left| \sum_{a \in \mathcal{A}} \frac{\sum_{a' \in \mathcal{A}} \sqrt{\kappa_1^\top(t_{a'}^{(c)} + \eta)}}{\sqrt{\kappa_1^\top(t_a^{(c)} + \eta)}} [t_a^{(c)}]_\pi - \sum_{a \in \mathcal{A}} \frac{\sum_{a' \in \mathcal{A}} \sqrt{\kappa_2^\top(t_{a'}^{(c)} + \eta)}}{\sqrt{\kappa_2^\top(t_a^{(c)} + \eta)}} [t_a^{(c)}]_\pi \right| \\
&= \sup_{c \in \mathcal{C}} \left| \sum_{a \in \mathcal{A}} \left[\frac{\sum_{a' \in \mathcal{A}} \sqrt{\kappa_1^\top(t_{a'}^{(c)} + \eta)} \sqrt{\kappa_2^\top(t_a^{(c)} + \eta)} - \sqrt{\kappa_2^\top(t_a^{(c)} + \eta)} \sqrt{\kappa_1^\top(t_a^{(c)} + \eta)}}{\sqrt{\kappa_1^\top(t_a^{(c)} + \eta)} \sqrt{\kappa_2^\top(t_a^{(c)} + \eta)}} [t_a^{(c)}]_\pi \right] \right| \\
&\leq \sup_{c \in \mathcal{C}} \sum_{a \in \mathcal{A}} \left[\frac{\sum_{a' \in \mathcal{A}} \left| \sqrt{\kappa_1^\top(t_{a'}^{(c)} + \eta)} \sqrt{\kappa_2^\top(t_a^{(c)} + \eta)} - \sqrt{\kappa_2^\top(t_{a'}^{(c)} + \eta)} \sqrt{\kappa_1^\top(t_a^{(c)} + \eta)} \right|}{\sqrt{\kappa_1^\top(t_a^{(c)} + \eta)} \sqrt{\kappa_2^\top(t_a^{(c)} + \eta)}} \right]. \tag{11}
\end{aligned}$$

Note that by triangular inequality

$$\begin{aligned}
& \left| \sqrt{\kappa_2^\top(t_a^{(c)} + \eta)} \sqrt{\kappa_1^\top(t_{a'}^{(c)} + \eta)} - \sqrt{\kappa_1^\top(t_a^{(c)} + \eta)} \sqrt{\kappa_2^\top(t_{a'}^{(c)} + \eta)} \right| \\
&\leq \left| \sqrt{\kappa_2^\top(t_a^{(c)} + \eta)} - \sqrt{\kappa_1^\top(t_a^{(c)} + \eta)} \right| \sqrt{\kappa_1^\top(t_{a'}^{(c)} + \eta)} \\
&\quad + \sqrt{\kappa_1^\top(t_a^{(c)} + \eta)} \left| \sqrt{\kappa_1^\top(t_{a'}^{(c)} + \eta)} - \sqrt{\kappa_2^\top(t_{a'}^{(c)} + \eta)} \right|.
\end{aligned}$$

Also note that

$$\begin{aligned}
\left| \sqrt{\kappa_2^\top(t_a^{(c)} + \eta)} - \sqrt{\kappa_1^\top(t_a^{(c)} + \eta)} \right| &= \frac{\left| \sum_{\pi \in \Pi} ([\kappa_1]_\pi - [\kappa_2]_\pi) (t_a^{(c)} + \eta)_\pi \right|}{\sqrt{\kappa_2^\top(t_a^{(c)} + \eta)} + \sqrt{\kappa_1^\top(t_a^{(c)} + \eta)}} \\
&\leq \frac{1}{2\sqrt{\eta\gamma_{\min}}} \|\kappa_2 - \kappa_1\|_1.
\end{aligned}$$

Therefore, (11) is bounded by $|\mathcal{A}|^2 \frac{1}{\eta\gamma_{\min}} \frac{1}{2\sqrt{\eta\gamma_{\min}}} \|\kappa_2 - \kappa_1\|_1$, so condition 2 is satisfied with $L = \frac{|\mathcal{A}|^2}{2(\eta\gamma_{\min})^{3/2}}$. Then, by Lemma D.5, with probability at least $1 - \delta$,

$$\begin{aligned}
& \sup_{(\lambda, \gamma) \in \gamma_{\max} \Delta_\Pi} \left| \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top(t_{a'}^{(c)} + \eta)}}{\sqrt{(\lambda \odot \gamma)^\top(t_a^{(c)} + \eta)}} (\gamma_\pi [t_a^{(c)}]_\pi) \right] \right. \\
& \quad \left. - \mathbb{E}_{c \sim \nu} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top(t_{a'}^{(c)} + \eta)}}{\sqrt{(\lambda \odot \gamma)^\top(t_a^{(c)} + \eta)}} (\gamma_\pi [t_a^{(c)}]_\pi) \right] \right| \\
&\leq \gamma_{\max} \left(\sqrt{\frac{|\mathcal{A}|^4 (1+\eta)\gamma_{\max}}{2\eta\gamma_{\min} m}} \log\left(\frac{2}{\delta}\right) + \frac{8|\mathcal{A}|^2\gamma_{\max}}{\sqrt{m}(\eta\gamma_{\min})^{3/2}} \sqrt{2k \log(3e|\Pi|/k)} \right).
\end{aligned}$$

□

E Proof of Theorem 3.3

We first write down Algorithm 3 in full detail in Algorithm 6. We aim to show that Algorithm 6 achieves the sample complexity lower bound. The two big goals here is to show that $\hat{\pi}_l \in S_l$ for all l , which shows that we get the optimal policy, and n_l achieves the sample complexity lower bound.

Algorithm 6 Full CODA Algorithm

Input: policies $\Pi = \{\pi : \mathcal{C} \rightarrow \mathcal{A}\}_\pi$, feature map $\phi : \mathcal{C} \times \mathcal{A} \rightarrow \mathbb{R}^d$, $\delta \in (0, 1)$, historical data $\mathcal{D} = \{\nu_s\}_s$

- 1: initiate $\hat{\pi}_0 \in \Pi$ arbitrarily, $\lambda_0 = \mathbf{e}_{\hat{\pi}_0}$, $\hat{\Delta}_0(\pi)$, γ_0 appropriately
- 2: **for** $l = 1, 2, \dots$ **do**
- 3: $\epsilon_l = 2^{-l}$, $\eta_l = C_1 \epsilon_l^2 |\mathcal{A}|^{-4}$, $\delta_l = \delta / (l^2 |\Pi|^2)$, T_l appropriately
- 4: $t_a^{(c)}(\pi') = \{\mathbf{1}\{\pi(c) = a, \pi'(c) \neq a\} + \mathbf{1}\{\pi(c) \neq a, \pi'(c) = a\}\}_{\pi \in \Pi} \in \mathbb{R}^\Pi$
- 5: Define $\gamma_{\min} := \frac{1}{3} \sqrt{\frac{\eta_l \log(1/\delta_l)}{n}}$, $\gamma_{\max} := \sqrt{\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n}}$
- 6: Define

$$h_l(\lambda, \gamma, n) = \sum_{\pi \in \Pi} \lambda_\pi \left(-\hat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \hat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_\pi n} \right) + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)}(\hat{\pi}_{l-1}) + \eta_l)} \right)^2 \right]. \quad (12)$$

- 7: Let $\lambda^l, \gamma^l, n_l = \text{FW-GD}(\Pi, |\mathcal{A}|, \hat{\pi}_{l-1}, \eta_l, T_l, \epsilon_l, \gamma_{\min}, \gamma_{\max})$. These are the solutions to

$$n_\ell := \min\{n \in \mathbb{N} : \max_{\lambda \in \Delta_\Pi} \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^{|\Pi|}} h_l(\lambda, \gamma, n) \leq \epsilon_\ell\} \quad (13)$$

- 8: Receive contexts $c_1, c_2, \dots, c_{n_l} \sim \nu$.
- 9: For each $c_s, s = 1, 2, \dots, n_l$, pull arms $a_s \sim p_{c_s}^{(\ell)}$ where $p_{c_s, a_s}^{(\ell)} \propto \sqrt{(\lambda^l \odot \gamma^l)^\top (t_{a_s}^{(c_s)}(\hat{\pi}_{l-1}) + \eta_l)}$, and observe rewards r_s where $t_{a_s}^{(c_s)}(\hat{\pi}_{l-1}) \in \mathbb{R}^{|\Pi|}$
- 10: For each $\pi \in \Pi$, define the IPS estimator

$$\hat{\Delta}_l^{\gamma^l}(\pi, \hat{\pi}_{l-1}) = \sum_{s=1}^{n_l} \frac{r_s}{p_{c_s, a_s}^{(\ell)} + [\gamma^l]_\pi} (\mathbf{1}\{\hat{\pi}_{l-1}(c_s) = a_s\} - \mathbf{1}\{\pi(c_s) = a_s\})$$

- 11: set

$$\hat{\pi}_l = \arg \min_{\pi \in \Pi} \hat{\Delta}_l^{\gamma^l}(\pi, \hat{\pi}_{l-1}) + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\frac{[\gamma^l]_\pi}{p_{c, a}^{(\ell)}} + \frac{[\gamma^l]_\pi}{p_{c, a'}}^{(\ell)} \right) \mathbf{1}\{\hat{\pi}_{l-1}(c) \neq \pi(c)\} \right] + \frac{\log(1/\delta_l)}{[\gamma^l]_\pi n_l}. \quad (14)$$

- 12: **end for**
- Output:** $\hat{\pi}_l$
-

Theorem E.1. *With probability at least $1 - \delta$, Algorithm 6 returns a policy $\hat{\pi}$ satisfying $V(\pi_*) - V(\hat{\pi}_\ell) \leq \epsilon$ in a number of samples not exceeding $O(\rho_{*, \epsilon} \log(|\Pi| \log_2(1/\Delta_\epsilon)/\delta) \log_2(1/\Delta_\epsilon))$ where $\Delta_\epsilon := \max\{\epsilon, \min_{\pi \in \Pi} V(\pi_*) - V(\pi)\}$.*

Proof. We first define our key events. Recall

$$\hat{\Delta}_l^{\gamma^l}(\pi, \hat{\pi}_{l-1}) = \sum_{s=1}^{n_l} \frac{r_s}{p_{c_s, a_s}^{(\ell)} + [\gamma^l]_\pi} (\mathbf{1}\{\hat{\pi}_{l-1}(c_s) = a_s\} - \mathbf{1}\{\pi(c_s) = a_s\})$$

and $\Delta(\pi, \pi') = V(\pi') - V(\pi)$. Define $w(\lambda, \gamma) \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{C}|}$ with

$$[w(\lambda, \gamma)]_{a, c} := \nu_c \cdot p_{c, a} = \nu_c \cdot \frac{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)}(\hat{\pi}_{l-1}) + \eta_l)}}{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)}(\hat{\pi}_{l-1}) + \eta_l)}}.$$

Then define the events

$$\mathcal{E}_l := \bigcap_{\pi, \pi' \in \Pi} \left\{ \left| \hat{\Delta}_l^{\gamma^l}(\pi, \pi') - \Delta(\pi, \pi') \right| \leq 2[\gamma^l]_\pi \|\phi_\pi - \phi_{\pi'}\|_{A(w(\lambda^l, \gamma^l))}^2 + \frac{2 \log(1/\delta_l)}{[\gamma^l]_\pi n_l} \right\},$$

and the good event $\mathcal{E} = \bigcap_{l=1}^\infty \mathcal{E}_l$. Lemma E.3 shows that \mathcal{E} happens with probability at least $1 - \delta$, and Lemma E.7 shows that under this event \mathcal{E} ,

$$n_l \lesssim \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\pi_*} - \phi_\pi\|_{A(w)}^2 \log(1/\delta_l)}{\epsilon_l^2 + \Delta(\pi, \pi_*)^2}.$$

Therefore, the total number of samples is no more than

$$\begin{aligned}
& \sum_{l=1}^{\log_2(1/\Delta_\epsilon)} \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\pi_*} - \phi_\pi\|_{A(w)}^2 \log(l^2 |\Pi|^2 / \delta)}{\epsilon_l^2 + \Delta(\pi, \pi_*)^2} \\
& \stackrel{(i)}{\leq} \sum_{l=1}^{\log_2(1/\Delta_\epsilon)} \min_{w \in \Omega} \max_{\pi \in \Pi \setminus \pi_*} \frac{2 \|\phi_{\pi_*} - \phi_\pi\|_{A(w)}^2 \log(l^2 |\Pi|^2 / \delta)}{\epsilon_l^2 + \Delta(\pi, \pi_*)^2} \\
& \stackrel{(ii)}{\leq} \sum_{l=1}^{\log_2(1/\Delta_\epsilon)} \min_{p^{(c)} \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \max_{\pi \in \Pi \setminus \pi_*} \frac{\mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{\pi_*(c)}} + \frac{1}{p_{\pi(c)}} \right) \mathbf{1}\{\pi_*(c) \neq \pi(c)\} \right] \log(l^2 |\Pi|^2 / \delta)}{\Delta(\pi, \pi_*)^2 + \epsilon_l^2} \\
& \lesssim \rho_{\star, \epsilon}(\Pi, \nu) \log(\log_2(1/\Delta_\epsilon) |\Pi| / \delta) \log_2(1/\Delta_\epsilon).
\end{aligned}$$

where (i) follows from the fact that π_* gives zero for the RHS, and (ii) follows from Lemma H.1. \square

In what follows, we will fill in the road map to the proof of Lemma E.3 and E.7. First, Lemma E.2 controls the estimation error of the gap and shows that $\mathbb{P}(\mathcal{E}_\ell) > 1 - \delta_\ell$, which leads to the high-probability of the good event \mathcal{E} (Lemma E.3). Lemma E.4 applies the duality machinery in Section G and controls the variance term. Lemma E.5 applies the result of Lemma E.4 and shows an upper bound for the difference between estimate gap and the true gap, which is a very similar result of Lemma C.3. Lemma E.6 is an important lemma showing the analytical solution of w given some λ and γ . With all of these results above, we get Lemma E.7 which gives the upper bound on the sample complexity.

