

---

# Incentivized Learning in Principal-Agent Bandit Games

---

Antoine Scheid<sup>1</sup> Daniil Tiapkin<sup>1,2</sup> Etienne Boursier<sup>3</sup> Aymeric Capitaine<sup>1</sup> Eric Moulines<sup>1</sup>  
Michael I. Jordan<sup>4,5</sup> El-Mahdi El-Mhamdi<sup>1</sup> Alain Durmus<sup>1</sup>

## Abstract

This work considers a repeated principal-agent bandit game, where the principal can only interact with her environment through the agent. The principal and the agent have misaligned objectives and the choice of action is only left to the agent. However, the principal can influence the agent’s decisions by offering incentives which add up to his rewards. The principal aims to iteratively learn an incentive policy to maximize her own total utility. This framework extends usual bandit problems and is motivated by several practical applications, such as healthcare or ecological taxation, where traditionally used mechanism design theories often overlook the learning aspect of the problem. We present nearly optimal (with respect to a horizon  $T$ ) learning algorithms for the principal’s regret in both multi-armed and linear contextual settings. Finally, we support our theoretical guarantees through numerical experiments.

## 1 Introduction

Decision-making under uncertainty is a ubiquitous feature of real-world applications of machine learning, arising in domains as diverse as recommendation systems (Li et al., 2010), healthcare (Yu et al., 2021), and agriculture (Evans et al., 2017). Multi-armed bandits provide a classical point of departure for decision-making under uncertainty in these settings (Thompson, 1933; Woodroffe, 1979; Lattimore & Szepesvári, 2020; Slivkins et al., 2019). The

---

<sup>1</sup>Centre de Mathématiques Appliquées – CNRS – École polytechnique – Institut Polytechnique de Paris – route de Saclay 91128 Palaiseau cedex <sup>2</sup>Université Paris-Saclay, CNRS, Laboratoire de mathématiques d’Orsay, 91405, Orsay, France <sup>3</sup>INRIA, Université Paris Saclay, LMO, Orsay, France <sup>4</sup>University of California, Berkeley <sup>5</sup>Inria, École Normale Supérieure, PSL Research University. Correspondence to: Antoine Scheid <antoine.scheid@polytechnique.edu>.

basic bandit solution involves an agent who learns which decisions yield high rewards via repeated experimentation. Real-world decision-making problems, however, often present challenges that are not addressed in this simple optimization framework. These include the challenge of scarcity when there are multiple decision-makers, issues of misaligned objectives, and problems arising from information asymmetries and signaling. The economics literature addresses these issues through the design of game-theoretic mechanisms, including auctions and contracts (see, e.g., Myerson, 1989; Laffont & Martimort, 2009), aiming to achieve favorable outcomes despite agents’ self-interest and limited information set. Unfortunately, the economics literature tends to neglect the learning aspect of the problem, often assuming that preferences, or distributions on preferences, are known a priori. Our work focuses on the blend of mechanism design and learning. We study a principal-agent model with information asymmetry and we develop a learning framework in which the principal aims to uncover the true preferences of the agent while optimizing her own gains.

Building on the work of Dogan et al. (2023a;b), we consider a repeated game between a principal and an agent, where, at each round, the principal proposes an incentive transfer associated with any action. The agent greedily chooses the action that maximizes the sum of his expected reward and the incentive. The goal of the principal is to learn an incentive policy which maximizes her own utility over time, taking into account both the rewards that she reaps and the costly incentives that she offers.

Our contributions are as follows:

- We present the Incentivized Principal-Agent Algorithm (IPA) framework, which comprises two steps. First, IPA estimates the minimal level of incentive needed to make the agent select any desired action. Subsequently, forming an upper estimate of these incentives, IPA uses a regret-minimization algorithm in a black-box fashion. The overall algorithm achieves both nearly optimal distribution-free and instance-dependent regret bounds.
- We extend IPA to the linear contextual bandit setting (see, e.g., Abe & Long, 1999; Auer, 2002; Dani et al., 2008), significantly broadening its applicability in various applications. Here Contextual IPA achieves

a  $\mathcal{O}(d\sqrt{T}\log(T))$  regret bound. We emphasize that Contextual IPA is the first known algorithm for incentivized learning in a contextual setting. Moreover it matches, up to logarithmic factors, the minimax lower bound  $\Omega(d\sqrt{T})$  for the easier problem of stochastic linear bandits (Rusmevichientong & Tsitsiklis, 2010).

## 2 Related Work

While classical work on bandit problems and reinforcement learning has predominantly focused on single-agent scenarios, many emerging applications require considering multiple agents. Recent literature has accordingly begun to study frameworks for learning in multi-agent multi-armed bandit settings (see, e.g., Boursier & Perchet, 2022).

Mansour et al. (2020) discuss how a social planner can simultaneously learn and influence self-interested agents’ decisions through Bayesian-Incentive Compatible (BIC) recommendations. The rationale behind this notion is that a BIC recommendation guarantees to each agent a maximal reward given the past, at any step. The social planner objective is to design BIC recommendations that maximize the global welfare. Mansour et al. (2020) propose an algorithm for solving this problem in both multi-armed and contextual bandit settings. Notably, their work turns any black-box bandit algorithm into a BIC algorithm. For this problem, Sellke & Slivkins (2021) show that Thompson sampling can be made BIC, with a sufficient number of initial observations. Hu et al. (2022) extend this work to the combinatorial bandits problem.

Another line of work due to Banihashem et al. (2023) and Simchowicz & Slivkins (2023) studies how a principal can provide recommendations to agents so that they explore all reachable states in a Markov Decision Process (MDP). To this end, the principal supplies the agents with a modified history, with the modifications carefully chosen to retain the agents’ trust. This line of work is closely related to the online Bayesian persuasion literature (see, e.g., Castiglioni et al., 2020; 2021), which dates to the seminal work of Kamenica & Gentzkow (2011). Online Bayesian persuasion consists of the principal sequentially influencing agents’ decision with signals in her own interest and has been extended to the online Information Design setting (Doval & Ely, 2020; Bernasconi et al., 2022).

In these works, the information asymmetry favors the principal, so that the principal can influence the agent’s decision at little or no cost. In our problem, the agent instead has perfect knowledge of the problem parameters and his action can only be influenced through utility transfers.

Ben-Porat et al. (2024) study a principal and an agent sharing a common Markov Decision Process (MDP) with different reward functions. Similarly to our setting, at each

step the action is chosen by an agent with full knowledge of the game. The objective of the principal is to minimize her cumulative regret under a constraint on the incentive budget. Despite extending our setup to MDPs, Ben-Porat et al. (2024) do not consider uncertainty in the principal’s side, turning the game into an optimization problem. Other leader/follower dynamics with Reinforcement Learning have for instance been studied in the works of Chen et al. (2023); Zhong et al. (2021) to find a Stackelberg equilibrium or an optimal policy for the leader.

Related issues arise in the study of dynamic pricing (Den Boer, 2015; Javanmard & Nazerzadeh, 2019; Mao et al., 2018; Golrezaei et al., 2023). Our work diverges from dynamic pricing in that in our case the principal not only faces uncertainty with respect to the agent’s utility, but also with respect to her own utility.

Some works study similar principal-agent games but with a specific focus on the achievable optimality of the contract Cohen et al. (2022) or a specific stochastic model for the agent’s behavior Conitzer & Garera (2006). The work of Zhu et al. (2022) interestingly studies online contract design but with a specific focus on the Stackelberg regret of the task, which provides results concerning the sample complexity.

Finally, our study is inspired by the work of Dogan et al. (2023b) to explore the principal’s learning mechanism within a principal-agent setting. They propose an  $\epsilon$ -Greedy algorithm with suboptimal regret guarantees. In particular, it suffers an exponential dependence in the number of actions. Dogan et al. (2023a) extend the work of (Dogan et al., 2023b), taking into account the presence of uncertainty on the agent’s side but with the same limitation. We consider the same setup and problem as introduced in (Dogan et al., 2023b;a). However, it corresponds to the only common ground between our work and theirs. First, their algorithm as well as their regret bound depends on a hyperparameter  $m$  that controls the level of exploration, as mentioned in Appendix B. If  $m$  is too small, there is not enough exploration and the regret scales poorly with respect to  $T$  whereas if  $m$  is too large, the convergence takes longer to occur.  $m$  needs to be carefully chosen and depends on the specific bandit instance considered, which makes things more complicated. Furthermore, their algorithm is harder to analyze and to implement in practical settings: a lot of difficulties come from a simultaneous exploration/exploitation of both the principal’s rewards and the agent’s ones. In contrast, the algorithmic solution that we provide for the multi-armed case is completely orthogonal, relying on a two phases approach and avoids fine-tuning of any kind. We provide both distribution free and instance-dependent regret bounds that nearly match the known lower bounds. Also note that we extend our approach to the non-trivial contextual case.

### 3 Multi-Armed Principal-Agent Learning

**Setup.** We consider a repeated principal-agent game. A contextual version of the game is introduced in Section 4. The action set for the agent (or set of arms) is fixed to be  $\mathcal{A} := [K] = \{1, \dots, K\}$ ,  $K \in \mathbb{N}^*$ . We assume that the agent's rewards,  $\mathbf{s} = (s_1, \dots, s_K) \in \mathbb{R}_+^K$ , are deterministic and that they are known to the agent and unknown to the principal.

For each action  $a \in [K]$ , the rewards of the principal are given by a random, i.i.d. sequence  $(X_a(t))_{t \in [T]}$ , where  $X_a(t) \sim \nu_a$  and  $\nu_a$  is the arm distribution. The distributions  $\{\nu_a\}_{a=1}^K$  are unknown to the principal and are learned as a consequence of the following principal-agent interaction.

At each round  $t \in [T]$ , where  $T$  is the game horizon, the principal proposes an incentive  $\pi(t) \in \mathbb{R}_+$  associated with an action  $a_t \in [K]$ . The agent then greedily chooses action  $A_t$  maximizing his utility:

$$A_t \in \operatorname{argmax}_{a \in [K]} \{s_a + \mathbb{1}_{a_t}(a)\pi(t)\}, \quad (1)$$

breaking ties arbitrarily.<sup>1</sup> The principal then observes the arm  $A_t$  selected by the agent, as well as her reward given by  $X_{A_t}(t)$ . The utility of the principal on the round is given by  $X_{A_t}(t) - \mathbb{1}_{a_t}(A_t)\pi(t)$ . For any  $a \in [K]$ , the principal's mean reward is  $\theta_a := \mathbb{E}[X_a(t)]$ . See Table 1 in Appendix A for a summary of the main definitions used in this section.

The sequence of incentives  $(a_t, \pi(t))_{t \in [T]}$  defines a sequence of actions  $(A_t)_{t \in [T]}$  chosen by the agent. The goal of the principal is to maximize her total utility. On a single round, she thus aims at proposing an optimal incentive  $\pi^{\text{opt}}$  on an arm  $a^{\text{opt}} \in [K]$ , which solves

$$\begin{aligned} & \text{maximize } \int x \nu_a(dx) - \pi \text{ over } \pi \in \mathbb{R}_+, a \in [K] \\ & \text{such that } a \in \operatorname{argmax}_{a' \in [K]} \{s_{a'} + \mathbb{1}_a(a')\pi\}. \end{aligned} \quad (2)$$

This is consistent with the conventional framework for utility in bandit problems, where we subtract the cost of incentives to the principal. Here, the principal's influence is exerted solely through the strategic use of incentives, carefully designed to guide the agent's behavior. We define  $\mu^* := \int x \nu_{a^{\text{opt}}}(dx) - \pi^{\text{opt}}$ . Maximizing the total utility of the principal over  $T$  rounds is equivalent to minimizing the expected regret, defined as

$$\mathfrak{R}(T) := T \mu^* - \sum_{t=1}^T \mathbb{E}[X_{A_t}(t) - \mathbb{1}_{a_t}(A_t)\pi(t)]. \quad (3)$$

**Remark.** In the prior work of Dogan et al. (2023b), incentives were defined as a vector of size  $K$ , where the

<sup>1</sup>Note that the related works (Ben-Porat et al., 2024; Simchowitz & Slivkins, 2023) assume a tie-breaking in favor of the principal, an assumption that we do not need here.

incentive associated with an action  $a \in [K]$  was denoted  $\pi_a$ . In our setting, since the goal of the principal is to make sure that the agent picks one prescribed action, it is enough to consider a restricted family of the form  $\pi_a = \mathbb{1}_{a_t}(a)\pi(t)$ , where  $(a_t, \pi(t))$  are incentives in the sense defined above.

We make the following assumption.

**H1.** For any  $a \in [K]$ ,  $s_a \in [0, 1]$ .

Neither the distributions  $\nu_a$  nor the preferences of the agent are known to the principal. Another difficulty arises from designing the magnitude of the incentive  $\pi(t)$ : if it is too small, the agent might not choose the arm  $a_t$  proposed by the principal whereas using an overly large amount leads the principal to overpay, decreasing her utility.

This trade-off also arises in dynamic pricing, where sellers must strike a balance between attractive pricing and profitability. For discussion of the results in that literature, see the comprehensive overview by Den Boer (2015). In addition, there are links between dynamic pricing and bandit problems (see, e.g., Javanmard & Nazerzadeh, 2019; Cai et al., 2023).

**Optimal incentives.** Before introducing IPA, we highlight a pivotal observation. For any given round  $t \geq 1$ , action  $a \in [K]$  and  $\varepsilon > 0$ , the principal can entice the agent to choose  $a$  by offering an incentive,  $\pi_a^{*,\varepsilon} \in \mathbb{R}_+$ , defined as:

$$\pi_a^{*,\varepsilon} = \max_{a' \in [K]} s_{a'} - s_a + \varepsilon. \quad (4)$$

With this incentive, it holds that for any  $a' \in [K]$ ,  $a' \neq a$ :

$$s_{a'} < s_a + \pi_a^{*,\varepsilon},$$

which ensures that the agent chooses  $A_t = a$ , given that action  $a$  yields a superior reward. Consequently,  $\pi_a^* := \lim_{\varepsilon \rightarrow 0} \pi_a^{*,\varepsilon} = \max_{a' \in [K]} s_{a'} - s_a$  represents the infimal incentive necessary to make arm  $a$  the agent's selection. Assuming  $\mathbf{s}$  is known to the principal, then using  $\pi_a^{*,\varepsilon}$  for any  $\varepsilon > 0$  across all arms  $a \in [K]$ , will provide an expected reward of  $\theta_a - \pi_a^*$  per arm, which can be found using a standard bandit algorithm. Lemma 1 allows us to define the regret in a more convenient way.

**Lemma 1.** For any  $T \in \mathbb{N}$ , the regret of any algorithm on our problem instance can be written as

$$\begin{aligned} \mathfrak{R}(T) = & T \max_{a \in [K]} \{ \theta_a + s_a - \max_{a' \in [K]} s_{a'} \} \\ & - \mathbb{E} \left[ \sum_{t=1}^T \{ \theta_{A_t} - \mathbb{1}_{a_t}(A_t)\pi(t) \} \right]. \end{aligned}$$

**Warm up: fixed horizon solution and regret analysis.** IPA separates the problem of learning optimal incentives  $\pi_a^*$  for each action  $a$ —a problem that can be solved efficiently

via binary search (see Algorithm 3)—from estimation of the principal’s expected reward  $(\theta_a)_{a \in [K]}$ , which is achieved using a standard multi-armed bandit algorithm. With a known horizon  $T$ , the algorithm unfolds in two stages. First, for each action  $a \in [K]$ , the principal devotes  $N_T := \lceil \log_2 T \rceil$  rounds of binary search per arm to estimate  $\pi_a^*$ , maintaining lower  $\underline{\pi}_a(N_T)$  and upper  $\bar{\pi}_a(N_T)$  bounds with  $\underline{\pi}_a(N_T) \leq \pi_a^* \leq \bar{\pi}_a(N_T)$ . Denoting  $\bar{\pi}_a := \bar{\pi}_a(\lceil \log_2 T \rceil)$  for simplicity, we compute the estimate

$$\hat{\pi}_a := \bar{\pi}_a(\lceil \log_2 T \rceil) + 1/T, \quad (5)$$

where  $1/T$  is added to avoid any tie-breaking situation. We show formally in Lemma 8 that

$$\hat{\pi}_a - 2/T \leq \pi_a^* < \hat{\pi}_a. \quad (6)$$

In the second phase, IPA then employs an arbitrary multi-armed bandit subroutine Alg in a black-box manner to learn  $\theta$ .

**Bandit instance.** For any distributions  $(\tilde{\nu}_a)_{a \in [K]}$  and sequence of i.i.d. random variables  $(Y_a(t))_{t \in \mathbb{N}^*}$ ,  $a \in [K]$ , with  $Y_a(t) \sim \tilde{\nu}_a$ , we define the history  $\mathcal{G}_t := (A_s, U_s, Y_{A_s}(s))_{s \leq t}$  where for any  $s \in \mathbb{N}^*$ ,  $(U_s)_{s \in \mathbb{N}^*}$  is a family of independent uniform random variables on  $[0, 1]$  allowing for randomization in the subroutine. Let Alg be a bandit algorithm, i.e.,  $\text{Alg}: (U_t, \mathcal{G}_{t-1}) \mapsto a_t^{\text{Rec}}$ . We define the expected regret of Alg as

$$R_{\text{Alg}}(T, \tilde{\nu}) := T \max_{a \in [K]} \mathbb{E}_{Y \sim \tilde{\nu}} [Y_a(1)] - \mathbb{E}_{Y \sim \tilde{\nu}} \left[ \sum_{t=1}^T Y_{\text{Alg}}(U_t, \mathcal{G}_{t-1})(t) \right].$$

After the binary search phase, for  $t > K \lceil \log_2 T \rceil$ , the principal plays Alg on her bandit instance driven by her own mean rewards  $(\theta_a)_{a \in [K]}$  and the approximated incentives  $(\hat{\pi}_a)_{a \in [K]}$ . Alg will be fed with a shifted history, defined for any  $t > K \lceil \log_2 T \rceil$  as

$$\tilde{\mathcal{H}}_t := (a_s^{\text{Rec}}, U_s, X_{A_s}(s) - \hat{\pi}_{a_s^{\text{Rec}}})_{s \in [K \lceil \log_2 T \rceil + 1, t]}, \quad (7)$$

with  $a_t^{\text{Rec}}$  the action recommended by Alg at time  $t$  and  $A_t$  the action pulled by the agent. At time  $t$ , IPA offers the incentive  $\hat{\pi}_{a_t^{\text{Rec}}}$  to the agent if he chooses action  $a_t^{\text{Rec}}$ . Equation (6) ensures that this incentive makes  $a_t^{\text{Rec}}$  strictly preferable to any other action for the agent and so  $a_t^{\text{Rec}}$  is eventually played. As can be seen in (7), the shift of each arm’s mean by  $\hat{\pi}_a$  is taken into account while Alg is learning. We also define the shifted distribution  $\rho_a^T$  for any  $a \in [K]$  as the distribution of  $X_a(1) - \hat{\pi}_a$ .

