Learning to Persuade on the Fly: Robustness Against Ignorance

You Zu¹, Krishnamurthy Iyer¹, and Haifeng Xu²

¹Industrial and Systems Engineering, University of Minnesota, Minneapolis, MN {zu000002,kriyer}@umn.edu

²Department of Computer Science, University of Virginia, Charlottesville, VA hx4ad@virginia.edu

February 23, 2021

Abstract

We study a repeated persuasion setting between a sender and a receiver, where at each time t, the sender observes a payoff-relevant state drawn independently and identically from an unknown prior distribution, and shares state information with the receiver, who then myopically chooses an action. As in the standard setting, the sender seeks to persuade the receiver into choosing actions that are aligned with the sender's preference by selectively sharing information about the state. However, in contrast to the standard models, the sender does not know the prior, and has to persuade while gradually learning the prior on the fly.

We study the sender's learning problem of making persuasive action recommendations to achieve low regret against the optimal persuasion mechanism with the knowledge of the prior distribution. Our main positive result is an algorithm that, with high probability, is persuasive across all rounds and achieves $O(\sqrt{T\log T})$ regret, where T is the horizon length. The core philosophy behind the design of our algorithm is to leverage robustness against the sender's ignorance of the prior. Intuitively, at each time our algorithm maintains a set of candidate priors, and chooses a persuasion scheme that is simultaneously persuasive for all of them. To demonstrate the effectiveness of our algorithm, we further prove that no algorithm can achieve regret better than $\Omega(\sqrt{T})$, even if the persuasiveness requirements were significantly relaxed. Therefore, our algorithm achieves optimal regret for the sender's learning problem up to terms logarithmic in T.

1 Introduction

Examples of online platforms recommending content or products to their users abound in online economy. For instance, a marketplace like Etsy recommends vintage items made by independent sellers, a styling service like Stitch Fix recommends clothing designs made by custom brands, and an online platform like YouTube recommends content generated by independent channels to its users. Often, the platform making such recommendations must balance the dual objectives of being persuasive (a.k.a., being obedient [Bergemann and Morris, 2016]), i.e., make recommendations that will be adopted by the users, as well as furthering the platform's goals, such as increased sales, fewer returns or more engaged users.

To make a concrete illustration, consider a platform that recommends content created by independent creators ("channels") to its users. New channels regularly join the platform, and the distribution of the quality and/or relevance of a channel's content is initially unknown to the platform.

The users of the platform seek to consume fresh and high-quality content, while the platform itself may have other goals, such as maximizing content consumption, that is not fully aligned with users' interests. For each new content from a channel, the platform observes its quality/relevance (perhaps after an initial exploration or through in-house reviewers) and decides whether or not to recommend the content to its users. If the channel's content quality distribution is known, the platform can reliably make such recommendations that optimize its own goals while maintaining user satisfaction. However, in reality, the platform typically does not know the quality of a new channel, and thus must learn to make such persuasive recommendations over time.

In this paper, we study the problem faced by such a platform learning to make recommendations over time. Formally, we study a repeated persuasion setting between a sender and a receiver, where at each time t, the sender shares information about a payoff-relevant state with the receiver. The state at each time t is drawn independently and identically from an unknown distribution, and subsequent to receiving information about it, the receiver (myopically) chooses an action from a finite set. The sender seeks to persuade the receiver into choosing actions that are aligned with her preference by selectively sharing information about the state.

In contrast to the standard persuasion setting, we focus on the case where neither the sender nor the receiver knows the distribution of the payoff relevant state. Instead, the sender learns this distribution over time by observing the state realizations. We adopt the assumption common in the literature on Bayesian persuasion that at each time period, prior to observing the realized state in that period, the sender commits to a signaling mechanism that maps each state to a possibly random action recommendation. Subsequent to the state observation, the sender recommends an action as per the chosen signaling mechanism.

One natural requirement is for the sender to make recommendations that the receiver will find optimal to follow, i.e., recommendations that are persuasive. In the case where the receiver knows the prior, this requirement is easily justified from the observation that a rational Bayesian receiver will update her beliefs after receiving the recommendation and choose an action that maximizes her expected utility w.r.t. to her posterior beliefs. However, even if the receiver does not know the prior, practical considerations such as building and maintaining a reputation can make persuasiveness important to the sender. This is especially so in our example settings, where the sender is a long-lived platform and the receiver corresponds to a stream of users [Rayo and Segal, 2010], who may be able to verify the quality of recommendations ex post.

A sender who simply recommends the receiver's best action at the realized state will certainly be persuasive, but may end up with a significant loss in utility when compared to her utility if she had known the prior. Thus, the sender seeks to make persuasive action recommendations that achieve low regret against the optimal signaling mechanism with the knowledge of the prior distribution.

The primary contribution of this work is an efficient algorithm that, with high probability, makes persuasive action recommendations and at the same time achieves vanishing average regret. The algorithm we propose proceeds by maintaining at each time a set of candidate priors, based on the observed state realizations in the past. To compensate for the ignorance of the true prior, the algorithm resorts to robustness: it chooses a signaling mechanism that is simultaneously persuasive for each of the candidate priors and maximizes the sender's utility. Due to this aspect of the algorithm, we name it the *Robustness against Ignorance* (\mathfrak{Rai}) algorithm.

Our main positive result, Theorem 2, establishes that for any persuasion setting satisfying certain regularity conditions, the \mathfrak{Rai} algorithm achieves $O(\sqrt{T\log T})$ regret with high probability, where T is the horizon length. To show this result, we define a quantity Gap that measures the sender's cost of robust persuasion. Formally, $\mathsf{Gap}(\mu,\mathcal{B})$ captures the loss in the sender's expected utility (under belief μ) from using a signaling mechanism that is persuasive for all beliefs in the set \mathcal{B} , as opposed to using one that is persuasive only for the belief μ . The crux of the proof argument lies at

showing, in Proposition 1, that the sender's cost of robust persuasion $\mathsf{Gap}(\mu, \mathcal{B})$ is at most linear in the radius of the set \mathcal{B} . This is achieved via an explicit construction of a signaling mechanism that is persuasive for all beliefs in \mathcal{B} and achieves sender's utility close to optimum.

We strengthen our contribution by proving in Theorem 3 a matching lower-bound (up to $\log T$ terms) for the regret of any algorithm that makes persuasive recommendations. In particular, we construct a persuasion instance for which no persuasive algorithm can achieve regret better than $\Omega(\sqrt{T})$. Furthermore, in Theorem 4, we show this lower bound holds even if the persuasiveness requirements on the algorithm were significantly relaxed. As a byproduct, our regret analysis also leads to useful insight about robust persuasion w.r.t. to persuasiveness. Specifically, to prove our lower bound result, we carefully craft the persuasion instance and use its geometry to prove a linear cost for robust persuasion; this instance thus serves as a lower bound example for robust persuasion, which may be of independent interest.

Our results contribute to the work on online learning that seeks to evaluate the value of knowing the underlying distributional parameters in settings with repeated interactions [Kleinberg and Leighton, 2003]. In particular, our results fully characterize the sender's value of knowing the prior for repeated persuasion. Our attempt of relaxing the known prior assumptions is also aligned with the prior-independent mechanism design literature Dhangwatnotai et al. [2015], Chawla et al. [2013].

1.1 Literature Survey

Our paper contributes to the burgeoning literature on Bayesian persuasion and information design in economics, operations research and computer science. We refer curious readers to [Kamenica and Gentzkow, 2011] and [Bergemann and Morris, 2019] for a general overview of the recent developments as well as a survey from the algorithmic perspective by Dughmi [2017].

Online learning & mechanism design. Our work subscribes to the recent line of work that studies the interplay of learning and mechanism design in incomplete-information settings, in the absence of common knowledge on the prior. We briefly discuss the ones closely related to our work.

Castiglioni et al. [2020] focus on persuasion setting with a commonly known prior distribution of the state but unknown receiver types chosen adversarially from a finite set. They show that effective learning, in this case, is computationally intractable but does admit $O(\sqrt{T})$ regret learning algorithm, after relaxing the computability constraint. Our model complements theirs by focusing on known receiver types but unknown state distributions in a stochastic setup. Moreover, we achieve a similar (and tight) regret bound through a computationally efficient algorithm. Also relevant to us the recent line of work on Bayesian exploration Kremer et al. [2014], Mansour et al. [2015, 2016], which is also motivated by online recommendation systems. Opposite to us, these models assume a commonly known prior but the realized state is unobservable and thus needs to be learned during the repeated interactions.

Dispensing with the common prior itself, Camara et al. [2020] study an adversarial online learning model where both a mechanism designer and the agent learn about the states over time. The agent is assumed to minimize her counterfactual (internal) regret in response to the mechanism designer's policies, which are assumed to be non-responsive to the agent's actions. The authors characterize the regret of the mechanism designer relative to the best-in-hindsight fixed mechanism. Similar to our work, the regret bounds require the characterization of a "cost of robustness" of the underlying design problem. While related, our model is stochastic rather than adversarial, and thus a prior exists. Our model is similar in spirit to the prior-independent mechanism design literature Dhangwatnotai et al. [2015], Chawla et al. [2013], however, our setup is different. Moreover, our algorithm is measured by the regret whereas approximation ratios are often adopted for prior-independent mechanism design.

Recent works by Hahn et al. [4019, 2020] study information design in online optimiza-

tion problems such as the secretary problem Hahn et al. [2019] and the prophet inequalities Hahn et al. [2020]. The propose constant-approximation persuasive schemes. These online optimization problems often take the adversarial approach, which is different from our stochastic setup and learning-focused tasks. Therefore, our results are not comparable.

Robust persuasion: The algorithm we propose relies crucially on robust persuasion due to the ignorance of the prior, and as a part of establishing the regret bounds for the algorithm, we quantify the sender's cost of robustness. Kosterina [2018] studies a persuasion setting in the absence of the common prior assumption. In particular, the sender has a known prior, whereas only the set in which the receiver's prior lies is known to the sender. Furthermore, the sender evaluates the expected utility under each signaling mechanism with respect to the worst-case prior of the receiver. Similarly, Hu and Weng [2020] study the problem of sender persuading a privately informed receiver, where the sender seeks to maximize her expected payoff under the worst-case information of the receiver. Finally, Dworczak and Pavan [2020] study a related setting and propose a lexicographic solution concept where the sender first identifies the signaling mechanisms that maximize her worst-case payoff, and then among them chooses the one that maximizes the expected utility under her conjectured prior. In contrast to these work, our model focuses on a setting with common, but unknown, prior, and where the receiver has no private information. Instead, our notion of robustness is with respect to this unknown (common) prior.

Our work also relates to several other lines of research. Since the persuasion problem can be posed as a linear program, our setting relates to online convex optimization [Mahdavi et al., 2011, Yu et al., 2017, Yu and Neely, 2020, Yuan and Lamperski, 2018]. As in our work, these authors consider an online convex optimization problem with unknown objective function and/or constraints, and study algorithms that minimize regret while at the same time ensuring low constraint violations. However, most of the work here seeks to bound the magnitude of the constraint violation, whereas our persuasiveness guarantees correspond to constraint satisfaction with high probability. Finally, by characterizing the persuasion problem as a Stackelberg game between the sender's choice of a signaling mechanism and the receiver's subsequent choice of an action, our work is related to the broader work on the characterization of regret in repeated Stackelberg settings [Balcan et al., 2015, Dong et al., 2018, Chen et al., 2019].

2 Model

2.1 Persuasion instance

Consider a persuasion setting with a single sender persuading a receiver sequentially over a time horizon of length T. At each time $t \in [T] = \{0, \dots, T-1\}$, a state $\omega_t \in \Omega$ is drawn independently and identically from a prior distribution $\mu^* \in \Delta(\Omega)$. We focus on the setting where Ω is a known finite set, however the distribution μ^* is unknown to both the sender and the receiver. To capture the sender's initial knowledge (before time t = 0) about the prior μ^* , we assume that the sender knows that μ^* lies in the set $\mathcal{B}_0 \subseteq \Delta(\Omega)$.

At each time $t \in [T]$, the realized state ω_t is observed solely by the sender, who then shares with the receiver an action recommendation² $a_t \in A$ (chosen according to a signaling algorithm, which we define below), where A is a finite set of available actions for the receiver. The receiver then chooses an action \hat{a}_t (not necessarily equal to a_t), whereupon she receives a utility $u(\omega_t, \hat{a}_t)$, whereas the sender receives utility $v(\omega_t, \hat{a}_t)$. We assume that the receiver chooses her actions myopically, in order

¹For any finite set X, let $\Delta(X)$ denote the set of all probability distributions over X.

²Invoking the standard revelation principle, it can be shown that this is without loss of generality.

to model settings where a single long-run sender is persuading a stream of receivers. Furthermore, while our baseline model assumes that the receiver's utility is homogeneous across time, our model and the results easily apply to the case with heterogeneous receiver utility, assuming the sender observes the receiver's type prior to persuasion.