Lemma E.2. *For any $l > 0$, $\pi, \pi' \in \Pi$, with probability at least $1 - \delta_l$,*

$$\left| \widehat{\Delta}_l^{\gamma^l}(\pi, \pi') - \Delta(\pi, \pi') \right| \leq 2[\gamma^l]_\pi \|\phi_\pi - \phi_{\pi'}\|_{A(w(\lambda^l, \gamma^l))}^2 + \frac{2 \log(1/\delta_l)}{[\gamma^l]_\pi n_l}.$$

Proof. Define

$$\widehat{V}_l^{\gamma^l}(\pi) := \sum_{s=1}^{n_l} \frac{r_s}{p_{c_s, a_s}^{(\ell)} + [\gamma^l]_\pi} \mathbf{1}\{\pi(c_s) = a_s\},$$

so that

$$\widehat{\Delta}_l^{\gamma^l}(\pi, \pi') = \widehat{V}_l^{\gamma^l}(\pi') - \widehat{V}_l^{\gamma^l}(\pi).$$

First, note that below.

$$\begin{aligned}
V(\pi) &= \mathbb{E}_{c \sim \nu} [r(c, \pi(c))] \\
&= \mathbb{E}_{c \sim \nu} \left[\mathbb{E}_{a \sim p_c^{(\ell)}} \left[r(c, a) \frac{\mathbf{1}\{\pi(c) = a\}}{p_{c, a}^{(\ell)}} \middle| c \right] \right] = \mathbb{E} \left[\frac{1}{t} \sum_{s=1}^t \frac{r_s}{p_{c_s, a_s}^{(\ell)}} \mathbf{1}\{\pi(c_s) = a_s\} \right].
\end{aligned}$$

Therefore,

$$\begin{aligned}
& \left| \mathbb{E} \left[\widehat{V}_l^{\gamma^l}(\pi) - \widehat{V}_l^{\gamma^l}(\pi') \right] - [V(\pi) - V(\pi')] \right| \\
& \leq \left| \mathbb{E} \left[\frac{1}{n_l} \sum_{s=1}^{n_l} \left(\frac{1}{p_{c_s, a_s}^{(\ell)} + [\gamma^l]_\pi} - \frac{1}{p_{c_s, a_s}^{(\ell)}} \right) (\mathbf{1}\{\pi(c_s) = a_s\} - \mathbf{1}\{\pi'(c_s) = a_s\}) \right] \right| \\
& = \left| \mathbb{E} \left[\frac{1}{n_l} \sum_{s=1}^{n_l} \frac{-[\gamma^l]_\pi}{p_{c_s, a_s}^{(\ell)} (p_{c_s, a_s}^{(\ell)} + [\gamma^l]_\pi)} (\mathbf{1}\{\pi(c_s) = a_s\} - \mathbf{1}\{\pi'(c_s) = a_s\}) \right] \right| \\
& \leq \mathbb{E} \left[\frac{1}{n_l} \sum_{s=1}^{n_l} \frac{[\gamma^l]_\pi (\mathbf{1}\{\pi'(c_s) = a_s, \pi(c_s) \neq a_s\} + \mathbf{1}\{\pi'(c_s) \neq a_s, \pi(c_s) = a_s\})}{p_{c_s, a_s}^{(\ell)} (p_{c_s, a_s}^{(\ell)} + [\gamma^l]_\pi)} \right] \\
& = [\gamma^l]_\pi \mathbb{E} \left[\frac{1}{p_{c, a}^{(\ell)} (p_{c, a}^{(\ell)} + [\gamma^l]_\pi)} \nu_c^2 [\phi_\pi - \phi_{\pi'}]_{a, c}^2 \right] \\
& = [\gamma^l]_\pi \sum_{c \in \mathcal{C}} \nu_c \sum_{a \in \mathcal{A}} \frac{p_{c, a}^{(\ell)}}{p_{c, a}^{(\ell)} \nu_c^2 (p_{c, a}^{(\ell)} + [\gamma^l]_\pi)} [\phi_\pi - \phi_{\pi'}]_{a, c}^2 \\
& \leq [\gamma^l]_\pi \|\phi_\pi - \phi_{\pi'}\|_{A(w(\lambda^l, \gamma^l))}^2
\end{aligned}$$

where the last inequality follows since $\nu_c p_{c, a}^{(\ell)} = [w(\lambda^l, \gamma^l)]_{a, c}$. Meanwhile, note that

$$\frac{r_s}{p_{c_s, a_s}^{(\ell)} + [\gamma^l]_\pi} (\mathbf{1}\{\pi(c_s) = a_s\} - \mathbf{1}\{\pi'(c_s) = a_s\}) \leq \frac{1}{[\gamma^l]_\pi},$$

and

$$\begin{aligned}
& \mathbb{E} \left[\left(\frac{r_s}{p_{c_s, a_s}^{(\ell)} + [\gamma^l]_\pi} (\mathbf{1}\{\pi(c_s) = a_s\} - \mathbf{1}\{\pi'(c_s) = a_s\}) \right)^2 \right] \\
& \leq \mathbb{E} \left[\frac{1}{(p_{c_s, a_s}^{(\ell)} + [\gamma^l]_\pi)^2} (\mathbf{1}\{\pi(c_s) = a_s\} - \mathbf{1}\{\pi'(c_s) = a_s\})^2 \right] \\
& = \mathbb{E} \left[\frac{1}{(p_{c_s, a_s}^{(\ell)} + [\gamma^l]_\pi)^2 \nu_c^2} [\phi_\pi - \phi_{\pi'}]_{a, c}^2 \right] \\
& \leq \|\phi_\pi - \phi_{\pi'}\|_{A(w(\lambda^l, \gamma^l))}^2
\end{aligned}$$

by a similar argument as before. Therefore, by Bernstein's inequality, we have with probability at least $1 - \delta$,

$$\left| \widehat{V}_l^{\gamma^l}(\pi) - \widehat{V}_l^{\gamma^l}(\pi') - \mathbb{E} \left[\widehat{V}_l^{\gamma^l}(\pi) - \widehat{V}_l^{\gamma^l}(\pi') \right] \right| \leq \sqrt{\|\phi_\pi - \phi_{\pi'}\|_{A(w(\lambda^l, \gamma^l))}^2 \frac{2 \log(1/\delta)}{n_l}} + \frac{\log(1/\delta)}{[\gamma^l]_\pi n_l}.$$

Combining this with the deviation on expectation gives us

$$\begin{aligned}
& \left| \widehat{\Delta}_l^{\gamma^l}(\pi, \pi') - \Delta(\pi, \pi') \right| \\
& \leq [\gamma^l]_\pi \|\phi_\pi - \phi_{\pi'}\|_{A(w(\lambda^l, \gamma^l))}^2 + \sqrt{\|\phi_\pi - \phi_{\pi'}\|_{A(w(\lambda^l, \gamma^l))}^2 \frac{2 \log(1/\delta)}{n_l}} + \frac{2 \log(1/\delta)}{[\gamma^l]_\pi n_l} \\
& \leq 2[\gamma^l]_\pi \|\phi_\pi - \phi_{\pi'}\|_{A(w(\lambda^l, \gamma^l))}^2 + \frac{4 \log(1/\delta)}{[\gamma^l]_\pi n_l}.
\end{aligned}$$

□

Lemma E.3. $\mathbb{P}(\mathcal{E}) \geq 1 - \delta$.

Proof. By Lemma E.2 and a union bound over all policies, we have

$$\mathbb{P}(\mathcal{E}_l \mid \mathcal{E}_{l-1}, \dots, \mathcal{E}_1) \geq 1 - \frac{\delta}{l^2}.$$

Since $\mathcal{E} = \bigcap_{l=0}^{\infty} \mathcal{E}_l$,

$$\begin{aligned} \mathbb{P}(\mathcal{E}^c) &= \mathbb{P}((\bigcap_{l=0}^{\infty} \mathcal{E}_l)^c) = \mathbb{P}(\bigcup_{l=0}^{\infty} \mathcal{E}_l^c) = \mathbb{P}(\bigcup_{l=0}^{\infty} (\mathcal{E}_l^c \setminus (\bigcup_{j<l} \mathcal{E}_j^c))) \\ &\leq \sum_{l=0}^{\infty} \mathbb{P}(\mathcal{E}_l^c \setminus (\bigcup_{j<l} \mathcal{E}_j^c)) \leq \sum_{l=0}^{\infty} \mathbb{P}(\mathcal{E}_l^c \mid (\bigcap_{j<l} \mathcal{E}_j)) \leq \sum_{l=0}^{\infty} \frac{\delta}{l^2} \leq \delta. \end{aligned}$$

Therefore, $\mathbb{P}(\mathcal{E}) \geq 1 - \delta$. \square

Lemma E.4. Under \mathcal{E} , we have for any $\pi \in \Pi$,

$$[\gamma^l]_{\pi} \|\phi_{\pi} - \phi_{\hat{\pi}_{l-1}}\|_{A(w(\lambda^l, \gamma^l))}^2 + \frac{\log(1/\delta_l)}{[\gamma^l]_{\pi} n_l} \leq \frac{1}{6} \epsilon_l + \frac{1}{64} \widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \hat{\pi}_{l-1}).$$

Proof. We know that the choice of n_l ensures

$$h_l(\lambda^l, \gamma^l, n_l) \leq \epsilon_l.$$

Also, by Theorem G.1 we have

$$\frac{1}{3} \epsilon_l \geq \max_{\pi \in \Pi} \left(-\frac{1}{8} \widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \hat{\pi}_{l-1}) + 8[\gamma^l]_{\pi} \|\phi_{\pi} - \phi_{\hat{\pi}_{l-1}}\|_{A(w(\lambda^l, \gamma^l))}^2 + \frac{8 \log(1/\delta_l)}{[\gamma^l]_{\pi} n_l} \right) - h_l(\lambda^l, \gamma^l, n_l).$$

Combining the above two displays gives us

$$\begin{aligned} \epsilon_l &\geq h_l(\lambda^l, \gamma^l, n_l) \\ &\geq \max_{\pi \in \Pi} \left(-\frac{1}{8} \widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \hat{\pi}_{l-1}) + 8[\gamma^l]_{\pi} \|\phi_{\hat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w(\lambda^l, \gamma^l))}^2 + \frac{8 \log(1/\delta_l)}{[\gamma^l]_{\pi} n_l} \right) - \frac{1}{3} \epsilon_l. \end{aligned}$$

Therefore, for any $\pi \in \Pi$,

$$[\gamma^l]_{\pi} \|\phi_{\pi} - \phi_{\hat{\pi}_{l-1}}\|_{A(w(\lambda^l, \gamma^l))}^2 + \frac{\log(1/\delta_l)}{[\gamma^l]_{\pi} n_l} \leq \frac{1}{6} \epsilon_l + \frac{1}{64} \widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \hat{\pi}_{l-1}).$$

\square

Lemma E.5. Under \mathcal{E} , for all $l \in \mathbb{N}$, the following holds:

1. $|\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \hat{\pi}_{l-1}) - \Delta(\pi, \pi_*)| \leq 2\epsilon_{l-1} + \frac{1}{4} \Delta(\pi, \pi_*)$.
2. $\hat{\pi}_l \in S_l := \{\pi \in \Pi : \Delta(\pi, \pi_*) \leq \epsilon_l\}$.

Proof. We prove this by induction. First, in round $l = 0$, this holds since our rewards are bounded by 1. Then, assume that in round $l - 1$, we have $\hat{\pi}_{l-1} \in S_{l-1}$ and

$$|\widehat{\Delta}_{l-2}^{\gamma^{l-2}}(\pi, \hat{\pi}_{l-2}) - \Delta(\pi, \pi_*)| \leq 2\epsilon_{l-2} + \frac{1}{4} \Delta(\pi, \pi_*).$$

Then, on round l ,

$$\begin{aligned} &|\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \hat{\pi}_{l-1}) - \Delta(\pi, \pi_*)| \\ &= |\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \hat{\pi}_{l-1}) - \Delta(\pi, \hat{\pi}_{l-1}) - \Delta(\hat{\pi}_{l-1}, \pi_*)| \\ &\leq 2[\gamma^{l-1}]_{\pi} \|\phi_{\pi} - \phi_{\hat{\pi}_{l-1}}\|_{A(w(\lambda^{l-1}, \gamma^{l-1}))}^2 + \frac{2 \log(1/\delta_{l-1})}{[\gamma^{l-1}]_{\pi} n_{l-1}} + \epsilon_{l-1} \\ &\hspace{15em} \text{(from event } \mathcal{E} \text{ and inductive hypothesis)} \\ &\leq \frac{2}{3} \epsilon_l + \frac{1}{64} \widehat{\Delta}_{l-2}^{\gamma^{l-2}}(\pi, \hat{\pi}_{l-2}) + \frac{1}{64} \widehat{\Delta}_{l-2}^{\gamma^{l-2}}(\hat{\pi}_{l-1}, \hat{\pi}_{l-2}) + \epsilon_{l-1} \hspace{5em} \text{(from Lemma E.4)} \\ &\leq \frac{5}{3} \epsilon_{l-1} + \frac{1}{64} \left(2\epsilon_{l-2} + \frac{5}{4} \Delta(\pi, \pi_*) + 2\epsilon_{l-2} + \frac{5}{4} \Delta(\hat{\pi}_{l-1}, \pi_*) \right) \hspace{2em} \text{(from inductive hypothesis)} \\ &\leq \frac{5}{3} \epsilon_{l-1} + \frac{1}{64} \left(2\epsilon_{l-2} + \frac{5}{4} \Delta(\pi, \pi_*) + 2\epsilon_{l-2} + \frac{5}{4} \epsilon_{l-1} \right) \\ &\leq 2\epsilon_{l-1} + \frac{1}{4} \Delta(\pi, \pi_*). \end{aligned}$$

Also,

$$\begin{aligned}
\Delta(\widehat{\pi}_l, \widehat{\pi}_{l-1}) &\leq \widehat{\Delta}_l^{\gamma^l}(\widehat{\pi}_l, \widehat{\pi}_{l-1}) + [\gamma^l]_{\widehat{\pi}_l} \|x_{\widehat{\pi}_l} - \phi_{\widehat{\pi}_{l-1}}\|_{A(w(\lambda^l, \gamma^l))^{-1}}^2 + \frac{\log(1/\delta_l)}{[\gamma^l]_{\widehat{\pi}_l} n_l} && \text{(from } \mathcal{E} \text{)} \\
&\leq \widehat{\Delta}_l^{\gamma^l}(\pi_*, \widehat{\pi}_{l-1}) + [\gamma^l]_{\pi_*} \|\phi_{\pi_*} - \phi_{\widehat{\pi}_{l-1}}\|_{A(w(\lambda^l, \gamma^l))^{-1}}^2 + \frac{\log(1/\delta_l)}{[\gamma^l]_{\pi_*} n_l} \\
&&& \text{(eqn (6), the minimum)} \\
&\leq \Delta(\pi_*, \widehat{\pi}_{l-1}) + 2[\gamma^l]_{\pi_*} \|\phi_{\pi_*} - \phi_{\widehat{\pi}_{l-1}}\|_{A(w(\lambda^l, \gamma^l))^{-1}}^2 + \frac{2\log(1/\delta_l)}{[\gamma^l]_{\pi_*} n_l} && \text{(from } \mathcal{E} \text{)} \\
&\leq \Delta(\pi_*, \widehat{\pi}_{l-1}) + \frac{1}{3}\epsilon_l + \frac{1}{32}\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi_*, \widehat{\pi}_{l-1}) && \text{(from Lemma E.4)} \\
&\leq \Delta(\pi_*, \widehat{\pi}_{l-1}) + \frac{1}{3}\epsilon_l + \frac{1}{32}\left(2\epsilon_{l-1} + \frac{5}{4}\Delta(\pi_*, \pi_*)\right). && \text{(from the above)}
\end{aligned}$$

Therefore,

$$\begin{aligned}
\Delta(\widehat{\pi}_l, \pi_*) &= \Delta(\widehat{\pi}_l, \widehat{\pi}_{l-1}) - \Delta(\pi_*, \widehat{\pi}_{l-1}) \\
&\leq \frac{1}{3}\epsilon_l + \frac{1}{16}2\epsilon_l \\
&\leq \epsilon_l
\end{aligned}$$

Therefore, $\Delta(\widehat{\pi}_l, \pi_*) \leq \epsilon_l$, so $\widehat{\pi}_l \in S_l$.

□

Lemma E.6. For any $\lambda \in \Delta_\Pi$, $\gamma \in \mathbb{R}^{|\Pi|}$, and $\pi' \in \Pi$, we have

$$\min_{w \in \Omega} \sum_{\pi \in \Pi} \lambda_\pi \gamma_\pi \|\phi_\pi - \phi_{\pi'}\|_{A(w)^{-1}}^2 = \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top t_a^{(c)}(\pi')} \right)^2 \right].$$

where $w_{a,c} = \nu_c p_a^{(c)}$ and $p_a^{(c)} \propto \sqrt{\sum_{\pi \in \Pi} \lambda_\pi \gamma_\pi (\mathbf{1}\{\pi'(c) = a, \pi(c) \neq a\} + \mathbf{1}\{\pi'(c) \neq a, \pi(c) = a\})}$ and \odot denotes element-wise multiplication.