**Theorem 1.** IPA run with any multi-armed bandit subroutine Alg has an overall regret  $\mathfrak{R}(T)$  such that

$$\mathfrak{R}(T) \leq 2 + (1 + \max_{a \in [K]} \{\theta_a\} - \min_{a \in [K]} \{\theta_a\})(1 + K \log_2 T)$$

---

**Algorithm 1** IPA
 

---

- 1: **Input:** Set of actions  $\mathcal{A} := [K]$ , time horizon  $T$ , subroutine Alg
  - 2: Compute  $\tilde{\mathcal{H}}_s := \emptyset$  for any  $s \leq K \lceil \log_2 T \rceil$
  - 3: **for**  $a \in [K]$  **do**
  - 4:   **# See Algorithm 3**
  - 5:    $\underline{\pi}_a, \bar{\pi}_a = \text{Binary Search}(a, \lceil \log_2 T \rceil, 0, 1)$
  - 6: **end for**
  - 7: For any action  $a \in [K]$ ,  $\hat{\pi}_a = \bar{\pi}_a + 1/T$
  - 8: **for**  $t = K \lceil \log_2 T \rceil + 1, \dots, T$  **do**
  - 9:   Sample  $U_t \sim \text{U}(0, 1)$  and get a recommended action by Alg,  $a_t^{\text{Rec}} = \text{Alg}(U_t, \tilde{\mathcal{H}}_{t-1})$
  - 10:   **Offer an incentive**  $\hat{\pi}_{a_t^{\text{Rec}}}$  on action  $a_t^{\text{Rec}}$  and nothing for any other action  $a' \in [K]$ ,  $a' \neq a$
  - 11:   **Observe**  $A_t, X_{A_t}(t)$  and compute history  $\tilde{\mathcal{H}}_t = (a_s^{\text{Rec}}, U_s, X_{A_s}(s) - \hat{\pi}_{a_s^{\text{Rec}}})_{s \in [K \lceil \log_2 T \rceil + 1, t]}$
  - 12: **end for**
- 

$$+ R_{\text{Alg}}(T - K \lceil \log_2 T \rceil, \rho^T),$$

where  $R_{\text{Alg}}$  stands for the regret induced by Alg on the shifted vanilla multi-armed bandit problem  $\rho^T$ .

The proof is postponed to Appendix C.

**Corollary 1.** Assume the principal’s reward distribution  $\nu_a$  for any action  $a \in [K]$  is 1-subgaussian. Then, IPA run with the bandit subroutine Alg = UCB has a regret bounded for any  $T \in \mathbb{N}$  as follows

$$\begin{aligned} \mathfrak{R}(T) &\leq 3 + 3 \sum_{a \in [K], \Delta_a^* > 0} \Delta_a^* \\ &\quad + (1 + \max_{a \in [K]} \{\theta_a\} - \min_{a \in [K]} \{\theta_a\})(1 + 9K \log_2 T) \\ &\quad + 8 \min \left\{ \sqrt{TK \log T}; \sum_{a \in [K], \Delta_a^* > 0} \frac{4 \log T}{\Delta_a^*} \right\}, \end{aligned}$$

where  $\Delta_a^* := \max_{a' \in [K]} \{\theta_{a'} + s_{a'}\} - (\theta_a + s_a)$  are the reward gaps.

Note that any black-box algorithm, not necessarily UCB, can be employed, yielding other concrete bounds in the corollary. We recover the usual multi-armed UCB bounds (both distribution-free and instance-dependent): this is why IPA achieves the bound provided in Corollary 1. For completeness, the UCB subroutine is given in Appendix E.

## 4 Contextual Principal-Agent Learning

In this section, we study the same interaction between a principal and an agent, but in a contextual setting (see, e.g., Abe & Long, 1999; Auer, 2002; Dani et al., 2008). We use the simplified model of stochastic linear bandits for both the

agent and the principal. Consider a set of possible actions in  $B(0, 1)$ , where  $B(0, 1)$  stands for the unit closed ball in  $\mathbb{R}^d$ , and a family of zero-mean distributions indexed by  $B(0, 1)$ ,  $(\tilde{\nu}_a)_{a \in B(0,1)}$  such that for any  $a \in B(0, 1)$ ,  $t \in [T]$ ,  $\eta_a(t) \sim \tilde{\nu}_a$ . The principal's reward is given by the sequence  $\{(X_a(t))_{a \in B(0,1)} : t \in [T]\}$  of independent random variables such that for any  $t \in [T]$ ,  $a \in B(0, 1)$ ,

$$X_a(t) := \langle \theta^*, a \rangle + \eta_a(t),$$

where  $\theta^*$  is unknown to the principal. The agent's reward function is defined as  $a \mapsto \langle s^*, a \rangle$ , where  $s^* \in \mathbb{R}^d$  is known to the agent and unknown to the principal. With this notation, the agent and the principal observe an action set  $\mathcal{A}_t \subseteq B(0, 1)$ , at each round  $t \geq 1$ . Note that this set is no longer stationary. The precise timeline is as follows. At each round, the principal proposes an incentive function,  $\kappa(t, \cdot) : \mathcal{A}_t \rightarrow \mathbb{R}_+$ , associating any action  $a \in \mathcal{A}_t \subseteq B(0, 1)$  with a transfer of incentives  $\kappa(t, a)$  from the principal to the agent. The principal chooses  $\kappa(t, \cdot)$  as a function with a finite support, which makes it upper semi-continuous. The agent then greedily chooses the action  $A_t$  as follows

$$A_t \in \operatorname{argmax}_{a \in \mathcal{A}_t} \{\langle s^*, a \rangle + \kappa(t, a)\}, \quad (8)$$

which is well-defined since  $\kappa(t, \cdot)$  is upper semi-continuous and  $\mathcal{A}_t$  satisfies the following assumption.

**H2.** For any  $t \geq 1$ ,  $\mathcal{A}_t$  is closed, therefore compact. Moreover,  $s^* \in B(0, 1)$  and  $\theta^* \in B(0, 1)$ .

The principal then observes the arm  $A_t$  selected by the agent, as well as her incurred reward given by  $X_{A_t}(t)$ . The utility of the principal on the round is given by  $X_{A_t}(t) - \kappa(t, A_t)$ . This defines, for a sequence of principal's incentive functions  $\{\kappa(t, \cdot), t \in [T]\}$ , the sequence of actions  $\{A_t : t \in [T]\}$  chosen by the agent. The goal of the principal is to maximize her total utility. On a single round  $t$ , she thus aims at proposing an optimal incentive function  $\kappa(t, \cdot)$  which solves

$$\begin{aligned} & \text{maximize } \langle \theta^*, a \rangle - \kappa(t, a) \text{ over } \kappa(t, \cdot) : \mathcal{A}_t \rightarrow \mathbb{R}_+, \\ & \text{such that } a \in \operatorname{argmax}_{a' \in \mathcal{A}_t} \{\langle s^*, a' \rangle + \kappa(t, a')\}. \end{aligned} \quad (9)$$

In addition, we define the optimal average reward at  $t$  as

$$\begin{aligned} \mu_t^* &:= \sup_{\kappa(t, \cdot) : \mathcal{A}_t \rightarrow \mathbb{R}_+} \{\langle \theta^*, a \rangle - \kappa(t, a)\} \\ & \text{such that } a \in \operatorname{argmax}_{a' \in \mathcal{A}_t} \{\langle s^*, a' \rangle + \kappa(t, a')\}. \end{aligned} \quad (10)$$

Maximizing the total utility of the principal over  $T$  rounds is equivalent to minimizing the expected regret, defined as

$$\mathfrak{R}(T) = \sum_{t=1}^T \mu_t^* - \mathbb{E} \left[ \sum_{t=1}^T (X_{A_t}(t) - \kappa(t, A_t)) \right]. \quad (11)$$

Similarly to Lemma 1, the following result provides an alternative definition for the regret.

**Lemma 2.** For any  $T \in \mathbb{N}$ , the regret of any algorithm on our contextual problem instance can be written as

$$\begin{aligned} \mathfrak{R}(T) &= \sum_{t=1}^T \max_{a \in \mathcal{A}_t} \{\langle \theta^* + s^*, a \rangle - \max_{a' \in \mathcal{A}_t} \langle s^*, a' \rangle\} \\ & \quad - \mathbb{E} \left[ \sum_{t=1}^T (\langle \theta^*, A_t \rangle - \kappa(t, A_t)) \right]. \end{aligned}$$

The proof is deferred to Appendix D.

**Design of the optimal incentives.** At any round  $t \geq 1$ , for the agent to necessarily choose action  $a \in \mathcal{A}_t$ , the principal can provide the agent with the incentive function  $\kappa_a^{*,\varepsilon}(t, a') := \mathbb{1}_a(a') \pi^{*,\varepsilon}(t, a)$ , where for any  $\varepsilon > 0$ ,

$$\pi^{*,\varepsilon}(t, a) := \max_{a'_t \in \mathcal{A}_t} \{\langle s^*, a'_t - a \rangle + \varepsilon\}.$$

Lemma 10 in Appendix D guarantees that this choice of  $\kappa_a^{*,\varepsilon}$  gives  $A_t = a$ , where  $A_t$  is defined in (8). Define

$$\begin{aligned} a_t^{\text{ag}} &:= \operatorname{argmax}_{a' \in \mathcal{A}_t} \langle s^*, a' \rangle, \\ \pi^*(t, a) &:= \langle s^*, a_t^{\text{ag}} \rangle - \langle s^*, a \rangle \end{aligned} \quad (12)$$

and  $\kappa_a^*(t, a') := \mathbb{1}_a(a') \pi^*(t, a)$ . As in the non-contextual case, taking  $\varepsilon \rightarrow 0$  makes the incentive function  $\kappa_a^*$  the infimal function that makes the choice of  $a$  strictly preferable to any other arm  $a' \in \mathcal{A}_t$  at time  $t$ .

Similarly to the multi-armed setting, we decompose the problem into two distinct components. First, we aim to estimate the agent's reward  $a \mapsto \langle s^*, a \rangle$  based on the observation of agent's selected actions given an appropriate choice of incentives. As discussed below, this can be achieved with a binary-search-like procedure. Second, once this function is accurately estimated, the principal can use a contextual bandit algorithm `CtxAlg` in a black-box manner to minimize her own regret with the estimated incentive function to determine the agent's behavior.

**Estimation of the agent's reward.** The approach that we propose is based on a sequence of confidence sets  $\{\mathcal{S}_t\}_{t \in [T]}$  that satisfy  $s^* \in \mathcal{S}_t$  for any  $t \in [T]$ . We construct the sequence  $(\mathcal{S}_t)_{t \in [T]}$  recursively such that their diameters decrease along the iterations. This is motivated by Lemma 3 which allows us to control the estimation error of  $\pi^*$  and relates it to the diameter of these sets. The proof is postponed to Appendix D.

**Lemma 3.** For any  $t \in [T]$  and closed subset  $\mathcal{S} \subset B(0, 1)$  with  $s^* \in \mathcal{S}$ , it holds, for any  $a \in \mathcal{A}_t$ ,  $|\max_{s \in \mathcal{S}, a' \in \mathcal{A}_t} \langle s, a' - a \rangle - \pi^*(t, a)| \leq 2 \operatorname{diam}(\mathcal{S}, \mathcal{A}_t)$  where  $\operatorname{diam}(\mathcal{S}, \mathcal{A}_t) := \max_{a' \in \mathcal{A}_t} \max_{s_1, s_2 \in \mathcal{S}} |\langle s_1 - s_2, a' \rangle|$ .

In the light of Lemma 3, we thus aim to build confidence sets  $\mathcal{S}_t$  with decreasing diameters such that  $s^* \in \mathcal{S}_t$  for any  $t$ . To

this end, the principal can offer an incentive function  $\kappa(t, \cdot)$  concentrated on a *single point*  $a \in \mathcal{A}_t$  as in the multi-armed case:  $\kappa(t, a') = \pi(t) \cdot \mathbb{1}_a(a')$  for  $\pi(t) \in \mathbb{R}_+$ . In this case, the principal receives the agent's choice as a feedback; by (8), either  $A_t = a_t$  or  $A_t = \operatorname{argmax}_{a' \in \mathcal{A}_t} \langle s^*, a' \rangle = a_t^{\text{ag}}$ . In addition,  $A_t = a_t$  is equivalent to the fact that

$$\langle s^*, a_t^{\text{ag}} - a_t \rangle \leq \pi(t).$$

The information  $A_t = a_t$  or  $A_t = a_t^{\text{ag}}$  can be used as *binary search feedback* in the direction  $a_t^{\text{ag}} - a_t$ , as follows. Given a current confidence set  $\mathcal{S}_t$  at time  $t$  it can be updated either as  $\mathcal{S}_{t+1} = \mathcal{S}_t \cap \{s: \langle s, a_t^{\text{ag}} - a_t \rangle \leq \pi(t)\}$  if  $A_t = a_t$  or  $\mathcal{S}_{t+1} = \mathcal{S}_t \cap \{s: \langle s, a_t^{\text{ag}} - a_t \rangle \geq \pi(t)\}$  otherwise.

However, since the action set  $\mathcal{A}_t$  is non-stationary, we cannot determine the action  $a_t^{\text{ag}}$  by just observing the first round. Consequently, the new set  $\mathcal{S}_{t+1}$  cannot be computed as previously in the non-contextual setting. This makes the *single-point* incentive functions not suited for an efficient learning of  $s^*$  over the iterations. Instead, at any time  $t$ , we seek for a form of binary-search feedback in the direction  $a_t^1 - a_t^2$  for any two arms  $a_t^1 \neq a_t^2 \in \mathcal{A}_t$ . As we will see, this can be achieved by considering an incentive function  $\kappa$  with support  $\{a_t^1, a_t^2\} \subseteq \mathcal{A}_t$ .

Indeed, an important remark is that the amount of incentive needed to make the agent play any particular action is bounded under **H2** since

$$\max_{a \in \mathcal{A}_t} \pi^*(t, a) = \max_{a \in \mathcal{A}_t} \langle s^*, a_t^{\text{ag}} - a \rangle \leq 2, \quad (13)$$

this bound being known by the principal. For any  $a \in B(0, 1)$ , this makes the incentive function  $a' \mapsto 3 \cdot \mathbb{1}_a(a')$  sufficient to ensure  $A_t = a$  from (8). The value 3 in the definition of the incentive function is chosen instead of 2 to avoid an arbitrary tie-breaking.

Consequently, under the choice  $\kappa(t, a_t^1) = 3$ ,  $\kappa(t, a_t^2) = 3 + \pi(t)$  for  $\pi(t) \geq 0$  and  $\kappa(t, a') = 0$  for any other arm  $a'$ , (13) guarantees that only  $a_t^1$  and  $a_t^2$  may be chosen by the agent, helping the principal to update her confidence set  $\mathcal{S}_t$  in a known direction. Specifically for such an incentive  $\kappa$ , the choice  $A_t = a_t^1$  reveals the following information on  $s^*$ :

$$\langle s^*, a_t^1 - a_t^2 \rangle \geq \pi(t),$$

that permits the definition of a *binary search-like feedback* in the direction  $a_t^1 - a_t^2$  and thus allows us to update the confidence set  $\mathcal{S}_t$  following  $\mathcal{S}_{t+1} = \mathcal{S}_t \cap \{s: \langle s, a_t^1 - a_t^2 \rangle \geq \pi(t)\}$  if  $A_t = a_t^1$  or  $\mathcal{S}_{t+1} = \mathcal{S}_t \cap \{s: \langle s, a_t^1 - a_t^2 \rangle \leq \pi(t)\}$  otherwise.

**Binary search.** This update turns our estimation of the optimal incentives  $\pi^*$  into a multidimensional binary search where the unknown quantity is the vector  $s^*$ . At each iteration  $t$ , a vector  $w_t$  from the unit sphere is given. Then,

the algorithm has to guess the value of  $\langle s^*, w_t \rangle$  using its previous observations. Finally, an oracle reveals as feedback whether the guess is above or below the true value  $\langle s^*, w_t \rangle$ , and the algorithm updates its observation history. In our case,  $w_t := (a_t^1 - a_t^2) / \|a_t^1 - a_t^2\|$  for  $a_t^1, a_t^2 \in \mathcal{A}_t$  and the resulting feedback is given through the agent picking either  $a_t^1$  or  $a_t^2$ . However, extending the binary search to the multidimensional case is non-trivial for two reasons.

**Direction of the multidimensional binary search.** In the contextual bandit setting, we cannot divide the horizon into two successive phases. Indeed, the principal cannot choose any binary search direction in  $\mathbb{R}^d$ , since  $w_t$  depends on the action set  $\mathcal{A}_t$  available at each iteration. For instance, action sets  $\mathcal{A}_t$  could be restricted to a small dimensional subspace of  $\mathbb{R}^d$  during the whole binary search procedure, so that the principal can only get a good estimate of  $s^*$  in this subspace. After this phase, received action sets could be totally different (e.g., in the orthogonal subspace or the whole of  $\mathbb{R}^d$ ) during the remainder of the game.