Without loss of generality, we assume that with $|\Omega| \geq 2$, $|A| \geq 2$ and $v(\omega, a) \in [0, 1]$ for all $\omega \in \Omega$ and $a \in A$. We refer to the tuple $\mathcal{I} = (\Omega, A, u, v, \mathcal{B}_0)$ with $u : \Omega \times A \to \mathbb{R}$ and $v : \Omega \times A \to [0, 1]$ as an *instance* of our problem.

2.2 Persuasion with known prior

Informally, given a persuasion instance \mathcal{I} , the sender's goal is to send action recommendations such that her long-run total expected utility is maximized. To formalize this goal, we begin by focusing on the setting where the sender and the receiver commonly know that the prior distribution $\mu^* = \mu \in \Delta(\Omega)$. In this setting, the sender's problem decouples across time periods, and standard results Kamenica and Gentzkow [2011], Dughmi and Xu [2019] imply that the sender's problem can be formulated as a linear program. To elaborate, for $\omega \in \Omega$ and $a \in A$, let $\sigma(\omega, a)$ denote the probability with which the sender recommends action a when the realized state is ω . We refer to $\sigma = (\sigma(\omega, a) : \omega \in \Omega, a \in A)$ as a signaling mechanism, and let $\mathcal{S} = \{\sigma : \sigma(\omega, \cdot) \in \Delta(A) \text{ for each } \omega \in \Omega\}$ denote the set of all signaling mechanisms. A signaling mechanism $\sigma \in \mathcal{S}$ is persuasive, if conditioned on receiving an action recommendation $\sigma \in \mathcal{S}$ is indeed optimal for the receiver to choose action σ . We denote the set of persuasive mechanisms by σ .

$$\mathsf{Pers}(\mu) \triangleq \left\{ \sigma \in \mathcal{S} : \sum_{\omega \in \Omega} \mu(\omega) \sigma(\omega, a) \left(u(\omega, a) - u(\omega, a') \right) \ge 0, \text{ for all } a, a' \in A \right\}. \tag{1}$$

Here, using Bayes' rule and assuming $\sum_{\omega \in \Omega} \mu(\omega) \sigma(\omega, a) > 0$, the receiver's posterior belief that the realized state is ω , upon receiving the recommendation $a \in A$, is given by $\frac{\mu(\omega)\sigma(\omega, a)}{\sum_{\omega' \in \Omega} \mu(\omega')\sigma(\omega', a)}$, and hence $\sum_{\omega \in \Omega} \left(\frac{\mu(\omega)\sigma(\omega, a)}{\sum_{\omega' \in \Omega} \mu(\omega')\sigma(\omega', a)}\right) u(\omega, a')$ denotes her expected utility of choosing action $a' \in A$ conditioned on receiving the recommendation $a \in A$. Thus, the inequality in the preceding definition³ captures the requirement that the receiver's expected utility for choosing action a' is at most the expected utility for choosing action $a \in A$. We note that $\mathsf{Pers}(\mu)$ is a non-empty convex polytope.

Given a persuasive signaling mechanism $\sigma \in \mathsf{Pers}(\mu)$, the receiver is incentivized to choose the recommended action and thus the sender's expected utility will be

$$V(\mu, \sigma) \triangleq \sum_{\omega \in \Omega} \sum_{a \in A} \mu(\omega) \sigma(\omega, a) v(\omega, a) \in [0, 1].$$
 (2)

Consequently, the sender's problem of selecting an optimal persuasive signaling mechanism that maximizes her expected utility can be formulated as the following linear program:

$$\mathsf{OPT}(\mu) \triangleq \max_{\sigma} V(\mu, \sigma)$$
, subject to $\sigma \in \mathsf{Pers}(\mu)$.

2.3 Persuasion with unknown prior

We now return to the setting where neither the sender nor the receiver knows the prior μ^* .

In general, the sender chooses at each time t an action recommendation a_t based on the complete history, namely the past state realizations, the past action recommendations by the sender as well

³If $\sum_{\omega \in \Omega} \mu(\omega) \sigma(\omega, a) = 0$, the inequality is trivially satisfied.

as the past actions chosen by the receiver. However, since the receiver does not know the prior, neither the past actions recommended by the sender nor the past actions chosen by the receiver carry any information about the prior beyond that contained in the state realizations. Thus, the history relevant to the sender consists solely of the state realizations until time t. To formalize this description, we first define the history h_t at the beginning of time t to be the state realizations prior to time t: $h_t = \bigcup_{\tau < t} \{(\tau, \omega_\tau)\}$ (with $h_0 = \emptyset$).

A signaling algorithm $\mathfrak{a} \equiv \mathfrak{a}(\mathcal{I})$ for the sender specifies, at each time $t \in [T]$, a signaling mechanism $\sigma^{\mathfrak{a}}[h_t] \in \mathcal{S}$ (we sometimes drop the superscript \mathfrak{a} when it is clear from the context). Implicitly, this reflects the assumption that the sender commits to the mechanism for sending recommendations at time t after observing history h_t but prior to observing the state realization ω_t at time t. Once the state ω_t is realized, the action recommendation a_t is drawn (independently) according to the distribution $\sigma[h_t](\omega_t, \cdot) \in \Delta(A)$. Thus, the probability that an action $a \in A$ is recommended is given by $\sigma[h_t](\omega_t, a)$, which, through an abuse of notation, we simplify to $\sigma(\omega_t, a; h_t)$.

We say a signaling algorithm \mathfrak{a} is β -persuasive for some $\beta \in [0, 1]$, if with probability at least $1 - \beta$, each $\sigma[h_t]$ is persuasive:

$$\inf_{\mu^* \in \mathcal{B}_0} \mathbf{P}_{\mu^*} \left(\sigma^{\mathfrak{a}}[h_t] \in \mathsf{Pers}(\mu^*), \text{ for each } t \in [T] \right) \ge 1 - \beta. \tag{3}$$

(Here, \mathbf{P}_{μ^*} represents the probability with respect to the (unknown) prior μ^* and any independent randomizations in the algorithm.) We simply call a 0-persuasive algorithm persuasive. Persuasiveness is a natural condition to impose on a signaling algorithm if the receiver knows the prior, as it ensures that the action recommendations will be accepted and implemented by the (Bayesian) receiver. However, even in settings where the receiver does not know the prior, imposing β -persuasiveness is still important as a means to guarantee that the recommendations will (mostly) be in the receiver's best interests. Moreover, it is easy to see that β -persuasive algorithms exist: the algorithm \mathfrak{Full} that always recommends $a_t \in \arg\max_{a \in A} u(\omega_t, a)$ after any history h_t is clearly 0-persuasive.

Given the preceding discussion, we henceforth assume that the receiver always accept the sender's recommendation, noting that this assumption is valid with probability at least $1 - \beta$ for a β -persuasive signaling algorithm. Thus, the sender's total utility from using a signaling algorithm \mathfrak{a} is given by

$$V_{\mathcal{I}}(\mathfrak{a},T) \triangleq \sum_{t \in [T]} v(\omega_t, a_t).$$

On the other hand, if the prior μ^* were known to the sender, her optimal total expected utility is given by $\mathsf{OPT}(\mu^*) \cdot T$. We measure the sender's regret from using the signaling algorithm \mathfrak{a} by

$$\mathsf{Reg}_{\mathcal{I}}(\mathfrak{a}, T, \mu^*) \triangleq \mathsf{OPT}(\mu^*) \cdot T - V_{\mathcal{I}}(\mathfrak{a}, T). \tag{4}$$

We are now ready to formalize the sender's learning problem. Begin by noticing that one must require the signaling algorithm $\mathfrak a$ to be β -persuasive, for some small β , for the regret defined in Eqn. (4) to be meaningful. At the same time, persuasiveness alone does not guarantee low regret: it is straightforward to show that the signaling algorithm \mathfrak{Full} , which is 0-persuasive, typically has $\Omega(T)$ regret. Thus, the central problem is to design, for any given instance \mathcal{I} , an algorithm $\mathfrak a$ that is β -persuasive for small β and simultaneously achieves vanishing average regret with high probability.

3 Bounding the optimal regret

Having described the learning problem faced by the sender, in this section, we present the signaling algorithm $\Re \mathfrak{a}$ and show that it is persuasive with 1 - o(1) probability and meanwhile achieves an

ALGORITHM 1: The Robustness Against Ignorance (Rai) algorithm

```
Input: Instance \mathcal{I}, Time horizon T

Parameters: \gamma_0 \in \mathcal{B}_0, \{\epsilon_t > 0 : t \in [T]\}

Output: a_t \in A for each t \in [T]

for t = 0 to T - 1 do

Choose any \sigma[h_t] \in \arg \max_{\sigma} \{V(\gamma_t, \sigma) : \sigma \in \mathsf{Pers}(\mathcal{B}_t)\};

Recommend a_t = a \in A with probability \sigma(\omega_t, a; h_t);

Update \gamma_{t+1}(\omega) \leftarrow \frac{1}{t+1} \sum_{\tau=0}^t \mathbf{I}\{\omega_\tau = \omega\} for each \omega \in \Omega;

Set \mathcal{B}_{t+1} \leftarrow \mathsf{B}_1(\gamma_{t+1}, \epsilon_{t+1});
end
```

average regret $O(\sqrt{\frac{\log T}{T}})$.

3.1 The Robustness Against Ignorance (Rai) Algorithm

Before describing the algorithm, we need some notation. First, for any set $\mathcal{B} \subseteq \Delta(\Omega)$, let $\mathsf{Pers}(\mathcal{B})$ denote the set of signaling mechanisms that are simultaneously persuasive under all priors $\mu \in \mathcal{B}$: $\mathsf{Pers}(\mathcal{B}) = \cap_{\mu \in \mathcal{B}} \mathsf{Pers}(\mu)$. We remark that for any non-empty set $\mathcal{B} \subseteq \Delta(\Omega)$, the set $\mathsf{Pers}(\mathcal{B})$ is convex since it is an intersection of convex sets $\mathsf{Pers}(\mu)$, and is non-empty since it contains the full-information signaling mechanism. Second, let $\mathsf{B}_1(\mu, \epsilon) \triangleq \{\mu' \in \Delta(\Omega) : \|\mu' - \mu\|_1 \leq \epsilon\}$ denote the (closed) ℓ_1 -ball of radius $\epsilon > 0$ at $\mu \in \Delta(\Omega)$.

The $\Re \mathfrak{a}$ is algorithm is formally described in Algorithm 1. At a high level, at each time $t \geq 0$, the algorithm maintains a set \mathcal{B}_t of candidates for the (unknown) prior μ^* and a particular estimate γ_t . It then selects a robustly persuasive signaling scheme that maximizes the sender utility w.r.t. to the currently estimated prior γ_t . Concretely, among signaling mechanisms that are persuasive for all beliefs $\mu \in \mathcal{B}_t$, $\Re \mathfrak{a}$ is selects the one that maximizes the sender's expected utility under γ_t . Finally, it makes an action recommendation a_t using this signaling mechanism, given the state realization ω_t .

From the intuitive description, it follows that to overcome the sender's ignorance of the prior, the algorithm seeks to be persuasive robustly against a conservative set of priors to maintain persuasiveness of the algorithm. In this goal, the parameters $\{\epsilon_t : t \in [T]\}$ determine how conservative the algorithm is in its persuasion: larger values of ϵ_t imply that the algorithm is more likely to be persuasive. (In particular, for $\epsilon_t > 2$, the algorithm chooses the full-information mechanism \mathfrak{Full} , and hence is persuasive with certainty.) Unsurprisingly, larger values of ϵ_t lead to larger regret, and hence the sender, in the choice of ϵ_t , must optimally trade-off the certainty of persuasiveness of the algorithm against its regret.

Our first result characterizes the optimal choice of the parameters, and shows that the algorithm $\Re \mathfrak{a}$ is efficient, and persuasive with high probability.

Theorem 1. For each $t \in [T]$, let $\epsilon_t = \min\{\sqrt{\frac{|\Omega|}{t}} \left(1 + \sqrt{\Phi \log T}\right), 2\}$ with $\Phi > 0$. Then, the \mathfrak{Rai} algorithm runs efficiently in polynomial time, and is β -persuasive with

$$\beta = \sup_{\mu^* \in \mathcal{B}_0} \mathbf{P}_{\mu^*} \left(\cap_{t \in [T]} \mathcal{B}_t \not\ni \mu^* \right) \le T^{1 - \frac{3\Phi\sqrt{\Omega}}{56}}.$$

In particular, for $\Phi > 20$, we have $\beta \leq T^{-0.5}$.