Proof. For any $\lambda \in \Delta_{\Pi}$,

$$\begin{aligned}
& \min_{w \in \Omega} \sum_{\pi \in \Pi} \lambda_{\pi} \gamma_{\pi} \|\phi_{\pi} - \phi_{\pi'}\|_{A(w)}^2 \\
&= \min_{w \in \Omega} \sum_{\pi \in \Pi} \sum_{a,c} \frac{\lambda_{\pi} \gamma_{\pi}}{w_{a,c}} (\phi_{\pi} - \phi_{\pi'})^{\top} e_{a,c} e_{a,c}^{\top} (\phi_{\pi} - \phi_{\pi'}) \\
&= \min_{p_1, \dots, p_{|c|} \in \Delta_{\mathcal{A}}} \sum_{\pi \in \Pi} \sum_{a,c} \frac{\lambda_{\pi} \gamma_{\pi}}{\nu_c p_{c,a}} (\phi_{\pi} - \phi_{\pi'})^{\top} e_{a,c} e_{a,c}^{\top} (\phi_{\pi} - \phi_{\pi'}) \\
&= \sum_c \min_{p_c \in \Delta_{\mathcal{A}}} \sum_a \sum_{\pi \in \Pi} \frac{\lambda_{\pi} \gamma_{\pi}}{\nu_c p_{c,a}} (\phi_{\pi} - \phi_{\pi'})^{\top} e_{a,c} e_{a,c}^{\top} (\phi_{\pi} - \phi_{\pi'}) \\
&= \sum_c \frac{1}{\nu_c} \min_{p_c \in \Delta_{\mathcal{A}}} \sum_a \frac{1}{p_{c,a}} \left(\sum_{\pi \in \Pi} \lambda_{\pi} \gamma_{\pi} (\phi_{\pi} - \phi_{\pi'})^{\top} e_{a,c} e_{a,c}^{\top} (\phi_{\pi} - \phi_{\pi'}) \right) \\
&= \sum_c \frac{1}{\nu_c} \left(\sum_{a \in \mathcal{A}} \sqrt{\sum_{\pi \in \Pi} \lambda_{\pi} \gamma_{\pi} (\phi_{\pi} - \phi_{\pi'})^{\top} e_{a,c} e_{a,c}^{\top} (\phi_{\pi} - \phi_{\pi'})} \right)^2 \\
&= \sum_c \frac{1}{\nu_c} \left(\sum_{a \in \mathcal{A}} \sqrt{\sum_{\pi \in \Pi} \lambda_{\pi} \gamma_{\pi} \nu_c^2 (\mathbf{1}\{\pi'(c) = a, \pi(c) \neq a\} + \mathbf{1}\{\pi'(c) \neq a, \pi(c) = a\})} \right)^2 \\
&= \sum_c \nu_c \left(\sum_{a \in \mathcal{A}} \sqrt{\sum_{\pi \in \Pi} \lambda_{\pi} \gamma_{\pi} (\mathbf{1}\{\pi'(c) = a, \pi(c) \neq a\} + \mathbf{1}\{\pi'(c) \neq a, \pi(c) = a\})} \right)^2 \\
&= \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top} t_a^{(c)}(\pi')} \right)^2 \right].
\end{aligned}$$

Note that the minimizer

$$\begin{aligned}
p_{c,a} &= \frac{\sqrt{\sum_{\pi \in \Pi} \lambda_{\pi} \gamma_{\pi} (\phi_{\pi} - \phi_{\pi'})^{\top} e_{a,c} e_{a,c}^{\top} (\phi_{\pi} - \phi_{\pi'})}}{\sum_{a'} \sqrt{\sum_{\pi \in \Pi} \lambda_{\pi} \gamma_{\pi} (\phi_{\pi} - \phi_{\pi'})^{\top} e_{a',c} e_{a',c}^{\top} (\phi_{\pi} - \phi_{\pi'})}} \\
&\propto \sqrt{\sum_{\pi \in \Pi} \lambda_{\pi} \gamma_{\pi} (\mathbf{1}\{\pi'(c) = a, \pi(c) \neq a\} + \mathbf{1}\{\pi'(c) \neq a, \pi(c) = a\})}.
\end{aligned}$$

□

Lemma E.7. Under \mathcal{E} , the choice for n_l in the algorithm satisfies

$$n_l \lesssim \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\pi_*} - \phi_{\pi}\|_{A(w)}^2 \log(1/\delta_l)}{\epsilon_l^2 + \Delta(\pi)^2}.$$

Proof.

$$\begin{aligned}
& h_l(\lambda^l, \gamma^l, n_l) \\
&= \sum_{\pi \in \Pi} [\lambda^l]_{\pi} \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{[\gamma^l]_{\pi} n} \right) + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda^l \odot \gamma^l)^{\top} (t_a^{(c)} + \eta_l)} \right)^2 \right] \\
&\leq \max_{\lambda \in \Delta_{\Pi}} \min_{\gamma} \sum_{\pi \in \Pi} \lambda_{\pi} \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} \right) + \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right)^2 \right] + \frac{1}{4} \epsilon_l \\
&\hspace{15em} \text{(by Theorem G.2, the saddle point argument)} \\
&\leq \max_{\lambda \in \Delta_{\Pi}} \min_{\gamma} \sum_{\pi \in \Pi} \lambda_{\pi} \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} \right) + \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top} t_a^{(c)}} \right)^2 \right] + \frac{1}{2} \epsilon_l \\
&\hspace{15em} \text{(by Lemma H.3, controlling the bias)} \\
&= \max_{\lambda \in \Delta_{\Pi}} \min_{w \in \Omega} \min_{\gamma \in \mathbb{R}_+^{|\Pi|}} \sum_{\pi \in \Pi} \lambda_{\pi} \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \gamma_{\pi} \|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w)-1}^2 + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} \right) + \frac{1}{2} \epsilon_l \\
&\hspace{15em} \text{(by Lemma E.6, the definition of } w \text{)} \\
&= \min_{w \in \Omega} \max_{\pi \in \Pi} \min_{\gamma > 0} -\frac{1}{8} \widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + 8\gamma \|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w)-1}^2 + 8 \frac{\log(1/\delta_l)}{\gamma n_l} + \frac{1}{2} \epsilon_l \\
&\hspace{15em} \text{(by Lemma G.17, the strong duality)} \\
&\leq \min_{w \in \Omega} \max_{\pi \in \Pi} \min_{\gamma} \left(-\frac{3}{32} \Delta(\pi, \pi_*) + 8\gamma \|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w)-1}^2 + 8 \frac{\log(1/\delta_l)}{\gamma n_l} \right) + \frac{3}{4} \epsilon_l \quad \text{(by Lemma E.5)} \\
&\leq \min_{w \in \Omega} \max_{\pi \in \Pi} \left(-\frac{3}{32} \Delta(\pi, \pi_*) + 16 \sqrt{\frac{\|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w)-1}^2 \log(1/\delta_l)}{n_l}} \right) + \frac{3}{4} \epsilon_l \\
&\leq \min_{w \in \Omega} \max_{\pi \in \Pi} \left(-\frac{3}{32} \Delta(\pi, \pi_*) + 16 \sqrt{\frac{\|\phi_{\pi_*} - \phi_{\pi}\|_{A(w)-1}^2 \log(1/\delta_l)}{n_l}} \right. \\
&\quad \left. + 16 \sqrt{\frac{\|\phi_{\pi_*} - \phi_{\widehat{\pi}_{l-1}}\|_{A(w)-1}^2 \log(1/\delta_l)}{n_l}} \right) + \frac{3}{4} \epsilon_l \\
&\leq \min_{w \in \Omega} \max_{\pi \in \Pi} \left(-\frac{3}{32} \Delta(\pi, \pi_*) + 16 \sqrt{\frac{\|\phi_{\pi_*} - \phi_{\pi}\|_{A(w)-1}^2 \log(1/\delta_l)}{n_l}} \right. \\
&\quad \left. + 16 \sqrt{\frac{\max_{\pi' \in S_{l-1}} \|\phi_{\pi_*} - \phi_{\pi'}\|_{A(w)-1}^2 \log(1/\delta_l)}{n_l}} \right) + \frac{3}{4} \epsilon_l.
\end{aligned}$$

which is less than ϵ_l whenever

$$n_l \gtrsim \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\pi_*} - \phi_{\pi}\|_{A(w)-1}^2 \log(1/\delta_l)}{\epsilon_l^2 + \Delta(\pi)^2}. \quad (15)$$

□

F Intuition for convergence of duality gap

It could seem mysterious that one could find a $\log(|\Pi|)$ -sparse and ϵ -good solution of the optimization problem $\max_{\lambda \in \Delta_{\Pi}} \min_{\gamma \in \mathbb{R}_+^{|\Pi|}} h_l(\lambda, \gamma)$. In this section, we aim to provide some intuition and a constructive proof of an easier case.

The existence of such a solution critically relies on the fact that for any fixed $\gamma \in [\gamma_{\min}, \gamma_{\max}]^{|\Pi|}$, we can find a $\log(|\Pi|)$ -sparse solution λ^t such that $\max_{\lambda \in \Delta_{\Pi}} g(\lambda, \gamma) - g(\lambda^t, \gamma) \leq \epsilon_l$. Also, if we consider $\min_{\gamma \in \mathbb{R}_+^{|\Pi|}} h_l(\lambda, \gamma)$ for a fixed λ , the gradient descent algorithm allows us to find a good solution of γ in arbitrary precision. In what follows, we provide an argument for convergence analysis of the unregularized objective h_l assuming L -Lipschitz gradient and we can solve γ exactly.

Suppose the primal and dual problems are defined as follows:

$$\mathcal{P}_l(w, \gamma, n) = \max_{\pi \in \Pi} \left[-\Delta(\pi) + \frac{\log(1/\delta_l)}{\gamma_\pi n} + \gamma_\pi \|\phi_{\pi_*} - \phi_\pi\|_{A(w)^{-1}}^2 \right]$$

$$h_l(\lambda, \gamma, n) = \sum_{\pi \in \Pi} \lambda_\pi \cdot \left(-\Delta(\pi) + \frac{\log(1/\delta_l)}{\gamma_\pi n} \right) + \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top t_a^{(c)}} \right)^2 \right],$$

then

$$\nabla_\lambda h_l(\lambda, \gamma, n) = -\Delta(\pi) + \frac{\log(1/\delta_l)}{\gamma_\pi n} + \gamma_\pi \|\phi_{\pi_*} - \phi_\pi\|_{A(w(\lambda, \gamma))^{-1}}^2.$$

Observe that

$$\max_{\pi \in \Pi} [\nabla_\lambda h_l(\lambda, \gamma)]_\pi = \mathcal{P}_l(w(\lambda, \gamma), \gamma, n).$$

Therefore, the Frank-Wolfe gap

$$\begin{aligned} g_t &= \langle \nabla_\lambda h_t(\lambda^t, \gamma^t, n), e_{\pi_t} - \lambda^t \rangle \\ &= \max_{\pi \in \Pi} [\nabla_\lambda h_l(\lambda^t, \gamma^t, n)]_\pi - h_l(\lambda^t, \gamma^t, n) \\ &= \mathcal{P}_l(w(\lambda^t, \gamma^t), \gamma^t, n) - h_l(\lambda^t, \gamma^t, n). \end{aligned}$$

Note that if we assume $\gamma^t = \arg \min_\gamma h_l(\lambda^t, \gamma, n)$, we have

$$\begin{aligned} &\mathcal{P}_l(w(\lambda^t, \gamma^t), \gamma^t, n) \\ &= \max_{\pi \in \Pi} \left[-\Delta(\pi) + \frac{\log(1/\delta_l)}{[\gamma^t]_\pi n} + [\gamma^t]_\pi \|\phi_{\pi_*} - \phi_\pi\|_{A(w(\lambda^t, \gamma^t))^{-1}}^2 \right] \\ &\geq \max_{\pi \in \Pi} \min_{\gamma} \left[-\Delta(\pi) + \frac{\log(1/\delta_l)}{\gamma_\pi n} + \gamma_\pi \|\phi_{\pi_*} - \phi_\pi\|_{A(w(\lambda^t, \gamma^t))^{-1}}^2 \right] \\ &\geq \min_{w \in \Omega} \max_{\pi \in \Pi} \min_{\gamma} \left[-\Delta(\pi) + \frac{\log(1/\delta_l)}{\gamma_\pi n} + \gamma_\pi \|\phi_{\pi_*} - \phi_\pi\|_{A(w)^{-1}}^2 \right] \\ &= \max_{\lambda \in \Delta_\Pi} \min_{w \in \Omega} \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^\Pi} \sum_{\pi \in \Pi} \lambda_\pi \left(-\Delta(\pi) + \frac{\log(1/\delta_l)}{\gamma_\pi n} + \gamma_\pi \|\phi_{\pi_*} - \phi_\pi\|_{A(w)^{-1}}^2 \right) \\ &= \max_{\lambda \in \Delta_\Pi} \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^\Pi} \sum_{\pi} \lambda_\pi \cdot \left(-\Delta(\pi) + \frac{\log(1/\delta_l)}{\gamma_\pi n} \right) + \mathbb{E}_c \left[\left(\sum_a \sqrt{(\lambda \odot \gamma)^\top t_a^{(c)}} \right)^2 \right] \\ &= h_l(\lambda^*, \gamma^*) \\ &\geq \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^\Pi} \sum_{\pi} [\lambda^t]_\pi \cdot \left(-\Delta(\pi) + \frac{\log(1/\delta_l)}{\gamma_\pi n} \right) + \mathbb{E}_c \left[\left(\sum_a \sqrt{(\lambda^t \odot \gamma)^\top t_a^{(c)}} \right)^2 \right] \\ &= h_l(\lambda^t, \gamma^t, n). \end{aligned}$$

Therefore, to show that $h_l(\lambda^*, \gamma^*, n) - h_l(\lambda^t, \gamma^t, n)$ is small, it is sufficient to show that $\mathcal{P}_l(w(\lambda^t, \gamma^t), \gamma^t, n) - h_l(\lambda^t, \gamma^t, n)$ is small, which corresponds to a small Frank-Wolfe gap. Then, we can use similar arguments in Lemmas G.4 and G.5 to show that the Frank-Wolfe gap is small.

G Convergence analysis of FW-GD

G.1 Statement of the convergence results

In this section, we will characterize the performance of Algorithm 6, a.k.a. Algorithm 3. Our goal is to show two results: the duality gap converges to zero, and our algorithm converges to the saddle point. It is known that Frank-Wolfe algorithm directly deals with the duality gap [32], so we will define our primal and dual problem in what follows. Since we are computing n_l via binning, in each

inner loop n is fixed. Then, we define our dual objective the same as (12) with the shorthand notation $h_l(\lambda, \gamma) := h_l(\lambda, \gamma, n)$. We formulate our primal objective as

$$\mathcal{P}_l(w(\lambda, \gamma), \gamma) := \max_{\pi \in \Pi} \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \gamma_\pi \|\phi_\pi - \phi_{\widehat{\pi}_{l-1}}\|_{A(w(\lambda, \gamma))^{-1}}^2 + \frac{\log(1/\delta_l)}{\gamma_\pi n} \right), \quad (16)$$

where $w(\lambda, \gamma) \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{C}|}$ such that

$$[w(\lambda, \gamma)]_{a,c} = \nu_c \cdot p_{c,a} = \nu_c \cdot \frac{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta)}}{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta)}}. \quad (17)$$

Then we will show those two results. First, Theorem G.1 bounds the duality gap of the primal and dual objective. Second, Theorem G.2 shows that Algorithm 3 converges to a saddle point.

Theorem G.1. *For any $l \in \mathbb{N}$, with the number of FW-GD iterations $T_l = O(L^2 \epsilon_l^{-2})$ where $L = |\mathcal{A}|^2 \frac{((1+\eta_l)\gamma_{\max})^{5/2}}{\eta_l^{3/2} \gamma_{\min}^2}$, we have*

$$|\mathcal{P}_l(w(\lambda^l, \gamma^l), \gamma^l) - h_l(\lambda^l, \gamma^l)| \leq \epsilon_l.$$

Moreover, T_l depends at most polynomially on $|\mathcal{A}|, \epsilon_l^{-1}, \log(1/\delta_l)$.

Proof. First, Lemma H.2 shows that for any λ, γ , and n , $h_l(\lambda, \gamma, n) = \langle \lambda, \nabla_\lambda h_l(\lambda, \gamma, n) \rangle$. Therefore, at some iteration t , the Frank-Wolfe gap

$$g_t = \langle \nabla_\lambda h_l(\lambda^t, \gamma^t), \mathbf{e}_{\pi_t} - \lambda^t \rangle = \max_{\pi \in \Pi} [\nabla_\lambda h_l(\lambda^t, \gamma^t)]_\pi - h_l(\lambda^t, \gamma^t).$$

Lemma G.6 shows that with a small choice of the regularization parameter the primal objective is close to the maximum component of the gradient, i.e. $|\mathcal{P}_l(w(\lambda^l, \gamma^l), \gamma^l) - \max_{\pi \in \Pi} [\nabla_\lambda h_l(\lambda^l, \gamma^l)]_\pi| \leq \frac{\epsilon_l}{2}$. Also, Lemma G.5 shows that if $t \geq L^2 \epsilon_l^{-2}$ is large enough, the Frank-Wolfe gap is bounded by ϵ_l . Combining these two lemmas, for $t \geq L^2 \epsilon_l^{-2}$, we have

$$\begin{aligned} & |\mathcal{P}_l(w(\lambda^l, \gamma^l), \gamma^l) - h_l(\lambda^l, \gamma^l)| \\ & \leq |\mathcal{P}_l(w(\lambda^l, \gamma^l), \gamma^l) - \max_{\pi \in \Pi} [\nabla_\lambda h_l(\lambda^l, \gamma^l)]_\pi| + |h_l(\lambda^l, \gamma^l) - \max_{\pi \in \Pi} [\nabla_\lambda h_l(\lambda^l, \gamma^l)]_\pi| \\ & \leq |\mathcal{P}_l(w(\lambda, \gamma), \gamma) - \max_{\pi \in \Pi} [\nabla_\lambda h_l(\lambda, \gamma)]_\pi| + g_t \\ & \leq \frac{\epsilon_l}{2} + \frac{\epsilon_l}{2} = \epsilon_l. \end{aligned}$$

Finally, we conclude that $T_l = \text{poly}(|\mathcal{A}|, \epsilon_l^{-1}, \log(1/\delta_l))$ since $\gamma_{\max} = O(|\mathcal{A}|^{-1} \eta_l^{-1/2})$, $\gamma_{\min} = O(\sqrt{\eta_l})$, and $\eta_l = O(|\mathcal{A}|^{-4} \epsilon_l^2)$ all depends polynomially on $|\mathcal{A}|$ and ϵ_l^{-1} . This shows Theorem G.1. \square

We now have the second main result of this section.