We solve the issue of constraint directions for the binary search by running it in an adaptive way, depending on the available action set at each time step and on the current level of estimation on this set. More precisely, at iteration  $t$ , the principal's estimate of the true value  $\langle s^*, w_t \rangle$  is  $\langle \hat{s}_t, w_t \rangle$ , where  $\hat{s}_t$  is defined as the centroid of  $\mathcal{S}_t$ :

$$\hat{s}_t := \frac{1}{\operatorname{Vol}(\mathcal{S}_t)} \int_{\mathcal{S}_t} x dx \quad \text{with} \quad \operatorname{Vol}(\mathcal{S}_t) = \int_{\mathcal{S}_t} dx.$$

Whenever  $|\langle \hat{s}_t, w_t \rangle - \langle s^*, w_t \rangle| < 1/T$ , the principal incurs a negligible cost to incentivize the agent to choose her desired action. Then, in this context, for any action  $a \in \mathcal{A}_t$  that the principal wants to play, she designs the incentive

$$\hat{\pi}(t, a) := \max_{a' \in \mathcal{A}_t} \langle \hat{s}_t, a' \rangle - \langle \hat{s}_t, a \rangle + 2/T,$$

$$\hat{\kappa}_a(t, a') = \mathbb{1}_a(a') \hat{\pi}(t, a) \quad \text{for any } a' \in \mathcal{A}_t.$$

To control the precision of the estimation  $\hat{\pi}(t, a)$  of  $\pi^*(t, a)$  for any  $a \in \mathcal{A}_t$ , Lemma 4 shows that it is sufficient to consider the event  $\mathcal{E}_t$ , defined as

$$\mathcal{E}_t := \left\{ \max_{a_t^1 \neq a_t^2 \in \mathcal{A}_t} \operatorname{diam} \left( \mathcal{S}_t, \frac{a_t^1 - a_t^2}{\|a_t^1 - a_t^2\|} \right) < \frac{1}{T} \right\}, \quad (14)$$

where we recall the definition of the projected diameter:  $\operatorname{diam}(\mathcal{S}_t, x) := \max_{s_1, s_2 \in \mathcal{S}_t} |\langle s_1 - s_2, x \rangle|$  for any  $x \in \mathbb{R}^d$ . When  $\mathcal{E}_t$  is false, the principal does not have a good characterization of the incentive function that she needs to provide and thus Contextual IPA runs a multidimensional binary search step, which is explained in the paragraph below. Otherwise, Contextual IPA runs a contextual bandit subroutine `CtxAlg` in a black-box manner on her bandit instance driven by the principal's own mean rewards  $\langle \theta^*, a \rangle$  and the approximated incentives  $\hat{\pi}(t, a)$  for any

$a \in \mathcal{A}_t$ . Lemma 4 guarantees that these approximations are upper estimates of  $\pi^*(t, a)$ . The principal proposes an incentive function  $\hat{\kappa}_{a_t^{\text{Rec}}}$  depending on the estimate to make the agent select the action  $a_t^{\text{Rec}}$  recommended by the bandit subroutine. Again, we do not impose any assumption on the tie-breaking, which can be arbitrary.

**Lemma 4.** *Consider  $t \in [T]$ ,  $\mathcal{A}_t \subseteq \text{B}(0, 1)$ ,  $\mathcal{S}_t \subseteq \text{B}(0, 1)$  such that  $\mathcal{E}_t$  defined in (14) is true. Then for any action  $a \in \mathcal{A}_t$ , we have:  $\pi^*(t, a) < \hat{\kappa}_a(t, a) \leq \pi^*(t, a) + 4/T$ .*

A corollary of Lemma 4 is that running Contextual IPA, under  $\mathcal{E}_t$ ,  $A_t = a_t^{\text{Rec}}$ .

**Issue of the diameter reduction.** We illustrate the challenge of the multidimensional constrained binary search on a very simple problem. At time  $t$ , we can only run a binary search step in one of the directions  $w_t = a - a'$  for  $a, a' \in \mathcal{A}_t$ . Suppose that we have two directions of interest,  $v_1, v_2$  in  $\mathbb{R}^d$ , such that we aim to decrease the diameter of  $\mathcal{S}_t$  in the direction of  $v_1$  or  $v_2$ . Even if we divide the diameter of  $\mathcal{S}_t$  in two in a direction  $w_t$ , which is always possible, this does not necessarily imply that the diameter of  $\mathcal{S}_t$  would reduce along any direction  $v_i$ , as illustrated on Figure 1.

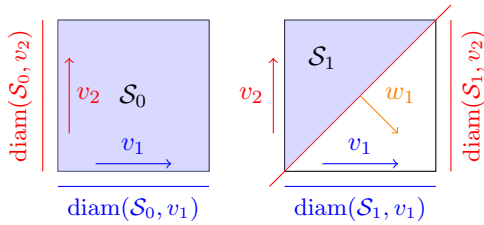


Figure 1: Illustration of a case where the volume  $\mathcal{S}_0$  is cut along a direction  $w_1$  to give a new confidence set  $\mathcal{S}_1$ ; while the diameter is not reduced along the directions  $v_1$  nor  $v_2$ .

An early attempt to tackle this multidimensional binary search problem with adversarial directions was presented by Cohen et al. (2020), who used ellipsoid methods. Here, we use the recent strategy proposed by Lobel et al. (2018) and their Projected Volume subroutine, which is described further in Appendix E.

**Non-stationarity of the reward shift.** For any rewards  $\{(Y_a(t))_{a \in \mathcal{A}_t} : t \in [T]\}$ , we define the history as  $\mathcal{G}_t := (\mathcal{A}_s, A_s, U_s, Y_{A_s}(s))_{s \leq t}$  where  $(U_s)_{s \in \mathbb{N}}$  is a family of independent uniform random variables on  $[0, 1]$  to allow randomization in the decision making. Let  $\text{CtxAlg}$  be a linear contextual bandit algorithm, i.e.,  $\text{CtxAlg} : (U_t, \mathcal{G}_{t-1}, \mathcal{A}_t) \mapsto a_t^{\text{Rec}} \in \mathcal{A}_t$ . When the principal is not running a binary search step, i.e., when  $\mathcal{E}_t$  is true, she plays the  $\text{CtxAlg}$  subroutine on her bandit instance. We define a subset  $I_t$  of all the iterations during which  $\text{CtxAlg}$  is run and a shifted history  $\tilde{\mathcal{H}}_t$  available at time  $t$  as

$$I_t := \{s \in [t] \text{ such that } \mathcal{E}_s \text{ is true}\}, \quad (15)$$

$$\text{and } \tilde{\mathcal{H}}_t := (\mathcal{A}_s, A_s, U_s, X_{a_s^{\text{Rec}}}(s) - \hat{\kappa}_{a_s^{\text{Rec}}}(s, A_s))_{s \in I_t}.$$

In our setup,  $\text{CtxAlg}$  will be fed in a black-box manner with this shifted history to issue recommendations  $a_t^{\text{Rec}}$ ,  $\text{CtxAlg} : (U_t, \tilde{\mathcal{H}}_{t-1}, \mathcal{A}_t) \mapsto a_t^{\text{Rec}}$ . At time  $t$ , for an action  $a_t^{\text{Rec}}$  recommended by  $\text{CtxAlg}$ , our meta-algorithm Contextual IPA proposes an incentive designed so that the agent eventually picks  $a_t^{\text{Rec}}$  (Lemma 4):  $A_t = a_t^{\text{Rec}}$ .

However, the last difference between the non-contextual and contextual cases is that the shift between  $\hat{\pi}(t, a_t)$  and  $\pi^*(t, a_t)$  is not constant anymore on bandit steps  $t \in I_T$ . This shift of rewards is interpreted as adversarial corruption (Bubeck & Slivkins, 2012; Lykouris et al., 2018).

At each round, taking into account this shift, the optimal average utility associated with action  $a \in \mathcal{A}_t$  for the principal is  $r^*(t, a) := \langle \theta^*, a \rangle - \pi^*(t, a)$ , while the principal can only estimate a *non-stationary* expected reward<sup>2</sup>  $r^*(t, a) + \varepsilon_t^{\text{corrupt}}$  with the *corruption level*  $\varepsilon_t^{\text{corrupt}}$  defined as

$$\varepsilon_t^{\text{corrupt}} := \pi^*(t, A_t) - \hat{\pi}(t, A_t) \quad (16)$$

and  $\varepsilon_{I_t}^{\text{corrupt}} := (\varepsilon_s^{\text{corrupt}})_{s \in I_t}$ . In this setup, we can define a corrupted regret as follows

$$R_{\text{CtxAlg}}^{\text{corrupt}}(I_T, \varepsilon_{I_T}^{\text{corrupt}}) = \mathbb{E} \left[ \sum_{t \in I_T} \max_{a \in \mathcal{A}_t} r^*(t, a) - r^*(t, a_t^{\text{Rec}}) \right], \quad (17)$$

where  $a_t^{\text{Rec}} = \text{CtxAlg}(U_t, \tilde{\mathcal{H}}_{t-1}(\varepsilon_{I_{t-1}}^{\text{corrupt}}), \mathcal{A}_t)$  and  $\tilde{\mathcal{H}}_t(\varepsilon_{I_t}^{\text{corrupt}}) = (\mathcal{A}_s, A_s, U_s, r^*(s, a_s^{\text{Rec}}) + \eta_{a_s^{\text{Rec}}}(s) + \varepsilon_s^{\text{corrupt}})_{s \in I_t}$ . Then, we aim to minimize the corrupted regret with  $\text{CtxAlg}$ , which is not possible using a naive linear contextual bandit algorithm.

**Regret analysis.** We split the regret into three components, each of them being bounded separately. One of these components comes from the bias in the estimation of the optimal incentives. Secondly, the principal incurs a cost due to the iterations of  $\text{CtxAlg}$  on the corrupted bandit instance. We use the results from He et al. (2022) with a known corruption level to bound this term. Finally, the last term follows from the multidimensional binary search steps used to estimate  $s^*$ . Lemma 5 allows us to bound the number of such steps; see Appendix D for a proof which builds on the work of Lobel et al. (2018).

**Lemma 5.** *Consider  $\mathcal{E}_t$  defined by (14) with  $(\mathcal{S}_t)_{t \in [T]}$  defined by Contextual IPA. Then it holds almost surely*

<sup>2</sup>Even if we were to feed the stochastic observations  $(X_{A_s}(s) - \hat{\pi}(t, A_s))_{s \leq t}$  at time  $t$ , past algorithmic decisions would depend on different observation distributions, making the direct use of classical regret bounds of the bandit subroutine impossible.

**Algorithm 2** Contextual IPA

---

```

1: Input: horizon  $T$ , subroutine  $\text{CtxAlg}$ ,  $\bar{\delta} = 1/16T^2d(d+1)^2$ 
2: Initialize:  $\mathcal{H}_0 = V_0 = \emptyset$ ,  $\mathcal{S}_0 = \{s \in \mathbb{R}^d : \|s\| \leq 1\}$ 
3: for  $t = 1, \dots, T$  do
4:   Observe available action set  $\mathcal{A}_t$ 
5:   if  $\mathcal{E}_t$  is FALSE then
6:     # Where  $\mathcal{E}_t$  is defined in (14)
7:      $\mathcal{S}_{t+1}, V_{t+1}$ 
8:     = Projected Volume( $T, \bar{\delta}, \mathcal{S}_t, V_t, a_t^1, a_t^2$ )
9:   else
10:    Compute  $\hat{s}_t$  as the centroid of  $\mathcal{S}_t$ 
11:    Sample  $U_t \sim \text{U}(0, 1)$ 
12:    Get  $a_t^{\text{Rec}} = \text{CtxAlg}(U_t, \tilde{\mathcal{H}}_t, \mathcal{A}_t)$ 
13:    Let  $\hat{\pi}(t, a_t^{\text{Rec}}) = \max_{a' \in \mathcal{A}_t} \langle \hat{s}_t, a' \rangle - \langle \hat{s}_t, a_t^{\text{Rec}} \rangle + \frac{2}{T}$ 
14:    Propose incentive function  $\hat{\kappa}_{a_t^{\text{Rec}}}(t, a) = \mathbb{1}_{a_t^{\text{Rec}}}(a) \cdot \hat{\pi}(t, a_t^{\text{Rec}})$ 
15:    Observe  $A_t$  as defined in (8),  $X_{A_t}$ 
16:    Update  $\tilde{\mathcal{H}}_t$  with  $(\mathcal{A}_t, A_t, U_t, X_{A_t} - \kappa(t, A_t))$ 
17:   end if
18: end for

```

---

that

$$\sum_{t \geq 1} \mathbb{1}_{\mathcal{E}_t^c} \leq 192 \cdot d \log(dT),$$

where  $\mathcal{E}_t$  is defined by (14).

**Theorem 2.** *If Contextual IPA is run with any linear contextual subroutine  $\text{CtxAlg}$ , then with the same constant 192 as in Lemma 5, the regret of Contextual IPA is bounded as*

$$\mathfrak{R}(T) \leq 2 + R_{\text{CtxAlg}}^{\text{corrupt}}(I_T, \varepsilon_{I_T}^{\text{corrupt}}) + 1344 d \log(dT).$$

We emphasize that our results still hold with any contextual linear bandit algorithm and that the overall regret is mostly driven by the term  $R_{\text{CtxAlg}}^{\text{corrupt}}$ : although the principal has to solve simultaneously a *pricing-like* problem and a *stochastic bandit* problem, she almost achieves linear bandit state-of-the-art regret.

The main difference between traditional and corrupted rewards in the bandit setting lies in the fact that, in the former case, rewards are typically assumed to be i.i.d., whereas in the latter case, they may be chosen adversarially. An algorithm is robust to corruptions if it yields regret guarantees for any possible reward corruption within a specific budget. This kind of problem was first considered by Bubeck & Slivkins (2012) and has been extensively studied since (see, e.g., Kapoor et al., 2019). In our setting, a bound for the corruption budget is available, thanks to Lemma 6 below.

**Lemma 6.** *Consider  $(\mathcal{E}_t)_{t \in [T]}$  defined by (14) with  $(\mathcal{S}_t)_{t \in [T]}$  defined by Contextual IPA. Let  $I_T$  and  $(\varepsilon_t^{\text{corrupt}})_{t \in [T]}$  as defined in (15) and (16) and  $t \in [T]$  such that  $\mathcal{E}_t$  is true. Then  $|\varepsilon_t^{\text{corrupt}}| \leq 4/T$  and  $\sum_{t \in I_T} |\varepsilon_t^{\text{corrupt}}| := C_{\text{corrupt}} \leq 4$ .*

With standard bandit assumptions, we can then consider a corruption robust algorithms, such as CW – OFUL from He et al. (2022).

**H 3.** *At each round  $t \geq 1$ , for any action  $a \in \mathcal{A}_t$ , the principal’s reward  $X_a(t)$  is  $\mathcal{H}_t$ -conditionally 1-subgaussian, i.e., for any  $\lambda \in \mathbb{R}$ , we have  $\mathbb{E}[e^{\lambda(X_a(t) - \mathbb{E}[X_a(t)])} | \mathcal{H}_{t-1}] \leq e^{\lambda^2/2}$ .*

**Corollary 2.** *Suppose that H3 is true. If Algorithm 2 is run with the subroutine  $\text{CtxAlg} := \text{CW} - \text{OFUL}$  proposed by He et al. (2022), the regularization parameter  $\lambda = 1$  and a confidence level  $\delta = 1/T$ , the following bound holds*

$$\mathfrak{R}(T) \leq 11 + 1344 d \log(dT) + C_{\text{CtxAlg}} d \sqrt{T} \log(T),$$

with  $C_{\text{CtxAlg}}$  being an universal constant.

As in the multi-armed setting, the obtained regret bounds are comparable to the achievable best performance in the standard bandit settings, where the principal does not need to estimate the agent’s parameters  $s^*$ .

## 5 Experiments

We illustrate our theoretical findings with experiments on a toy example and compare IPA with the Principal’s  $\varepsilon$ -Greedy algorithm of Dogan et al. (2023b). We also compare with a UCB Oracle baseline that runs UCB on the shifted bandit instance with arm means  $\mu_a := \theta_a - \pi_a^*$  and no principal-agent consideration. This baseline corresponds to the case where the principal knows the agent’s reward vector and therefore only has to consider a bandit algorithm. Experimental details can be found in Appendix B. We observe in Figure 2 that the Principal  $\varepsilon$ -Greedy Algorithm from Dogan et al. (2023b) exhibits suboptimal performance. Additionally, another issue arises from its computational complexity, requiring an optimization step at every round. In comparison, IPA yields a regret nearly equal to the one of Oracle UCB, illustrating that the cost of estimating the agent’s preferences, obtained from binary search, is negligible for IPA.

## 6 Lower Bounds

For the sake of clarity, we stick to the multi-armed case of Section 3 in this section. A simple observation yields that

$$\mathfrak{R}(T) \geq \mathbb{E} \left[ \sum_{t=1}^T \mu^* - \mu_{A_t} \right],$$



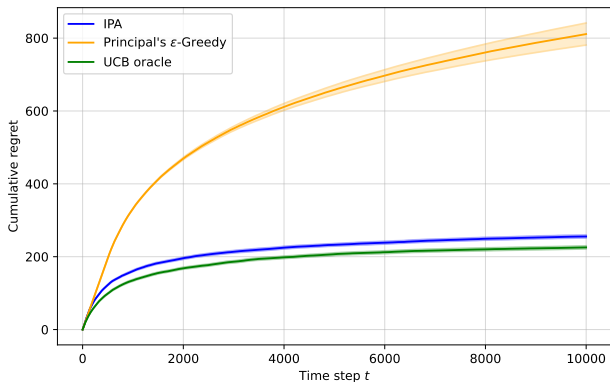


Figure 2: Cumulative regret for different algorithms on a 5 arms instance.

where  $\mu_a = \theta_a - \pi_a^*$  and  $\mu^* = \max_{a \in [K]} \mu_a$ . Even if the principal was to know the optimal incentives  $(\pi_a^*)_{a \in [K]}$ , she would still face a bandit instance with arm means  $\mu_a$ . From there, we can directly extend standard lower bounds from the bandit literature to our setting (Lai & Robbins, 1985; Burnetas & Katehakis, 1996).