To see the efficiency of the $\Re \mathfrak{a}$ algorithm, note that at each time t the algorithm has to solve the optimization problem $\max_{\sigma} \{V(\gamma_t, \sigma) : \sigma \in \mathsf{Pers}(\mathcal{B}_t)\}$. Since $\mathcal{B}_t = \mathsf{B}_1(\gamma_t, \epsilon_t)$ is an ℓ_1 -ball of radius ϵ_t ,

it is a convex polyhedron with at most $|\Omega| \cdot (|\Omega| - 1)$ vertices.⁴ By the linearity of the obedience constraints and the convexity of \mathcal{B}_t , it follows that $\mathsf{Pers}(\mathcal{B}_t)$ is obtained by imposing the obedience constraints at priors corresponding to each of these vertices. Since there are $O(|\Omega| + |A|^2)$ obedience constraints for each prior, we obtain that the optimization problem is a polynomially-sized linear program, and hence can be solved efficiently.

The proof of the persuasiveness of \mathfrak{Rai} follows by showing that the empirical distribution γ_t concentrates around the unknown prior μ^* with high probability. Since, after any history h_t , the signaling mechanism $\sigma[h_t]$ chosen by the algorithm is persuasive for all priors in an ℓ_1 -ball around γ_t , we deduce that it is persuasive under μ^* as well. To show the concentration result, we use a concentration inequality for independent random vectors in a Banach space [Foucart and Rauhut, 2013]; the full proof is provided in Appendix A.2.

We observe that to get strong persuasiveness guarantees, the choice of ϵ_t in the preceding theorem requires the knowledge of the time horizon T. However, applying the standard doubling tricks [Besson and Kaufmann, 2018], one can convert our algorithm to an *anytime* version that has the same regret upper bound guarantee, at the cost of a weakened persuasiveness guarantee, where the persuasiveness β is weakened to a constant arbitrarily close to 0.

3.2 Regret bound

Given the persuasiveness of the algorithm, we now devote the rest of this section to proving the regret upper bound, under minor regularity conditions on the problem instance.

To state the regularity conditions under which we obtain our regret bound, we need a definition. For each action $a \in A$, let \mathcal{P}_a denote the set of beliefs for which action a is optimal for the receiver:

$$\mathcal{P}_a \triangleq \left\{ \mu \in \Delta(\Omega) : \mathbf{E}_{\mu} \left[u(\omega, a) \right] \geq \mathbf{E}_{\mu} \left[u(\omega, a') \right], \text{ for all } a' \in A \right\}.$$

It is without loss of generality to assume that for each $a \in A$, the set \mathcal{P}_a is non-empty.⁵ Our primary regularity condition further requires that each such set has a non-empty (relative) interior.

Regularity Conditions. The instance \mathcal{I} satisfies the following conditions:

- 1. There exists d > 0 such that for each $a \in A$, the set \mathcal{P}_a contains an ℓ_1 -ball of size d. Let D > 0 denote the largest value of d for which the preceding is true, and let $\eta_a \in \mathcal{P}_a$ be such that $\mathsf{B}_1(\eta_a, D) \subseteq \mathcal{P}_a$.
- 2. There exists a $p_0 > 0$ such that for all $\mu \in \mathcal{B}_0$ we have $\min_{\omega} \mu(\omega) \geq p_0 > 0$.

These regularity conditions are essential to ensure the possibility of successful learning; it is not hard to see that if some \mathcal{P}_a has zero measure, then the sender cannot hope to persuasively recommend action a without complete certainty of the prior. The first condition ensures that such degeneracies do not arise. The second condition is technical and is made primarily to ensure the potency of the first condition: without it, the sets \mathcal{P}_a may satisfy the first condition in $\Delta(\Omega)$, while failing to satisfy it relative to the subset $\Delta(\{\omega : \mu^*(\omega) > 0\})$. Taken together, these regularity conditions serve to avoid pathologies, and henceforth we restrict our attention only to those instances satisfying these regularity conditions.

With these regularity conditions in place, our main positive result establishes a regret upper bound for the $\Re \mathfrak{a}$ i algorithm. We observe that while p_0 appears in our regret bound, it is not required by the $\Re \mathfrak{a}$ i algorithm for its operation.

⁴These vertices are all of the form $\gamma_t + \frac{\epsilon_t}{2} (e_\omega - e_{\omega'})$, where e_ω is the belief that puts all its weight on ω .

⁵This is because the receiver can never be persuaded to play an action $a \in A$ for which \mathcal{P}_a is empty, and hence such an action can be dropped from A.

Theorem 2. For $t \in [T]$, let $\epsilon_t = \min\{\sqrt{\frac{|\Omega|}{t}}(1+\sqrt{\Phi \log T}), 2\}$ with $\Phi > 0$. Then, for all $\mu^* \in \mathcal{B}_0$, with probability at least $1-T^{1-\frac{3\Phi\sqrt{\Omega}}{56}}-T^{-8\Phi|\Omega|}$, the \mathfrak{Rai} algorithm satisfies

$$\mathrm{Reg}_{\mathcal{I}}(\mathfrak{Rai},\mu^*,T) \leq 2\left(\frac{20}{p_0^2D}+1\right)\left(1+\sqrt{|\Omega|T}(1+2\sqrt{\Phi\log T})\right).$$

In particular, the regret is of order $O\left(\frac{\sqrt{\Omega}}{p_0^2 D} \sqrt{T \log T}\right)$ with high probability.

To obtain the preceding regret bound, we start by bounding a quantity that measures the loss in the sender's expected utility for being robustly persuasive for a subset of priors close to each other. Formally, for $\mathcal{B} \subseteq \mathcal{B}_0$ and $\mu \in \mathcal{B}$, define

$$\mathsf{Gap}(\mu, \mathcal{B}) \triangleq \sup_{\sigma \in \mathsf{Pers}(\mu)} V(\mu, \sigma) - \sup_{\sigma \in \mathsf{Pers}(\mathcal{B})} V(\mu, \sigma). \tag{5}$$

Note that $\mathsf{Gap}(\mu, \mathcal{B})$ captures the difference in the sender's expected utility between using the optimal persuasive signaling mechanism for prior $\mu \in \mathcal{B}$ and using the optimal signaling mechanism that is persuasive for all priors $\mu' \in \mathcal{B}$.

The following proposition provides a bound on $\mathsf{Gap}(\mu,\mathsf{B}_1(\mu,\epsilon))$, and is the central result that underpins our proof of Theorem 2.

Proposition 1. For any
$$\mu \in \mathcal{B}_0$$
 and for all $\epsilon \geq 0$, we have $\mathsf{Gap}(\mu, \mathsf{B}_1(\mu, \epsilon)) \leq \left(\frac{4}{p_0^2 D}\right) \epsilon$.

The proof of the proposition is obtained through an explicit construction of a signaling mechanism $\widehat{\sigma}$ that is persuasive for all priors in the set $\mathsf{B}_1(\mu,\epsilon)$. To do this, we first use the geometry of the instance to split the prior μ into a convex combination of beliefs that either fully reveal the state, or are well-situated in the interior of the sets \mathcal{P}_a . (It is here that we make use of the two regularity assumptions.) We then construct a signaling mechanism $\widehat{\sigma}$ that induces, under prior μ , the aforementioned beliefs as posteriors, and show that the induced posteriors under a prior μ' close to μ are themselves close to the posteriors under μ (hence within the sets \mathcal{P}_a) or are fully revealing. This proves the persuasiveness of $\widehat{\sigma}$ for all priors μ' close to μ . The bound on $\mathsf{Gap}(\mu,\mathsf{B}_1(\mu,\epsilon))$ then follows by showing that the sender's payoff under $\widehat{\sigma}$ for prior μ is close to the payoff under the optimal signaling mechanism that is persuasive under μ .

Before presenting the proof of Proposition 1, we briefly sketch the proof of Theorem 2. Using the definition (4), in Lemma 2 we show the following bound on the regret:

$$\begin{split} \operatorname{Reg}_{\mathcal{I}}(\mathfrak{Rai}, \mu^*, T) &\leq \sum_{t \in [T]} \operatorname{Gap}(\mu^*, \operatorname{B}_1(\mu^*, \|\mu^* - \gamma_t\|_1)) + \sum_{t \in [T]} \operatorname{Gap}(\gamma_t, \operatorname{B}_1(\gamma_t, \epsilon_t)) \\ &+ \sum_{t \in [T]} \|\mu^* - \gamma_t\|_1 + \sum_{t \in [T]} \mathbf{E}_{\mu^*}[v(\omega_t, a_t)|h_t] - v(\omega_t, a_t). \end{split}$$

On the event $\{\mu^* \in \cap_{t \in [T]} \mathcal{B}_t\}$, we have $\|\mu^* - \gamma_t\|_1 \leq \epsilon_t$. Together with Proposition 1, we obtain that the first three terms are of order $\sum_{t \in [T]} \epsilon_t = O(\sqrt{T \log T})$. The final term is also of the same order due to a simple application of the Azuma-Hoeffding inequality. The complete proof is provided in Appendix A.1.

We end this section with the proof of Proposition 1.

Proof of Proposition 1. Observe that for $\epsilon > \frac{p_0^2 D}{4}$, we have $\frac{4\epsilon}{p_0^2 D} > 1$, and hence the specified bound is trivial. Hence, hereafter, we assume $\epsilon \leq \frac{p_0^2 D}{4}$.

To begin, let $\sigma \in \arg\max_{\sigma' \in \mathsf{Pers}(\mu)} V(\mu, \sigma')$ denote the optimal signaling mechanism under the prior μ . Let $A_+ = \{a \in A : \sum_{\omega \in \Omega} \sigma(\omega, a) > 0\}$ denote the set of all actions that are recommended with positive probability under σ . For each $a \in A_+$, let μ_a denote the receiver's posterior belief (under signaling mechanism σ) upon receiving the action recommendation a. Note that since σ is persuasive under μ , we must have $\mu_a \in \mathcal{P}_a$. By the splitting lemma [Aumann et al., 1995], it then follows that μ can be written as a convex combination $\sum_{a \in A_+} w_a \mu_a$ of $\{\mu_a : a \in A_+\}$, where $w_a \in [0,1]$ is given by $w_a = \sum_{\omega \in \Omega} \mu(\omega) \sigma(\omega,a)$.

We next explicitly construct a signaling mechanism $\hat{\sigma}$. To simplify the proof argument, the signaling mechanism $\hat{\sigma}$ we construct is not a *straightforward* mechanism, in the sense that it reveals more than just action recommendations for signals in S. Using revelation principle, one can construct an equivalent straightforward mechanism $\bar{\sigma}$ by *coalescing* [Anunrojwong et al., 2020] signals with the same best response for the signal. We omit the details of this reduction. We start with some definitions that are needed to construct the signaling mechanism $\hat{\sigma}$.

Let $\eta_a \in \mathcal{P}_a$ be such that $\mathsf{B}_1(\eta_a, D) \subseteq \mathcal{P}_a$. For $\delta = \frac{2\epsilon}{p_0 D} \in [0, 1]$, define $\xi_a = (1 - \delta)\mu_a + \delta\eta_a \in \mathcal{P}_a$ for each $a \in A_+$ and let $\xi = \sum_{a \in A_+} w_a \xi_a$. Furthermore, since $\mu_a \in \mathcal{P}_a$ and $\mathsf{B}_1(\eta_a, D) \subseteq \mathcal{P}_a$, the convexity of the set \mathcal{P}_a implies that $\mathsf{B}_1(\xi_a, \delta D) \subseteq \mathcal{P}_a$.

convexity of the set \mathcal{P}_a implies that $\mathsf{B}_1(\xi_a, \delta D) \subseteq \mathcal{P}_a$. Since $\mu \in \mathcal{B}_0 \subseteq \mathsf{relint}(\Delta(\Omega))$, we have $\frac{1}{1-\rho}(\mu-\rho\xi) \in \Delta(\Omega)$ for all small enough $\rho > 0$. Let $\bar{\rho} \triangleq \sup \left\{ \rho \in [0,1] : \frac{1}{1-\rho}(\mu-\rho\xi) \in \Delta(\Omega) \right\}$ be the largest such value in [0,1], and define χ as

$$\chi \triangleq \begin{cases} \frac{1}{1-\bar{\rho}} \left(\mu - \bar{\rho} \xi \right), & \text{if } \bar{\rho} < 1; \\ \mu, & \text{if } \bar{\rho} = 1. \end{cases}$$

Then, we obtain $\mu = \bar{\rho}\xi + (1 - \bar{\rho})\chi$. Furthermore, if $\bar{\rho} < 1$, we have

$$\bar{\rho} = \frac{\|\chi - \mu\|_1}{\|\chi - \mu\|_1 + \|\mu - \xi\|_1} \ge \frac{p_0}{p_0 + \delta},$$

where the inequality follows from $\|\mu - \xi\|_1 \leq \sum_{a \in A_+} w_a \|\mu_a - \xi_a\|_1 = \delta \sum_{a \in A_+} w_a \|\eta_a - \xi_a\|_1 \leq 2\delta$ and from the fact that χ lies in the boundary of $\Delta(\Omega)$, which implies $\|\chi - \mu\|_1 \geq 2 \min_{\omega} \mu(\omega) \geq 2p_0$.