Theorem G.2. *For any l , with $T_l = \text{poly}(|\mathcal{A}|, \epsilon_l^{-1}, \log(1/\delta_l))$ and the size of the history $\mathcal{D} \geq \text{poly}(|\mathcal{A}|, \epsilon^{-1}, \log(1/\delta), \log(|\Pi|))$, Algorithm 3 converges to a saddle point, i.e.*

$$\left| \max_{\lambda \in \Delta_\Pi} \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^\Pi} h_l(\lambda, \gamma) - h_l(\lambda^l, \gamma^l) \right| \leq \epsilon_l.$$

Proof. Note that

$$\begin{aligned}
& \mathcal{P}_l(w(\lambda^l, \gamma^l), \gamma^l) \\
&= \max_{\pi \in \Pi} \left[-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{[\gamma^l]_{\pi} n} + [\gamma^l]_{\pi} \|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w(\lambda^l, \gamma^l))^{-1}}^2 \right] \\
&\geq \max_{\pi \in \Pi} \min_{\gamma} \left[-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} + \gamma_{\pi} \|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w(\lambda^l, \gamma^l))^{-1}}^2 \right] \\
&\geq \min_{w \in \Omega} \max_{\pi \in \Pi} \min_{\gamma} \left[-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} + \gamma_{\pi} \|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w)^{-1}}^2 \right] \\
&= \max_{\lambda \in \Delta_{\Pi}} \min_{w \in \Omega} \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^{\Pi}} \sum_{\pi \in \Pi} \lambda_{\pi} \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} + \gamma_{\pi} \|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w)^{-1}}^2 \right) \\
&\hspace{15em} \text{(by Lemma G.17, strong duality)} \\
&= \max_{\lambda \in \Delta_{\Pi}} \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^{\Pi}} \sum_{\pi} \lambda_{\pi} \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} \right) \\
&\quad + \mathbb{E}_{c \sim \nu} \left[\left(\sum_a \sqrt{(\lambda \odot \gamma)^{\top} t_a^{(c)}} \right)^2 \right] \hspace{5em} \text{(by Lemma E.6)} \\
&\geq \max_{\lambda \in \Delta_{\Pi}} \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^{\Pi}} \sum_{\pi \in \Pi} \lambda_{\pi} \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} \right) \\
&\quad + \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right)^2 \right] - \frac{1}{2} \epsilon_l \hspace{5em} \text{(by Lemma H.3)} \\
&\geq \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^{\Pi}} \sum_{\pi} [\lambda^l]_{\pi} \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} \right) \\
&\quad + \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda^l \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right)^2 \right] - \frac{1}{2} \epsilon_l \\
&\geq \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^{\Pi}} \sum_{\pi} [\lambda^l]_{\pi} \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} \right) \\
&\quad + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda^l \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right)^2 \right] - \frac{3}{4} \epsilon_l \\
&\hspace{15em} \text{(by Lemma D.6, controlling the history)} \\
&\geq \sum_{\pi} [\lambda^l]_{\pi} \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{[\gamma^l]_{\pi} n} \right) \\
&\quad + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda^l \odot \gamma^l)^{\top} (t_a^{(c)} + \eta_l)} \right)^2 \right] - \epsilon_l \\
&\hspace{15em} \text{(by Lemma G.7, the GD convergence)} \\
&= h_l(\lambda^l, \gamma^l) - \epsilon_l.
\end{aligned}$$

In other words,

$$\mathcal{P}_l(w(\lambda^l, \gamma^l), \gamma^l) \geq \max_{\lambda \in \Delta_{\Pi}} \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^{\Pi}} h_l(\lambda, \gamma) \geq h_l(\lambda^l, \gamma^l) - \epsilon_l.$$

On the other hand, by Theorem G.1, we have $\mathcal{P}_l(w(\lambda^l, \gamma^l), \gamma^l) \leq h_l(\lambda^l, \gamma^l) + \epsilon_l$. Therefore, we have

$$\max_{\lambda \in \Delta_{\Pi}} \min_{\gamma \in [\gamma_{\min}, \gamma_{\max}]^{\Pi}} h_l(\lambda, \gamma) \in [h_l(\lambda^l, \gamma^l) - \epsilon_l, h_l(\lambda^l, \gamma^l) + \epsilon_l]$$

and so we have our result. \square

G.2 Technical proofs

G.2.1 Guarantees on γ

We first provides some guarantees of γ and the convergence of the GD subroutine.

Lemma G.3. *Consider a fixed n . Let $\gamma^* = \arg \min_{\gamma} h_l(\lambda, \gamma, n)$. Then we have for all i ,*

$$[\gamma^*]_i \in \left[\frac{1}{3} \sqrt{\frac{\eta_l \log(1/\delta_l)}{n}}, \min \left\{ \sqrt{\frac{\log(1/\delta_l)}{2n\mathbb{E}_c[\mathbf{1}\{\pi(c) \neq \pi^*(c)\}]}} , \sqrt{\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n}} \right\} \right].$$

Proof.

$$\begin{aligned} & [\nabla_{\gamma} h_l(\lambda, \gamma)]_{\pi} \\ &= \mathbb{E}_c \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right) \cdot \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_{\pi} ([t_{a'}^{(c)}]_{\pi} + \eta_l)}{\sqrt{(\lambda \odot \gamma)^{\top} (t_{a'}^{(c)} + \eta_l)}} \right) \right] - \frac{\lambda_{\pi} \log(1/\delta_l)}{\gamma_{\pi}^2 n} \\ &\geq \mathbb{E}_c \left[\left(\sum_{a \in \mathcal{A}} \sqrt{\lambda_{\pi} ([t_a^{(c)}]_{\pi} + \eta_l)} \right)^2 \right] - \frac{\lambda_{\pi} \log(1/\delta_l)}{\gamma_{\pi}^2 n} \\ &\geq |\mathcal{A}|^2 \eta_l \lambda_{\pi} + 2\lambda_{\pi} \mathbb{E}_c[\mathbf{1}\{\pi(c) \neq \pi^*(c)\}] - \frac{\lambda_{\pi} \log(1/\delta_l)}{\gamma_{\pi}^2 n}, \end{aligned}$$

where the first to second line follows from Cauchy-Schwartz - $(\sum_a x_a) \sum_a \left(\frac{y_a}{x_a}\right) \geq (\sum_a \sqrt{y_a})^2$.

We first solve $\frac{\lambda_{\pi} \log(1/\delta_l)}{\gamma_{\pi}^2 n} < |\mathcal{A}|^2 \eta_l \lambda_{\pi}$ and get $\gamma_{\pi} > \sqrt{\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n}}$. We also solve $\frac{\lambda_{\pi} \log(1/\delta_l)}{\gamma_{\pi}^2 n} < 2\lambda_{\pi} \mathbb{E}_c[\mathbf{1}\{\pi(c) \neq \pi^*(c)\}]$ and get $\gamma_{\pi} < \sqrt{\frac{\log(1/\delta_l)}{2n\mathbb{E}_c[\mathbf{1}\{\pi(c) \neq \pi^*(c)\}]}}$. Therefore, the π th component of the gradient is always positive whenever $\gamma_{\pi} > \min \left\{ \sqrt{\frac{\log(1/\delta_l)}{2n\mathbb{E}_c[\mathbf{1}\{\pi(c) \neq \pi^*(c)\}]}} , \sqrt{\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n}} \right\}$. Therefore, the minimum γ should have $\gamma_{\pi} \leq \min \left\{ \sqrt{\frac{\log(1/\delta_l)}{2n\mathbb{E}_c[\mathbf{1}\{\pi(c) \neq \pi^*(c)\}]}} , \sqrt{\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n}} \right\}$. On the other hand, let $s = \arg \min_{\pi} \gamma_{\pi}$. Then,

$$\eta_l \gamma_s \leq (\lambda \odot \gamma)^{\top} (t_a^{(c)} + \eta_l) = \left(\lambda \odot (t_a^{(c)} + \eta_l) \right)^{\top} \gamma \leq \left\| \lambda \odot (t_a^{(c)} + \eta_l) \right\|_1 \cdot \|\gamma\|_{\infty}.$$

Then

$$\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \leq \sum_{a \in \mathcal{A}} \sqrt{\left\| \lambda \odot (t_a^{(c)} + \eta_l) \right\|_1} \cdot \sqrt{\|\gamma\|_{\infty}}.$$

Note that

$$\begin{aligned} \left(\sum_{a \in \mathcal{A}} \sqrt{\left\| \lambda \odot (t_a^{(c)} + \eta_l) \right\|_1} \right)^2 &= \left(\sum_{a \in \mathcal{A}} \sqrt{\lambda^{\top} (t_a^{(c)} + \eta_l)} \right)^2 \\ &\leq \left(\sum_{a \in \mathcal{A}} \lambda^{\top} (t_a^{(c)} + \eta_l) \right) |\mathcal{A}| \\ &\leq |\mathcal{A}|(1 + \eta_l). \end{aligned}$$

Since for any π , $\sum_{a' \in \mathcal{A}} [t_{a'}^{(c)}]_{\pi} \leq 2$, so

$$[\nabla_{\gamma} h_l(\lambda, \gamma)]_{\pi} \leq \sqrt{|\mathcal{A}|(1 + \eta_l) \|\gamma\|_{\infty}} \cdot \frac{(2 + \eta_l) \lambda_{\pi}}{\sqrt{\eta_l \gamma_s}} - \frac{\lambda_{\pi} \log(1/\delta_l)}{\gamma_{\pi}^2 n}.$$

Let $\pi = s$, then by the fact that $\|\gamma\|_{\infty} \leq \sqrt{\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n}}$, we have

$$[\nabla_{\gamma} h_l(\lambda, \gamma)]_s \leq \sqrt{|\mathcal{A}|(1 + \eta_l)} \left(\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n} \right)^{1/4} \cdot \frac{(2 + \eta_l) \lambda_s}{\sqrt{\eta_l \gamma_s}} - \frac{\lambda_s \log(1/\delta_l)}{\gamma_s^2 n}.$$

We solve $\sqrt{|\mathcal{A}|(1+\eta_l)} \left(\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n} \right)^{1/4} \cdot \frac{(2+\eta_l)\lambda_s}{\sqrt{\eta_l \gamma_s}} - \frac{\lambda_s \log(1/\delta_l)}{\gamma_s^2 n} < 0$. Then we get

$$\gamma_s < (1+\eta_l)^{-1/3} (2+\eta_l)^{-2/3} \sqrt{\frac{\eta_l \log(1/\delta_l)}{n}}.$$

Since $(1+\eta_l)^{-1/3} (2+\eta_l)^{-2/3} > \frac{1}{3}$ whenever $\eta_l \leq 1$, the s th component of the gradient is negative whenever $\gamma_s < \frac{1}{3} \sqrt{\frac{\eta_l \log(1/\delta_l)}{n}}$. Therefore, $\min_{\pi} \gamma_{\pi} \geq \frac{1}{3} \sqrt{\frac{\eta_l \log(1/\delta_l)}{n}}$. \square

G.2.2 Convergence of Frank-Wolfe gap

Lemma G.4 and G.5 shows that the Frank-Wolfe gap is small. The proof technique follows from the general Frank-Wolfe analysis.

Lemma G.4. For any $\xi \in [0, 1]$, any t , with $L = |\mathcal{A}|^2 \frac{((1+\eta_l)\gamma_{\max})^{5/2}}{\eta_l^{3/2} \gamma_{\min}^2}$, we have $h_l(\lambda^{t+1}, \gamma^{t+1}) \geq h_l(\lambda^t, \gamma^t) + \xi g_t - \frac{1}{2} \xi^2 L - \kappa_t$.

Proof. By L -Lipschitz gradient condition of $-h_\ell$ in λ given in Lemma G.12 we have

$$-h_l(\lambda^{t+1}, \gamma^{t+1}) \leq -h_l(\lambda^t, \gamma^{t+1}) - \langle \nabla_{\lambda} h_l(\lambda^t, \gamma^{t+1}), \lambda^{t+1} - \lambda^t \rangle + \frac{L}{2} \|\lambda^{t+1} - \lambda^t\|_1^2.$$

Therefore,

$$h_l(\lambda^{t+1}, \gamma^{t+1}) \geq h_l(\lambda^t, \gamma^{t+1}) + \langle \nabla_{\lambda} h_l(\lambda^t, \gamma^{t+1}), \lambda^{t+1} - \lambda^t \rangle - \frac{L}{2} \|\lambda^{t+1} - \lambda^t\|_1^2.$$

Plugging in $\lambda^{t+1} = (1 - \beta_t)\lambda^t + \beta_t \mathbf{e}_{\pi_t}$ as in line 8 of Algorithm 4, we have

$$\begin{aligned} & h_l((1 - \beta_t)\lambda^t + \beta_t \mathbf{e}_{\pi_t}, \gamma^{t+1}) \\ & \geq h_l(\lambda^t, \gamma^{t+1}) + \langle \nabla_{\lambda} h_l(\lambda^t, \gamma^{t+1}), (1 - \beta_t)\lambda^t + \beta_t \mathbf{e}_{\pi_t} - \lambda^t \rangle - \frac{L}{2} \|(1 + \beta_t)\lambda^t - \beta_t \mathbf{e}_{\pi_t} - \lambda^t\|_1^2 \\ & = h_l(\lambda^t, \gamma^{t+1}) + \beta_t \langle \nabla_{\lambda} h_l(\lambda^t, \gamma^{t+1}), \mathbf{e}_{\pi_t} - \lambda^t \rangle - \frac{L\beta_t^2}{2} \|\mathbf{e}_{\pi_t} - \lambda^t\|_1^2 \\ & = h_l(\lambda^t, \gamma^{t+1}) + \beta_t g_t - \frac{L\beta_t^2}{2} \|\mathbf{e}_{\pi_t} - \lambda^t\|_1^2. \end{aligned}$$

Choose $\beta_t := \arg \max_{\xi \in [0, 1]} \{\xi g_t - \frac{\xi^2 L}{2} \|\mathbf{e}_{\pi_t} - \lambda^t\|_1^2\}$. Plugging in this expression gives us

$$\begin{aligned} h_l(\lambda^{t+1}, \gamma^{t+1}) & \geq h_l(\lambda^t, \gamma^{t+1}) + \beta_t \langle \nabla_{\lambda} h_l(\lambda^t, \gamma^{t+1}), \mathbf{e}_{\pi_t} - \lambda^t \rangle - \frac{L\beta_t^2}{2} \|\mathbf{e}_{\pi_t} - \lambda^t\|_1^2 \\ & = h_l(\lambda^t, \gamma^{t+1}) + \max_{\xi \in [0, 1]} \left\{ \xi g_t - \frac{\xi^2 L}{2} \|\mathbf{e}_{\pi_t} - \lambda^t\|_1^2 \right\} \\ & \geq h_l(\lambda^t, \gamma^{t+1}) + \xi g_t - \frac{\xi^2 L}{2} \end{aligned}$$

for any $\xi \in [0, 1]$ since $\|\mathbf{e}_{\pi_t} - \lambda^t\|_1^2 \leq 1$. Also, by construction of γ^{t+1} and Lemma G.7, we have

$$h_l(\lambda^t, \gamma^{t+1}) \geq \min_{\gamma} h_l(\lambda^t, \gamma) \geq h_l(\lambda^t, \gamma^t) - \kappa_t.$$

Therefore, our result follows. \square

Lemma G.5. We have for any t , with $L = |\mathcal{A}|^2 \frac{((1+\eta_l)\gamma_{\max})^{5/2}}{\eta_l^{3/2} \gamma_{\min}^2}$, $\min_{i \in [1, t]} g_i \leq \frac{L}{\sqrt{t+1}}$.