**Proposition 1.** *Let  $\mathcal{D}$  be a class of distributions. Consider the multi-armed case of Section 3 and a policy satisfying for any instance  $\nu \in \mathcal{D}^K$  and  $\alpha > 0$ ,  $\mathfrak{R}(T) = o(T^\alpha)$ . Then, for any  $\nu \in \mathcal{D}^K$ ,*

$$\liminf_{T \rightarrow \infty} \frac{\mathfrak{R}(T)}{\log T} \geq \sum_{a, \mu_a < \mu^*} \frac{\mu^* - \mu_a}{\text{KL}_{\text{inf}}(\nu_a - \pi_a^*, \mu^*, \mathcal{D})},$$

where denoting by KL the Kullback-Leibler divergence,

$$\begin{aligned} & \text{KL}_{\text{inf}}(\rho, \mu^*, \mathcal{D}) \\ & := \inf \{ \text{KL}(\rho, \rho') : \rho' \in \mathcal{D}, \int x \rho'(dx) > \mu^* \}, \end{aligned}$$

The complete proof is postponed to Appendix F. Proposition 1 states that IPA yields a nearly optimal regret. Similar arguments can be made in the contextual setting.

## 7 Conclusion and Possible Extensions

This paper presents two novel algorithms called IPA and Contextual IPA, tackling generalizations of both multi-armed and contextual bandits that account for principal-agent interactions. By decoupling the learning of the agent and the estimation of the principal’s parameters, we are able to obtain a nearly optimal algorithm, improving over the previous work of Dogan et al. (2023b). Overall, we obtain an efficient principal-agent bandit framework that allows us to take into account an interaction between a principal and an agent with misaligned interests in a bandit environment. There are various possible extensions of our work, among which considering strategic behaviors for

repeated interactions with a single agent or uncertainty on the agent’s side.

**Information rent.** Again, we consider the multi-armed case for the sake of clarity. We assumed that the agent is always greedy and therefore chooses at time  $t$  an action following

$$A_t = \operatorname{argmax}_{a \in \mathcal{A}_t} \{s_a + \mathbb{1}_{a_t}(a) \pi(t)\}. \quad (18)$$

However, nothing prevents the agent from lying and choosing another action instead of (18). The maximal total welfare that can be extracted at each round is  $\max_{a \in \mathcal{A}} \{s_a + \theta_a\}$ . In our setting, with a trustful agent, this reward was shared between the two actors with an average reward  $\max_{a' \in \mathcal{A}} s_{a'}$  for the agent and  $\max_{a \in \mathcal{A}} \{s_a + \theta_a\} - \max_{a' \in \mathcal{A}} s_{a'}$  for the principal. However, as it is exposed by Dogan et al. (2023b, Section 4), the agent could play with a malicious policy and choose  $A_t$  as if he had different  $(s_a)_{a \in [K]}$ . In that case, he can extract an individual reward  $\max_{a \in \mathcal{A}} \{s_a + \theta_a\} - \min_{a' \in \mathcal{A}} \theta_{a'}$ , while letting a reward  $\min_{a' \in [K]} \theta_{a'}$  to the principal. In that case, the agent exploits his *information rent* to increase his profits. Against such adversarial and powerful agents, the principal cannot do more than play and learn with the  $(s_a)_{a \in [K]}$  announced by the agent. However, this situation is not an issue with myopic agents who act greedily since each of them tries to maximize his own instantaneous reward and eventually select  $A_t$  from (18). This situation is encountered in many applications, where each agent has a single round interaction with the principal for the whole game.

**Learning agents.** A possible extension would be to incorporate uncertainty on the agent’s side and consider learning agents (see, e.g., Dogan et al., 2023a). However again, when considering single round interactions with the agents, each agent myopically maximizes his reward *a priori*. Consequently, the agent policy is stationary and would be driven by  $s_a$ , the expected beliefs on the action rewards. The single interaction model is already well suited for numerous real world applications. In the case of repeated interactions between the principal and a single learning agent, it becomes much more complex, as this agent can both learn his true rewards on the run while trying to influence future actions of the principal with his own choices. Restricting the agent’s policy to a specific set might then be necessary, as done by Dogan et al. (2023a) with Agent’s  $\epsilon$ -Greedy strategy. The major learning difficulty from the principal side would then come from the non-stationarity of the agents decisions and could be handled using non-stationary bandits algorithms (see, e.g., Gittins, 1979; Lattimore & Szepesvári, 2020).

## Acknowledgements

Funded by the European Union (ERC, Ocean, 101071601). Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council Executive Agency. Neither the European Union nor the granting authority can be held responsible for them. The work of DT has been supported by the Paris Île-de-France Région in the framework of DIM AI4IDF.

## Impact Statement

This paper aims to advance the field of principal-agent problems, particularly in the online setting. We believe that this study could lead to several positive real-world outcomes, though we do not feel it is necessary to highlight any specific consequences here. However, we did mention a few relevant domains of interest in the introduction.

## References

- Abe, N. and Long, P. M. Associative reinforcement learning using linear probabilistic concepts. In *ICML*, pp. 3–11, 1999.
- Auer, P. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422, 2002.
- Banihashem, K., Hajiaghayi, M., Shin, S., and Slivkins, A. Bandit social learning: Exploration under myopic behavior. *arXiv preprint arXiv:2302.07425*, 2023.
- Ben-Porat, O., Mansour, Y., Moshkovitz, M., and Taitler, B. Principal-agent reward shaping in mdps. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 9502–9510, 2024.
- Bernasconi, M., Castiglioni, M., Marchesi, A., Gatti, N., and Trovò, F. Sequential information design: Learning to persuade in the dark. *Advances in Neural Information Processing Systems*, 35:15917–15928, 2022.
- Bertsimas, D. and Vempala, S. Solving convex programs by random walks. *Journal of the ACM (JACM)*, 51(4): 540–556, 2004.
- Boursier, E. and Perchet, V. A survey on multi-player bandits. *arXiv preprint arXiv:2211.16275*, 2022.
- Bubeck, S. and Slivkins, A. The best of both worlds: Stochastic and adversarial bandits. In *Conference on Learning Theory*, pp. 42–1. JMLR Workshop and Conference Proceedings, 2012.
- Burnetas, A. N. and Katehakis, M. N. Optimal adaptive policies for sequential allocation problems. *Advances in Applied Mathematics*, 17(2):122–142, 1996.
- Cai, J., Chen, R., Wainwright, M. J., and Zhao, L. Doubly high-dimensional contextual bandits: An interpretable model for joint assortment-pricing. *arXiv preprint arXiv:2309.08634*, 2023.
- Castiglioni, M., Celli, A., Marchesi, A., and Gatti, N. Online bayesian persuasion. *Advances in Neural Information Processing Systems*, 33:16188–16198, 2020.
- Castiglioni, M., Marchesi, A., Celli, A., and Gatti, N. Multi-receiver online bayesian persuasion. In *International Conference on Machine Learning*, pp. 1314–1323. PMLR, 2021.
- Chen, S., Wang, M., and Yang, Z. Actions speak what you want: Provably sample-efficient reinforcement learning of the quantal stackelberg equilibrium from strategic feedbacks. *arXiv preprint arXiv:2307.14085*, 2023.
- Cohen, A., Deligkas, A., and Koren, M. Learning approximately optimal contracts. In *International Symposium on Algorithmic Game Theory*, pp. 331–346. Springer, 2022.
- Cohen, M. C., Lobel, I., and Paes Leme, R. Feature-based dynamic pricing. *Management Science*, 66(11):4921–4943, 2020.
- Conitzer, V. and Garera, N. Learning algorithms for online principal-agent problems (and selling goods online). In *Proceedings of the 23rd international conference on Machine learning*, pp. 209–216, 2006.
- Dani, V., Hayes, T. P., and Kakade, S. M. Stochastic linear optimization under bandit feedback. 2008.
- Den Boer, A. V. Dynamic pricing and learning: historical origins, current research, and new directions. *Surveys in operations research and management science*, 20(1): 1–18, 2015.
- Dogan, I., Shen, Z.-J. M., and Aswani, A. Estimating and incentivizing imperfect-knowledge agents with hidden rewards. *arXiv preprint arXiv:2308.06717*, 2023a.
- Dogan, I., Shen, Z.-J. M., and Aswani, A. Repeated principal-agent games with unobserved agent rewards and perfect-knowledge agents. *arXiv preprint arXiv:2304.07407*, 2023b.
- Doval, L. and Ely, J. C. Sequential information design. *Econometrica*, 88(6):2575–2608, 2020.

- Evans, K. J., Terhorst, A., and Kang, B. H. From data to decisions: helping crop producers build their actionable knowledge. *Critical reviews in plant sciences*, 36(2): 71–88, 2017.
- Gittins, J. C. Bandit processes and dynamic allocation indices. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 41(2):148–164, 1979.
- Golrezaei, N., Jaillet, P., and Liang, J. C. N. Incentive-aware contextual pricing with non-parametric market noise. In *International Conference on Artificial Intelligence and Statistics*, pp. 9331–9361. PMLR, 2023.
- Grötschel, M., Lovász, L., and Schrijver, A. *Geometric algorithms and combinatorial optimization*, volume 2. Springer Science & Business Media, 2012.
- He, J., Zhou, D., Zhang, T., and Gu, Q. Nearly optimal algorithms for linear contextual bandits with adversarial corruptions. *Advances in Neural Information Processing Systems*, 35:34614–34625, 2022.
- Hu, X., Ngo, D., Slivkins, A., and Wu, S. Z. Incentivizing combinatorial bandit exploration. *Advances in Neural Information Processing Systems*, 35:37173–37183, 2022.
- Javanmard, A. and Nazerzadeh, H. Dynamic pricing in high-dimensions. *The Journal of Machine Learning Research*, 20(1):315–363, 2019.
- Kamenica, E. and Gentzkow, M. Bayesian persuasion. *American Economic Review*, 101(6):2590–2615, 2011.
- Kapoor, S., Patel, K. K., and Kar, P. Corruption-tolerant bandit learning. *Machine Learning*, 108(4):687–715, 2019.
- Laffont, J.-J. and Martimort, D. The theory of incentives: the principal-agent model. In *The Theory of Incentives*. Princeton University Press, 2009.
- Lai, T. L. and Robbins, H. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1): 4–22, 1985.
- Lattimore, T. and Szepesvári, C. *Bandit Algorithms*. Cambridge University Press, 2020.
- Li, L., Chu, W., Langford, J., and Schapire, R. E. A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th International Conference on the World Wide Web*, pp. 661–670, 2010.
- Lobel, I., Leme, R. P., and Vladu, A. Multidimensional binary search for contextual decision-making. *Operations Research*, 66(5):1346–1361, 2018.
- Lykouris, T., Mirrokni, V., and Paes Leme, R. Stochastic bandits robust to adversarial corruptions. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 114–122, 2018.
- Mansour, Y., Slivkins, A., and Syrgkanis, V. Bayesian incentive-compatible bandit exploration. *Operations Research*, 68(4):1132–1161, 2020.
- Mao, J., Leme, R., and Schneider, J. Contextual pricing for lipschitz buyers. *Advances in Neural Information Processing Systems*, 31, 2018.
- Myerson, R. B. *Mechanism design*. Springer, 1989.
- Rademacher, L. A. Approximating the centroid is hard. In *Proceedings of the Twenty-Third Annual Symposium on Computational Geometry*, pp. 302–305, 2007.
- Rusmevichientong, P. and Tsitsiklis, J. N. Linearly parameterized bandits. *Mathematics of Operations Research*, 35(2):395–411, 2010.
- Sellke, M. and Slivkins, A. The price of incentivizing exploration: A characterization via thompson sampling and sample complexity. In *Proceedings of the 22nd ACM Conference on Economics and Computation*, pp. 795–796, 2021.
- Simchowitz, M. and Slivkins, A. Exploration and incentives in reinforcement learning. *Operations Research*, 2023.
- Slivkins, A. et al. Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning*, 12(1-2):1–286, 2019.
- Smith, D. J. and Vamanamurthy, M. K. How small is a unit ball? *Mathematics Magazine*, 62(2):101–107, 1989.
- Thompson, W. R. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.
- Woodroffe, M. A one-armed bandit problem with a concomitant variable. *Journal of the American Statistical Association*, 74(368):799–806, 1979.
- Yu, C., Liu, J., Nemati, S., and Yin, G. Reinforcement learning in healthcare: A survey. *ACM Computing Surveys (CSUR)*, 55(1):1–36, 2021.
- Zhong, H., Yang, Z., Wang, Z., and Jordan, M. I. Can reinforcement learning find stackelberg-nash equilibria in general-sum markov games with myopic followers? *arXiv preprint arXiv:2112.13521*, 2021.
- Zhu, B., Bates, S., Yang, Z., Wang, Y., Jiao, J., and Jordan, M. I. The sample complexity of online contract design. *arXiv preprint arXiv:2211.05732*, 2022.

**A Notation**

$\mathcal{A} := [K]$	Set of possible arms.
$T$	Horizon.
$N_T := \log_2 T$	Number of steps dedicated to the binary search on each arm in IPA.
$a_t$	Arm on which the principal offers an incentive.
$\pi(t)$	Amount of incentive offered by the principal on action $a_t$ .
$A_t$	Arm chosen by the agent, maximizing his utility, known by everyone.
$s_a + \mathbb{1}_{a_t}(a)\pi(t)$	Agent's utility for action $a$ .
$\nu_a$	Principal's reward distribution for action $a$ .
$X_a(t) - \mathbb{1}_{a_t}(a)\pi(t)$	Principal's utility for action $a$ .
$\mu_a$	Principal's expected utility for action $a$ , using the optimal incentive $\pi_a^*$ .
$\mu^*$	Maximal expected utility for the principal.
$\theta_a := \mathbb{E}[X_a(1)]$	Principal's expected reward.
$\pi_a^*$	Infimum amount of incentives to be offered on action $a_t = a$ to make the agent choose it.
$\tilde{\mathcal{H}}$	Shifted history used to feed Alg.
$R_{\text{Alg}}(T)$	Regret of the subroutine Alg on a horizon $T$ .
$\mathfrak{R}(T)$	Overall regret of IPA on a horizon $T$ .

Table 1: Notations used in Section 3.

$T$	Horizon.
$B(0, 1)$	Unit ball in $\mathbb{R}^d$ , $d \geq 1$ .
$\mathcal{A}_t \subseteq B(0, 1)$	Action set at time $t$ among which the agent selects $A_t$ .
$\pi(t)$	Amount of incentive offered by the principal on some action.
$\eta_a(t)$	Noise distribution of the principal's reward associated with action $a$ at time $t$ .
$r^*(t, a) = \langle \theta^*, a \rangle + \eta_a(t) - \pi^*(t, a)$	Utility collected by the principal on action $a$ at time $t$ if the optimal amount of incentive is used.
$\mu_t^*$	Principal's maximal expected utility at time $t$ .
$\pi^*(t, a)$	Infimal amount of incentives to be offered on action $a$ with $\kappa(t, a') = \mathbb{1}_a(a')\pi_a^*$ so that the agent eventually chooses action $a$ .
$\hat{\pi}(t, a)$	Principal's estimation of $\pi^*(\cdot, a)$
$\kappa(t, \cdot): B(0, 1) \rightarrow \mathbb{R}_+$	Incentive function, associating each action with some amount of incentives.
$\mathbf{s}^* \in B(0, 1)$	Agent's true reward vector.
$\mathcal{S}_t \ni \mathbf{s}^*$	Principal's confidence set for $\mathbf{s}^*$ at time $t$ .
$\langle \mathbf{s}^*, a \rangle + \kappa(t, a)$	Agent's utility for action $a$ .
$a_t^{\text{ag}} = \operatorname{argmax}_{a \in \mathcal{A}_t} \langle \mathbf{s}^*, a \rangle$	Agent's optimal action with null incentives at time $t$ .
$a_t^{\text{Rec}}$	Action recommended by CtxAlg at time $t$ .
$\langle \theta^*, a \rangle + \eta_a(t) - \kappa(t, a)$	Principal's utility for action $a$ .
$\mu_t^*$	Maximal expected utility for the principal at time $t$ .
$\varepsilon_t^{\text{corrupt}} = \pi^*(t, a) - \hat{\pi}(t, a)$	Shift between the optimal incentives and the estimated ones.
$C_{\text{corrupt}}$	Total corruption budget due to the shift between $\hat{\pi}(t, a)$ and $\pi^*(t, a)$ over the rounds.
$\mathcal{E}_t$	Event being true if the diameter of $\mathcal{S}_t$ projected $\mathcal{A}_t$ is small than $1/T$ .
$I_t$	Rounds up to time $t$ during which $\mathcal{E}_s$ , $s \leq t$ is true.
$\tilde{\mathcal{H}}$	Shifted history used to feed Alg.
$R_{\text{CtxAlg}}(T)$	Regret of the subroutine CtxAlg on a horizon $T$ .
$\mathfrak{R}(T)$	Overall regret of Contextual IPA on a horizon $T$ .

Table 2: Notation used in Section 4.

## B Experimental details

We ran the experiments in Figure 2 for a horizon  $T = 10\,000$  on an average of 100 runs on a five arms bandit. We plotted the standard error across the different runs. The expected rewards for the principal ( $\theta$ ) and the agent (s) are given in Table 3. The principal’s rewards  $X_a(t)$  are drawn from an i.i.d. distribution  $X_a(t) \sim \mathcal{N}(\theta_a, 1)$  for any  $a \in [K], t \in [T]$ . We also run an oracle UCB instance with rewards following a Gaussian distribution  $\mathcal{N}(\mu_a, 1)$  where for any  $a \in [K]$ ,  $\mu_a := \theta_a + s_a - \max_{a' \in [K]} s_{a'}$ , as if a UCB algorithm was run with the full knowledge of the optimal incentives and was learning his own mean rewards ( $\mu$ ), taking into account these incentives. The mean rewards  $\mu$  are also given in Table 3. We observe that the additional exploration steps needed to learn the optimal incentives in IPA are not very costly compared to the regret achieved by the UCB oracle.