With the preceding definitions in place, we are now ready to construct the mechanism $\widehat{\sigma}$. Let a_{ω} be a best response for the receiver at state $\omega \in \Omega$, and let $S = \{(\omega, a_{\omega}) \in \Omega \times A : \chi(\omega) > 0\}$. Consider the signaling mechanism $\widehat{\sigma}$, with the set of signals $A_+ \cup S$, defined as follows: for each $\omega \in \Omega$, let

$$\widehat{\sigma}(\omega, s) \triangleq \begin{cases} \overline{\rho} \frac{w_a \xi_a(\omega)}{\mu(\omega)}, & \text{for } s \in A_+; \\ (1 - \overline{\rho}) \frac{\chi(\omega)}{\mu(\omega)}, & \text{for } s = (\omega, a_\omega) \in S; \\ 0, & \text{otherwise.} \end{cases}$$
(6)

We now show that the signaling mechanism $\widehat{\sigma}$ is persuasive for all priors in $\mathsf{B}_1(\mu,\epsilon)$, in the sense that for all signals $s \in A_+$ it is optimal for the receiver to play s, and for all signals $s = (\omega, a_\omega) \in S$, it is optimal for the receiver to play a_ω . To see this, for any $\gamma \in \mathsf{B}_1(\mu,\epsilon)$, let $\gamma(\cdot|s)$ denote the receiver's posterior under signaling mechanism $\widehat{\sigma}$ upon receiving the signal $s \in A_+ \cup S$. For $s = (\omega, a_\omega) \in S$, we have $\gamma(\cdot|s) = e_\omega$, where e_ω is the belief that puts all its weight on $\omega \in \Omega$. Thus, upon receiving the signal $s = (\omega, a_\omega)$ it is optimal for the receiver with prior γ to take action a_ω . Thus, it only remains to show that signals $s = a \in A_+$ are persuasive.

For $a \in A_+$, we have for $\omega \in \Omega$,

$$\mu(\omega|a) = \frac{\mu(\omega)\widehat{\sigma}(\omega, a)}{\sum_{\omega' \in \Omega} \mu(\omega')\widehat{\sigma}(\omega', a)} = \xi_a(\omega)$$

$$\gamma(\omega|a) = \frac{\gamma(\omega)\widehat{\sigma}(\omega, a)}{\sum_{\omega' \in \Omega} \gamma(\omega')\widehat{\sigma}(\omega', a)} = \frac{\gamma(\omega)}{\mu(\omega)} \cdot \frac{\xi_a(\omega)}{\sum_{\omega' \in \Omega} \frac{\gamma(\omega')\xi_a(\omega')}{\mu(\omega')}}.$$

Then, using triangle inequality and some algebra, we obtain

$$\begin{split} \|\gamma(\cdot|a) - \mu(\cdot|a)\|_1 &= \sum_{\omega \in \Omega} |\gamma(\omega|a) - \xi_a(\omega)| \\ &\leq \sum_{\omega \in \Omega} \left| \gamma(\omega|a) - \frac{\gamma(\omega)}{\mu(\omega)} \cdot \xi_a(\omega) \right| + \sum_{\omega \in \Omega} \left| \frac{\gamma(\omega)}{\mu(\omega)} \cdot \xi_a(\omega) - \xi_a(\omega) \right| \\ &\leq 2 \cdot \sup_{\omega \in \Omega} \frac{\xi_a(\omega)}{\mu(\omega)} \cdot \|\gamma - \mu\|_1 \\ &\leq \frac{2\epsilon}{p_0}, \end{split}$$

where in the final inequality, we have used $\min_{\omega} \mu(\omega) \geq p_0$ to get $\sup_{\omega \in \Omega} \frac{\xi_a(\omega)}{\mu(\omega)} \leq \frac{1}{p_0}$. Since $\mu(\cdot|a) = \xi_a$, this implies that $\gamma(\cdot|a) \in \mathsf{B}_1\left(\xi_a, \frac{2\epsilon}{p_0}\right) = \mathsf{B}_1(\xi_a, \delta D) \subseteq \mathcal{P}_a$. Thus, the signal $a \in A_+$ is persuasive for the prior $\gamma \in \mathsf{B}_1(\mu, \epsilon)$. Taken together, we obtain that the signaling mechanism $\widehat{\sigma}$ is persuasive for all $\gamma \in \mathsf{B}_1(\mu, \epsilon)$.

The persuasiveness of $\hat{\sigma}$ for all $\gamma \in \mathsf{B}_1(\mu, \epsilon)$ implies that

$$\begin{split} \sup_{\sigma' \in \mathsf{Pers}(\mathsf{B}_1(\mu, \epsilon))} V(\mu, \sigma') &\geq V(\mu, \widehat{\sigma}) \\ &= \sum_{\omega \in \Omega} \sum_{a \in A_+} \mu(\omega) \widehat{\sigma}(\omega, a) v(\omega, a) + \sum_{\omega \in \Omega} \sum_{s \in S} \mu(\omega) \widehat{\sigma}(\omega, s) v(\omega, a_\omega) \\ &\geq \sum_{\omega \in \Omega} \sum_{a \in A_+} \bar{\rho} w_a \xi_a(\omega) v(\omega, a) \\ &= \bar{\rho} \sum_{\omega \in \Omega} \sum_{a \in A_+} w_a \left((1 - \delta) \mu_a(\omega) + \delta \eta_a(\omega) \right) v(\omega, a) \\ &\geq \bar{\rho} (1 - \delta) \sum_{\omega \in \Omega} \sum_{a \in A_+} w_a \mu_a(\omega) v(\omega, a) \\ &= \bar{\rho} (1 - \delta) \mathsf{OPT}(\mu). \end{split}$$

Thus, we obtain

$$\begin{split} \mathsf{Gap}(\mu,\mathsf{B}_1(\mu,\epsilon)) &= \mathsf{OPT}(\mu) - \sup_{\sigma' \in \mathsf{Pers}(\mathsf{B}_1(\mu,\epsilon))} V(\mu,\sigma') \\ &\leq \left(1 - \bar{\rho}(1-\delta)\right) \mathsf{OPT}(\mu) \\ &\leq \left(\frac{4}{p_\sigma^2 D}\right) \epsilon, \end{split}$$

where the final inequality follows from $\bar{\rho} \geq \frac{p_0}{p_0 + \delta}$, $\delta = \frac{2\epsilon}{p_0 D}$ and $\mathsf{OPT}(\mu) \leq 1$.

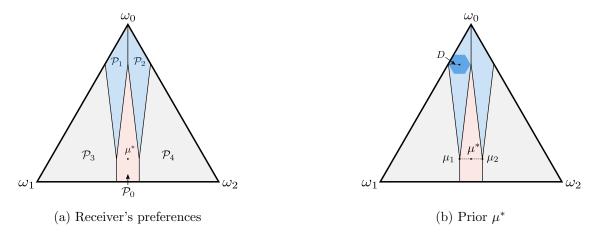


Figure 1: The persuasion instance \mathcal{I}_0 .

4 Lower Bounding the Regret

In this section, we show that the our regret upper bound in Theorem 2 are essentially tight with respect to the parameter D, T (up to a lower order $\sqrt{\log T}$ factor). We also show that the inverse polynomial dependence on p_0 , the smallest probability of states, is necessary though the exact order of the dependence on p_0 is left as an interesting open question.

Theorem 3. There exists an instance \mathcal{I}_0 , a prior $\mu^* \in \mathcal{B}_0$, and $T_0 > 0$ such that for any $T \geq T_0$ and any β_T -persuasive algorithm \mathfrak{a} the following holds with probability at least $\frac{1}{2} - 2\beta_T$:

$$\mathrm{Reg}_{\mathcal{I}}(\mathfrak{a},T,\mu^*) = T \cdot \mathrm{OPT}(\mu^*) - \sum_{t \in [T]} v(\omega_t,a_t) \geq \frac{\sqrt{T}}{32Dp_0}.$$

The persuasion instance \mathcal{I}_0 in the preceding theorem is carefully crafted to result in a substantial loss to the sender for being robustly persuasive. We begin by providing a geometric overview, and the underlying intuition, behind the crafting of this persuasion instance.

In the persuasion instance \mathcal{I}_0 , there are three states $\Omega = \{\omega_0, \omega_1, \omega_2\}$ and five actions $A = \{a_0, a_1, a_2, a_3, a_4\}$ for the receiver. At a high level, the receiver's preference can be illustrated as in Fig. 1a, which depicts the receiver's optimal action for any belief in the simplex. The regions \mathcal{P}_i in the figure correspond to the set of beliefs that induce action $a_i \in A$ as the receiver's best response. The instance is crafted in a way such that the sets \mathcal{P}_1 and \mathcal{P}_2 that induce actions a_1 and a_2 respectively are symmetric and extremely narrow with the width controlled by an ℓ_1 -ball⁶ of radius D contained, as depicted in Fig. 1b. (For completeness, the receiver's utility is listed explicitly in Table 1.) The sender seeks to persuade the receiver into choosing one of actions a_1 and a_2 (regardless of the state); all other actions are strictly worse for the sender. (Formally, we set $v(\omega, a) = 1$ if $a \in \{a_1, a_2\}$ and 0 otherwise, for all ω .) The sender's initial knowledge regarding the prior is captured by the set $\mathcal{B}_0 = \{\mu \in \Delta(\Omega) : \min_{\omega} \mu \geq p_0\}$, while the prior of interest is $\mu^* = (p_0, \frac{1-p_0}{2}, \frac{1-p_0}{2})$, corresponding to the midpoint of the tips of the sets \mathcal{P}_i , as shown in Fig. 1b. We focus on the setting where the instance parameters D and p_0 satisfy $Dp_0 < 1/64$.

The proof of Theorem 3 relies crucially on the following proposition, which shows that in the instance \mathcal{I}_0 , it is costly to require the signaling mechanism to be robustly persuasive for a set of priors around μ^* . It also implies that the bound on $\mathsf{Gap}(\cdot)$ obtained in Proposition 1 is almost tight,

⁶Since $|\Omega| = 3$, the ℓ_1 -ball here is an hexagon.

Table 1: Receiver's utility in instance \mathcal{I}_0 , with $u(\omega, a_0)$ normalized to 0 for all $\omega \in \Omega$.

	a_1	a_2	a_3	a_4
ω_0	$2D^2$	$2D^2$	$-2D(1-p_0-2D)$	$-2D(1-p_0-2D)$
ω_1	$(1-2D)(1-D)-p_0$	$(D+1)(2D-1) + p_0$	$2(1 - p_0 - 2D)(1 - D)$	$-2(1 - p_0 - 2D)(D+1)$
ω_2	$(D+1)(2D-1) + p_0$	$(1-2D)(1-D)-p_0$	$-2(1 - p_0 - 2D)(D+1)$	$2(1 - p_0 - 2D)(1 - D)$

up to a factor of $1/p_0$. We remark that this result also serves as a worst-case (lower bound) example for robust persuasion, which may be of independent interest.

Proposition 2 (The Cost of Robustness). For the instance \mathcal{I}_0 , we have $\mathsf{OPT}(\mu^*) = 1$. Furthermore, for all $\epsilon \in (0, D)$, we have

$$\mathsf{Gap}(\mu^*,\mathsf{Pers}(\mu^*,\bar{\mu}_1,\bar{\mu}_2)) \geq \frac{\epsilon}{8Dn_0},$$

where $\bar{\mu}_1 = \mu^* + \frac{\epsilon}{2}(e_1 - e_2)$, $\bar{\mu}_2 = \mu^* + \frac{\epsilon}{2}(e_2 - e_1)$, where the belief e_i puts all its weight on ω_i .

We defer the rigorous algebraic proof of the proposition to Appendix C.1 and present a brief sketch using a geometric argument here. Since the prior μ can be written as a convex combination $\mu = (\mu_1 + \mu_2)/2$, where μ_1 and μ_2 are the tips of region \mathcal{P}_1 and \mathcal{P}_2 respectively (see Fig. 1b), by the splitting lemma [Aumann et al., 1995], it follows that the optimal signaling mechanism sends signals induces posterior beliefs μ_1 and μ_2 leading to receiver's choice of a_1 and a_2 respectively. Since the sender can always persuade the receiver to choose one of her preferred actions, we obtain $\mathsf{OPT}(\mu^*) = 1$.

On the other hand, for a signaling mechanism to be robustly persuasive for all priors ϵ -close to the prior μ^* for sufficiently small ϵ , the posteriors for the sender's preferred actions a_1, a_2 induced by the signaling mechanism have to be shifted up significantly in the narrow region. Such a large discrepancy ultimately forces the sender to suffer a substantial loss in the expected utility.

Armed with this proposition, we are now ready to present the proof of Theorem 3.