Proof. With Lemma G.4, we have

$$h_l(\lambda^{t+1}, \gamma^{t+1}, n_r) \geq h_l(\lambda^t, \gamma^t, n_r) + \xi g_t - \frac{1}{2} \xi^2 L - \kappa_t.$$

Plugging in the choice $\xi = \min\{\frac{g_t}{L}, 1\}$, we have $h_l(\lambda^{t+1}, \gamma^{t+1}, n_r) \geq h_l(\lambda^t, \gamma^t, n_r) + \frac{g_t}{2} \min\{\frac{g_t}{L}, 1\} - \kappa_t$. Summing this up from 0 to t gives us

$$\begin{aligned} h_l(\lambda^{t+1}, \gamma^{t+1}, n_r) - h_l(\lambda_0, \gamma_0, n_r) &\geq \sum_{i=0}^t \frac{g_i}{2} \min\{\frac{g_i}{L}, 1\} - \delta_i \\ &\geq (t+1)g_t^* \min\{\frac{g_t^*}{L}, 1\} - \sum_{i=0}^t \delta_i. \end{aligned}$$

where $g_t^* = \min_{i=0, \dots, t} g_i$. Then, as long as $\sum_{i=0}^t \delta_i \leq \epsilon_l$, by the fact that $h_l(\lambda^{t+1}, \gamma^{t+1}) - h_l(\lambda_0, \gamma_0) \leq \max_{\lambda \in \Delta_{\Pi}} \min_{\gamma} h_l(\lambda, \gamma) - h_l(\lambda_0, \gamma_0) < \infty$. Therefore, we have $\min_{i \in [1, t]} g_i \leq \frac{L}{\sqrt{t+1}}$. \square

G.2.3 Connect the Frank-Wolfe gap to the duality gap

Lemma G.6 shows that the primal objective is approximately the maximum component of the gradient of the dual objective, which simplifies our Frank-Wolfe gap expression.

Lemma G.6. Consider some $\lambda \in \Delta_{\Pi}$, $\gamma \in \mathbb{R}_+^{|\Pi|}$, and $n \in \mathbb{N}$. For $\eta_l < |\mathcal{A}|^{-4} \epsilon_l^2$, we have $|\mathcal{P}_l(w(\lambda^l, \gamma^l), \gamma^l) - \max_{\pi \in \Pi} [\nabla_{\lambda} h_l(\lambda^l, \gamma^l)]_{\pi}| \leq \epsilon_l$.

Proof. Observe that for any $\pi, \pi' \in \Pi$ and any γ ,

$$\begin{aligned} &\gamma_{\pi} \|\phi_{\pi'} - \phi_{\pi}\|_{A(w(\lambda, \gamma))^{-1}}^2 \\ &= \gamma_{\pi} \sum_{a, c} \frac{\nu_c^2}{[w(\lambda, \gamma)]_{a, c}} (\mathbf{1}\{\pi'(c) = a, \pi(c) \neq a\} + \mathbf{1}\{\pi'(c) \neq a, \pi(c) = a\}) \\ &= \gamma_{\pi} \sum_c \nu_c \sum_a \left(\frac{\nu_c}{[w(\lambda, \gamma)]_{a, c}} (\mathbf{1}\{\pi'(c) = a, \pi(c) \neq a\} + \mathbf{1}\{\pi'(c) \neq a, \pi(c) = a\}) \right) \\ &= \gamma_{\pi} \mathbb{E}_{c \sim \nu} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top} (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)}} (\mathbf{1}\{\pi'(c) = a, \pi(c) \neq a\} + \mathbf{1}\{\pi'(c) \neq a, \pi(c) = a\}) \right] \\ &= \mathbb{E}_{c \sim \nu} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top} (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)}} (\gamma_{\pi} [t_a^{(c)}]_{\pi}) \right]. \end{aligned}$$

Therefore,

$$\begin{aligned} &\mathcal{P}_l(w(\lambda^l, \gamma^l), \gamma^l) \\ &= \max_{\pi \in \Pi} \left\{ -\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi) + [\gamma^l]_{\pi} \|\phi_{\pi} - \phi_{\widehat{\pi}_{l-1}}\|_{A(w(\lambda^l, \gamma^l))^{-1}}^2 + \frac{\log(1/\delta_l)}{[\gamma^l]_{\pi} n} \right\} \\ &= \max_{\pi \in \Pi} \left\{ -\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi) + \mathbb{E}_{c \sim \nu} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda^l \odot \gamma^l)^{\top} (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda^l \odot \gamma^l)^{\top} (t_a^{(c)} + \eta_l)}} ([\gamma^l]_{\pi} [t_a^{(c)}]_{\pi}) \right] + \frac{\log(1/\delta_l)}{[\gamma^l]_{\pi} n} \right\}. \end{aligned}$$

Lemma D.7 guarantees that we could replace the expectation over context to history of contexts $\nu_{\mathcal{D}}$ without incurring much error. In particular, for a sufficiently large history \mathcal{D} , it guarantees

$$\begin{aligned} &\max_{\pi \in \Pi} \left| \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda^l \odot \gamma^l)^{\top} (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda^l \odot \gamma^l)^{\top} (t_a^{(c)} + \eta_l)}} ([\gamma^l]_{\pi} [t_a^{(c)}]_{\pi}) \right] \right. \\ &\quad \left. - \mathbb{E}_{c \sim \nu} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda^l \odot \gamma^l)^{\top} (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda^l \odot \gamma^l)^{\top} (t_a^{(c)} + \eta_l)}} ([\gamma^l]_{\pi} [t_a^{(c)}]_{\pi}) \right] \right| \leq \frac{\epsilon_l}{2}. \end{aligned}$$

On the other hand,

$$\begin{aligned} & \max_{\pi \in \Pi} \left\{ -\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi) + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda^l \odot \gamma^l)^\top (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda^l \odot \gamma^l)^\top (t_a^{(c)} + \eta_l)}} ([\gamma^l]_\pi [t_a^{(c)}]_\pi) \right] + \frac{\log(1/\delta_l)}{[\gamma^l]_\pi n} \right\} \\ &= \max_{\pi \in \Pi} \left\{ [\nabla_\lambda h_l(\lambda^l, \gamma^l)]_\pi - \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda^l \odot \gamma^l)^\top (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda^l \odot \gamma^l)^\top (t_a^{(c)} + \eta_l)}} [\gamma^l]_\pi \eta_l \right] \right\}. \end{aligned}$$

Note that when $\gamma_\pi \in [\gamma_{\min}, \gamma_{\max}]$,

$$\mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)}} \gamma_\pi \eta_l \right] \in \left[0, |\mathcal{A}|^2 \sqrt{\frac{\gamma_{\max}(1 + \eta_l)}{\gamma_{\min} \eta_l}} \gamma_{\max} \eta_l \right].$$

Therefore, for $\eta_l < |\mathcal{A}|^{-4} \epsilon_l^2$,

$$\left| \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_a \frac{\sum_{a' \in \mathcal{A}} \sqrt{(\lambda^l \odot \gamma^l)^\top (t_{a'}^{(c)} + \eta_l)}}{\sqrt{(\lambda^l \odot \gamma^l)^\top (t_a^{(c)} + \eta_l)}} [\gamma^l]_\pi \eta_l \right] \right| \leq \frac{\epsilon_l}{2}.$$

Therefore, we have our results. \square

G.3 Convergence of gradient descent

In this subsection we show convergence for gradient descent.

Algorithm 7 GD

Input: λ^t, n, κ_t

1: define $\iota^t = \epsilon_l^3 t^{-3} |\mathcal{A}|^{-6}$

2: clip λ and define $\tilde{\lambda} = \text{clip}(\lambda, \iota^t)$

3: run gradient descent of on γ for $h_l(\tilde{\lambda}, \gamma, n)$ over $\text{supp}(\tilde{\lambda})$ and output γ^t

Output: γ^t

We will first state the main result of this section.

Lemma G.7. *With the number of iterations $T = O\left(\frac{L_\gamma}{\iota^t} + \frac{1}{\kappa_t \iota^t}\right)$ with $L_\gamma = |\mathcal{A}|^2 \frac{((1+\eta_l)\gamma_{\max})^{3/2}}{\eta_l^{3/2} \gamma_{\min}^2} + \frac{2\log(1/\delta_l)}{n\gamma_{\min}^3}$, we have $h_l(\lambda, \gamma^t, n) - \min_\gamma h_l(\lambda, \gamma, n) \leq \kappa_t$.*

Proof sketch. Lemma G.9 shows that this clipping does not affect the function value that much. Since we do not assume our function to be convex for γ , we will show that the stationary point is unique and the gradient is strictly positive around the stationary point. Lemma G.14 first shows that our function is locally strongly convex around any stationary point. In particular, if we are at a point where the L_1 norm of the gradient is less than λ_{\min} , we are locally strongly convex. Lemma G.13 shows our gradient is Lipschitz with respect to the L_1 norm. Then, Lemma G.8 then shows that the gradient descent algorithm converges to a stationary point. It is the classical argument for gradient descent algorithm on non-convex objectives [22].

Lemma G.8. *For any K , with $L_\gamma = |\mathcal{A}|^2 \frac{((1+\eta_l)\gamma_{\max})^{3/2}}{\eta_l^{3/2} \gamma_{\min}^2} + \frac{2\log(1/\delta_l)}{n\gamma_{\min}^3}$,*

$$\min_{k \leq K} \|\nabla_\gamma h_l(\lambda, \gamma_k, n)\|_1^2 \leq 2L_\gamma \frac{h_l(\lambda, \gamma_0, n) - \min_\gamma h_l(\lambda, \gamma, n)}{K}.$$

With this lemma, we have for a sufficiently large K , the minimum gradient can be made arbitrarily small. In particular, for $K \geq L_\gamma \lambda_{\min}^{-1}$ we have that the minimum gradient has L_1 -norm less than

λ_{\min} , and thus we are in a neighborhood of our stationary point by Lemma G.15. After that, it takes $O(\frac{1}{\kappa_t \lambda_{\min}})$ steps to converge to a point whose value is at most κ_t away from the value of the stationary point. The results in [30] coupled with Lemma G.14 ensure that our stationary point is unique. Intuitively, if we have two locally strongly convex stationary points, there must be a "hill" between them, which also corresponds to a stationary point, but we have shown that all stationary points must be "holes" due to local strong convexity, so the stationary point has to be unique. Thanks to the clipping, we can lower bound λ_{\min} by ι_t , so the total number of steps is $\frac{L}{\lambda_{\min}} + \frac{1}{\kappa_t \lambda_{\min}} = \frac{L}{\iota_t} + \frac{1}{\kappa_t \iota_t}$ which matches the result in Lemma G.7. \square

Lemma G.9. *For some iterate t , let $\iota_t = \epsilon_t^3 t^{-3} |\mathcal{A}|^{-6}$ and denote $\tilde{\lambda} := \text{clip}(\lambda, \iota_t)$ where $[\text{clip}(\lambda, \epsilon)]_\pi := \lambda_\pi \mathbf{1}\{\lambda_\pi \geq \epsilon\}$. Then, for any γ , we have*

$$\left| h_l(\tilde{\lambda}, \gamma, n) - h_l(\lambda, \gamma, n) \right| \leq \kappa_t.$$

Proof. For the first term in h_l , in the case where $\lambda_\pi \geq \iota_t$, $h_l(\lambda, \gamma, n) = h_l(\tilde{\lambda}, \gamma, n)$. When $0 < \lambda_\pi < \iota_t$. We see that

$$\sum_{\pi \in \Pi, \lambda_\pi < \iota_t} \lambda_\pi \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_\pi n} \right) < t\epsilon \left(\frac{1}{\gamma_{\min}} + \frac{1}{\gamma_{\min}} \right) = \frac{2t\iota_t}{\gamma_{\min}}.$$

Then we focus on the expectation part of $h_l(\lambda, \gamma, n)$. Note that

$$\begin{aligned} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} &= \sqrt{\sum_{\pi, \lambda_\pi \geq \iota_t} \lambda_\pi \gamma_\pi [t_a^{(c)} + \eta_l]_\pi + \sum_{\pi, \lambda_\pi < \iota_t} \lambda_\pi \gamma_\pi [t_a^{(c)} + \eta_l]_\pi} \\ &= \sqrt{(\tilde{\lambda} \odot \gamma)^\top (t_a^{(c)} + \eta_l) + \sum_{\pi, \lambda_\pi < \iota_t} \lambda_\pi \gamma_\pi [t_a^{(c)} + \eta_l]_\pi} \\ &\leq \sqrt{(\tilde{\lambda} \odot \gamma)^\top (t_a^{(c)} + \eta_l) + t\iota_t \gamma_{\max}} \\ &\leq \sqrt{(\tilde{\lambda} \odot \gamma)^\top (t_a^{(c)} + \eta_l) + \sqrt{t\iota_t \gamma_{\max}}}. \end{aligned}$$

Therefore,

$$\begin{aligned} &\mathbb{E} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right)^2 \right] - \mathbb{E} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\tilde{\lambda} \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right)^2 \right] \\ &= \mathbb{E} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} + \sqrt{(\tilde{\lambda} \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right) \right. \\ &\quad \left. \left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} - \sqrt{(\tilde{\lambda} \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right) \right] \\ &\leq |\mathcal{A}| \sqrt{\gamma_{\max}} |\mathcal{A}| \sqrt{t\iota_t \gamma_{\max}} \\ &= |\mathcal{A}|^2 \gamma_{\max} \sqrt{t\iota_t}. \end{aligned}$$

Combining two displays above and plugging in γ_{\min} and γ_{\max} gives

$$\begin{aligned} \left| h_l(\tilde{\lambda}, \gamma, n) - h_l(\lambda, \gamma, n) \right| &\leq \frac{2t\iota_t}{\gamma_{\min}} + |\mathcal{A}| \sqrt{\frac{t\iota_t}{\eta_l}} \\ &= \frac{2t\iota_t |\mathcal{A}| \epsilon_l^{-1}}{\sqrt{\eta_l}} + |\mathcal{A}| \sqrt{\frac{t\iota_t}{\eta_l}}. \end{aligned}$$

Let RHS be κ_t and solve for ι_t we get $\iota_t \leq \min\{\frac{\sqrt{\eta_l} \kappa_t \epsilon_l}{2t|\mathcal{A}|}, \frac{\eta_l \kappa_t}{|\mathcal{A}|^2 t}\}$. Plugging in $\eta_l = |\mathcal{A}|^{-4} \epsilon_l^2$ gives the result. \square

Lemma G.10. Suppose γ^t satisfies that $h_l(\tilde{\lambda}, \gamma^t, n) - \min_{\gamma} h_l(\tilde{\lambda}, \gamma, n) \leq \kappa_t$, then we also have $h_l(\lambda, \gamma^t, n) - \min_{\gamma} h_l(\lambda, \gamma, n) \leq \kappa_t$, i.e. γ^t satisfies the desired property.

Proof. Let $\tilde{\gamma}_* = \arg \min_{\gamma} h_l(\tilde{\lambda}, \gamma, n)$ and $\gamma_* = \arg \min_{\gamma} h_l(\lambda, \gamma, n)$. The result follows from applying Lemma G.9 twice on $h_l(\tilde{\lambda}, \gamma^t, n)$ and $h_l(\tilde{\lambda}, \gamma_*, n)$. In particular,

$$\begin{aligned}
h_l(\lambda, \gamma^t, n) &\leq h_l(\tilde{\lambda}, \gamma^t, n) + \kappa_t && \text{(Lemma G.9)} \\
&\leq h_l(\tilde{\lambda}, \tilde{\gamma}_*, n) + 2\kappa_t && \text{(convergence of GD)} \\
&\leq h_l(\tilde{\lambda}, \gamma_*, n) + 2\kappa_t && \text{(minimality of } \tilde{\gamma}_*) \\
&\leq h_l(\lambda, \gamma_*, n) + 3\kappa_t && \text{(Lemma G.9)} \\
&= \min_{\gamma} h_l(\lambda, \gamma, n) + 3\kappa_t.
\end{aligned}$$

□

G.4 Guarantees for strong concavity and local strong convexity

The following series of lemmas show that our optimization problem is strongly concave in λ and local strongly convex around the minimum γ , as well as explicitly constructing the Lipschitz constants. These serve as the conditions for convergence of the Frank-Wolfe and gradient descent algorithms.