For the Principal’s  $\varepsilon$ -Greedy algorithm, we use the hyperparameters  $\alpha = 1$  and  $m = 500$ . The hyperparameter  $m$  controls the number of exploration steps. We ran the Principal’s  $\varepsilon$ -Greedy algorithm on the same bandit setting for different values  $m = 30, m = 100, m = 200, m = 300, m = 400, m = 500, m = 600, m = 800, m = 1000, m = 2000, m = 5000, m = 10\,000$ . Below  $m = 500$ , the algorithm does not explore enough and incurs a linear regret on some runs, consequently yielding a poor mean regret, whereas above  $m = 500$ , the algorithm explores excessively, leading to a higher regret due to overexploration. We ran the same experiments on longer horizons  $T = 100\,000$  and  $T = 1\,000\,000$  and the algorithms exhibited the same behavior. In practice, the tuning of the  $\varepsilon$ -Greedy algorithm depends on the reward gaps and is not common to use. This is why another advantage of IPA compared to the Principal’s  $\varepsilon$ -Greedy algorithm of Dogan et al. (2023b) lies in the fact that it does not need any tuning of hyperparameters, leading to a better use in practice, on potentially broader bandit instances.

s	0.64	0.99	0.73	0.61	0.59
$\theta$	0.30	0.24	0.88	0.07	0.65
$\mu$	-0.05	0.24	0.62	-0.31	0.25

Table 3: Experimental parameters for Figure 2.

We did not run Contextual IPA in a contextual bandit setting because it is quite tedious to implement, due to the use of the Projected Volume subroutine from the work of Lobel et al. (2018). Even though they obtain an excellent regret bound, the computations raise specific challenges. The first issue is the computation of the centroid which is known to be a #P-hard problem (Rademacher, 2007). However, it can be solved through an approximation of the centroid, which is computable in polynomial time (see, Bertsimas & Vempala, 2004, Lemma 5 and Theorem 12). A second issue is finding directions along which the set  $\mathcal{S}_t$  has a small diameter, which is needed to compute the set  $V_t$ . It is solved by Lobel et al. (2018) with an ellipsoidal approximation  $E$  of  $\mathcal{S}_t$  such that  $E \subseteq \mathcal{S}_t \subseteq \alpha E$  with  $\alpha > 1$ , since such an ellipsoid can be computed in polynomial time, (see, Grötschel et al., 2012, Corollary 4.6.9). Such a variation of the Projected Volume subroutine is presented in the work of Lobel et al. (2018, Section 9.3). It is shown that one can achieve polynomial time computations with still the same regret bound for the multidimensional binary search steps (Lobel et al., 2018, Theorem 9.4). This line of work needs to be explored for implementing Contextual IPA in practice, which is feasible but still requires a significant amount of work.

## C Regret Bound for Non-Contextual Setting

**Notations.** We define  $\bar{\pi}_a(t) \in \mathbb{R}_+$  as the upper estimate and  $\underline{\pi}_a(t) \in \mathbb{R}_+$  as the lower estimate of  $\pi_a^*$  after  $t$  rounds of binary search on arm  $a$ . For any  $t \in [T]$  and  $a \in [K]$ , we define  $\pi_a^{\text{mid}}(t) := (\bar{\pi}_a(t) + \underline{\pi}_a(t))/2$ . We define  $N_T := \lceil \log_2 T \rceil$  as the number of binary search steps per arm and  $\hat{\pi}_a$  is the estimated incentive to make the agent choose action  $a$  after  $N_T$  steps of binary search:  $\hat{\pi}_a = \bar{\pi}_a(N_T) + 1/T$ . Since our problem is stationary, we write  $a^{\text{pr}} := \operatorname{argmax}_{a \in [K]} \{\theta_a - \pi_a^*\}$  for the optimal action that the principal could aim to play at each round.

**Lemma 1.** For any  $T \in \mathbb{N}$ , the regret of any algorithm on our problem instance can be written as

$$\mathfrak{R}(T) = T \max_{a \in [K]} \{ \theta_a + s_a - \max_{a' \in [K]} s_{a'} \} - \mathbb{E} \left[ \sum_{t=1}^T \{ \theta_{A_t} - \mathbb{1}_{a_t}(A_t) \pi(t) \} \right].$$

*Proof of Lemma 1.* Recall that the regret is defined in (3) as  $\mathfrak{R}(T) := T\mu^* - \sum_{t=1}^T \mathbb{E}[X_{A_t}(t) - \mathbb{1}_{a_t}(A_t)\pi(t)]$ , where  $\mu^* = \sup_{a \in [K], \pi \in \mathbb{R}_+} \mathbb{E}_\nu[X_a(1)] - \pi$ , such that  $a \in \arg\max_{a' \in [K]} \{s_{a'} + \pi\}$ . Note that we can write  $\mu^* = \sup_{a \in [K], \pi \in \mathbb{R}_+} \{\theta_a - \mathbb{1}_{\tilde{A}}(a, \pi)\pi\}$ , where  $\tilde{A} := \{(a, \pi) : s_a + \pi \geq \max_{a'} s_{a'} + \mathbb{1}_a(a')\pi\}$ . First note that if  $(a, \pi) \in \tilde{A}$ , then  $s_a - s_{a'} \geq -\pi$  for any  $a' \in [K]$ , which implies by definition of the optimal incentives (4) that  $\pi_a^* \leq \pi$ . Consequently,

$$\mu^* = \max_{a \in [K]} \{\mathbb{E}_\nu[X_{a^{\text{opt}}}(1)] - \pi_a^*\} = \max_{a \in [K]} \{\theta_a - \max_{a' \in [K]} \{s_{a'}\} + s_a\},$$

hence our result about the regret.  $\square$

**Lemma 7.** Assume **H1** and that we run Algorithm 3 for an action  $a \in [K]$  and a number of binary searches  $N_T \in \mathbb{N}$ . Then, for any  $t \in [N_T]$ ,  $0 \leq \underline{\pi}_a(t) \leq \pi_a^{\text{mid}}(t) \leq \bar{\pi}_a(t) \leq 1$ .

*Proof.* The proof is by induction on  $t \in [N_T]$ . For  $t = 0$ , it is defined by definition. Then suppose that it holds true for  $t \geq 0$ . Note that line 6 in Algorithm 3 can be written as

$$\begin{aligned} \bar{\pi}_a(t+1) &= \mathbb{1}_a(A_t)\pi_a^{\text{mid}}(t) + (1 - \mathbb{1}_a(A_t))\bar{\pi}_a(t) \\ \underline{\pi}_a(t+1) &= (1 - \mathbb{1}_a(A_t))\pi_a^{\text{mid}}(t) + \mathbb{1}_a(A_t)\underline{\pi}_a(t), \end{aligned} \quad (19)$$

which completes the proof by applying the induction hypothesis.  $\square$

**Lemma 8.** Assume **H1** and that we run Algorithm 3 for an action  $a \in [K]$  and a number of binary searches  $N_T \in \mathbb{N}$ . Then, for any  $t \in [N_T]$ ,

$$\pi_a^* \in [\underline{\pi}_a(t), \bar{\pi}_a(t)] \text{ and } |\bar{\pi}_a(t) - \underline{\pi}_a(t)| \leq 1/2^t.$$

*Proof.* The proof is by induction on  $t$ . The case  $t = 0$  is trivial by the initialization of Algorithm 3 and **H1**.

Suppose that the statement holds for  $t$ . Note that  $\mathbb{1}(\{A_t = a\}) = \mathbb{1}(\{\pi_a^{\text{mid}}(t) \geq \pi_a^*\})$ , therefore using (19), we obtain by using the induction hypothesis that

$$\pi_a^* \in [\underline{\pi}_a(t+1), \bar{\pi}_a(t+1)], \quad \bar{\pi}_a(t+1) - \underline{\pi}_a(t+1) = \frac{\bar{\pi}_a(t) - \underline{\pi}_a(t)}{2},$$

which completes the proof.  $\square$

*Proof of Theorem 1.* Recall that Lemma 1 implies that

$$\mathfrak{R}(T) = \mathbb{E} \left[ \sum_{t=1}^T \max_{a \in [K]} \{\theta_a - \pi_a^*\} - (X_{A_t}(t) - \mathbb{1}_{a_t}(A_t)\pi(t)) \right] \quad (20)$$

We decompose the regret between the  $K \lceil \log_2 T \rceil = KN_T$  first steps during which we run the Binary Search Subroutine and all the subsequent ones

$$\begin{aligned} \mathfrak{R}(T) &= (\mathbf{A}) + (\mathbf{B}) \\ (\mathbf{A}) &= \mathbb{E} \left[ \sum_{t=1}^{K \lceil \log_2 T \rceil} \max_{a \in [K]} \{\theta_a - \pi_a^*\} - \{X_{A_t}(t) - \mathbb{1}_{a_t}(A_t)\pi(t)\} \right] \\ (\mathbf{B}) &= \mathbb{E} \left[ \sum_{t=K \lceil \log_2 T \rceil + 1}^T \max_{a \in [K]} \{\theta_a - \pi_a^*\} - \{X_{A_t}(t) - \mathbb{1}_{a_t}(A_t)\pi(t)\} \right]. \end{aligned}$$

We separate the analysis of the regret, bounding independently two terms in the right-hand side of the previous decomposition. Since  $\pi(t)$  is always equal to  $\pi_a^{\text{mid}}(t)$  for some  $t \in [N_T]$ ,  $a \in [K]$ , during the binary search phase, we use Lemma 7 to bound  $\pi(t)$  by 1 for any  $t \leq K \lceil \log_2 T \rceil$  in **(A)**, giving

$$(\mathbf{A}) = \mathbb{E} \left[ \sum_{t=1}^{K \lceil \log_2 T \rceil} \max_{a \in [K]} \{\theta_a - X_{A_t}(t) + \mathbb{1}_{a_t}(A_t)(t)\pi(t) - \pi_{a^{\text{pr}}}^*\} \right] \quad (21)$$

$$\leq \sum_{t=1}^{K \lceil \log_2 T \rceil} (1 + \max_{a \in [K]} \{\theta_a\} - \min_{a \in [K]} \{\theta_a\}) \leq (1 + \max_{a \in [K]} \{\theta_a\} - \min_{a \in [K]} \{\theta_a\}) K (1 + \log_2 T). \quad (22)$$

At the end of the binary search phase and for all the subsequent rounds  $t > K \lceil \log_2 T \rceil$ , Alg recommends an action  $a_t^{\text{Rec}}$  and the principal proposes the incentive  $\pi(t) = \hat{\pi}_{a_t^{\text{Rec}}} = \bar{\pi}_{a_t^{\text{Rec}}}(\lceil \log_2 T \rceil) + 1/T$  on action  $a_t^{\text{Rec}}$  to make the agent choose it. Lemma 8 ensures that after  $\lceil \log_2 T \rceil$  rounds of binary search on action  $a \in [K]$ , we have

$$\underline{\pi}_a(\lceil \log_2 T \rceil) \leq \pi_a^* \leq \bar{\pi}_a(\lceil \log_2 T \rceil) \quad \text{and} \quad \bar{\pi}_a(\lceil \log_2 T \rceil) - \underline{\pi}_a(\lceil \log_2 T \rceil) \leq 1/2^{\lceil \log_2 T \rceil} \leq 1/T.$$

Therefore,

$$\pi_a^* < \hat{\pi}_a \quad \text{and} \quad \hat{\pi}_a - \pi_a^* \leq 2/T.$$

For the agent, the utility associated with action  $a_t^{\text{Rec}}$  is  $\mathfrak{s}_{a_t^{\text{Rec}}} + \hat{\pi}_{a_t^{\text{Rec}}} > \mathfrak{s}_{a_t^{\text{Rec}}} + \pi_{a_t^{\text{Rec}}}^*$ , which guarantees that he eventually selects  $a_t^{\text{Rec}}$  at time  $t$  because of (1) and (4). It ensures that for any  $t > K \lceil \log_2 T \rceil$ ,  $A_t = a_t^{\text{Rec}}$ .

To compute these recommendations, Alg is fed at any time  $t \in \mathbb{N}$  with the shifted history defined in (7):  $\tilde{\mathcal{H}}_t = (a_s^{\text{Rec}}, U_s, X_{a_s^{\text{Rec}}}(s) - \hat{\pi}_{a_s^{\text{Rec}}})_{s \in [K \lceil \log_2 T \rceil + 1, t]}$ . Recall that we defined the shifted distribution  $\rho_a^T$  for any  $a \in [K]$  as the distribution of  $X_a(1) - \hat{\pi}_a$ . For any  $t > K \lceil \log_2 T \rceil$ ,  $a_t^{\text{Rec}} = \text{Alg}(U_t, \tilde{\mathcal{H}}_{t-1})$  and we define  $Y_a(t) \sim \rho_a^T$  for any  $t \in [K \lceil \log_2 T \rceil + 1, T]$ ,  $a \in [K]$ . In this setup, the regret of Alg after  $\tau$  subsequent steps is defined as

$$R_{\text{Alg}}(\tau, \rho^T) = \tau \max_{a \in [K]} \mathbb{E}_{\rho^T} [Y_a(K \lceil \log_2 T \rceil + 1)] - \mathbb{E} \left[ \sum_{s=K \lceil \log_2 T \rceil + 1}^{K \lceil \log_2 T \rceil + \tau} Y_{\text{Alg}(U_s, \tilde{\mathcal{H}}_{s-1})}(s) \right].$$

Consequently, since  $a_t = a_t^{\text{Rec}} = A_t$

$$\begin{aligned} \text{(B)} &= \mathbb{E} \left[ \sum_{t=K \lceil \log_2 T \rceil + 1}^T \max_{a \in [K]} \{\theta_a - \pi_a^*\} - (X_{a_t^{\text{Rec}}}(t) - \hat{\pi}_{a_t^{\text{Rec}}}) \right] \\ &\leq \mathbb{E} \left[ \sum_{t=K \lceil \log_2 T \rceil + 1}^T \max_{a \in [K]} \left\{ \theta_a - \hat{\pi}_a - (X_{a_t^{\text{Rec}}}(t) - \hat{\pi}_{a_t^{\text{Rec}}}) \right\} + \max_{a' \in [K]} \{\hat{\pi}_{a'} - \pi_{a'}^*\} \right] \\ &= \mathbb{E} \left[ \sum_{t=K \lceil \log_2 T \rceil + 1}^T \max_{a \in [K]} \{\theta_a - \hat{\pi}_a\} - (X_{a_t^{\text{Rec}}}(t) - \hat{\pi}_{a_t^{\text{Rec}}}) \right] + \mathbb{E} \left[ \sum_{t=K \lceil \log_2 T \rceil + 1}^T \max_{a' \in [K]} \{\hat{\pi}_{a'} - \pi_{a'}^*\} \right] \\ &= (T - K \lceil \log_2 T \rceil) \max_{a \in [K]} \mathbb{E}[Y_a(K \lceil \log_2 T \rceil + 1)] - \mathbb{E} \left[ \sum_{t=K \lceil \log_2 T \rceil + 1}^T (X_{a_t^{\text{Rec}}}(t) - \hat{\pi}_{a_t^{\text{Rec}}}) \right] \\ &\quad + (T - K \lceil \log_2 T \rceil) \max_{a' \in [K]} \{\hat{\pi}_{a'} - \pi_{a'}^*\} \\ &\leq R_{\text{Alg}}(T - K \lceil \log_2 T \rceil, (\rho_a^T)_{a \in [K]}) + 2. \end{aligned}$$

Plugging (A) and (B) together finally gives a bound for the regret

$$\mathfrak{R}(T) \leq 2 + (1 + \max_{a \in [K]} \{\theta_a\} - \min_{a \in [K]} \{\theta_a\})(1 + K \log_2 T) + R_{\text{Alg}}(T - K \lceil \log_2 T \rceil, \rho^T).$$

□

*Proof of Corollary 1.* The case  $T \leq 9K$  is trivial. Assume then that  $T \geq 9K$ . Note that after  $\lceil \log_2 T \rceil$  rounds of binary search,  $\pi_a^* \leq \bar{\pi}_a \leq \pi_a^* + 1/T$ . We define  $\Delta_a^* := \max_{a' \in [K]} \{\theta_{a'} - \pi_{a'}^*\} - (\theta_a - \pi_a^*) = \max_{a' \in [K]} \{\theta_{a'} + \mathfrak{s}_{a'}\} - (\theta_a + \mathfrak{s}_a)$  and  $\tilde{\Delta}_a := \max_{a' \in [K]} \{\theta_{a'} - \hat{\pi}_{a'}\} - (\theta_a - \hat{\pi}_a) = \max_{a' \in [K]} \{\theta_{a'} - \bar{\pi}_{a'}\} - (\theta_a - \bar{\pi}_a)$  using  $\hat{\pi}_a = \bar{\pi}_a + 1/T$ . Since  $\pi_a^* \leq \bar{\pi}_a \leq \pi_a^* + 1/T$ , we have  $|\Delta_a^* - \tilde{\Delta}_a| \leq 2/T$ .

Using the results about UCB algorithm that can be found in (Lattimore & Szepesvári, 2020, Theorems 7.1 and 7.2), since Alg is run in a black-box manner on a shifted bandit instance  $\rho^T$  for  $T - K \lceil \log_2 T \rceil$  rounds with reward gaps  $\tilde{\Delta}_a$ , we have

$$\begin{aligned}
 R_{\text{Alg}}(T - K \lceil \log_2 T \rceil, \rho^T) &\leq 3 \sum_{\tilde{\Delta}_a > 0} \tilde{\Delta}_a \\
 &\quad + 8 \min \left\{ \sqrt{(T - K \lceil \log_2 T \rceil) K \log(T - K \lceil \log_2 T \rceil)} ; \sum_{\tilde{\Delta}_a > 0} \frac{2 \log(T - K \lceil \log_2 T \rceil)}{\tilde{\Delta}_a} \right\} \\
 &\leq 3 \sum_{a \in [K], \Delta_a^* > 0} \left( \Delta_a^* + \frac{2}{T} \right) + 8 \min \left\{ \sqrt{TK \log T} ; \sum_{a \in [K], \tilde{\Delta}_a > 0} \frac{2 \log T}{\tilde{\Delta}_a} \right\} \\
 &\leq 1 + 3 \sum_{a \in [K], \Delta_a^* > 0} \Delta_a^* + 8 \min \left\{ \sqrt{TK \log T} ; \sum_{a \in [K], \tilde{\Delta}_a > 0} \frac{2 \log T}{\tilde{\Delta}_a} \right\}, \tag{23}
 \end{aligned}$$

where the last line holds because of  $T \geq 9K > 6K$ .