Proof of Theorem 3. For a prior $\mu \in \mathcal{B}_0$, define the event $\mathcal{E}_T(\mu)$ as

$$\mathcal{E}_T(\mu) = \{h_T : \sigma^{\mathfrak{a}}[h_t] \in \mathsf{Pers}(\mu), \text{ for each } t \in [T]\}.$$

In words, under the event $\mathcal{E}_T(\mu)$, the signaling mechanisms $\sigma^{\mathfrak{a}}[h_t]$ chosen by the algorithm \mathfrak{a} after any history $h_t \in \mathcal{E}_T(\mu)$ is persuasive for the prior μ . Since the algorithm \mathfrak{a} is β_T -persuasive, we obtain

$$\mathbf{P}_{\mu}\left(\mathcal{E}_{T}(\mu)\right) \geq 1 - \beta_{T}, \text{ for all } \mu \in \mathcal{B}_{0}.$$

Fix an $\epsilon \in (0, \frac{1-3p_0}{2})$ to be chosen later, and consider the priors $\bar{\mu}_0 = \mu^* = (p_0, \frac{1-p_0}{2}, \frac{1-p_0}{2})$ and $\bar{\mu}_1 = \mu^* + \frac{\epsilon}{2} (e_1 - e_2)$ and $\bar{\mu}_2 = \mu^* + \frac{\epsilon}{2} (e_2 - e_1)$, where e_j is the belief that puts all its weight on state ω_j for $j \in \{1, 2\}$. Observe that for each $i \in \{0, 1, 2\}$ and for all $\epsilon \in (0, \frac{1-3p_0}{2})$, we have $\bar{\mu}_i \in \mathcal{B}_0$ and hence $\mathbf{P}_{\bar{\mu}_i}(\mathcal{E}_T(\bar{\mu}_i)) \geq 1 - \beta_T$.

Now, on the event $\mathcal{E}_T(\bar{\mu}_0) \cap \mathcal{E}_T(\bar{\mu}_1) \cap \mathcal{E}_T(\bar{\mu}_2)$, the signaling mechanisms $\sigma^{\mathfrak{a}}[h_t]$ chosen by the algorithm after any history h_t is persuasive for all the priors $\bar{\mu}_i$, i = 0, 1, 2. This implies that on this event, we have

$$T\cdot \mathsf{OPT}(\bar{\mu}_0) - \sum_{t\in [T]} V(\bar{\mu}_0, \sigma^{\mathfrak{a}}[h_t]) \geq T\cdot \mathsf{Gap}(\bar{\mu}_0, \{\bar{\mu}_0, \bar{\mu}_1, \bar{\mu}_2\}) \geq \frac{\epsilon T}{8Dp_0},$$

where the first inequality follows from the definition 5 of $\mathsf{Gap}(\cdot)$, and the second inequality follows from Proposition 2.

Now, we have

$$2 |\mathbf{P}_{\bar{\mu}_0} (\mathcal{E}_T(\bar{\mu}_1)) - \mathbf{P}_{\bar{\mu}_1} (\mathcal{E}_T(\bar{\mu}_1))|^2 \le \sum_{t \in [T]} \mathsf{KL} (\bar{\mu}_0 || \bar{\mu}_1)$$

$$= \frac{1 - p_0}{2} \log \left(\frac{(1 - p_0)^2}{(1 - p_0)^2 - \epsilon^2} \right) T$$

$$= \frac{1 - p_0}{2} \log \left(1 + \frac{\epsilon^2}{(1 - p_0)^2 - \epsilon^2} \right) T$$

$$\le \frac{1 - p_0}{2} \left(\frac{\epsilon^2}{(1 - p_0)^2 - \epsilon^2} \right) T,$$

where the first inequality is the Pinsker's inequality, and the first equality is from the definition of the Kullback-Leibler divergence, and the final inequality follows from $\log(1+x) \le x$ for $x \ge 0$. Thus, for $\epsilon < \frac{1-p_0}{2}$, we obtain

$$2 |\mathbf{P}_{\bar{\mu}_0} (\mathcal{E}_T(\bar{\mu}_1)) - \mathbf{P}_{\bar{\mu}_1} (\mathcal{E}_T(\bar{\mu}_1))|^2 \le \frac{2\epsilon^2 T}{3(1-p_0)} \le \epsilon^2 T,$$

where we have used $p_0 \leq \frac{1}{|\Omega|} = \frac{1}{3}$ in the final inequality. Thus, we obtain that

$$\mathbf{P}_{\bar{\mu}_0}\left(\mathcal{E}_T(\bar{\mu}_1)\right) \ge \mathbf{P}_{\bar{\mu}_0}\left(\mathcal{E}_T(\bar{\mu}_0)\right) - |\mathbf{P}_{\bar{\mu}_0}\left(\mathcal{E}_T(\bar{\mu}_1)\right) - \mathbf{P}_{\bar{\mu}_1}\left(\mathcal{E}_T(\bar{\mu}_1)\right)|$$

$$\ge 1 - \beta_T - \epsilon \sqrt{\frac{T}{2}}.$$

By the same argument, we obtain $\mathbf{P}_{\bar{\mu}_0}\left(\mathcal{E}_T(\bar{\mu}_2)\right) \geq 1 - \beta_T - \epsilon \sqrt{\frac{T}{2}}$.

By the linearity of the obedience constraints, we obtain that if $\sigma \in \mathsf{Pers}(\bar{\mu}_1) \cap \mathsf{Pers}(\bar{\mu}_2)$, then $\sigma \in \mathsf{Pers}(\bar{\mu}_0)$. Thus, we have $\mathcal{E}_T(\bar{\mu}_1) \cap \mathcal{E}_T(\bar{\mu}_2) \subseteq \mathcal{E}_T(\bar{\mu}_0)$, and hence

$$\mathbf{P}_{\bar{\mu}_0}(\mathcal{E}_T(\mu) \cap \mathcal{E}_T(\bar{\mu}_1) \cap \mathcal{E}_T(\bar{\mu}_2)) = \mathbf{P}_{\bar{\mu}_0}(\mathcal{E}_T(\bar{\mu}_1) \cap \mathcal{E}_T(\bar{\mu}_2))$$

$$\geq \mathbf{P}_{\bar{\mu}_0}(\mathcal{E}_T(\bar{\mu}_1)) + \mathbf{P}_{\bar{\mu}_0}(\mathcal{E}_T(\bar{\mu}_2)) - 1$$

$$\geq 1 - 2\beta_T - \epsilon \sqrt{2T}.$$

Finally, by the Azuma-Hoeffding inequality, we obtain

$$\mathbf{P}_{\bar{\mu}_0} \left(\sum_{t \in [T]} V(\bar{\mu}_0, \sigma^{\mathfrak{a}}[h_t]) - \sum_{t \in [T]} v(\omega_t, a_t) < -\sqrt{T} \right) < e^{-1/2}.$$

Taken together, we obtain that with probability at least $2-2\beta_T - \epsilon\sqrt{2T} - e^{-1/2}$, we have

$$\mathrm{Reg}_{\mathcal{I}}(\mathfrak{a},T,\bar{\mu}_0) = T \cdot \mathrm{OPT}(\bar{\mu}_0) - \sum_{t \in [T]} v(\omega_t,a_t) \geq \frac{\epsilon T}{8Dp_0} - \sqrt{T}.$$

For $T \ge T_0 = \frac{1}{(1-3p_0)^2}$, choosing $\epsilon = \frac{1}{2\sqrt{T}} \le \frac{1-3p_0}{2}$, we obtain, with probability at least $\frac{1}{2} - 2\beta_T$,

$$\mathrm{Reg}_{\mathcal{I}}(\mathfrak{a},T,\bar{\mu}_0) = T \cdot \mathrm{OPT}(\bar{\mu}_0) - \sum_{t \in [T]} v(\omega_t,a_t) \geq \sqrt{T} \left(\frac{1}{16Dp_0} - 1\right) \geq \frac{\sqrt{T}}{32Dp_0},$$

for
$$Dp_0 < 1/32$$
.

We end this section by observing that the lower bound in Theorem 3 extends also to signaling algorithms that satisfy a much weaker persuasiveness requirement. Specifically, we say an algorithm \mathfrak{a} is weakly β -persuasive if during T rounds of persuasion, the expected number of rounds that the algorithm is not persuasive is at most βT . Clearly, any β -persuasive algorithm is also weakly β -persuasive. To see how this definition is weaker, imagine an algorithm that is never persuasive at time 1 but always persuasive at all other times. Then it is a 1-persuasive algorithm (meaning it is never persuasive across all rounds) but is weakly 1/T-persuasive. Despite the much weaker restriction, the following theorem shows that weakly β -persuasive algorithms still cannot guarantee a significantly better regret (i.e., by more than logarithmic terms) than the \mathfrak{Rai} algorithm with a stronger persuasiveness guarantee.

Theorem 4. There exists an instance \mathcal{I}_0 , a prior $\mu^* \in \mathcal{B}_0$, and $T_0 > 0$ such that for any $T \geq T_0$ and for any algorithm \mathfrak{a} that is weakly β_T -persuasive, the following holds with probability at least $\frac{1}{2} - 8\beta_T$:

$$\sum_{t \in [T]} \left(\mathsf{OPT}(\mu^*) - v(\omega_t, a_t) \right) \mathbf{I}\{\sigma^{\mathfrak{a}}[h_t] \in \mathsf{Pers}(\mu^*)\} \geq \frac{\sqrt{T}}{64Dp_0},$$

We remark that in the preceding theorem, the regret is measured only over those time periods where the signaling mechanism chosen by the algorithm is persuasive under μ^* . An immediate implication of the result is that for any algorithm, the maximum of the expected number of unpersuasive recommendations and the regret for the persuasive recommendations is $\Omega(\sqrt{T})$. The proof, which uses the same instance \mathcal{I}_0 but considers different events, is provided in Appendix C.2.

5 Conclusion

We studied a repeated Bayesian persuasion problem where the prior distribution of payoff-relevant states is unknown to the sender. The sender learns this distribution from observing state realizations while making recommendations to the receiver. We propose the \mathfrak{Rai} algorithm which persuades robustly and achieves $O(\sqrt{T \log T})$ regret against the optimal signaling mechanism under the knowledge of the prior. To match this upper-bound, we construct a persuasion instance for which no persuasive algorithm achieves regret better than $\Omega(\sqrt{T})$. Taken together, our work precisely characterizes the value of knowing the prior distribution in repeated persuasion.

While in our analysis we have assumed that the receiver's utility is fixed across time periods, our model and the analysis can be easily extended to accommodate heterogeneous receivers, as long as the sender observes the receiver's type prior to making the recommendation, and the cost of robustness Gap can be uniformly bounded across different receiver types. More interesting is the setting where the sender must persuade a receiver with an unknown type. In such a setting, assuming the sender cannot elicit the receiver's type prior to making the recommendation, the sender makes a menu of action recommendations (one for each receiver type). It can be shown the complete information problem in this setting corresponds to public persuasion of a group of receivers with no externality, which is known to be a computationally hard linear program with exponentially many constraints [Dughmi and Xu, 2017]. Consequently, our algorithm ceases to be computationally efficient. Nevertheless, our results imply that the algorithm continues to maintain the $O(\sqrt{T \log T})$ regret bound.

References

- Jerry Anunrojwong, Krishnamurthy Iyer, and David Lingenbrink. Persuading risk-conscious agents: A geometric approach. Available at SSRN 3386273, 2020.
- Robert J Aumann, Michael Maschler, and Richard E Stearns. Repeated games with incomplete information. MIT press, 1995.
- Maria-Florina Balcan, Avrim Blum, Nika Haghtalab, and Ariel D Procaccia. Commitment without regrets: Online learning in stackelberg security games. In *Proceedings of the sixteenth ACM conference on economics and computation*, pages 61–78, 2015.
- Dirk Bergemann and Stephen Morris. Bayes correlated equilibrium and the comparison of information structures in games. *Theoretical Economics*, 11(2):487–522, 2016.
- Dirk Bergemann and Stephen Morris. Information design: A unified perspective. *Journal of Economic Literature*, 57(1):44–95, 2019.
- Lilian Besson and Emilie Kaufmann. What doubling tricks can and can't do for multi-armed bandits. arXiv preprint arXiv:1803.06971, 2018.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. Concentration inequalities: A nonasymptotic theory of independence. Oxford university press, 2013.
- Modibo Camara, Jason Hartline, and Aleck Johnsen. Mechanisms for a no-regret agent: Beyond the common prior. arXiv preprint arXiv:2009.05518, 2020.
- Matteo Castiglioni, Andrea Celli, Alberto Marchesi, and Nicola Gatti. Online bayesian persuasion. Advances in Neural Information Processing Systems, 33, 2020.
- Shuchi Chawla, Jason D Hartline, David Malec, and Balasubramanian Sivan. Prior-independent mechanisms for scheduling. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 51–60, 2013.
- Yiling Chen, Yang Liu, and Chara Podimata. Learning strategy-aware linear classifiers. arXiv preprint arXiv:1911.04004, 2019.
- Peerapong Dhangwatnotai, Tim Roughgarden, and Qiqi Yan. Revenue maximization with a single sample. Games and Economic Behavior, 91:318–333, 2015.
- Jinshuo Dong, Aaron Roth, Zachary Schutzman, Bo Waggoner, and Zhiwei Steven Wu. Strategic classification from revealed preferences. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 55–70, 2018.
- Shaddin Dughmi. Algorithmic information structure design: a survey. ACM SIGecom Exchanges, 15(2):2–24, 2017.
- Shaddin Dughmi and Haifeng Xu. Algorithmic persuasion with no externalities. In *Proceedings of the 2017 ACM Conference on Economics and Computation*, pages 351–368, 2017.
- Shaddin Dughmi and Haifeng Xu. Algorithmic bayesian persuasion. SIAM Journal on Computing, (0):STOC16-68, 2019.