Lemma G.11. $h_l(\lambda, \gamma, n)$ is a concave function of λ .

Proof. Note that

$$\mathbb{E} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right)^2 \right] = \mathbb{E} \left[\sum_{a \in \mathcal{A}} \sum_{a' \in \mathcal{A}} \sqrt{(t_{a'}^{(c)} + \eta_l)^\top (\lambda \odot \gamma) (\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right].$$

we know that $\lambda \mapsto (t_{a'}^{(c)} + \eta_l)^\top (\lambda \odot \gamma)$ and $\lambda \mapsto (\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)$ are concave, the square root function is concave and non-decreasing, and sum of concave functions is concave. Therefore, $h_l(\lambda, \gamma, n)$ is concave in λ by property of concave functions. □

Lemma G.12. Consider some λ, γ and n . For any $\lambda_1, \lambda_2 \in \Delta_{\Pi}$, with $L = |\mathcal{A}|^2 \frac{((1+\eta_l)\gamma_{\max})^{5/2}}{\eta_l^{3/2} \gamma_{\min}^2}$,

$$f(\lambda_2, \gamma, n) \leq f(\lambda_1, \gamma, n) + \nabla_{\lambda} f(\lambda_1, \gamma, n)^\top (\lambda_2 - \lambda_1) + L \|\lambda_2 - \lambda_1\|_1^2,$$

where $f(\lambda, \gamma, n)$ could be either $h_l(\lambda, \gamma, n)$ or $-h_l(\lambda, \gamma, n)$.

Proof. The proof for the negative case is exactly the same as the positive case, so we focus on $f(\lambda, \gamma, n) = h_l(\lambda, \gamma, n)$. We take the gradient of h_l with respect to λ and get

$$\begin{aligned}
[\nabla_{\lambda} h_l(\lambda, \gamma, n)]_{\pi} &= -\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} \\
&\quad + \mathbb{E}_{c \sim \nu_D} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right) \left(\sum_{a' \in \mathcal{A}} \frac{\gamma_{\pi} (t_{a'}^{(c)} + \eta_l)_{\pi}}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}} \right) \right].
\end{aligned}$$

By Lemma H.2, for any $\lambda \in \Delta_{\Pi}$, we have $\langle \lambda, \nabla_{\lambda} h_l(\lambda, \gamma, n) \rangle = h_l(\lambda, \gamma, n)$. If we use the shortcut $f(\lambda) := h_l(\lambda, \gamma, n)$, we have

$$f(\lambda_2) - f(\lambda_1) - \nabla_{\lambda} f(\lambda_1)^\top (\lambda_2 - \lambda_1) = f(\lambda_2) - \nabla_{\lambda} f(\lambda_1)^\top \lambda_2 = (\nabla f(\lambda_2) - \nabla f(\lambda_1))^\top \lambda_2.$$

Note that

$$\begin{aligned}
& (\nabla_{\lambda} f(\lambda_2) - \nabla_{\lambda} f(\lambda_1))^{\top} \lambda_2 \\
&= \sum_{\pi \in \Pi} [\lambda_2]_{\pi} \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right) \left(\sum_{a' \in \mathcal{A}} \frac{\gamma_{\pi} \cdot (t_{a'}^{(c)} + \eta_l)_{\pi}}{\sqrt{(\lambda_2 \odot \gamma)^{\top} (t_{a'}^{(c)} + \eta_l)}} \right) \right. \\
&\quad \left. - \left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right) \left(\sum_{a' \in \mathcal{A}} \frac{\gamma_{\pi} \cdot (t_{a'}^{(c)} + \eta_l)_{\pi}}{\sqrt{(\lambda_1 \odot \gamma)^{\top} (t_{a'}^{(c)} + \eta_l)}} \right) \right] \\
&= \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_{a' \in \mathcal{A}} (\lambda_2 \odot \gamma)^{\top} (t_{a'}^{(c)} + \eta_l) \right. \\
&\quad \cdot \left. \sum_{a \in \mathcal{A}} \frac{\sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} - \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)}}{\sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)}} \right] \\
&\leq \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_{a' \in \mathcal{A}} (\lambda_2 \odot \gamma)^{\top} (t_{a'}^{(c)} + \eta_l) \right. \\
&\quad \cdot \left. \sum_{a \in \mathcal{A}} \frac{\left| \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} - \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right|}{\sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)}} \right] \\
&\leq \sum_{a' \in \mathcal{A}} \frac{(1 + \eta_l) \gamma_{\max}}{\eta_l \gamma_{\min}} \cdot \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\sum_{a \in \mathcal{A}} \left| \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right. \right. \\
&\quad \left. \left. - \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right| \right] \tag{18}
\end{aligned}$$

Note that by triangular inequality

$$\begin{aligned}
& \left| \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} - \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right| \\
&\leq \left| \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} - \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right| \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \\
&\quad + \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \left| \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} - \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right|.
\end{aligned}$$

Also note that

$$\begin{aligned}
& \left| \sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} - \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} \right| \\
&= \frac{\left| \sum_{\pi \in \Pi} ((\lambda_2)_{\pi} - (\lambda_1)_{\pi}) \gamma_{\pi} (t_a^{(c)} + \eta_l)_{\pi} \right|}{\sqrt{(\lambda_2 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)} + \sqrt{(\lambda_1 \odot \gamma)^{\top} (t_a^{(c)} + \eta_l)}} \\
&\leq \frac{(1 + \eta_l) \gamma_{\max}}{2\sqrt{\eta_l} \gamma_{\min}} \|\lambda_2 - \lambda_1\|_1,
\end{aligned}$$

so (18) is bounded by

$$\begin{aligned}
& \sum_{a' \in \mathcal{A}} \frac{(1 + \eta_l) \gamma_{\max}}{\eta_l \gamma_{\min}} \cdot \left(\sum_{a \in \mathcal{A}} 2 \cdot \frac{(1 + \eta_l) \gamma_{\max}}{2\sqrt{\eta_l} \gamma_{\min}} \|\lambda_2 - \lambda_1\|_1 \sqrt{(1 + \eta_l) \gamma_{\max}} \right) \\
&= |\mathcal{A}|^2 \frac{((1 + \eta_l) \gamma_{\max})^{5/2}}{\eta_l^{3/2} \gamma_{\min}^2} \|\lambda_2 - \lambda_1\|_1.
\end{aligned}$$

□

Lemma G.13. Consider some λ and n . For any $\gamma_1, \gamma_2 \in \Delta_\Pi$, with $L_\gamma = |\mathcal{A}|^2 \frac{((1+\eta_l)\gamma_{\max})^{3/2}}{\eta_l^{3/2}\gamma_{\min}^2} + \frac{2 \log(1/\delta_l)}{n\gamma_{\min}^3}$,

$$h_l(\lambda, \gamma_2, n) \leq h_l(\lambda, \gamma_1, n) + \nabla_\gamma h_l(\lambda, \gamma_1, n)^\top (\gamma_2 - \gamma_1) + L_\gamma \|\gamma_2 - \gamma_1\|_1^2.$$

Proof.

$$[\nabla_\gamma h_l(\lambda, \gamma)]_\pi = \mathbb{E}_c \left[\left(\sum_{a' \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)} \right) \cdot \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_\pi ([t_{a'}^{(c)}]_\pi + \eta_l)}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}} \right) \right] - \frac{\lambda_\pi \log(1/\delta_l)}{\gamma_\pi^2 n}.$$

Then we have similar to the proof of Lemma G.12, for any γ we have $h_l(\lambda, \gamma, n) - \nabla_\gamma h_l(\lambda, \gamma, n)^\top \gamma = 2 \sum_\pi \frac{\lambda_\pi \log(1/\delta_l)}{\gamma_\pi^2 n}$, so

$$\begin{aligned} & h_l(\lambda, \gamma_2, n) - h_l(\lambda, \gamma_1, n) - \nabla_\gamma h_l(\lambda, \gamma_1, n)^\top (\gamma_2 - \gamma_1) \\ &= 2 \sum_\pi \frac{\lambda_\pi \log(1/\delta_l)}{[\gamma_2]_\pi^2 n} - 2 \sum_\pi \frac{\lambda_\pi \log(1/\delta_l)}{[\gamma_1]_\pi^2 n} + (\nabla_\gamma h_l(\lambda, \gamma_2, n) - \nabla_\gamma h_l(\lambda, \gamma_1, n))^\top \gamma_2. \end{aligned}$$

First, we can follow similar techniques in the proof of Lemma G.12 to bound the second part and get

$$\begin{aligned} & (\nabla_\gamma h_l(\lambda, \gamma_2, n) - \nabla_\gamma h_l(\lambda, \gamma_1, n))^\top \gamma_2 \\ & \leq \sum_{a' \in \mathcal{A}} (\lambda \odot \gamma_2)^\top (t_{a'}^{(c)} + \eta_l) \\ & \quad \cdot \mathbb{E}_{c \sim \nu_D} \left\{ \sum_{a \in \mathcal{A}} \left[\frac{1}{\sqrt{(\lambda \odot \gamma_2)^\top (t_{a'}^{(c)} + \eta_l)} \sqrt{(\lambda \odot \gamma_1)^\top (t_{a'}^{(c)} + \eta_l)}} \right. \right. \\ & \quad \left. \left. \cdot \left| \sqrt{(\lambda \odot \gamma_1)^\top (t_{a'}^{(c)} + \eta_l)} \sqrt{(\lambda \odot \gamma_2)^\top (t_{a'}^{(c)} + \eta_l)} - \sqrt{(\lambda \odot \gamma_2)^\top (t_{a'}^{(c)} + \eta_l)} \sqrt{(\lambda \odot \gamma_1)^\top (t_{a'}^{(c)} + \eta_l)} \right| \right] \right\} \\ & \leq \sum_{a' \in \mathcal{A}} \frac{(1 + \eta_l)\gamma_{\max}}{\eta_l \gamma_{\min}} \cdot \mathbb{E}_{c \sim \nu_D} \left[\sum_{a \in \mathcal{A}} \left| \sqrt{(\lambda \odot \gamma_2)^\top (t_{a'}^{(c)} + \eta_l)} \sqrt{(\lambda \odot \gamma_1)^\top (t_{a'}^{(c)} + \eta_l)} \right. \right. \\ & \quad \left. \left. - \sqrt{(\lambda \odot \gamma_1)^\top (t_{a'}^{(c)} + \eta_l)} \sqrt{(\lambda \odot \gamma_2)^\top (t_{a'}^{(c)} + \eta_l)} \right| \right]. \end{aligned}$$

Also, note that

$$\begin{aligned} & \left| \sqrt{(\lambda \odot \gamma_2)^\top (t_{a'}^{(c)} + \eta_l)} - \sqrt{(\lambda \odot \gamma_1)^\top (t_{a'}^{(c)} + \eta_l)} \right| \\ &= \frac{\left| \sum_{\pi \in \Pi} (\lambda_\pi ([\gamma_2]_\pi - [\gamma_1]_\pi) (t_{a'}^{(c)})_\pi) \right|}{\sqrt{(\lambda \odot \gamma_2)^\top (t_{a'}^{(c)} + \eta_l)} + \sqrt{(\lambda \odot \gamma_1)^\top (t_{a'}^{(c)} + \eta_l)}} \\ &\leq \frac{1}{2\sqrt{\eta_l \gamma_{\min}}} \|\gamma_2 - \gamma_1\|_1^2, \end{aligned}$$

Therefore, similarly we can bound

$$\begin{aligned} & \left| \sqrt{(\lambda \odot \gamma_2)^\top (t_{a'}^{(c)} + \eta_l)} \sqrt{(\lambda \odot \gamma_1)^\top (t_{a'}^{(c)} + \eta_l)} - \sqrt{(\lambda \odot \gamma_1)^\top (t_{a'}^{(c)} + \eta_l)} \sqrt{(\lambda \odot \gamma_2)^\top (t_{a'}^{(c)} + \eta_l)} \right| \\ & \leq \frac{\sqrt{(1 + \eta_l)\gamma_{\max}}}{2\sqrt{\eta_l \gamma_{\min}}} \|\gamma_2 - \gamma_1\|_1^2. \end{aligned}$$

For the second term,

$$\begin{aligned} & 2 \sum_\pi \frac{\lambda_\pi \log(1/\delta_l)}{[\gamma_2]_\pi^2 n} - 2 \sum_\pi \frac{\lambda_\pi \log(1/\delta_l)}{[\gamma_1]_\pi^2 n} \\ &= \frac{2 \log(1/\delta_l)}{n} \sum_\pi \lambda_\pi \frac{[\gamma_1]_\pi^2 - [\gamma_2]_\pi^2}{[\gamma_1]_\pi^2 [\gamma_2]_\pi^2} \\ &\leq \frac{2 \log(1/\delta_l)}{n\gamma_{\min}^3} \|\gamma_2 - \gamma_1\|_1^2. \end{aligned}$$

Therefore, we have the result stated above. \square

Lemma G.14. Consider some fixed $\lambda \in \Delta_\Pi$ and n . Assume γ_* is a stationary point of $h_l(\lambda, \gamma, n)$, then $h_l(\lambda, \gamma, n)$ is locally strongly convex at γ_* , i.e. for $L_{\text{hess}} = \frac{\lambda_{\min} \log(1/\delta_l)}{\gamma_{\max}^3 n}$, there exists $\epsilon > 0$ such that for all $\gamma \in B_\epsilon(\gamma_*)$, $h_l(\lambda, \gamma, n) \geq h_l(\lambda, \gamma_*, n) + \frac{L_{\text{hess}}}{2} \|\gamma - \gamma_*\|^2$.

Proof. Since λ and n are fixed, we use the shortcut $g(\gamma) := h_l(\lambda, \gamma, n)$ in the proof. Denote the Hessian of g as M . We aim to show that the Hessian $M \succeq L_{\text{hess}} I$ at γ_* . First, since γ_* is a stationary point, $\nabla_\gamma g(\gamma_*) = 0$, and so for any i ,

$$\sum_{c \in \mathcal{D}} \nu_{c\mathcal{D}} \left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right) \cdot \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i [t_{a'}^{(c)}]_i + \eta_l}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}} \right) = \frac{\lambda_i \log(1/\delta_l)}{\gamma_i^2 n}. \quad (19)$$

Also, we have for $i \neq j$,

$$\begin{aligned} \frac{\partial^2 g(\gamma)}{\partial \gamma_i \partial \gamma_j} &= \sum_{c \in \mathcal{D}} \nu_{c\mathcal{D}} \left(\sum_{a' \in \mathcal{A}} \frac{1}{2} \frac{\lambda_i [t_{a'}^{(c)} + \eta_l]_i}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}} \right) \cdot \left(\sum_{a \in \mathcal{A}} \frac{\lambda_j [t_a^{(c)} + \eta_l]_j}{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)}} \right) \\ &\quad + \left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right) \cdot \left(\sum_{a' \in \mathcal{A}} -\frac{1}{2} \cdot \frac{\lambda_i \lambda_j [t_{a'}^{(c)} + \eta_l]_i [t_{a'}^{(c)} + \eta_l]_j}{((\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l))^{3/2}} \right). \end{aligned}$$

And

$$\begin{aligned} \frac{\partial^2 g(\gamma)}{\partial \gamma_i^2} &= \frac{2\lambda_i \log(1/\delta_l)}{\gamma_i^3 n} + \sum_{c \in \mathcal{D}} \nu_{c\mathcal{D}} \frac{1}{2} \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i [t_{a'}^{(c)} + \eta_l]_i}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}} \right)^2 \\ &\quad - \frac{1}{2} \left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right) \cdot \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i^2 [t_{a'}^{(c)} + \eta_l]_i^2}{((\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l))^{3/2}} \right). \end{aligned}$$

Then, for any vector $\mu \in \mathbb{R}^{|\Pi|}$ with $\|\mu\| = 1$, we have

$$\begin{aligned} \mu^\top M \mu &= \sum_i \sum_j \mu_i \mu_j M_{ij} = \sum_i \mu_i^2 M_{ii} + \sum_{i \neq j} \mu_i \mu_j M_{ij} \\ &= \sum_i \mu_i^2 \frac{2\lambda_i \log(1/\delta_l)}{\gamma_i^3 n} \end{aligned} \quad (20)$$

$$\begin{aligned} &+ \sum_c \nu_c \sum_i \sum_j \mu_i \mu_j \frac{1}{2} \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i [t_{a'}^{(c)} + \eta_l]_i}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}} \right) \cdot \left(\sum_{a \in \mathcal{A}} \frac{\lambda_j [t_a^{(c)} + \eta_l]_j}{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)}} \right) \\ &+ \mu_i \mu_j \left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right) \cdot \left(\sum_{a' \in \mathcal{A}} -\frac{1}{2} \cdot \frac{\lambda_i \lambda_j [t_{a'}^{(c)} + \eta_l]_i [t_{a'}^{(c)} + \eta_l]_j}{((\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l))^{3/2}} \right). \end{aligned} \quad (21)$$