We now analyse the sum  $\sum_{a \in [K], \tilde{\Delta}_a > 0} 2 \log T / \tilde{\Delta}_a$  and consider two cases: either there exists  $\tilde{a} \in [K]$  such that  $\Delta_{\tilde{a}}^* \leq 4/T$  or not.

*First case:* if there exists  $\tilde{a} \in [K]$  such that  $\Delta_{\tilde{a}}^* \leq 4/T$ , since  $T > 9K$ , we have for such an action  $\tilde{a}$ :  $2 \log T / \Delta_{\tilde{a}}^* \geq T \log T / 2 > \sqrt{TK \log T}$  as well as  $\tilde{\Delta}_{\tilde{a}} \leq \Delta_{\tilde{a}}^* + 2/T \leq 6/T$  which is equivalent to  $2 \log T / \tilde{\Delta}_{\tilde{a}} \geq T \log T / 3 \geq \sqrt{TK \log T}$ . Consequently,

$$\min \left\{ \sqrt{TK \log T} ; \sum_{a \in [K], \tilde{\Delta}_a > 0} \frac{2 \log T}{\tilde{\Delta}_a} \right\} = \sqrt{TK \log T} = \min \left\{ \sqrt{TK \log T} ; \sum_{a \in [K], \Delta_a^* > 0} \frac{2 \log T}{\Delta_a^*} \right\}. \tag{24}$$

*Second case:* for any  $a \in [K]$ ,  $\Delta_a^* > 4/T$ . Therefore  $\tilde{\Delta}_a \geq \Delta_a^* - 2/T > \Delta_a^* - \Delta_a^*/2 = \Delta_a^*/2$ . Consequently

$$\begin{aligned}
 R_{\text{Alg}}(T - K \lceil \log_2 T \rceil, \rho^T) &\leq 1 + 3 \sum_{a \in [K], \Delta_a^* > 0} \Delta_a^* + 8 \min \left\{ \sqrt{TK \log T} ; \sum_{a \in [K], \tilde{\Delta}_a > 0} \frac{2 \log T}{\tilde{\Delta}_a} \right\} \\
 &\leq 1 + 3 \sum_{a \in [K], \Delta_a^* > 0} \Delta_a^* + 8 \min \left\{ \sqrt{TK \log T} ; \sum_{a \in [K], \Delta_a^* > 0} \frac{4 \log T}{\Delta_a^*} \right\}. \tag{25}
 \end{aligned}$$

Finally, combining (24) and (25) in (23) completes the proof.  $\square$

## D Regret Bound for the Contextual Setting

For the whole section, we define  $\hat{a}_t^{\text{ag}} := \operatorname{argmax}_{a \in \mathcal{A}_t} \langle \hat{s}_t, a \rangle$  and recall that  $a_t^{\text{ag}} = \operatorname{argmax}_{a \in \mathcal{A}_t} \langle s^*, a \rangle$ .

### D.1 Technical lemmas

**Lemma 2.** For any  $T \in \mathbb{N}$ , the regret of any algorithm on our contextual problem instance can be written as

$$\begin{aligned}
 \mathfrak{R}(T) &= \sum_{t=1}^T \max_{a \in \mathcal{A}_t} \{ \langle \theta^* + s^*, a \rangle - \max_{a' \in \mathcal{A}_t} \langle s^*, a' \rangle \} \\
 &\quad - \mathbb{E} \left[ \sum_{t=1}^T (\langle \theta^*, A_t \rangle - \kappa(t, A_t)) \right].
 \end{aligned}$$

*Proof of Lemma 2.* Recall that the regret is defined in (11) as  $\mathfrak{R}(T) = \sum_{t=1}^T \mu_t^* - \mathbb{E} \left[ \sum_{t=1}^T (X_{A_t}(t) - \kappa(t, A_t)) \right]$ , where  $\mu_t^* := \sup_{\kappa(t, \cdot): \mathbb{R}^d \rightarrow \mathbb{R}_+} \{ \langle \theta^*, a \rangle - \kappa(t, a) \}$  such that  $a \in \operatorname{argmax}_{a' \in \mathcal{A}_t} \{ \langle s^*, a' \rangle + \kappa(t, a') \}$ . Note that we can write



$\mu_t^* = \sup_{a \in \mathcal{A}_t, \pi \in \mathbb{R}_+} \{\langle \theta^*, a \rangle - \mathbb{1}_{\tilde{\mathcal{A}}_t}(a, \pi)\pi\}$  where  $\tilde{\mathcal{A}}_t := \{(a, \pi) : \langle s^*, a \rangle + \pi \geq \max_{a' \in \mathcal{A}_t} \langle s^*, a' \rangle + \mathbb{1}_a(a')\pi\}$ . First note that if  $(a, \pi) \in \tilde{\mathcal{A}}_t$ , then  $\langle s^*, a - a' \rangle \geq -\pi$  for any  $a' \in \mathcal{A}_t$ , which implies by definition of the optimal incentives (12) that  $\pi^*(t, a) \leq \pi$ . Consequently

$$\mu_t^* = \max_{a \in \mathcal{A}_t} \{\langle \theta^*, a \rangle - \pi^*(t, a)\} = \max_{a \in \mathcal{A}_t} \{\langle \theta^*, a \rangle - \max_{a' \in \mathcal{A}_t} \{\langle s^*, a' \rangle\} + \langle s^*, a \rangle\},$$

hence our result about the regret.  $\square$

**Lemma 3.** For any  $t \in [T]$  and closed subset  $\mathcal{S} \subset B(0, 1)$  with  $s^* \in \mathcal{S}$ , it holds, for any  $a \in \mathcal{A}_t$ ,  $|\max_{s \in \mathcal{S}, a' \in \mathcal{A}_t} \langle s, a' - a \rangle - \pi^*(t, a)| \leq 2 \text{diam}(\mathcal{S}, \mathcal{A}_t)$  where  $\text{diam}(\mathcal{S}, \mathcal{A}_t) := \max_{a' \in \mathcal{A}_t} \max_{s_1, s_2 \in \mathcal{S}} |\langle s_1 - s_2, a' \rangle|$ .

*Proof of Lemma 3.* For any  $t \in [T]$ ,  $\mathcal{S} \in B(0, 1)$  with  $s^* \in \mathcal{S}$ ,  $\mathcal{A}_t \in B(0, 1)$  and  $a \in \mathcal{A}_t$ , recall that we defined  $a_t^{\text{ag}} = \arg\max_{a' \in \mathcal{A}_t} \langle s^*, a' \rangle$  and  $\pi^*(t, a) = \langle s^*, a_t^{\text{ag}} \rangle - \langle s^*, a \rangle$ . Consequently, defining for any  $s \in \mathcal{S}$ ,  $a_t^s := \arg\max_{a' \in \mathcal{A}_t} \langle s, a' \rangle$  (the compactness of both  $\mathcal{S}$  and  $\mathcal{A}_t$  as well as the continuity of the applications that we consider guarantee the existence of such an argmax), we have, since  $\langle s^*, a_t^{\text{ag}} \rangle \geq \langle s^*, a_t^s \rangle$  for any  $s \in \mathcal{S}$  and associated  $a_t^s$ ,

$$\begin{aligned} \max_{s \in \mathcal{S}, a' \in \mathcal{A}_t} \langle s, a' - a \rangle - \pi^*(t, a) &= \max_{s \in \mathcal{S}} \max_{a' \in \mathcal{A}_t} \{\langle s, a' - a \rangle - \langle s^*, a_t^{\text{ag}} - a \rangle\} \\ &= \max_{s \in \mathcal{S}} \{\langle s, a_t^s - a \rangle - \langle s^*, a_t^{\text{ag}} - a \rangle\} \\ &\leq \max_{s \in \mathcal{S}} \{\langle s, a_t^s - a \rangle - \langle s^*, a_t^s - a \rangle\} \\ &\leq \max_{s \in \mathcal{S}} |\langle s - s^*, a_t^s \rangle| + \max_{s \in \mathcal{S}} |\langle s - s^*, a \rangle| \\ &\leq 2 \text{diam}(\mathcal{S}, \mathcal{A}_t). \end{aligned}$$

Similarly, we have

$$\begin{aligned} \pi^*(t, a) - \max_{s \in \mathcal{S}, a' \in \mathcal{A}_t} \langle s, a' - a \rangle &\leq \langle s^*, a_t^{\text{ag}} - a \rangle - \max_{s \in \mathcal{S}} \langle s, a_t^{\text{ag}} - a \rangle \\ &\leq \max_{s \in \mathcal{S}} |\langle s^* - s, a_t^{\text{ag}} \rangle| + \max_{s \in \mathcal{S}} |\langle s^* - s, a \rangle| \\ &\leq 2 \text{diam}(\mathcal{S}, \mathcal{A}_t), \end{aligned}$$

and the proof follows.  $\square$

**Lemma 4.** Consider  $t \in [T]$ ,  $\mathcal{A}_t \subseteq B(0, 1)$ ,  $\mathcal{S}_t \subseteq B(0, 1)$  such that  $\mathcal{E}_t$  defined in (14) is true. Then for any action  $a \in \mathcal{A}_t$ , we have:  $\pi^*(t, a) < \hat{\kappa}_a(t, a) \leq \pi^*(t, a) + 4/T$ .

*Proof.* The proof is similar to the proof of Lemma 3. Consider  $t \geq 1$ ,  $\mathcal{A}_t \subseteq B(0, 1)$ ,  $\mathcal{S}_t \subseteq B(0, 1)$  such that  $\mathcal{E}_t$  holds,  $a_t \in \mathcal{A}_t$ . Then  $1/T > \max_{a_t^1 \neq a_t^2 \in \mathcal{A}_t} \text{diam}(\mathcal{S}_t, (a_t^1 - a_t^2)/\|a_t^1 - a_t^2\|)$ . Recall that we defined  $a_t^{\text{ag}} = \arg\max_{a \in \mathcal{A}_t} \langle s^*, a \rangle$  and  $\hat{a}_t^{\text{ag}} = \arg\max_{a \in \mathcal{A}_t} \langle \hat{s}_t, a \rangle$ ,  $\pi^*(t, a_t) = \langle s^*, a_t^{\text{ag}} \rangle - \langle s^*, a_t \rangle$  and  $\hat{\pi}(t, a_t) = \langle \hat{s}_t, \hat{a}_t^{\text{ag}} \rangle - \langle \hat{s}_t, a_t \rangle + 2/T$ , giving

$$\hat{\pi}(t, a_t) \geq \langle \hat{s}_t, a_t^{\text{ag}} \rangle - \langle \hat{s}_t, a_t \rangle + 2/T.$$

Therefore, under  $\mathcal{E}_t$ , it holds

$$\begin{aligned} \hat{\pi}(t, a_t) - \pi^*(t, a_t) &\geq \langle \hat{s}_t, a_t^{\text{ag}} \rangle - \langle \hat{s}_t, a_t \rangle + 2/T - \langle s^*, a_t^{\text{ag}} \rangle + \langle s^*, a_t \rangle \\ &= \langle \hat{s}_t - s^*, a_t^{\text{ag}} - a_t \rangle + 2/T \\ &> -\text{diam}\left(\mathcal{S}_t, \frac{a_t^{\text{ag}} - a_t}{\|a_t^{\text{ag}} - a_t\|}\right) \underbrace{\|a_t^{\text{ag}} - a_t\|}_{\leq 2} + 2 \max_{a_t^1 \neq a_t^2 \in \mathcal{A}_t} \text{diam}\left(\mathcal{S}_t, \frac{a_t^1 - a_t^2}{\|a_t^1 - a_t^2\|}\right) \geq 0. \end{aligned} \quad (26)$$

Similarly, since  $\langle s^*, a_t^{\text{ag}} \rangle \geq \langle s^*, \hat{a}_t^{\text{ag}} \rangle$ , under  $\mathcal{E}_t$ , we have

$$\begin{aligned} \hat{\pi}(t, a_t) - \pi^*(t, a_t) &= \langle \hat{s}_t, \hat{a}_t^{\text{ag}} - a_t \rangle + 2/T - \langle s^*, a_t^{\text{ag}} - a_t \rangle \\ &\leq \langle \hat{s}_t, \hat{a}_t^{\text{ag}} - a_t \rangle - \langle s^*, a_t^{\text{ag}} - a_t \rangle + 2/T = \langle \hat{s}_t - s^*, \hat{a}_t^{\text{ag}} - a_t \rangle + 2/T \\ &\leq \|\hat{a}_t^{\text{ag}} - a_t\| \cdot \text{diam}\left(\mathcal{S}_t, \frac{\hat{a}_t^{\text{ag}} - a_t}{\|\hat{a}_t^{\text{ag}} - a_t\|}\right) + 2/T \leq 4/T. \end{aligned} \quad (27)$$

Combining (26) and (27) with the definition of  $\hat{\kappa}$ , we obtain the result.  $\square$

**Lemma 9.** Let  $t \in [T]$  such that  $\mathcal{E}_t$  does not hold. Then  $\text{Vol}(\Pi_{V_t^\perp} \mathcal{S}_t) \geq \bar{\delta}^{2d}/d^{2d}$ , where  $\bar{\delta} = 1/16T^2d(d+1)^2$  and  $V_t$  is defined in (31).

*Proof.* Since  $\mathcal{E}_t$  does not hold, there exists a direction  $u \in \mathcal{A}_t \subseteq B(0,1)$  such that  $\text{diam}(\mathcal{S}_t, u) \geq 1/T \geq \bar{\delta}$ . For any  $u \in V_t^\perp$ ,  $\text{diam}(\mathcal{S}_t, u) \geq \bar{\delta}$ . Lemma 6.3 of Lobel et al. (2018, Section 6) guarantees that  $\Pi_{V_t^\perp}(\mathcal{S}_t)$  contains a  $k$ -dimensional ball of radius  $\bar{\delta}/k$ , where  $k := \dim(V_t^\perp)$ . Therefore,  $\text{Vol}(\Pi_{V_t^\perp}(\mathcal{S}_t)) \geq \bar{\delta}^k \pi^{k/2}/k^k \Gamma(k/2 + 1)$ , where  $\Gamma$  stands for Euler's Gamma function (Smith & Vamanamurthy, 1989). Therefore, we have by definition of Euler's Gamma function:  $\text{Vol}(\Pi_{V_t^\perp}(\mathcal{S}_t)) \geq \bar{\delta}^k \pi^{k/2}/k^k k^k \geq \bar{\delta}^{2d}/d^{2d}$ .  $\square$

**Lemma 5.** Consider  $\mathcal{E}_t$  defined by (14) with  $(\mathcal{S}_t)_{t \in [T]}$  defined by Contextual IPA. Then it holds almost surely that

$$\sum_{t \geq 1} \mathbb{1}_{\mathcal{E}_t^c} \leq 192 \cdot d \log(dT),$$

where  $\mathcal{E}_t$  is defined by (14).

*Proof of Lemma 5.* This proof follows the same line as the proof of the main theorem of Lobel et al. (2018, Section 7). Let  $t \in [T]$  such that  $\mathcal{E}_t$  does not hold. We define  $w_t$  as

$$w_t := \operatorname{argmax} \left\{ \text{diam} \left( \mathcal{S}_t, \frac{a_t^1 - a_t^2}{\|a_t^1 - a_t^2\|} \right) : \frac{a_t^1 - a_t^2}{\|a_t^1 - a_t^2\|} \text{ such that } a_t^1 \neq a_t^2 \in \mathcal{A}_t \right\}, \quad (28)$$

where  $\mathcal{S}_t$  is defined in (30). Our goal is to bound the number of steps for which the diameter of  $\mathcal{S}_t$  in the direction  $w_t$  is strictly superior to  $1/T$ . Define

$$b_t := \mathbb{1}\{\text{diam}(\text{Cyl}(\mathcal{S}_t, V_t^\perp), w_t) \geq 1/T\},$$

where  $V_t$  is defined in (31). Let  $\Pi_E$  denote the orthogonal projection onto the subspace  $E \subset \mathbb{R}^d$  and  $\text{Vol}(\Pi_E \mathcal{S}) = \mu_E(\Pi_E \mathcal{S})$ , where  $\mu_E$  is the Lebesgue measure of  $\Pi_E \mathcal{S}$ , well-defined for  $\Pi_E \mathcal{S}$  being a convex body. Setting  $\bar{\delta} = T^{-2}/16d(d+1)^2$ , we can apply the projected Grünbaum lemma of Lobel et al. (2018, Lemma 7.1) to obtain that

$$\text{Vol}(\Pi_{V_t^\perp} \mathcal{S}_{t+1}) \leq (1 - e^{-2})^{b_t} \text{Vol}(\Pi_{V_t^\perp} \mathcal{S}_t).$$

By definition of  $V_t$  in (31), we have for any  $u \in V_t$ ,  $\text{diam}(\mathcal{S}_t, u) \geq \bar{\delta}$ . Therefore, by definition of  $\mathcal{S}_{t+1}$  in (30), the directional Grünbaum Theorem (Lobel et al., 2018, Theorem 5.3) guarantees that we have:  $\text{diam}(\mathcal{S}_{t+1}, u) \geq \bar{\delta}/(d+1)$ . Note that for  $\mathcal{S}_t$  being a convex body, Lemma 11 ensures that  $\mathcal{S}_{t+1}$  remains a convex body.

If  $J_{t+1} = 0$ , where  $J_t$  is defined in (32), then  $V_{t+1}^\perp = V_t^\perp$ , and we have  $\text{Vol}(\Pi_{V_{t+1}^\perp} \mathcal{S}_{t+1}) = \text{Vol}(\Pi_{V_t^\perp} \mathcal{S}_{t+1})$ .