- Piotr Dworczak and Alessandro Pavan. Preparing for the worst but hoping for the best: Robust (bayesian) persuasion. 2020.
- Simon Foucart and Holger Rauhut. A Mathematical Introduction to Compressive Sensing. 2013. doi: 10.1007/978-0-8176-4948-7.
- Niklas Hahn, Martin Hoefer, and Rann Smorodinsky. The secretary recommendation problem. arXiv preprint arXiv:1907.04252, 2019.
- Niklas Hahn, Martin Hoefer, and Rann Smorodinsky. Prophet inequalities for bayesian persuasion. In *Proc. 29th Int. Joint Conf. Artif. Intell.(IJCAI)*, pages 175–181, 2020.
- Ju Hu and Xi Weng. Robust persuasion of a privately informed receiver. *Economic Theory*, pages 1–45, 2020.
- Emir Kamenica and Matthew Gentzkow. Bayesian persuasion. *American Economic Review*, 101(6): 2590–2615, 2011.
- Robert Kleinberg and Tom Leighton. The value of knowing a demand curve: Bounds on regret for online posted-price auctions. In *Proceedings of the 44th Annual IEEE Symposium on Foundations of Computer Science*, FOCS '03, page 594, USA, 2003. IEEE Computer Society. ISBN 0769520405.
- Svetlana Kosterina. Persuasion with unknown beliefs. Work. Pap., Princeton Univ., Princeton, NJ, 2018.
- Ilan Kremer, Yishay Mansour, and Motty Perry. Implementing the "wisdom of the crowd". *Journal of Political Economy*, 122(5):988–1012, 2014.
- Mehrdad Mahdavi, Rong Jin, and Tianbao Yang. Trading regret for efficiency: Online convex optimization with long term constraints. *CoRR*, abs/1111.6082, 2011. URL http://arxiv.org/abs/1111.6082.
- Yishay Mansour, Aleksandrs Slivkins, and Vasilis Syrgkanis. Bayesian incentive-compatible bandit exploration. In *Proceedings of the Sixteenth ACM Conference on Economics and Computation*, pages 565–582, 2015.
- Yishay Mansour, Aleksandrs Slivkins, Vasilis Syrgkanis, and Zhiwei Steven Wu. Bayesian exploration: Incentivizing exploration in bayesian games. In *Proceedings of the 2016 ACM Conference on Economics and Computation*, pages 661–661, 2016.
- Luis Rayo and Ilya Segal. Optimal information disclosure. *Journal of political Economy*, 118(5): 949–987, 2010.
- Hao Yu and Michael J. Neely. A low complexity algorithm with o(\(\hat{a}^*\)st) regret and o(1) constraint violations for online convex optimization with long term constraints. *Journal of Machine Learning Research*, 21(1):1-24, 2020. URL http://jmlr.org/papers/v21/16-494.html.
- Hao Yu, Michael J. Neely, and Xiaohan Wei. Online convex optimization with stochastic constraints, 2017.
- Jianjun Yuan and Andrew G. Lamperski. Online convex optimization for cumulative constraints. CoRR, abs/1802.06472, 2018. URL http://arxiv.org/abs/1802.06472.

Proofs from Section 3

This section provides the proof of our main theorems in Section 3. Throughout, we use the same notation as in the main text.

Proof of Theorem 2 A.1

In this section, we provide the proof of Theorem 2. In the process, we also state and prove several helper lemmas used in the proof.

Proof of Theorem 2. From Lemma 2, on the event $\{\mu^* \in \cap_{t \in [T]} \mathcal{B}_t\}$, we have

$$\begin{split} \operatorname{Reg}_{\mathcal{I}}(\mathfrak{Rai}, \mu^*, T) &\leq \left(\frac{20}{p_0^2 D} + 1\right) \sum_{t \in [T]} \epsilon_t + \sum_{t \in [T]} \mathbf{E}_{\mu^*}[v(\omega_t, a_t) | h_t] - v(\omega_t, a_t) \\ &\leq \left(\frac{20}{p_0^2 D} + 1\right) \left(2 + \sum_{t=1}^{T-1} \sqrt{\frac{|\Omega|}{t}} (1 + \sqrt{\Phi \log T})\right) \\ &+ \sum_{t \in [T]} \mathbf{E}_{\mu^*}[v(\omega_t, a_t) | h_t] - v(\omega_t, a_t) \\ &\leq 2 \left(\frac{20}{p_0^2 D} + 1\right) \left(1 + \sqrt{|\Omega| T} (1 + \sqrt{\Phi \log T})\right) \\ &+ \sum_{t \in [T]} \mathbf{E}_{\mu^*}[v(\omega_t, a_t) | h_t] - v(\omega_t, a_t), \end{split}$$

where in the final inequality, we have used the fact that $\sum_{t=1}^{T-1} 1/\sqrt{t} \leq 2\sqrt{T}$. From Theorem 1, we have $\mathbf{P}_{\mu}\left(\bigcap_{t\in[T]}\mathcal{B}_{t}\not\ni\mu\right)\leq T^{1-\frac{3\Phi\sqrt{\Omega}}{56}}$. Furthermore, from Lemma 7, we have for $\alpha > 0$,

$$\mathbf{P}_{\mu^*} \left(\sum_{t \in [T]} \mathbf{E}_{\mu^*} [v(\omega_t, a_t) | h_t] - v(\omega_t, a_t) \ge \sqrt{\alpha T \log T} \right) < \frac{1}{T^{2\alpha}}.$$

After choosing $\alpha = 4\Phi |\Omega|$ and taking the union bound, we obtain with probability at least $1 - T^{1 - \frac{3\Phi\sqrt{\Omega}}{56}} - T^{-8\Phi|\Omega|}$, we have

$$\operatorname{Reg}_{\mathcal{I}}(\mathfrak{Rai},\mu^*,T) \leq 2\left(\frac{20}{p_0^2D}+1\right)\left(1+\sqrt{|\Omega|T}(1+2\sqrt{\Phi\log T})\right). \hspace{1cm} \Box$$

Lemma 1. The Rai algorithm satisfies

$$\begin{split} \operatorname{Reg}_{\mathcal{I}}(\mathfrak{Rai}, \mu^*, T) &= \sum_{t \in [T]} (\operatorname{OPT}(\mu^*) - \operatorname{OPT}(\gamma_t)) + \sum_{t \in [T]} \operatorname{Gap}(\gamma_t, \mathcal{B}_t) \\ &+ \sum_{t \in [T]} \left(V(\gamma_t, \sigma[h_t]) - V(\mu^*, \sigma[h_t]) \right) \\ &+ \sum_{t \in [T]} \mathbf{E}_{\mu^*}[v(\omega_t, a_t) | h_t] - v(\omega_t, a_t). \end{split}$$

Proof. We have

$$\begin{split} \mathsf{Reg}_{\mathcal{I}}(\mathfrak{Rai}, \mu^*, T) &= \mathsf{OPT}(\mu^*) \cdot T - \sum_{t \in [T]} v(\omega_t, a_t) \\ &= \mathsf{OPT}(\mu^*) \cdot T - \sum_{t \in [T]} \mathbf{E}_{\mu^*}[v(\omega_t, a_t) | h_t] + \left(\sum_{t \in [T]} \mathbf{E}_{\mu^*}[v(\omega_t, a_t) | h_t] - v(\omega_t, a_t) \right). \end{split}$$

Now, note that for $t \in T$,

$$\begin{split} \mathbf{E}_{\mu^*}[v(\omega_t, a_t)|h_t] &= V(\mu^*, \sigma[h_t]) \\ &= V(\gamma_t, \sigma[h_t]) + (V(\mu^*, \sigma[h_t]) - V(\gamma_t, \sigma[h_t])) \\ &= \sup_{\sigma \in \mathsf{Pers}(\mathcal{B}_t)} V(\gamma_t, \sigma) + (V(\mu^*, \sigma[h_t]) - V(\gamma_t, \sigma[h_t])) \,. \end{split}$$

Thus, we have

$$\operatorname{Reg}_{\mathcal{I}}(\mathfrak{Rai}, \mu^{*}, T) = \operatorname{OPT}(\mu^{*}) \cdot T - \sum_{t \in [T]} \mathbf{E}_{\mu^{*}}[v(\omega_{t}, a_{t})|h_{t}] + \left(\sum_{t \in [T]} \mathbf{E}_{\mu^{*}}[v(\omega_{t}, a_{t})|h_{t}] - v(\omega_{t}, a_{t})\right) \\
= \operatorname{OPT}(\mu^{*}) \cdot T - \sum_{t \in [T]} V(\mu^{*}, \sigma[h_{t}]) + \left(\sum_{t \in [T]} \mathbf{E}_{\mu^{*}}[v(\omega_{t}, a_{t})|h_{t}] - v(\omega_{t}, a_{t})\right) \\
= \sum_{t \in [T]} \left(\operatorname{OPT}(\mu^{*}) - V(\mu^{*}, \sigma[h_{t}])\right) + \left(\sum_{t \in [T]} \mathbf{E}_{\mu^{*}}[v(\omega_{t}, a_{t})|h_{t}] - v(\omega_{t}, a_{t})\right). \tag{7}$$

Finally, note that

$$\begin{split} \mathsf{OPT}(\mu^*) - V(\mu^*, \sigma[h_t]) &= \mathsf{OPT}(\mu^*) - V(\gamma_t, \sigma[h_t]) + V(\gamma_t, \sigma[h_t]) - V(\mu^*, \sigma[h_t]) \\ &= (\mathsf{OPT}(\mu^*) - \mathsf{OPT}(\gamma_t)) + (\mathsf{OPT}(\gamma_t) - V(\gamma_t, \sigma[h_t])) \\ &+ (V(\gamma_t, \sigma[h_t]) - V(\mu^*, \sigma[h_t])) \\ &= (\mathsf{OPT}(\mu^*) - \mathsf{OPT}(\gamma_t)) + \mathsf{Gap}(\gamma_t, \mathcal{B}_t) \\ &+ (V(\gamma_t, \sigma[h_t]) - V(\mu^*, \sigma[h_t])) \,, \end{split}$$

where in the final equality, we have used the fact that $\mathsf{OPT}(\gamma_t) - V(\gamma_t, \sigma[h_t]) = \mathsf{Gap}(\gamma_t, \mathcal{B}_t)$. Substituting the preceding expression into (7) yields the lemma statement.