In what follows, we will first show that

$$\begin{aligned} &\sum_i \mu_i^2 \frac{\lambda_i \log(1/\delta_l)}{\gamma_i^3 n} - \sum_c \nu_c \sum_i \sum_j \mu_i \mu_j \left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right) \\ &\quad \cdot \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i \lambda_j [t_{a'}^{(c)} + \eta_l]_i [t_{a'}^{(c)} + \eta_l]_j}{((\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l))^{3/2}} \right) \geq 0. \end{aligned} \quad (22)$$

By equation 19, the LHS of (21) simplifies to

$$\begin{aligned} & \sum_c \nu_c \sum_i \mu_i^2 \frac{1}{\gamma_i} \left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta)} \right) \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i [t_{a'}^{(c)} + \eta]_i}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta)}} \right) \\ & - \sum_c \nu_c \sum_i \sum_j \mu_i \mu_j \left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta)} \right) \cdot \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i \lambda_j [t_{a'}^{(c)} + \eta]_i [t_{a'}^{(c)} + \eta]_j}{((\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta))^{3/2}} \right). \end{aligned}$$

Therefore, it is sufficient to show that

$$\sum_i \mu_i^2 \frac{1}{\gamma_i} \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i [t_{a'}^{(c)} + \eta]_i}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta)}} \right) - \sum_i \sum_j \mu_i \mu_j \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i \lambda_j [t_{a'}^{(c)} + \eta]_i [t_{a'}^{(c)} + \eta]_j}{((\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta))^{3/2}} \right) \geq 0.$$

Consider some $a' \in \mathcal{A}$. The LHS of the above simplifies to

$$\begin{aligned} & \sum_i \mu_i^2 \frac{1}{\gamma_i} \frac{\lambda_i [t_{a'}^{(c)} + \eta]_i}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta)}} - \sum_i \sum_j \mu_i \mu_j \frac{\lambda_i \lambda_j [t_{a'}^{(c)} + \eta]_i [t_{a'}^{(c)} + \eta]_j}{((\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta))^{3/2}} \\ & = \frac{1}{((\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta))^{3/2}} \left(\sum_i \frac{\mu_i^2}{\gamma_i} \lambda_i [t_{a'}^{(c)} + \eta]_i \left(\sum_j \lambda_j \gamma_j [t_{a'}^{(c)} + \eta]_j \right) \right. \\ & \quad \left. - \sum_i \sum_j \mu_i \mu_j \lambda_i \lambda_j [t_{a'}^{(c)} + \eta]_i [t_{a'}^{(c)} + \eta]_j \right) \\ & = \frac{1}{((\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta))^{3/2}} \left(\sum_i \sum_j \gamma_i^{-1} \left(\mu_i^2 \lambda_i [t_{a'}^{(c)} + \eta]_i \lambda_j \gamma_j [t_{a'}^{(c)} + \eta]_j \right. \right. \\ & \quad \left. \left. - \mu_i \mu_j \lambda_i \lambda_j \gamma_i [t_{a'}^{(c)} + \eta]_i [t_{a'}^{(c)} + \eta]_j \right) \right). \end{aligned}$$

Each summand is

$$\begin{aligned} & \gamma_i^{-1} \left(\mu_i^2 \lambda_i [t_{a'}^{(c)} + \eta]_i \lambda_j \gamma_j [t_{a'}^{(c)} + \eta]_j - \mu_i \mu_j \lambda_i \lambda_j \gamma_i [t_{a'}^{(c)} + \eta]_i [t_{a'}^{(c)} + \eta]_j \right) \\ & = \gamma_i^{-1} \mu_i \lambda_i \lambda_j [t_{a'}^{(c)} + \eta]_i [t_{a'}^{(c)} + \eta]_j (\mu_i \gamma_j - \mu_j \gamma_i) \\ & = \gamma_i^{-1} \gamma_j^{-1} \lambda_i \lambda_j [t_{a'}^{(c)} + \eta]_i [t_{a'}^{(c)} + \eta]_j (\mu_i \gamma_j) (\mu_i \gamma_j - \mu_j \gamma_i). \end{aligned}$$

Exchanging subscripts of i and j , we have

$$\gamma_j^{-1} \gamma_i^{-1} \lambda_j \lambda_i [t_{a'}^{(c)} + \eta]_j [t_{a'}^{(c)} + \eta]_i (\mu_j \gamma_i) (\mu_j \gamma_i - \mu_i \gamma_j).$$

The sum of these two terms is

$$\gamma_i^{-1} \gamma_j^{-1} \lambda_i \lambda_j [t_{a'}^{(c)} + \eta]_i [t_{a'}^{(c)} + \eta]_j (\mu_i \gamma_j - \mu_j \gamma_i)^2 \geq 0.$$

Therefore, we proved equation (22). We will show next that

$$\begin{aligned} & \sum_i \mu_i^2 \frac{\lambda_i \log(1/\delta_i)}{\gamma_i^3 n} + \sum_c \nu_c \sum_i \sum_j \mu_i \mu_j \frac{1}{2} \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i [t_{a'}^{(c)} + \eta]_i}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta)}} \right) \\ & \cdot \left(\sum_{a \in \mathcal{A}} \frac{\lambda_j [t_a^{(c)} + \eta]_j}{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta)}} \right) \geq 0. \end{aligned} \tag{23}$$

By similar calculation, we can obtain that the above simplifies to

$$\sum_c \nu_c \sum_i \mu_i \gamma_i^{-1} \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i [t_{a'}^{(c)} + \eta]_i}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta)}} \right) \cdot \left\{ \mu_i \sum_{a \in \mathcal{A}} \frac{\sum_j \lambda_j \gamma_j [t_a^{(c)} + \eta]_j}{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta)}} + \mu_j \gamma_i \sum_{a \in \mathcal{A}} \frac{\sum_j \lambda_j [t_a^{(c)} + \eta]_j}{\sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta)}} \right\}.$$

We can show that the sum of the above is positive by similar techniques for showing (22). Plugging equation 22 and 23 in equation 21, we have that

$$\mu^\top M \mu \geq \sum_i \mu_i^2 \frac{\lambda_i \log(1/\delta_i)}{\gamma_i^3 n} \geq \frac{\lambda_{\min} \log(1/\delta_i)}{\gamma_{\max}^3 n},$$

so the Hessian is positive-definite. \square

Note that the minimum eigenvalue of the Hessian at the stationary point is $\frac{\lambda_{\min} \log(1/\delta_i)}{\gamma_{\max}^3 n} > 0$, we can extend the result in Lemma G.14 to α -stationary points, where $\alpha < \frac{\lambda_{\min} \log(1/\delta_i)}{\gamma_{\max}^3 n}$, and still maintain local strong convexity.

Lemma G.15. *Consider some fixed $\lambda \in \Delta_\Pi$ and n . Assume γ_α is an α -stationary point of $h_l(\lambda, \gamma, n)$, where $\alpha = \frac{\lambda_{\min} \log(1/\delta_i)}{2\gamma_{\max}^3 n}$, then $h_l(\lambda, \gamma, n)$ is locally strongly convex at γ_α , i.e. for $L_{\text{hess}} = \frac{\lambda_{\min} \log(1/\delta_i)}{2\gamma_{\max}^3 n}$, there exists $\epsilon > 0$ such that for all $\gamma \in B_\epsilon(\gamma_\alpha)$, $h_l(\lambda, \gamma, n) \geq h_l(\lambda, \gamma_\alpha, n) + \frac{L_{\text{hess}}}{2} \|\gamma - \gamma_\alpha\|^2$.*

Proof. The proof follows almost identically from that of Lemma G.14. Note that the α -stationary point ensures that $\|\nabla_\gamma h_l(\lambda, \gamma)\|_1 \leq \alpha$, so equation 19 is rewritten as

$$\sum_i \left| \sum_{c \in \mathcal{D}} \nu_{c\mathcal{D}} \left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta)} \right) \cdot \left(\sum_{a' \in \mathcal{A}} \frac{\lambda_i ([t_{a'}^{(c)}]_i + \eta)}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta)}} \right) - \frac{\lambda_i \log(1/\delta_i)}{\gamma_i^2 n} \right| \leq \alpha. \quad (24)$$

Therefore, for any μ we can still use the same trick and get

$$\mu^\top M \mu \geq \sum_i \mu_i^2 \frac{\lambda_i \log(1/\delta_i)}{\gamma_i^3 n} - \alpha \geq \frac{\lambda_{\min} \log(1/\delta_i)}{2\gamma_{\max}^3 n},$$

so our result follows. \square

G.5 Proof of strong duality

In this section, we would like to show that strong duality holds. We first show that the primal problem is convex for w .

Lemma G.16. *The primal problem (12) is convex for w .*

Proof. Note that the primal problem could be written as

$$\min_{w \in \Omega} c \quad \text{s.t. } \forall \pi \in \Pi, -\Delta(\pi) + \sqrt{\frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{n}} \leq c.$$

Therefore, we consider the function $f(w) := -\Delta(\pi) + \sqrt{\frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{n}}$ for some $\pi \in \Pi$. Note that to show that $f(w) = -\Delta(\pi) + \sqrt{\frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{n}}$ is convex for w , it is equivalent to show that

$g(w) := \sqrt{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}$ is convex for w . Note that

$$\begin{aligned} g(w) &= \sqrt{\sum_{a,c} \nu_c^2 w_{a,c}^{-1} (\mathbf{1}\{\pi(c) = a, \pi_*(c) \neq a\} + \mathbf{1}\{\pi(c) \neq a, \pi_*(c) = a\})} \\ &= \sqrt{\sum_{a,c, t_a^{(c)}=1} \nu_c^2 w_{a,c}^{-1}}. \end{aligned}$$

So restricting to a, c such that $t_a^{(c)} = 1$

$$\frac{\partial g(w)}{\partial w_{a,c}} = \frac{1}{2\sqrt{\sum_{a,c, t_a^{(c)}=1} \nu_c^2 w_{a,c}^{-1}}} \cdot (-\nu_c^2 w_{a,c}^{-2}),$$

and

$$\begin{aligned} \frac{\partial^2 g(w)}{\partial w_{a,c}^2} &= -\frac{1}{4\left(\sum_{a,c, t_a^{(c)}=1} \nu_c^2 w_{a,c}^{-1}\right)^{3/2}} \cdot (-\nu_c^2 w_{a,c}^{-2} \cdot -\nu_c^2 w_{a,c}^{-2}) + \frac{1}{\sqrt{\sum_{a,c, t_a^{(c)}=1} \nu_c^2 w_{a,c}^{-1}}} \cdot \nu_c^2 w_{a,c}^{-3} \\ \frac{\partial^2 g(w)}{\partial w_{a_1, c_1} \partial w_{a_2, c_2}} &= -\frac{1}{4\left(\sum_{a,c, t_a^{(c)}=1} \nu_c^2 w_{a,c}^{-1}\right)^{3/2}} \cdot (-\nu_{c_1}^2 w_{a_1, c_1}^{-2} \cdot -\nu_{c_2}^2 w_{a_2, c_2}^{-2}) \end{aligned}$$

Denote the Hessian as M . Then, for any vector $\mu \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{C}|}$ with $\|\mu\|_2 = 1$, we have

$$\begin{aligned} \mu^\top M \mu &= -\frac{1}{4} \sum_{a,c, t_a^{(c)}=1} \sum_{a', c', t_{a'}^{(c')}=1} \mu_{a,c} \mu_{a', c'} \left(\sum_{a,c, t_a^{(c)}=1} \nu_c^2 w_{a,c}^{-1} \right)^{-3/2} \nu_c^2 \nu_{c'}^2 w_{a,c}^{-2} w_{a', c'}^{-2} \\ &\quad + \sum_{a,c, t_a^{(c)}=1} \mu_{a,c}^2 \nu_c^2 w_{a,c}^{-3} \left(\sum_{a,c, t_a^{(c)}=1} \nu_c^2 w_{a,c}^{-1} \right)^{-1/2}. \end{aligned}$$

To show that this is nonnegative, it is equivalent to show that

$$-\frac{1}{4} \sum_{a,c, t_a^{(c)}=1} \sum_{a', c', t_{a'}^{(c')}=1} \mu_{a,c} \mu_{a', c'} \nu_c^2 \nu_{c'}^2 w_{a,c}^{-2} w_{a', c'}^{-2} + \sum_{a,c, t_a^{(c)}=1} \mu_{a,c}^2 \nu_c^2 w_{a,c}^{-3} \left(\sum_{a', c', t_{a'}^{(c')}=1} \nu_{c'}^2 w_{a', c'}^{-1} \right) \geq 0,$$

which is equivalent to show that

$$\sum_{a,c, t_a^{(c)}=1} \sum_{a', c', t_{a'}^{(c')}=1} -\mu_{a,c} \mu_{a', c'} \nu_c^2 \nu_{c'}^2 w_{a,c}^{-2} w_{a', c'}^{-2} + \mu_{a,c}^2 \nu_c^2 w_{a,c}^{-3} \nu_{c'}^2 w_{a', c'}^{-1} \geq 0. \quad (25)$$

Note that

$$\begin{aligned} &-\mu_{a,c} \mu_{a', c'} \nu_c^2 \nu_{c'}^2 w_{a,c}^{-2} w_{a', c'}^{-2} + \mu_{a,c}^2 \nu_c^2 w_{a,c}^{-3} \nu_{c'}^2 w_{a', c'}^{-1} \\ &= \mu_{a,c} w_{a,c}^{-3} w_{a', c'}^{-2} \nu_c^2 \nu_{c'}^2 (\mu_{a,c} w_{a', c'} - \mu_{a', c'} w_{a,c}) \\ &= w_{a,c}^{-3} w_{a', c'}^{-3} \nu_c^2 \nu_{c'}^2 (\mu_{a,c} w_{a', c'}) (\mu_{a,c} w_{a', c'} - \mu_{a', c'} w_{a,c}). \end{aligned}$$

Then, exchanging the label of a and a' , we also get a term like

$$w_{a', c'}^{-3} w_{a,c}^{-3} \nu_{c'}^2 \nu_c^2 (\mu_{a', c'} w_{a,c}) (\mu_{a', c'} w_{a,c} - \mu_{a,c} w_{a', c'}).$$

The sum of these two terms is

$$\begin{aligned} &w_{a', c'}^{-3} w_{a,c}^{-3} \nu_{c'}^2 \nu_c^2 (\mu_{a', c'} w_{a,c}) (\mu_{a', c'} w_{a,c} - \mu_{a,c} w_{a', c'}) \\ &+ w_{a,c}^{-3} w_{a', c'}^{-3} \nu_c^2 \nu_{c'}^2 (\mu_{a,c} w_{a', c'}) (\mu_{a,c} w_{a', c'} - \mu_{a', c'} w_{a,c}) \\ &= w_{a', c'}^{-3} w_{a,c}^{-3} \nu_c^2 \nu_{c'}^2 (\mu_{a', c'} w_{a,c} - \mu_{a,c} w_{a', c'}) (\mu_{a', c'} w_{a,c} - \mu_{a,c} w_{a', c'}) \\ &= w_{a', c'}^{-3} w_{a,c}^{-3} \nu_{c'}^2 \nu_c^2 (\mu_{a', c'} w_{a,c} - \mu_{a,c} w_{a', c'})^2 \geq 0. \end{aligned}$$

Therefore, equation 25 becomes

$$\begin{aligned}
& \sum_{a,c,t_a^{(c)}=1} \sum_{\substack{a',c' \\ t_{a'}^{(c')}=1 \\ (a',c') > (a,c)}} (w_{a',c'}^{-3} w_{a,c}^{-3} \nu_c^2 \nu_c^2 (\mu_{a',c'} w_{a,c}) (\mu_{a',c'} w_{a,c} - \mu_{a,c} w_{a',c'})) \\
& + w_{a,c}^{-3} w_{a',c'}^{-3} \nu_c^2 \nu_c^2 (\mu_{a,c} w_{a',c'}) (\mu_{a,c} w_{a',c'} - \mu_{a',c'} w_{a,c}) \\
& = \sum_{a,c,t_a^{(c)}=1} \sum_{\substack{a',c' \\ t_{a'}^{(c')}=1 \\ (a',c') > (a,c)}} w_{a',c'}^{-3} w_{a,c}^{-3} \nu_c^2 \nu_c^2 (\mu_{a',c'} w_{a,c} - \mu_{a,c} w_{a',c'})^2 \geq 0.
\end{aligned}$$

Since the above holds for any vector μ , the Hessian is positive-semidefinite, and so the function $g(w)$ is convex for w . \square

Lemma G.17. *In the optimization problem 12, the strong duality holds, i.e.*

$$\min_{w \in \Omega} \max_{\pi \in \Pi} \left(-\Delta(\pi) + \sqrt{\frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{n}} \right) = \max_{\lambda \in \Delta_\Pi} \min_{w \in \Omega} \sum_{\pi \in \Pi} \lambda_\pi \left(-\Delta(\pi) + \sqrt{\frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{n}} \right).$$

Proof. By Lemma G.16, the primal problem is convex for w , so it is left to check the KKT conditions. Note that the lagrangian is

$$\mathcal{L}(w, \lambda, c) = c + \sum_{\pi \in \Pi} \lambda_\pi \cdot \left(-\Delta(\pi) + \sqrt{\frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{n}} - c \right).$$

Let $h_\pi(w) = -\Delta(\pi) + \sqrt{\frac{\|\phi_\pi - \phi_{\pi_*}\|_{A(w)^{-1}}^2}{n}} - c$. At an optimal solution w^* and λ^* , we would like to show that

$$\sum_{\pi \in \Pi} \lambda_\pi^* h_\pi(w^*) = 0.$$

We prove this by contradiction. If there is some π such that $\lambda_\pi > 0$ and $h_\pi(w^*) < 0$. Then we could find another $\lambda' \in \Delta_\Pi$ that places zero mass on this π and thus get a larger objective, so we get a contradiction. The other conditions follow from the optimality of w^* and λ^* . \square

H Useful lemmas

In this section, we state several algebraic facts of our function, which serves as the key to derive convergence as well as complexity.