Otherwise, let  $i \in [J_{t+1} - 1]$  and  $v \in V_t^{(i), \perp}$  such that  $V_t^{(i+1)} = V_t^{(i)} \cup \{v\}$ . We have  $V_t^{(i), \perp} \cap \text{span}(v)^\perp = \{x \in V_t^{(i), \perp} : v^\top x = 0\} \subseteq V_t^{(i+1), \perp} \subseteq V_t^{(i), \perp}$  and  $\dim(V_t^{(i+1), \perp}) = \dim(V_t^{(i), \perp}) - 1$ . Then, applying the Cylindrification Lemma from Lobel et al. (2018, Lemma 6.1) we obtain

$$\text{Vol}(\Pi_{V_t^{(i+1), \perp}} \mathcal{S}_{t+1}) \leq \frac{d(d+1)^2}{\bar{\delta}} \text{Vol}(\Pi_{V_t^{(i), \perp}} \mathcal{S}_{t+1}).$$

If  $J_{t+1} = r$ , the volume can blow up by at most  $(d(d+1)^2/\bar{\delta})^r$ . In particular, since the initial volume is bounded by  $\text{Vol}(B(0,1)) \leq 8\pi^2/15 \leq 6$  (Smith & Vamanamurthy, 1989), then by Lemma 9, we obtain:  $\bar{\delta}^{2d}/d^{2d} \leq \text{Vol}(\Pi_{L_t} \mathcal{S}_t) \leq 6 \cdot (d(d+1)^2/\bar{\delta})^d \cdot (1 - 1/e^2)^{\sum_{t=1}^T b_t}$ . Therefore, applying the logarithm function, we obtain

$$\sum_{t=1}^T b_t \leq -1/\log(1 - e^{-2})(\log 6 + 2d \log(16) + 5d \log(d) + 6d \log(d+1) + 4d \log(T)),$$

giving, since  $-1/\log(1 - e^{-2}) < 16$ :  $\sum_{t=1}^T \mathbb{1}\{\text{diam}(\text{Cyl}(\mathcal{S}_t, L_t), w_t) \geq 1/T\} \leq 192d \log(dT)$ . Therefore, lemma 12 ensures that  $\text{diam}(\mathcal{S}_t, w_t) \leq \text{diam}(\text{Cyl}(\mathcal{S}_t, V_t^\perp), w_t)$  and we get

$$\sum_{t=1}^T \mathbb{1}\left\{ \max_{s \in \mathcal{S}_t} |\langle s, w_t \rangle| \geq 1/T \right\} \leq 192 d \log(dT),$$

which concludes the proof by definition of  $w_t$  in (28) and  $\mathcal{E}_t$  in (14).  $\square$

**Lemma 6.** Consider  $(\mathcal{E}_t)_{t \in [T]}$  defined by (14) with  $(\mathcal{S}_t)_{t \in [T]}$  defined by Contextual IPA. Let  $I_T$  and  $(\varepsilon_t^{\text{corrupt}})_{t \in [T]}$  as defined in (15) and (16) and  $t \in [T]$  such that  $\mathcal{E}_t$  is true. Then  $|\varepsilon_t^{\text{corrupt}}| \leq 4/T$  and  $\sum_{t \in I_T} |\varepsilon_t^{\text{corrupt}}| := C_{\text{corrupt}} \leq 4$ .

*Proof of Lemma 6.* Consider  $t \in I_T$ ,  $a \in \mathcal{A}_t$ : by (15),  $\mathcal{E}_t$  is true. Consider  $\varepsilon_t^{\text{corrupt}}$  as defined in (16)

Using Lemma 4 gives  $|\varepsilon_t^{\text{corrupt}}| \leq 4/T$ , and summing over all the iterations  $t \in I_T$ , since  $|I_T| \leq T$ , we have  $\sum_{t \in I_T} |\varepsilon_t^{\text{corrupt}}| \leq 4 \cdot 1/T \cdot |I_T| \leq 4$ .  $\square$

**Lemma 10.** For any  $t \in [T]$ ,  $\varepsilon > 0$  and action  $a \in \mathcal{A}_t$ , set  $\pi^{*,\varepsilon}(t, a) := \max_{a_t^{\text{ag}} \in \mathcal{A}_t} \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle - \langle \mathbf{s}^*, a \rangle + \varepsilon$ , and define the incentive function  $\kappa_a^{*,\varepsilon}(t, a') = \mathbb{1}_a(a') \pi^{*,\varepsilon}(t, a)$  for any  $a' \in \mathcal{A}_t$ . Then  $A_t = a$ , where  $A_t$  is defined by (8).

*Proof of Lemma 10.* Note that for any  $a' \in \mathcal{A}_t$ ,  $a' \neq a$ ,  $\langle \mathbf{s}^*, a' \rangle + \kappa_a^{*,\varepsilon}(t, a') < \langle \mathbf{s}^*, a \rangle + \kappa_a^{*,\varepsilon}(t, a)$  and, as a result,  $A_t = a$ .  $\square$

## D.2 Proof of Theorem 2

*Proof of Theorem 2.* Denote for any  $t \in [T]$ ,  $a_t^{\text{pr}} := \operatorname{argmax}_{a \in \mathcal{A}_t} \{\langle \theta^*, a \rangle - \pi^*(t, a)\} = \operatorname{argmax}_{a \in \mathcal{A}_t} \langle \theta^* + \mathbf{s}^*, a \rangle$ ,  $a_t^{\text{ag}} := \operatorname{argmax}_{a \in \mathcal{A}_t} \langle \mathbf{s}^*, a \rangle$  and  $\hat{a}_t^{\text{ag}} = \operatorname{argmax}_{a \in \mathcal{A}_t} \langle \hat{\mathbf{s}}_t, a \rangle$ . Note that if  $\mathcal{E}_t$  holds, for any  $a_t^1, a_t^2 \in \mathcal{A}_t$ ,  $a_t^1 \neq a_t^2$ , **H2** gives that if

$$\operatorname{diam}\left(\mathcal{S}_t, \frac{a_t^1 - a_t^2}{\|a_t^1 - a_t^2\|}\right) < \frac{1}{T}, \text{ then } \max_{s \in \mathcal{S}_t} \langle s, a_t^1 - a_t^2 \rangle = \max_{s \in \mathcal{S}_t} \left\langle s, \frac{a_t^1 - a_t^2}{\|a_t^1 - a_t^2\|} \right\rangle \underbrace{\|a_t^1 - a_t^2\|}_{\leq 2} < 2 \cdot \frac{1}{T},$$

and therefore  $\operatorname{diam}(\mathcal{S}_t, a_t^1 - a_t^2) < 2/T$ .

If  $\mathcal{E}_t$  does not hold, the incentive function proposed in Algorithm 5 is given by  $\kappa(t, a) = 3 \cdot \mathbb{1}_{a_t^1}(a) + (3 + \langle \hat{\mathbf{s}}_t, a_t^1 - a_t^2 \rangle) \cdot \mathbb{1}_{a_t^2}(a)$ . Therefore:  $\kappa(t, a_t^1) = 3$ ,  $\kappa(t, a_t^2) = 3 + \langle \hat{\mathbf{s}}_t, a_t^1 - a_t^2 \rangle \leq 5$ . When  $\mathcal{E}_t$  holds, the incentive function proposed in Contextual IPA is defined as  $\kappa(t, a) = \mathbb{1}_{a_t^{\text{Rec}}}(a) \hat{\pi}(t, a_t^{\text{Rec}}) \leq 2$ .

For any  $t \in [T]$ , we define the instantaneous regret  $\operatorname{reg}_t$  at  $t$  as

$$\operatorname{reg}_t := \mu_t^* - (X_{A_t}(t) - \kappa(t, A_t)),$$

where  $\mu_t^*$  is defined in (10) and we decompose the regret into two terms, making use of the Cauchy-Schwarz inequality as well as **H2** to obtain

$$\begin{aligned} \mathfrak{R}(T) &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \operatorname{reg}_t + \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t^c\}} \operatorname{reg}_t \right] \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \operatorname{reg}_t \right] + \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t^c\}} (\max_{a \in \mathcal{A}_t} \langle \theta^* + \mathbf{s}^*, a \rangle - \max_{a' \in \mathcal{A}_t} \langle \mathbf{s}^*, a' \rangle - \langle \theta^*, A_t \rangle + \kappa(t, A_t)) \right] \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \operatorname{reg}_t \right] + \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t^c\}} (\underbrace{\max_{a \in \mathcal{A}_t} \langle \theta^*, a \rangle}_{\leq 1} + \underbrace{\langle \mathbf{s}^*, a - a_t^{\text{ag}} \rangle}_{\leq 0} - \underbrace{\langle \theta^*, A_t \rangle}_{\leq 1} + \underbrace{\kappa(t, A_t)}_{\leq 5}) \right] \\ &\leq \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \operatorname{reg}_t \right] + 7 \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t^c\}} \right]. \end{aligned}$$

Using Lemma 5, we can bound the second term

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t^c\}} \right] \leq 192 d \log(dT).$$

Now we bound the first term. Working on steps  $t$  such that  $\mathcal{E}_t$  is true with incentive function  $\hat{\kappa}_{a_t^{\text{Rec}}}(t, \cdot) = \mathbb{1}_{a_t^{\text{Rec}}}(\cdot) \hat{\pi}(t, a_t^{\text{Rec}})$ , Lemma 4 guarantees that  $A_t = a_t^{\text{Rec}}$ , and we have

$$\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \operatorname{reg}_t \right]$$

$$\begin{aligned}
 &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( \max_{a \in \mathcal{A}_t} \{\langle \theta^* + \mathbf{s}^*, a \rangle - \max_{a' \in \mathcal{A}_t} \langle \mathbf{s}^*, a' \rangle\} - \left\{ \langle \theta^*, a_t^{\text{Rec}} \rangle - \left( \langle \hat{\mathbf{s}}_t, \hat{a}_t^{\text{ag}} \rangle - \langle \hat{\mathbf{s}}_t, a_t^{\text{Rec}} \rangle + \frac{2}{T} \right) \right\} \right) \right] \\
 &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( \langle \theta^*, a_t^{\text{pr}} \rangle + \langle \mathbf{s}^*, a_t^{\text{pr}} \rangle - \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle - \langle \theta^*, a_t^{\text{Rec}} \rangle + \langle \hat{\mathbf{s}}_t, \hat{a}_t^{\text{ag}} \rangle - \langle \hat{\mathbf{s}}_t, a_t^{\text{Rec}} \rangle \right) + \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \frac{2}{T} \right] \\
 &\leq \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( \langle \theta^*, a_t^{\text{pr}} \rangle + \langle \mathbf{s}^*, a_t^{\text{pr}} \rangle - \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle - \langle \theta^*, a_t^{\text{Rec}} \rangle \right) \right] + 2T \frac{1}{T} \\
 &\quad + \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( -\langle \mathbf{s}^*, a_t^{\text{Rec}} \rangle + \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle + \langle \hat{\mathbf{s}}_t, \hat{a}_t^{\text{ag}} \rangle - \langle \hat{\mathbf{s}}_t, a_t^{\text{Rec}} \rangle + \langle \mathbf{s}^*, a_t^{\text{Rec}} \rangle - \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle \right) \right] \\
 &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( \langle \theta^*, a_t^{\text{pr}} \rangle + \langle \mathbf{s}^*, a_t^{\text{pr}} \rangle - \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle - \langle \theta^*, a_t^{\text{Rec}} \rangle - \langle \mathbf{s}^*, a_t^{\text{Rec}} \rangle + \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle \right) \right] \\
 &\quad + \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( \langle \hat{\mathbf{s}}_t, \hat{a}_t^{\text{ag}} \rangle + \langle \mathbf{s}^* - \hat{\mathbf{s}}_t, a_t^{\text{Rec}} \rangle - \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle \right) \right] + 2.
 \end{aligned}$$

Since  $\langle \mathbf{s}^*, a_t^{\text{ag}} \rangle \geq \langle \mathbf{s}^*, \hat{a}_t^{\text{ag}} \rangle$ , we have  $-\langle \mathbf{s}^*, a_t^{\text{ag}} \rangle \leq -\langle \mathbf{s}^*, \hat{a}_t^{\text{ag}} \rangle$ . Therefore

$$\begin{aligned}
 &\mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( \langle \hat{\mathbf{s}}_t, \hat{a}_t^{\text{ag}} \rangle + \langle \mathbf{s}^* - \hat{\mathbf{s}}_t, a_t^{\text{Rec}} \rangle - \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle \right) \right] \\
 &\leq \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( \langle \hat{\mathbf{s}}_t, \hat{a}_t^{\text{ag}} \rangle + \langle \mathbf{s}^* - \hat{\mathbf{s}}_t, a_t^{\text{Rec}} \rangle - \langle \mathbf{s}^*, \hat{a}_t^{\text{ag}} \rangle \right) \right] \\
 &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( \langle \mathbf{s}^* - \hat{\mathbf{s}}_t, a_t^{\text{Rec}} \rangle - \langle \mathbf{s}^* - \hat{\mathbf{s}}_t, \hat{a}_t^{\text{ag}} \rangle \right) \right] \\
 &= \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( \langle \mathbf{s}^* - \hat{\mathbf{s}}_t, a_t^{\text{Rec}} - \hat{a}_t^{\text{ag}} \rangle \right) \right] \\
 &\leq \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \mathbb{1}_{\{a_t^{\text{Rec}} \neq \hat{a}_t^{\text{ag}}\}} \underbrace{\|a_t^{\text{Rec}} - \hat{a}_t^{\text{ag}}\|}_{\leq 2} \text{diam} \left( \mathcal{S}_t, \frac{a_t^{\text{Rec}} - \hat{a}_t^{\text{ag}}}{\|a_t^{\text{Rec}} - \hat{a}_t^{\text{ag}}\|} \right) < 2T \cdot \frac{1}{T} = 2,
 \end{aligned}$$

and plugging this inequality gives

$$\begin{aligned}
 \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \text{reg}_t \right] &\leq \mathbb{E} \left[ \sum_{t=1}^T \mathbb{1}_{\{\mathcal{E}_t\}} \left( \langle \theta^* + \mathbf{s}^*, a_t^{\text{pr}} \rangle - \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle - \langle \theta^* + \mathbf{s}^*, a_t^{\text{Rec}} \rangle + \langle \mathbf{s}^*, a_t^{\text{ag}} \rangle \right) \right] + 4 \\
 &= \mathbb{E} \left[ \sum_{t \in I_T} \max_{a \in \mathcal{A}_t} \{\langle \theta^* + \mathbf{s}^*, a \rangle - \max_{a' \in \mathcal{A}_t} \langle \mathbf{s}^*, a' \rangle\} \right] - \mathbb{E} \left[ \sum_{t \in I_T} \langle \theta^* + \mathbf{s}^*, a_t^{\text{Rec}} \rangle - \max_{a \in \mathcal{A}_t} \langle \mathbf{s}^*, a \rangle \right] + 4 \\
 &= R_{\text{CtxAlg}}^{\text{corrupt}}(I_T, \varepsilon_{I_T}^{\text{corrupt}}) + 4,
 \end{aligned}$$

where  $R_{\text{Contextual IPA}}^{\text{corrupt}}$  is defined in (17) with  $\pi^*(t, a) = \max_{a' \in \mathcal{A}_t} \langle \mathbf{s}^*, a' \rangle - \langle \mathbf{s}, a \rangle$ . Plugging all the terms together gives the following upper-bound for the regret:

$$\mathfrak{R}(T) \leq 1344 d \log(dT) + 4 + R_{\text{CtxAlg}}^{\text{corrupt}}(I_T, \varepsilon_{I_T}^{\text{corrupt}}).$$

□

*Proof of Corollary 2.* Here, Lemma 6 guarantees that  $C_{\text{corrupt}} = 4$ . With the setup of He et al. (2022), we take  $L = 1, S = 1, R = 1$ , and  $\alpha = \sqrt{d}/4$ .

Choose  $\lambda = 1$  as a regularization parameter and  $\delta = 1/T$  as a confidence level. Using CW – OFUL proposed in He et al. (2022) as a subroutine robust to corruption in the stochastic linear case, since our subroutine is fed with the same reward  $r'$  as in their model while  $R_{\text{CtxAlg}}^{\text{corrupt}}(I_T, \varepsilon_{I_T}^{\text{corrupt}})$  is defined compared to the true reward  $r$ , we can use the result provided in He et al. (2022, Theorem 4.2) to bound  $R_{\text{CtxAlg}}^{\text{corrupt}}(I_T, \varepsilon_{I_T}^{\text{corrupt}})$  with probability  $1 - 1/T$  and use the fact that our instantaneous regret is always bounded by 7 as it is shown in the proof of Theorem 2 to get in expectation for some universal constant  $B > 0$

$$\begin{aligned} R_{\text{CtxAlg}}^{\text{corrupt}}(I_T, \varepsilon_{I_T}^{\text{corrupt}}) &\leq (1 - \delta)B \left( 2d\sqrt{T} \log\left(\frac{1+T}{\delta}\right) + \sqrt{d\lambda T} \sqrt{\log(1+T)} + 4d\sqrt{\log\left(\frac{1+T}{\delta}\right)^3} \right) + 7T\delta \\ &\leq \left(1 - \frac{1}{T}\right)B \left( 2d\sqrt{T} \log\left(\frac{1+T}{\frac{1}{T}}\right) + \sqrt{dT} \sqrt{\log\left(\frac{1+T}{\frac{1}{T}}\right)} + 4d\sqrt{\log\left(\frac{1+T}{\frac{1}{T}}\right)^3} \right) + 7T\frac{1}{T} \\ &\leq \left(1 - \frac{1}{T}\right)B \left( 2d\sqrt{T} \log(T+T^2) + \sqrt{dT} \sqrt{\log(T+T^2)} + 4d(\log(T+T^2))^{\frac{3}{2}} \right) + 7, \end{aligned}$$

therefore, since  $C_{\text{corrupt}} \leq 4$ , there exists a constant  $C_{\text{CtxAlg}}$  such that

$$R_{\text{CtxAlg}}^{\text{corrupt}}(I_T, \varepsilon_{I_T}^{\text{corrupt}}) \leq 7 + C_{\text{CtxAlg}} d\sqrt{T} \log T.$$

Finally plugging this term in the bound from Theorem 2 and integrating the 3 factor in constant  $B$  gives the result

$$\mathfrak{R}(T) \leq 11 + 1344 d \log(dT) + C_{\text{CtxAlg}} d\sqrt{T} \log T.$$

□

## E Algorithms

### E.1 UCB subroutine

We present the Binary search subroutine and the UCB algorithm that we use as a subroutine in IPA as formulated in Lattimore & Szepesvári (2020, Algorithm 3).