Lemma 2. On the event $\{\mu^* \in \cap_{t \in [T]} \mathcal{B}_t\}$, we have

$$\mathrm{Reg}_{\mathcal{I}}(\mathfrak{Rai},\mu^*,T) \leq \left(\frac{20}{p_0^2D}+1\right)\sum_{t\in[T]}\epsilon_t + \sum_{t\in[T]}\mathbf{E}_{\mu^*}[v(\omega_t,a_t)|h_t] - v(\omega_t,a_t).$$

Proof. In Lemma 1, we obtain the following expression for the regret:

$$\begin{split} \operatorname{Reg}_{\mathcal{I}}(\mathfrak{Rai}, \mu^*, T) &= \sum_{t \in [T]} \left(\operatorname{OPT}(\mu^*) - \operatorname{OPT}(\gamma_t) \right) + \sum_{t \in [T]} \operatorname{Gap}(\gamma_t, \mathcal{B}_t) \\ &+ \sum_{t \in [T]} \left(V(\gamma_t, \sigma[h_t]) - V(\mu^*, \sigma[h_t]) \right) + \sum_{t \in [T]} \mathbf{E}_{\mu^*}[v(\omega_t, a_t) | h_t] - v(\omega_t, a_t). \end{split}$$

Now, from Lemma 4, we obtain

$$\begin{split} \mathsf{OPT}(\mu^*) - \mathsf{OPT}(\gamma_t) & \leq \mathsf{Gap}(\mu^*, \mathsf{B}_1(\mu^*, \|\mu^* - \gamma_t\|_1) + \frac{1}{2} \cdot \|\mu^* - \gamma_t\|_1 \\ & \leq \left(\frac{4}{p_0^2 D}\right) \cdot \|\mu^* - \gamma_t\|_1 + \frac{1}{2} \cdot \|\mu^* - \gamma_t\|_1, \end{split}$$

where the second inequality follows from Proposition 1. Furthermore, from Lemma 5, we have

$$V(\gamma_t, \sigma[h_t]) - V(\mu^*, \sigma[h_t]) \le \frac{1}{2} \|\mu^* - \gamma_t\|_1$$

Taken together, we have

$$\begin{split} \operatorname{Reg}_{\mathcal{I}}(\mathfrak{Rai}, \mu^*, T) &\leq \sum_{t \in [T]} \operatorname{Gap}(\gamma_t, \mathcal{B}_t) + \left(\frac{4}{p_0^2 D} + 1\right) \sum_{t \in [T]} \|\mu^* - \gamma_t\|_1 \\ &+ \sum_{t \in [T]} \mathbf{E}_{\mu^*}[v(\omega_t, a_t) | h_t] - v(\omega_t, a_t). \end{split}$$

Finally, in Lemma 3, we show that on the event $\{\mu^* \in \mathcal{B}_t\}$, we have $\mathsf{Gap}(\gamma_t, \mathcal{B}_t) \leq \left(\frac{16}{p_0^2 D}\right) \epsilon_t$. Thus, on the event $\{\mu^* \in \cap_{t \in [T]} \mathcal{B}_t\}$, we obtain

$$\operatorname{Reg}_{\mathcal{I}}(\mathfrak{Rai}, \mu^*, T) \leq \left(\frac{20}{p_0^2 D} + 1\right) \sum_{t \in [T]} \epsilon_t + \sum_{t \in [T]} \mathbf{E}_{\mu^*}[v(\omega_t, a_t) | h_t] - v(\omega_t, a_t). \quad \Box$$

Lemma 3. For $t \in [T]$, on the event $\{\mu^* \in \mathcal{B}_t\}$, we have

$$\mathsf{Gap}(\gamma_t, \mathcal{B}_t) \leq \left(\frac{16}{p_0^2 D}\right) \epsilon_t.$$

Proof. On the event $\{\mu^* \in \mathcal{B}_t\}$, we have

$$\gamma_t(\omega) \ge \mu(\omega) - \|\gamma_t - \mu^*\|_1$$

$$\ge p_0 - \epsilon_t.$$

Thus, for $\epsilon_t < \frac{p_0}{2}$, we have $\min_{\omega} \gamma_t(\omega) \ge \frac{p_0}{2}$. Using the same argument as in Proposition 1, we then obtain

$$\operatorname{\mathsf{Gap}}(\gamma_t,\mathcal{B}_t) = \operatorname{\mathsf{Gap}}(\gamma_t,\mathsf{B}_1(\gamma_t,\epsilon_t)) \leq \left(\frac{4}{D\min_{\omega}\gamma_t(\omega)^2}\right)\epsilon_t \leq \left(\frac{16}{p_0^2D}\right)\epsilon_t.$$

For $\epsilon_t > p_0/2$, the bound holds trivially since $16\epsilon_t/p_0^2 D > 1$.

Lemma 4. For any $\mu_1, \mu_2 \in \Delta(\Omega)$, we have

$$\mathsf{OPT}(\mu_1) - \mathsf{OPT}(\mu_2) \leq \mathsf{Gap}(\mu_1, \mathsf{B}_1(\mu_1, \|\mu_1 - \mu_2\|_1)) + \frac{1}{2} \cdot \|\mu_1 - \mu_2\|_1.$$

Proof. Fix $\mu_1, \mu_2 \in \Delta(\Omega)$. For $i \in \{1, 2\}$, let $\sigma_i \in \arg\max_{\sigma' \in \mathsf{Pers}(\mu_i)} V(\mu_i, \sigma')$. By definition, we have $\mathsf{OPT}(\mu_i) = V(\mu_i, \sigma_i)$.

Next, among all signaling mechanisms that are persuasive for all $\mu \in \mathsf{B}_1(\mu_1, \|\mu_1 - \mu_2\|_1)$, let σ_3 maximize $V(\mu_1, \sigma)$. Since σ_3 is persuasive for μ_2 , we have $\mathsf{OPT}(\mu_2) = V(\mu_2, \sigma_2) \geq V(\mu_2, \sigma_3)$. Thus, we have

$$\begin{split} \mathsf{OPT}(\mu_1) - \mathsf{OPT}(\mu_2) &= V(\mu_1, \sigma_1) - V(\mu_2, \sigma_2) \\ &\leq V(\mu_1, \sigma_1) - V(\mu_2, \sigma_3) \\ &= V(\mu_1, \sigma_1) - V(\mu_1, \sigma_3) + V(\mu_1, \sigma_3) - V(\mu_2, \sigma_3) \\ &\leq \mathsf{Gap}(\mu_1, \mathsf{B}_1(\mu_1, \|\mu_1 - \mu_2\|_1)) + \frac{1}{2} \cdot \|\mu_1 - \mu_2\|_1. \end{split}$$

Here, the inequality follows from the definition of $Gap(\cdot)$, and from Lemma 5.

Lemma 5. For any $\mu_1, \mu_2 \in \Delta(\Omega)$ and any signaling mechanism σ , we have

$$|V(\mu_1, \sigma) - V(\mu_2, \sigma)| \le \frac{1}{2} \cdot ||\mu_1 - \mu_2||_1.$$

Proof. Fix $\mu_1, \mu_2 \in \Delta(\Omega)$. For any signaling mechanism σ that is persuasive under μ_1 , we have for any $x \in \mathbb{R}$,

$$|V(\mu_1, \sigma) - V(\mu_2, \sigma)| = \left| \sum_{\omega \in \Omega} (\mu_1(\omega) - \mu_2(\omega)) \left(\sum_{a \in A} \sigma(\omega, a) v(\omega, a) - x \right) \right|$$

$$\leq \|\mu_1 - \mu_2\|_1 \cdot \sup_{\omega \in \Omega} \left| \sum_{a \in A} \sigma(\omega, a) v(\omega, a) - x \right|,$$

where we have used the Hölder's inequality in the last line. Optimizing over x and using the fact that the sender's valuations lie in [0,1] yields the result.

A.2 Proof of Theorem 1

Proof of Theorem 1. If $\mu^* \in \mathcal{B}_t$ for each $t \in [T]$, then since $\sigma[h_t]$ is persuasive under all priors in \mathcal{B}_t , we deduce that $\sigma[h_t]$ is persuasive under prior μ^* for all $t \in [T]$. Thus, we obtain that the \mathfrak{Rai} -algorithm is β -persuasive for

$$\beta = \sup_{\mu^* \in \mathcal{B}_0} \mathbf{P}_{\mu^*} \left(\cap_{t \in [T]} \mathcal{B}_t \not\ni \mu^* \right).$$

Now, for any $\mu \in \mathcal{B}_0$, using the union bound we get

$$\mathbf{P}_{\mu} \left(\cap_{t \in [T]} \mathcal{B}_{t} \not\ni \mu \right) = \mathbf{P}_{\mu} \left(\cup_{t \in [T]} \mathcal{B}_{t}^{c} \ni \mu \right)$$

$$\leq \sum_{t \in [T]} \mathbf{P}_{\mu} \left(\mathcal{B}_{t}^{c} \ni \mu \right)$$

$$= \sum_{t \in [T]} \mathbf{P}_{\mu} \left(\| \gamma_{t} - \mu \|_{1} > \epsilon_{t} \right)$$

$$= \sum_{t \in [T]} \mathbf{P}_{\mu} \left(\| \gamma_{t} - \mu \|_{1} > \sqrt{\frac{|\Omega|}{t}} \left(1 + \sqrt{\Phi \log T} \right) \right).$$

For $t < \frac{1}{4}\Phi \log T$, we have

$$\sqrt{\frac{|\Omega|}{t}} \left(1 + \sqrt{\Phi \log T} \right) > 2\sqrt{|\Omega|} \left(1 + \frac{1}{\sqrt{\Phi \log T}} \right) \ge 2.$$

Hence, $\mathbf{P}_{\mu}\left(\|\gamma_t - \mu\|_1 > \sqrt{\frac{|\Omega|}{t}}\left(1 + \sqrt{\Phi \log T}\right)\right) = 0$. On the other hand, for $t \geq \frac{1}{4}\Phi \log T$, we have $\sqrt{\Phi \log T} \leq 2\sqrt{t}$, and hence from Lemma 6, we obtain

$$\sum_{t \ge \frac{1}{4}\Phi \log T} \mathbf{P}_{\mu} \left(\|\gamma_t - \mu\|_1 > \sqrt{\frac{|\Omega|}{t}} \left(1 + \sqrt{\Phi \log T} \right) \right) \le \sum_{t \ge \frac{1}{4}\Phi \log T} \exp \left(-\frac{3\Phi \log T \sqrt{\Omega}}{56} \right)$$

$$\le T^{-\frac{3\Phi\sqrt{\Omega}}{56}} \left(T - \frac{\Phi \log T}{4} \right)$$

$$\le T^{1 - \frac{3\Phi\sqrt{\Omega}}{56}}.$$

Setting $\Phi > 20$ implies the final term is at most $T^{-0.5}$.

B Concentration Inequalities

In this section, we provide some concentration inequalities used in the proofs of our main results. We use the same notation as in the main text. The following lemma provides a bound on the ℓ_1 -norm of the deviation of the empirical distribution from its mean.

Lemma 6. For each $t \in [T]$, and for any $\mu \in \Delta(\Omega)$, we have for all $0 < \Phi_t \le 2\sqrt{t}$,

$$\mathbf{P}_{\mu}\left(\left\|\gamma_{t}-\mu\right\|_{1} \geq \sqrt{\frac{|\Omega|}{t}}\left(1+\Phi_{t}\right)\right) \leq \exp\left(-\frac{3\Phi_{t}^{2}\sqrt{\Omega}}{56}\right)\mathbf{I}\left\{\sqrt{\frac{|\Omega|}{t}}\left(1+\Phi_{t}\right) \leq 2\right\}.$$

Proof. Let $X_t \in \{0,1\}^{|\Omega|}$ denote the random variable with $X_t(\omega) = \mathbf{I}\{\omega_t = \omega\}$, and define $Y_t = X_t - \mathbf{E}_{\mu}[X_t]$. Let $Z_t = \|\sum_{\tau \in [t]} Y_{\tau}\|_1$. Since $\|Y_t\|_1 \leq \|X_t - \mathbf{E}_{\mu}[X_t]\|_1 \leq 2$ for each $t \in [T]$, by Foucart and Rauhut [2013, Corollary 8.46], we obtain for each $t \in [T]$,

$$\mathbf{P}_{\mu}\left(Z_{t} \geq \mathbf{E}_{\mu}[Z_{t}] + s\right) \leq \exp\left(-\frac{3s^{2}}{4\left(6t + 6\mathbf{E}_{\mu}[Z_{t}] + s\right)}\right).$$

Next, letting $Z_{t,\omega} = |\sum_{\tau \in [t]} Y_{\tau}(\omega)|$ for $\omega \in \Omega$, we obtain

$$\mathbf{E}_{\mu}[Z_{t}] = \sum_{\omega \in \Omega} \mathbf{E}_{\mu}[Z_{t,\omega}]$$

$$= \sum_{\omega \in \Omega} \mathbf{E}_{\mu}[\sqrt{Z_{t,\omega}^{2}}]$$

$$\leq \sum_{\omega \in \Omega} \sqrt{\mathbf{E}_{\mu}[Z_{t,\omega}^{2}]}$$

$$= \sum_{\omega \in \Omega} \sqrt{\sum_{\tau \in [t]} \mathbf{Var}_{\mu}[Y_{\tau}(\omega)]}$$

$$= \sqrt{t} \cdot \sum_{\omega \in \Omega} \sqrt{\mu(\omega)(1 - \mu(\omega))}$$

$$\leq \sqrt{|\Omega|t},$$

where the first inequality follows from Jensen's inequality, and the third equality follows from the fact that, since $\mathbf{E}_{\mu}[Y_t(\omega)] = 0$, we have $\mathbf{E}[Z_{t,\omega}^2] = \sum_{\tau \in [t]} \mathbf{Var}_{\mu}[Y_{\tau}(\omega)]$. The final step follows from a straightforward optimization. Thus, we obtain

$$\mathbf{P}_{\mu}\left(Z_{t} \geq \sqrt{|\Omega|t} + s\right) \leq \exp\left(-\frac{3s^{2}}{4\left(6t + 6\sqrt{|\Omega|t} + s\right)}\right).$$

Choosing $s = \Phi_t \sqrt{|\Omega|t}$ for $0 < \Phi_t \le 2\sqrt{t}$, and noting that $Z_t = t \|\gamma_t - \mu\|_1$, we obtain

$$\begin{aligned} \mathbf{P}_{\mu} \left(\| \gamma_t - \mu \|_1 &\geq \sqrt{\frac{|\Omega|}{t}} \left(1 + \Phi_t \right) \right) \leq \exp \left(-\frac{3\Phi_t^2 |\Omega| t}{4 \left(6t + 6\sqrt{|\Omega|t} + \Phi_t \sqrt{|\Omega|t} \right)} \right) \\ &\leq \exp \left(-\frac{3\Phi_t^2 \sqrt{|\Omega|}}{4 \left(12 + \Phi_t / \sqrt{t} \right)} \right) \\ &\leq \exp \left(-\frac{3\Phi_t^2 \sqrt{\Omega}}{56} \right). \end{aligned}$$

The lemma statement then follows after noticing that for all $t \in [T]$, we have $\|\gamma_t - \mu\|_1 \leq \|\gamma_t\|_1 + \|\mu\|_1 \leq 2$.