Lemma H.1. *For any l ,*

$$\min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\hat{\pi}_{l-1}} - \phi_\pi\|_{A(w)^{-1}}^2}{\Delta(\pi)^2} = \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \frac{\mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{c, \hat{\pi}_{l-1}(c)}} + \frac{1}{p_{c, \pi(c)}} \right) \mathbf{1}\{\hat{\pi}_{l-1}(c) \neq \pi(c)\} \right]}{\Delta(\pi)^2}.$$

Proof. Let $w_{a,c} = \nu_c p_{c,a}$ for some $p_c \in \Delta_{\mathcal{A}}$. Then, for any $\pi \in \Pi$,

$$\begin{aligned}
& \frac{1}{\Delta(\pi)^2} \|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w)^{-1}}^2 \\
&= \frac{1}{\Delta(\pi)^2} \sum_{a,c} \frac{\nu_c^2}{w_{a,c}} (\mathbf{1}\{\widehat{\pi}_{l-1}(c) = a, \pi(c) \neq a\} + \mathbf{1}\{\widehat{\pi}_{l-1}(c) \neq a, \pi(c) = a\}) \\
&= \frac{1}{\Delta(\pi)^2} \sum_{a,c} \frac{\nu_c}{p_{c,a}} (\mathbf{1}\{\widehat{\pi}_{l-1}(c) = a, \pi(c) \neq a\} + \mathbf{1}\{\widehat{\pi}_{l-1}(c) \neq a, \pi(c) = a\}) \\
&= \frac{1}{\Delta(\pi)^2} \sum_c \nu_c \left(\frac{1}{p_{c,\widehat{\pi}_{l-1}(c)}} + \frac{1}{p_{c,\pi(c)}} \right) \mathbf{1}\{\widehat{\pi}_{l-1}(c) \neq \pi(c)\} \\
&= \frac{1}{\Delta(\pi)^2} \mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{c,\widehat{\pi}_{l-1}(c)}} + \frac{1}{p_{c,\pi(c)}} \right) \mathbf{1}\{\widehat{\pi}_{l-1}(c) \neq \pi(c)\} \right].
\end{aligned}$$

Therefore,

$$\min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_{\widehat{\pi}_{l-1}} - \phi_{\pi}\|_{A(w)^{-1}}^2}{\Delta(\pi)^2} = \min_{p_c \in \Delta_{\mathcal{A}}, \forall c \in \mathcal{C}} \frac{\mathbb{E}_{c \sim \nu} \left[\left(\frac{1}{p_{c,\widehat{\pi}_{l-1}(c)}} + \frac{1}{p_{c,\pi(c)}} \right) \mathbf{1}\{\widehat{\pi}_{l-1}(c) \neq \pi(c)\} \right]}{\Delta(\pi)^2}.$$

□

Lemma H.2. For any l , any $\lambda \in \Delta_{\Pi}$, $\gamma > 0$, and any n , we have $h_l(\lambda, \gamma, n) = \langle \lambda, \nabla_{\lambda} h_l(\lambda, \gamma, n) \rangle$.

Proof. We first compute

$$\begin{aligned}
[\nabla_{\lambda} h_l(\lambda, \gamma, n)]_{\pi} &= -\widehat{\Delta}_{l-1}^{\gamma_{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_{\pi} n} \\
&\quad + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top}(t_a^{(c)} + \eta_l)} \right) \left(\sum_{a' \in \mathcal{A}} \frac{\gamma_{\pi}(t_{a'}^{(c)} + \eta_l)_{\pi}}{\sqrt{(\lambda \odot \gamma)^{\top}(t_{a'}^{(c)} + \eta_l)}} \right) \right].
\end{aligned}$$

Then, by the fact that

$$\begin{aligned}
& \sum_{\pi \in \Pi} \lambda_{\pi} \cdot \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top}(t_a^{(c)} + \eta_l)} \right) \left(\sum_{a' \in \mathcal{A}} \frac{\gamma_{\pi}(t_{a'}^{(c)} + \eta_l)_{\pi}}{\sqrt{(\lambda \odot \gamma)^{\top}(t_{a'}^{(c)} + \eta_l)}} \right) \right] \\
&= \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top}(t_a^{(c)} + \eta_l)} \right) \left(\sum_{a' \in \mathcal{A}} \frac{(\lambda \odot \gamma)^{\top}(t_{a'}^{(c)} + \eta_l)}{\sqrt{(\lambda \odot \gamma)^{\top}(t_{a'}^{(c)} + \eta_l)}} \right) \right] \\
&= \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^{\top}(t_a^{(c)} + \eta_l)} \right)^2 \right],
\end{aligned}$$

we have

$$\begin{aligned}
& \langle \lambda, \nabla_\lambda h_l(\lambda, \gamma, n) \rangle \\
&= \sum_{\pi \in \Pi} \lambda_\pi [\nabla_\lambda h_l(\lambda, \gamma, n)]_\pi \\
&= \sum_{\pi \in \Pi} \lambda_\pi \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_\pi n} \right) \\
&\quad + \sum_{\pi \in \Pi} \lambda_\pi \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right) \left(\sum_{a' \in \mathcal{A}} \frac{\gamma_\pi (t_{a'}^{(c)} + \eta_l)_\pi}{\sqrt{(\lambda \odot \gamma)^\top (t_{a'}^{(c)} + \eta_l)}} \right) \right] \\
&= \sum_{\pi \in \Pi} \lambda_\pi \cdot \left(-\widehat{\Delta}_{l-1}^{\gamma^{l-1}}(\pi, \widehat{\pi}_{l-1}) + \frac{\log(1/\delta_l)}{\gamma_\pi n} \right) + \mathbb{E}_{c \sim \nu_{\mathcal{D}}} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right)^2 \right] \\
&= h_l(\lambda, \gamma, n).
\end{aligned}$$

□

Lemma H.3. For any $\lambda \in \Delta_\Pi$ and $\gamma \in \left[0, \min \left\{ \sqrt{\frac{\log(1/\delta_l)}{2n_l \mathbb{E}_c[\mathbf{1}\{\pi(c) \neq \pi^*(c)\}]}} , \sqrt{\frac{\log(1/\delta_l)}{|\mathcal{A}|^2 \eta_l n_l}} \right\} \right]^\Pi$, with $\eta_l = |\mathcal{A}|^{-4} \epsilon_l^2$, we have

$$0 \leq \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right)^2 \right] - \mathbb{E}_{c \sim \nu} \left[\left(\sum_a \sqrt{(\lambda \odot \gamma)^\top t_a^{(c)}} \right)^2 \right] \leq \epsilon_l.$$

Proof. The first inequality is clear since $\eta_l > 0$ and $\lambda_\pi, \gamma_\pi \geq 0$ for all $\pi \in \Pi$, so we focus on the upper bound. Note that

$$\begin{aligned}
& \mathbb{E}_{c \sim \nu} \left[\left(\sum_{a \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l)} \right)^2 \right] - \mathbb{E}_c \left[\left(\sum_a \sqrt{(\lambda \odot \gamma)^\top t_a^{(c)}} \right)^2 \right] \\
&= \mathbb{E}_{c \sim \nu} \left[\sum_{a \in \mathcal{A}} (\lambda \odot \gamma)^\top (t_a^{(c)} + \eta_l) + \sum_{a_1 \in \mathcal{A}} \sum_{a_2 \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a_1}^{(c)} + \eta_l)} (t_{a_2}^{(c)} + \eta_l)^\top (\lambda \odot \gamma) \right] \\
&\quad - \mathbb{E}_{c \sim \nu} \left[\sum_{a \in \mathcal{A}} (\lambda \odot \gamma)^\top t_a^{(c)} + \sum_{a_1 \in \mathcal{A}} \sum_{a_2 \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top t_{a_1}^{(c)}} t_{a_2}^{(c)\top} (\lambda \odot \gamma) \right]. \tag{26}
\end{aligned}$$

Note that

$$\begin{aligned}
& \mathbb{E}_{c \sim \nu} \left[\sum_{a_1 \in \mathcal{A}} \sum_{a_2 \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top (t_{a_1}^{(c)} + \eta_l)} (t_{a_2}^{(c)} + \eta_l)^\top (\lambda \odot \gamma) \right] \\
&= \mathbb{E}_{c \sim \nu} \left[\sum_{a_1 \in \mathcal{A}} \sum_{a_2 \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top t_{a_1}^{(c)}} (t_{a_2}^{(c)})^\top (\lambda \odot \gamma) + \eta_l \lambda^\top \gamma (\lambda \odot \gamma)^\top (t_{a_1}^{(c)} + t_{a_2}^{(c)}) + \eta_l^2 (\lambda^\top \gamma)^2 \right] \\
&\leq \mathbb{E}_{c \sim \nu} \left[\sum_{a_1 \in \mathcal{A}} \sum_{a_2 \in \mathcal{A}} \sqrt{(\lambda \odot \gamma)^\top t_{a_1}^{(c)}} (t_{a_2}^{(c)})^\top (\lambda \odot \gamma) \right] \\
&\quad + 2|\mathcal{A}| \mathbb{E}_{c \sim \nu} \left[\sum_{a \in \mathcal{A}} \sqrt{\eta_l \lambda^\top \gamma (\lambda \odot \gamma)^\top t_a^{(c)}} \right] + |\mathcal{A}|^2 \eta_l \lambda^\top \gamma.
\end{aligned}$$

Then (26) is upper bounded by

$$\begin{aligned}
& \mathbb{E}_{c \sim \nu} \left[\sum_{a \in \mathcal{A}} \eta_l \lambda^\top \gamma \right] + 2|\mathcal{A}| \mathbb{E}_{c \sim \nu} \left[\sum_{a \in \mathcal{A}} \sqrt{\eta_l \lambda^\top \gamma (\lambda \odot \gamma)^\top t_a^{(c)}} \right] + |\mathcal{A}|^2 \eta_l \lambda^\top \gamma \\
&= |\mathcal{A}| \eta_l \lambda^\top \gamma + |\mathcal{A}|^2 \eta_l \lambda^\top \gamma + 2|\mathcal{A}| \sqrt{\eta_l \lambda^\top \gamma} \mathbb{E}_{c \sim \nu} \left[\sum_{a \in \mathcal{A}} \sqrt{\sum_{\pi \in \Pi} \lambda_\pi \gamma_\pi [t_a^{(c)}]_\pi} \right] \\
&= |\mathcal{A}| \eta_l \lambda^\top \gamma + |\mathcal{A}|^2 \eta_l \lambda^\top \gamma + 2|\mathcal{A}|^2 \sqrt{\eta_l \lambda^\top \gamma} \mathbb{E}_{c \sim \nu} \left[\sum_{a \in \mathcal{A}} \frac{1}{|\mathcal{A}|} \sqrt{\sum_{\pi \in \Pi} \lambda_\pi \gamma_\pi [t_a^{(c)}]_\pi} \right] \\
&= |\mathcal{A}| \eta_l \lambda^\top \gamma + |\mathcal{A}|^2 \eta_l \lambda^\top \gamma + 2|\mathcal{A}|^2 \sqrt{\eta_l \lambda^\top \gamma} \mathbb{E}_{c \sim \nu} \left[\mathbb{E}_{a \sim \mu} \left[\sqrt{\sum_{\pi \in \Pi} \lambda_\pi \gamma_\pi [t_a^{(c)}]_\pi} \right] \right] \\
&\leq |\mathcal{A}| \eta_l \lambda^\top \gamma + |\mathcal{A}|^2 \eta_l \lambda^\top \gamma + 2|\mathcal{A}|^2 \sqrt{\eta_l \lambda^\top \gamma} \sqrt{\sum_{\pi \in \Pi} \lambda_\pi \gamma_\pi \frac{1}{|\mathcal{A}|} \mathbb{E}_{c \sim \nu} \left[\sum_{a \in \mathcal{A}} [t_a^{(c)}]_\pi \right]} \\
&= |\mathcal{A}| \eta_l \lambda^\top \gamma + |\mathcal{A}|^2 \eta_l \lambda^\top \gamma + 2|\mathcal{A}|^2 \sqrt{\eta_l \lambda^\top \gamma} \sqrt{\sum_{\pi \in \Pi} \lambda_\pi \gamma_\pi \frac{1}{|\mathcal{A}|} 2 \cdot \mathbb{E}_{c \sim \nu} [\mathbf{1}\{\pi(c) \neq \pi^*(c)\}]}. \quad (27)
\end{aligned}$$

Since $\gamma_\pi \leq \sqrt{\frac{\log(1/\delta_l)}{2n_l \mathbb{E}_c[\mathbf{1}\{\pi(c) \neq \pi^*(c)\}]}}$, $\gamma_\pi \mathbb{E}_{c \sim \nu}[\mathbf{1}\{\pi(c) \neq \pi^*(c)\}] \leq \sqrt{\frac{\mathbb{E}_c[\mathbf{1}\{\pi(c) \neq \pi^*(c)\}] \log(1/\delta_l)}{2n_l}} \leq \sqrt{\frac{\log(1/\delta_l)}{2n_l}}$. We know from the lower bound argument that

$$n_l \gtrsim \min_{w \in \Omega} \max_{\pi \in \Pi} \frac{\|\phi_\pi - \phi_{\pi^*}\|_{A(w)}^2}{\Delta(\pi)^2 + \epsilon_l^2} \log(1/\delta_l) \geq \epsilon_l^{-1} \log(1/\delta_l),$$

so $\sqrt{\frac{\log(1/\delta_l)}{2n_l}} \lesssim \sqrt{\epsilon_l}$. Therefore, (27) is upper bounded by

$$(|\mathcal{A}| + |\mathcal{A}|^2) \eta_l \lambda^\top \gamma + 2|\mathcal{A}|^{3/2} \sqrt{\epsilon_l \eta_l \lambda^\top \gamma}. \quad (28)$$

Since $\eta_l \lambda^\top \gamma \leq \eta_l \gamma_{\max} = \sqrt{\frac{\eta_l \log(1/\delta_l)}{|\mathcal{A}|^2 n_l}} \leq \sqrt{\eta_l} \frac{1}{|\mathcal{A}|}$. Plugging this as well as $\eta_l \leq |\mathcal{A}|^{-4} \epsilon_l^2$ in equation 28 gives that the bias is upper bounded by ϵ_l . \square