---

#### Algorithm 3 Binary Search Subroutine

---

- 1: **Input:** action  $a, N_T$
  - 2: **Initialize:**  $\underline{\pi}_a(0), \bar{\pi}_a(0) = 0, 1$
  - 3: **for**  $t = 1, \dots, N_T$  **do**
  - 4:    $\pi_a^{\text{mid}}(t-1) = \frac{\bar{\pi}_a(t-1) + \underline{\pi}_a(t-1)}{2}$
  - 5:   Propose incentive  $\pi_a^{\text{mid}}(t-1)$  on arm  $a$
  - 6:   **If**  $A_{t-1} = a$  **then**  $\bar{\pi}_a(t) = \pi_a^{\text{mid}}(t-1)$  and  $\underline{\pi}_a(t) = \underline{\pi}_a(t-1)$  **else**  $\underline{\pi}_a(t) = \pi_a^{\text{mid}}(t)$  and  $\bar{\pi}_a(t) = \bar{\pi}_a(t-1)$
  - 7: **end for**
  - 8: **Return**  $\underline{\pi}_a(N_T), \bar{\pi}_a(N_T)$
- 

### E.2 Projected volume algorithm

We present the Projected volume algorithm from Lobel et al. (2018) that we use as a subroutine in Contextual IPA. For any horizon  $T$  and  $0 < \bar{\delta} < T^{-2}/(16d(d+1)^2)$ , this algorithm defines recursively a sequence  $(\mathcal{S}_t, \mathbf{V}_t)_{t \in [T]}$  such that  $(\mathcal{S}_t)_{t \in [T]}$  is a sequence of decreasing subsets (for the inclusion) of  $B(0, 1)$  including  $\mathbf{s}^*$  and  $(\mathbf{V}_t)_{t \in [T]}$  is a sequence of increasing sets of  $\mathbb{R}^d$  containing orthogonal directions  $\{v_i\}_{i \in [n]}$  along which the principal has a good knowledge of  $\langle \mathbf{s}^*, v_i \rangle$ .

The main ingredient in Lobel et al. (2018) allowing low regret is *cylindrification*. Given a compact convex set  $\mathcal{S} \subseteq \mathbb{R}^d$ ,  $\mathbf{V} = \{v_1, \dots, v_n\}$ , and  $\Pi_{\mathbf{V}^\perp}(\mathcal{S})$  being the orthogonal projection of  $\mathcal{S}$  onto  $\mathbf{V}^\perp$ , we define the cylindrification of  $\mathcal{S}$  on  $\mathbf{V}$  as

$$\text{Cyl}(\mathcal{S}, \mathbf{V}) := \Pi_L(\mathcal{S}) + \Pi_{\text{span}(v_1)}(\mathcal{S}) + \dots + \Pi_{\text{span}(v_n)}(\mathcal{S})$$

**Algorithm 4** UCB Subroutine

- 1: **Input:** Set of arms  $K$ , horizon  $T$
- 2: **Initialize:** For any arm  $a \in [K]$ , set  $\hat{\mu}_a := 0$ ,  $T_a := 0$
- 3: **for**  $1 \leq t \leq K$ : **do**
- 4:   Pull arm  $a = t$
- 5:   Update  $\hat{\mu}_a = X_a(t)$ ,  $T_a(t) = 1$
- 6: **end for**
- 7: **for**  $t \geq K + 1$  **do**
- 8:   Pull arm  $a_{\max} \in \operatorname{argmax}_{a \in [K]} \left\{ \hat{\mu}_a(t-1) + 2\sqrt{\frac{\log T}{T_a(t-1)}} \right\}$
- 9:   Update  $T_a(t) = T_a(t-1) + 1$ ,  $\hat{\mu}_a(t) = \frac{1}{T_a(t)}(T_a(t-1)\hat{\mu}_a(t-1) + X_a(t))$
- 10: **end for**

$$= \left\{ x + \sum_{i=1}^n y_i v_i : x \in \Pi_L(\mathcal{S}), \min_{s \in \mathcal{S}} \langle s, v_i \rangle \leq y_i \leq \max_{s \in \mathcal{S}} \langle s, v_i \rangle \right\}.$$

At iteration  $t$ , given  $(\mathcal{S}_t, \mathbf{V}_t)$ , we define the estimate  $\hat{s}_t$  of  $s^*$  as the centroid of  $\operatorname{Cyl}(\mathcal{S}_t, \mathbf{V}_t)$ :

$$\hat{s}_t := \frac{1}{\operatorname{Vol}(\operatorname{Cyl}(\mathcal{S}_t, \mathbf{V}_t))} \int_{\operatorname{Cyl}(\mathcal{S}_t, \mathbf{V}_t)} x dx, \quad (29)$$

Note that the iterative construction of  $\mathcal{S}_t$  described below together with Lemma 11 guarantees that  $\mathcal{S}_t$  is always a convex body, making  $\operatorname{Vol}(\mathcal{S}_t)$  well-defined as the Lebesgue measure of  $\mathcal{S}_t$  in dimension  $d$  and  $\operatorname{Vol}(\mathcal{S}_t) > 0$ . Combined with Lemma 12, it guarantees that  $\operatorname{Vol}(\operatorname{Cyl}(\mathcal{S}_t, \mathbf{V}_t))$  is well-defined, and  $\operatorname{Vol}(\operatorname{Cyl}(\mathcal{S}_t, \mathbf{V}_t)) > 0$ . Therefore, (29) is well-defined.

At iteration  $t$ , given  $(\mathcal{S}_t, \mathbf{V}_t)$  and two actions  $a_t^1, a_t^2 \in \mathcal{A}_t$ , recall that we defined  $w_t$  as

$$w_t := \operatorname{argmax} \left\{ \operatorname{diam} \left( \mathcal{S}_t, \frac{a_t^1 - a_t^2}{\|a_t^1 - a_t^2\|} \right) : \frac{a_t^1 - a_t^2}{\|a_t^1 - a_t^2\|} \text{ such that } a_t^1 \neq a_t^2 \in \mathcal{A}_t \right\}.$$

Then, the principal offers the incentive function  $\kappa(t, a) = 5 \cdot \mathbb{1}_{a_t^1}(a) + (5 + \langle \hat{s}_t, a_t^1 - a_t^2 \rangle) \cdot \mathbb{1}_{a_t^2}(a)$ . Recall that  $\kappa(t, a)$  is always bounded by 5 and that  $a_t^1, a_t^2$  are chosen such that  $\langle \hat{s}_t, a_t^1 - a_t^2 \rangle \geq 0$ , ensuring that we either have  $A_t = a_t^1$  or  $A_t = a_t^2$ .

If  $A_t = a_t^1$ , it means that  $\langle s^*, w_t \rangle \geq 0$  and we update  $\mathcal{S}_{t+1} = \mathcal{S}_t \cap \{s \mid \langle s, w_t \rangle \geq x_t\}$ . Otherwise, if  $A_t = a_t^2$ , it means that  $\langle s^*, w_t \rangle \leq 0$  and we update  $\mathcal{S}_{t+1} = \mathcal{S}_t \cap \{s \mid \langle s, w_t \rangle \leq x_t\}$ . This defines the subset  $\mathcal{S}_{t+1}$  such that  $s^* \in \mathcal{S}_{t+1}$ :

$$\mathcal{S}_{t+1} = \mathcal{S}_t \cap (\mathbb{1}_{a_t^1}(A_t) \{s \mid \langle s, w_t \rangle \geq x_t\} + \mathbb{1}_{a_t^2}(A_t) \{s \mid \langle s, w_t \rangle \leq x_t\}) \quad (30)$$

Regarding the subspace  $\mathbf{V}_{t+1}$ , we consider  $\mathbf{V}_{t+1} = \mathbf{V}_{t+1}^{J_{t+1}}$  where  $(\mathbf{V}_{t+1}^i)_{i \in \mathbb{N}}$  is defined as

$$\begin{aligned} \mathbf{V}_{t+1}^{(0)} &:= \mathbf{V}_t \\ \mathbf{V}_{t+1}^{(i+1)} &:= \begin{cases} \mathbf{V}_{t+1}^{(i)} \cup \{v\} & \text{if } \exists v \in (\mathbf{V}_{t+1}^{(i)})^\perp : \operatorname{diam}(\mathcal{S}_{t+1}, v) \leq \bar{\delta} \\ \mathbf{V}_{t+1}^{(i)} & \text{otherwise} \end{cases} \end{aligned} \quad (31)$$

$$J_{t+1} := \min\{i : \text{it does not exist } v \in (\mathbf{V}_{t+1}^{(i)})^\perp \text{ such that } \operatorname{diam}(\mathcal{S}_{t+1}, v) \leq \bar{\delta}\}, \quad (32)$$

which exists since if it does not exist  $i$  such that it does not exist  $v \in (\mathbf{V}_{t+1}^{(i)})^\perp$  such that  $\operatorname{diam}(\mathcal{S}_{t+1}, v) \geq \bar{\delta}$ , we would have  $\dim(\operatorname{span}(\mathbf{V}_t^{(i+1)})) = \dim(\operatorname{span}(\mathbf{V}_t^{(i)})) + 1$  for any  $i \in \mathbb{N}$ , which would imply a contradiction.

In what follows, we provide technical results ensuring that  $\mathcal{S}_t$  is a convex body for any  $t \in [T]$  and  $\mathcal{S}_t \subseteq \operatorname{Cyl}(\mathcal{S}_t, \mathbf{V}_t)$ , which implies that  $\operatorname{Cyl}(\mathcal{S}_t, \mathbf{V}_t)$  has non-empty interior.

**Lemma 11.** *Let  $\mathcal{S}$  be a convex body in  $\mathbb{R}^d$  and  $s$  be a point in the interior of  $\mathcal{S}$ :  $s \in \mathring{\mathcal{S}}$ . Let  $\mathbf{G}$  be the half-space defined by  $\mathbf{G} := \{x \in \mathbb{R}^d : \langle h^*, x \rangle \geq 0\}$ , for some  $h^* \in \mathbb{R}^d$ . Suppose that  $s \in \mathbf{G}$ . Then  $\mathcal{S} \cap \mathbf{G}$  is a convex body.*

*Proof.* Note that we only need to show that  $\mathcal{S} \cap \mathcal{G}$  has non-empty interior since the intersection of two compact convex sets is compact and convex.

The only case that we consider is  $s \in \mathcal{H}$  where  $\mathcal{H}$  is the hyperplane defined by  $\mathcal{H} := \{x \in \mathbb{R}^d : \langle h^*, x \rangle = 0\}$ . The other case simply follows from the fact that the intersection of two open sets is also open.

Since  $\mathcal{S}$  is a convex body and  $s \in \overset{\circ}{\mathcal{S}}$ , there exists a ball  $B(s, r)$  centered in  $s$  with  $r > 0$ , such that  $B(s, r) \subseteq \overset{\circ}{\mathcal{S}}$ . For any  $x \in B(s, r) \cap \mathcal{G}$  is equivalent to  $\langle h^*, x - s \rangle \geq 0$  since  $\langle h^*, s \rangle = 0$  and  $\|x - s\| \leq r$ .

Now we define  $y_0 = s + rh^*/(2\|h^*\|)$  and consider the ball  $B(y_0, r/2)$ . For any  $y \in B(y_0, r/2)$ , we can write  $y = y_0 + \tilde{y}$  with  $\|\tilde{y}\| \leq r/2$ . Using  $\langle h^*, s \rangle = 0$ , we have  $\|y - s\| \leq \|y_0 - s\| + \|\tilde{y}\| \leq r$  and  $\langle h^*, y \rangle = \langle h^*, y_0 \rangle + \langle h^*, \tilde{y} \rangle$  with  $\langle h^*, y_0 \rangle = r/2$  and  $\langle h^*, \tilde{y} \rangle \geq -\|h^*\| \cdot \|\tilde{y}\| \geq -1 \cdot r/2 = -r/2$ . Therefore  $\langle h^*, y \rangle \geq 0$  and we obtain  $y \in B(s, r) \cap \mathcal{G}$ , which gives  $B(y_0, r/2) \subseteq B(s, r) \cap \mathcal{G}$ . □

**Lemma 12.** *Given a convex body  $\mathcal{S} \subseteq \mathbb{R}^d$  and a set of orthonormal vectors  $\mathcal{V} = \{v_1, \dots, v_n\} \subseteq \mathbb{R}^d$ , let  $\Pi_{\mathcal{V}^\perp}(\mathcal{S})$  be the projection of  $\mathcal{S}$  on the subspace  $\mathcal{V}^\perp = \{x \in \mathbb{R}^d \mid \langle x, v_i \rangle = 0\}$ . Define the cylindrification of  $\mathcal{S}$  onto  $\mathcal{V}$  as*

$$\begin{aligned} \text{Cyl}(\mathcal{S}, \mathcal{V}) &:= \Pi_{\mathcal{V}^\perp}(\mathcal{S}) + \Pi_{\text{span}(v_1)}(\mathcal{S}) + \dots + \Pi_{\text{span}(v_n)}(\mathcal{S}) \\ &= \left\{ x + \sum_{i=1}^n y_i v_i \mid x \in \Pi_{\mathcal{V}^\perp}(\mathcal{S}), \min_{s \in \mathcal{S}} \langle s, v_i \rangle \leq y_i \leq \max_{s \in \mathcal{S}} \langle s, v_i \rangle \right\}. \end{aligned}$$

Then it holds that  $\mathcal{S} \subseteq \text{Cyl}(\mathcal{S}, \mathcal{V})$ .

*Proof.* Define  $\Pi_{\mathcal{V}^\perp}$  as the orthogonal projector on  $\mathcal{V}^\perp$ . Then we have for any  $s \in \mathcal{S}$

$$s = \Pi_{\mathcal{V}^\perp} s + (I - \Pi_{\mathcal{V}^\perp})s,$$

where  $(I - \Pi_{\mathcal{V}^\perp})$  is an orthogonal projector on the space  $\text{span}(\{v_1, \dots, v_n\})$ , thus  $(I - \Pi_{\mathcal{V}^\perp})s = \sum_{i=1}^n y_i v_i$  for  $y_i = \langle s, v_i \rangle$ . This decomposition allows us to conclude. □

---

### Algorithm 5 Projected Volume

---

- 1: **Input:**  $T, \bar{\delta}, \mathcal{S}_t$  such that  $\text{diam}(\mathcal{S}_t) \geq \bar{\delta}, \mathcal{V}_t, a_t^1, a_t^2$
  - 2: Compute  $\text{Cyl}(\mathcal{S}_t, \mathcal{V}_t)$  and its centroid  $\hat{s}_t, w_t = (a_t^1 - a_t^2)/\|a_t^1 - a_t^2\|, x_t = \langle \hat{s}_t, w_t \rangle, \bar{\delta} \in (0, 1/16d(d+1)^2 T^2)$
  - 3: Propose an incentive function  $\kappa(t, a) = 5 \cdot \mathbb{1}_{a_t^1}(a) + (5 + \langle \hat{s}_t, a_t^1 - a_t^2 \rangle) \cdot \mathbb{1}_{a_t^2}(a)$
  - 4: **if**  $A_t = a_t^1$  **then**
  - 5:      $\mathcal{S}_{t+1} = \mathcal{S}_t \cap \{s \mid \langle s, w_t \rangle \geq x_t\}$
  - 6: **else**
  - 7:      $\mathcal{S}_{t+1} = \mathcal{S}_t \cap \{s \mid \langle s, w_t \rangle \leq x_t\}$
  - 8: **end if**
  - 9: Let  $\mathcal{V}_{t+1} = \mathcal{V}_t$
  - 10: **if**  $\exists v \perp \mathcal{V}_{t+1}$  such that  $\text{diam}(\mathcal{S}_{t+1}, v) \leq \bar{\delta}$  **then**
  - 11:     add  $v$  to  $\mathcal{V}_t$
  - 12:     Repeat this step as many times as necessary
  - 13: **end if**
  - 14: **Output:**  $\mathcal{S}_{t+1}, \mathcal{V}_{t+1}$ .
- 

## F Lower Bound

*Proof of Proposition 1.* Suppose that the principal was to know  $(s_a)_{a \in [K]}$ . For any incentive  $\pi(t)$  offered on action  $a_t$  at round  $t$ , the agent selects his action following (1):  $A_t = \arg\max_{a \in [K]} s_a + \mathbb{1}_{a_t}(a)\pi(t)$ . The principal's expected reward is

$$\theta_{A_t} - \mathbb{1}_{a_t}(A_t)\pi(t) \leq \theta_{A_t} - \pi_{A_t}^* = \mu_a,$$

by definition of  $\pi_a^* := \max_{a' \in [K]} s_{a'} - s_a$  as the infimal amount of incentive to be offered on action  $a$  to make the agent choose it. Consequently, we have

$$\begin{aligned} \mathfrak{R}(T) &= T \mu^* - \sum_{t=1}^T \mathbb{E}[\theta_{A_t} - \mathbb{1}_{a_t}(A_t)\pi(t)] \\ &\geq \mathbb{E} \left[ \sum_{t=1}^T \mu^* - (\theta_{A_t} - \pi_{A_t}^*) \right] \\ &\geq \mathbb{E} \left[ \sum_{t=1}^T \mu^* - \mu_{A_t} \right]. \end{aligned}$$

Assuming the principal knows  $(s_a)_{a \in [K]}$ , observing  $X_a(t)$  is equivalent to observing  $X_a(t) - \pi_a^*$ . Using the result of Burnetas & Katehakis (1996) (see, e.g., Lattimore & Szepesvári, 2020, Theorem 16.2.), it then comes

$$\liminf_{T \rightarrow \infty} \frac{\mathfrak{R}(T)}{\log T} \geq \sum_{a, \mu_a < \mu^*} \frac{\mu^* - \mu_a}{\text{KL}_{\text{inf}}(\nu_a - \pi_a^*, \mu^*, \mathcal{D})}.$$

□