The following lemma is a standard application of the Azuma-Hoeffding inequality Boucheron et al. [2013], and used to obtain bounds on the regret.

Lemma 7. For all $\mu^* \in \mathcal{B}_0$ and $\alpha > 0$, we have

$$\mathbf{P}_{\mu^*} \left(\sum_{t \in [T]} \mathbf{E}_{\mu^*} [v(\omega_t, a_t) | h_t] - v(\omega_t, a_t) \ge \sqrt{\alpha T \log T} \right) < \frac{1}{T^{2\alpha}}.$$

Proof. For $t \in [T]$, let $X_t \triangleq \mathbf{E}_{\mu^*}[v(\omega_t, a_t)|h_t] - v(\omega_t, a_t)$. Observe that $\mathbf{E}_{\mu^*}[X_t|h_t] = 0$ and $|X_t| \leq 1$. Thus the sequence $\{X_t : t \in [T]\}$ is a bounded martingale difference sequence. Hence, from Azuma-Hoeffding Boucheron et al. [2013], we obtain for $z \geq 0$,

$$\mathbf{P}_{\mu^*} \left(\sum_{t \in [T]} \mathbf{E}_{\mu^*} [v(\omega_t, a_t) | h_t] - v(\omega_t, a_t) \ge z \right) < \exp\left(-\frac{2z^2}{T}\right).$$

Choosing $z = \sqrt{\alpha T \log T}$ for $\alpha \ge 0$ yields the lemma statement.

C Proofs from Section 4

This section provides the proof of our main results in Section 4.

C.1 Proof of Proposition 2

In this section, we provide the proof of Proposition 2.

Proof of Proposition 2. It is straightforward to verify that the following signaling mechanism $\sigma^* \in \text{Pers}(\mu^*)$ optimizes the sender's expected utility among all mechanisms in $\text{Pers}(\mu^*)$:

$$\sigma^*(\omega_0, a_1) = \sigma^*(\omega_0, a_2) = \frac{1}{2},$$

$$\sigma^*(\omega_1, a_1) = \sigma^*(\omega_2, a_2) = \frac{1}{2} + \frac{D}{2(1 - p_0)},$$

$$\sigma^*(\omega_1, a_2) = \sigma^*(\omega_2, a_1) = \frac{1}{2} - \frac{D}{2(1 - p_0)},$$

$$\sigma^*(\omega, a) = 0, \text{ otherwise.}$$

Since the action recommendations are always in $\{a_1, a_2\}$, we obtain $\mathsf{OPT}(\mu^*) = 1$.

By the linearity of obedience constraints and $\mu^* = (\bar{\mu}_1 + \bar{\mu}_2)/2$, it follows that $\mathsf{Pers}(\{\mu^*, \bar{\mu}_1, \bar{\mu}_2\})$ can be obtained by imposing the obedience constraints at priors $\bar{\mu}_1$ and $\bar{\mu}_2$. The optimization problem $\max_{\sigma} \{V(\mu^*, \sigma) : \sigma \in \mathsf{Pers}(\{\bar{\mu}_1, \bar{\mu}_2\})\}$ can be solved to obtain the following optimal signaling mechanism:

$$\hat{\sigma}(\omega_0, a_1) = \hat{\sigma}(\omega_0, a_2) = \frac{1}{2}$$

$$\hat{\sigma}(\omega_1, a_1) = \hat{\sigma}(\omega_2, a_2) = \frac{X}{Z}$$

$$\hat{\sigma}(\omega_1, a_2) = \hat{\sigma}(\omega_2, a_1) = \frac{Y}{Z}$$

$$\hat{\sigma}(\omega_1, a_3) = \hat{\sigma}(\omega_2, a_4) = 1 - \hat{\sigma}(\omega_1, a_1) - \hat{\sigma}(\omega_1, a_2)$$

$$\hat{\sigma}(\omega, a) = 0, \quad \text{otherwise,}$$

where

$$X = 2p_0(1 - p_0 - \epsilon)(1 - p_0 + D)D^2 + p_0(1 - p_0 + \epsilon)(1 - p_0 - D - 2D^2)$$

$$Y = p_0(1 - p_0 - \epsilon)(1 - p_0 - 3D + 2D^2) + 2p_0(1 - p_0 + \epsilon)(1 - p_0 - D)(1 - 2D)D^2$$

$$Z = (1 - p_0 + \epsilon)^2(1 - p_0 - D)(1 - 2D)(1 - p_0 - D - 2D^2)$$

$$- (1 - p_0 - \epsilon)^2(1 - p_0 + D)(1 - p_0 - 3D + 2D^2)$$

The difference in the sender's expected utility between using the optimal persuasive signaling mechanism for prior $\mu^* \in \mathcal{B}$ and using the optimal signaling mechanism that is persuasive for all priors in $\{\mu^*, \bar{\mu}_1, \bar{\mu}_2\}$ is given by

$$\begin{split} \mathsf{Gap}(\mu^*, \mathsf{Pers}(\mu^*, \bar{\mu}_1, \bar{\mu}_2)) &= V(\mu^*, \sigma^*) - V(\mu^*, \hat{\sigma}) \\ &\geq \frac{\epsilon}{2} \frac{1/2 + Dp_0(1 + \epsilon/2 - Dp_0 - D)}{Dp_0 + \epsilon} \\ &\geq \frac{\epsilon}{8Dp_0}. \end{split}$$

C.2 Proof of Theorem 4

Proof of Theorem 4. The proof uses the same instance \mathcal{I}_0 as in the proof of Theorem 3, but uses a different construction for the random events. Specifically, for a prior $\mu \in \mathcal{B}_0$, let $Y_t(\mu)$ indicate

whether the signaling mechanism $\sigma^{\mathfrak{a}}[h_t]$ chosen by the algorithm \mathfrak{a} after history h_t is persuasive for prior μ :

$$Y_t(\mu) \triangleq \mathbf{I} \left\{ \sigma^{\mathfrak{a}}[h_t] \in \mathsf{Pers}(\mu) \right\}.$$

Since \mathfrak{a} is weakly β_T -persuasive, we have $\mathbf{E}_{\mu}[\sum_{t\in[T]}Y_t(\mu)]\geq T(1-\beta_T)$ for all $\mu\in\mathcal{B}_0$. Define the event

$$\mathcal{F}_T(\mu) \triangleq \left\{ h_T : \sum_{t \in [T]} Y_t(\mu) \ge \frac{3T}{4} \right\}.$$

Thus, on the event $\mathcal{F}_T(\mu)$ the signaling mechanism chosen by the algorithm \mathfrak{a} is persuasive for at least 3T/4 time periods. Since $\mathbf{E}_{\mu}[\sum_{t\in[T]}Y_T(\mu)] \leq T\mathbf{P}_{\mu}(\mathcal{F}_T(\mu)) + \frac{3T}{4}(1-\mathbf{P}_{\mu}(\mathcal{F}_T(\mu)))$, we have $\mathbf{P}_{\mu}(\mathcal{F}_T(\mu)) \geq 1 - 4\beta_T$ for all $\mu \in \mathcal{B}_0$.

For $\epsilon \in (0, \frac{1-3p_0}{2})$ to be chosen later, let $\bar{\mu}_0, \bar{\mu}_1, \bar{\mu}_2$ be defined as in the proof of Theorem 3. For each $i \in \{0, 1, 2\}$ and for all $\epsilon \in (0, \frac{1-3p_0}{2})$, we have $\bar{\mu}_i \in \mathcal{B}_0$ and hence $\mathbf{P}_{\mu}(\mathcal{F}_T(\bar{\mu}_i)) \geq 1 - 4\beta_T$.

On the event $\mathcal{F}_T(\bar{\mu}_1) \cap \mathcal{F}_T(\bar{\mu}_2)$, we have

$$\sum_{t \in [T]} Y_t(\bar{\mu}_1) Y_t(\bar{\mu}_2) \ge \sum_{t \in [T]} Y_t(\bar{\mu}_1) + \sum_{t \in [T]} Y_t(\bar{\mu}_1) - T \ge \frac{T}{2},$$

where the first inequality follows from $ab \geq a+b-1$ for $a,b \in [0,1]$. Thus, on at least T/2 time periods, the signaling mechanism $\sigma[h_t]$ is persuasive for both priors $\bar{\mu}_1$ and $\bar{\mu}_2$, and hence also for prior $\bar{\mu}_0$ because of the linearity of the obedience constraints. On each of these time periods, we have $\mathsf{OPT}(\bar{\mu}_0) - V(\bar{\mu}_0, \sigma^{\mathfrak{a}}[h_t]) \geq \mathsf{Gap}(\bar{\mu}_0, \{\bar{\mu}_0, \bar{\mu}_1, \bar{\mu}_2\}) \geq \frac{\epsilon}{8Dp_0}$. Thus, we obtain, on the event $\mathcal{F}_T(\bar{\mu}_1) \cap \mathcal{F}_T(\bar{\mu}_2)$,

$$\sum_{t\in[T]} \left(\mathsf{OPT}(\bar{\mu}_0) - V(\bar{\mu}_0, \sigma^{\mathfrak{a}}[h_t])\right) Y_t(\bar{\mu}_1) Y_t(\bar{\mu}_2) \geq \frac{\epsilon T}{16Dp_0}.$$

Now, by the same argument as in the proof of Theorem 3, using the Pinsker's inequality, we obtain $\mathbf{P}_{\bar{\mu}_0} (\mathcal{F}_T(\bar{\mu}_1) \cap \mathcal{F}_T(\bar{\mu}_2)) \geq 1 - 8\beta_T - \epsilon \sqrt{2T}$. Moreover, since $Y_t(\bar{\mu}_0) \in \{0, 1\}$ is h_t -measurable, by the Azuma-Hoeffding inequality Boucheron et al. [2013], we have

$$\mathbf{P}_{\bar{\mu}_0} \left(\sum_{t \in [T]} \left(V(\bar{\mu}_i, \sigma^{\mathfrak{a}}[h_t]) - v(\omega_t, a_t) \right) Y_t(\bar{\mu}_0) < -\sqrt{T} \right) < e^{-1/2}.$$

Thus, with probability at least $2 - 8\beta_T - \epsilon\sqrt{2T} - e^{-1/2}$, we have

$$\begin{split} \sum_{t \in [T]} \left(\mathsf{OPT}(\bar{\mu}_0) - v(\omega_t, a_t) \right) \right) Y_t(\bar{\mu}_0) &\geq \sum_{t \in [T]} \left(\mathsf{OPT}(\bar{\mu}_0) - V(\bar{\mu}_0, \sigma^{\mathfrak{a}}[h_t]) \right) Y_t(\bar{\mu}_0) - \sqrt{T} \\ &\geq \sum_{t \in [T]} \left(\mathsf{OPT}(\bar{\mu}_0) - V(\bar{\mu}_0, \sigma^{\mathfrak{a}}[h_t]) \right) Y_t(\bar{\mu}_1) Y_t(\bar{\mu}_2) - \sqrt{T} \\ &\geq \frac{\epsilon T}{16Dp_0} - \sqrt{T}. \end{split}$$

For $T \geq T_0 = \frac{1}{(1-3p_0)^2}$, choosing $\epsilon = \frac{1}{2\sqrt{T}} \leq \frac{1-3p_0}{2}$, we obtain with probability at least $\frac{1}{2} - 8\beta_T$,

$$\sum_{t \in [T]} \left(\mathsf{OPT}(\bar{\mu}_0) - v(\omega_t, a_t) \right)) \, \mathbf{I} \{ \sigma^{\mathfrak{a}}[h_t] \in \mathsf{Pers}(\bar{\mu}_0) \} \geq \sqrt{T} \left(\frac{1}{32Dp_0} - 1 \right) \geq \frac{\sqrt{T}}{64Dp_0},$$

for
$$Dp_0 < 1/64$$
.