

# The Price of Incentivizing Exploration: A Characterization via Thompson Sampling and Sample Complexity\*

Mark Sellke<sup>†</sup>  
Stanford University  
msellke@stanford.edu

Aleksandrs Slivkins  
Microsoft Research NYC  
slivkins@microsoft.com

First version: February 2020

This version: June 2022

## Abstract

We consider *incentivized exploration*: a version of multi-armed bandits where the choice of arms is controlled by self-interested agents, and the algorithm can only issue recommendations. The algorithm controls the flow of information, and the information asymmetry can incentivize the agents to explore. Prior work achieves optimal regret rates up to multiplicative factors that become arbitrarily large depending on the Bayesian priors, and scale exponentially in the number of arms. A more basic problem of sampling each arm once runs into similar factors.

We focus on the *price of incentives*: the loss in performance, broadly construed, incurred for the sake of incentive-compatibility. We prove that Thompson Sampling, a standard bandit algorithm, is incentive-compatible if initialized with sufficiently many data points. The performance loss due to incentives is therefore limited to the initial rounds when these data points are collected. The problem is largely reduced to that of sample complexity: how many rounds are needed? We address this question, providing matching upper and lower bounds and instantiating them in various corollaries. Typically, the optimal sample complexity is polynomial in the number of arms and exponential in the “strength of beliefs”.

---

\*An extended abstract of this paper appeared at *ACM EC 2021* (ACM Symp. on Economics and Computation).

This is the *full version* for this extended abstract. Compared to the initial version from Feb’20, the versions since Feb’21 contain several new extensions: Sections 6.3 and 6.5, the lower bounds in Section 6.4, and Section 7. The monotonicity result for Thompson Sampling (Section 4.2) appears since Jun’22.

We are grateful to the anonymous referees of *ACM EC 2021* and *Operations Research* for thoughtful comments and suggestions. We thank Xinyan Hu and Dung Daniel Ngo for comments on the manuscript.

<sup>†</sup>This work was supported by an NSF GRFP and a Stanford Graduate Fellowship.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Related work</b>	<b>7</b>
<b>3</b>	<b>Preliminaries</b>	<b>8</b>
<b>4</b>	<b>Incentivized Exploration via Thompson Sampling</b>	<b>10</b>
4.1	Proofs	12
4.2	Monotonicity of Thompson Sampling	14
<b>5</b>	<b>Collecting Initial Samples</b>	<b>14</b>
<b>6</b>	<b>Sample Complexity of Incentivized Exploration</b>	<b>20</b>
6.1	Lower bound	20
6.2	Polynomially matching upper/lower bounds	21
6.3	Dependence on the number of arms	21
6.4	Canonical priors	22
6.5	One well-known arm	23
<b>7</b>	<b>Extensions via Improved Algorithms</b>	<b>24</b>
<b>8</b>	<b>Explorability Characterization</b>	<b>25</b>
	<b>References</b>	<b>26</b>
<b>A</b>	<b>Tools from Probability</b>	<b>29</b>
A.1	Fortuin-Kasteleyn-Ginibre (FKG) inequality for correlation	29
A.2	Bayesian concentration: proof of Lemma 3.1	29
A.3	Stochastic and MLR Domination	30
A.4	Tails of sub-Gaussian distributions	31
<b>B</b>	<b>Initial Sampling: proofs for Section 5</b>	<b>32</b>
B.1	Existence of a Suitable Policy: proof of Lemma 5.3	32
B.2	A suitable policy for bootstrapping: proof of Lemma 5.4	35
B.3	BIC property for bootstrapping: proof of Lemma 5.5	35
<b>C</b>	<b>Sample Complexity for Arbitrary Priors (for Section 6)</b>	<b>36</b>
C.1	Necessity of the Non-Degeneracy Assumption	38
<b>D</b>	<b>Sample Complexity for Truncated Gaussians and Beta priors</b>	<b>39</b>
<b>E</b>	<b>Extension: A more efficient version of <code>ExponentialExploration</code></b>	<b>45</b>
<b>F</b>	<b>Extension: Improved Algorithm for “Easy” problem Instances</b>	<b>45</b>
<b>G</b>	<b>Extension: Efficient Computation for Beta Priors</b>	<b>47</b>

# 1 Introduction

Consider an online platform where users need to choose between some actions (*e.g.*, products or experiences) of initially unknown quality, and can jointly learn which actions are better. The users collectively face the tradeoff between *exploring* various actions so as to acquire new information, and *exploiting* this information to choose better actions. A benevolent dictator controlling the users would run an algorithm to resolve this tradeoff so as to maximize social welfare. The online platform may wish to coordinate the users in a similar way. However, each user is a self-interested agent making her own choices, and her incentives are heavily skewed in favor of exploitation. This is because she suffers full costs of her exploration, whereas its benefits are spread among many agents in the future. Misaligned incentives can lead to under-exploration, whereby better alternatives are explored very slowly or not at all if they are unappealing initially. These issues are common in online platforms that present recommendations and ratings based on user feedback, which are ubiquitous in a variety of domains: movies, restaurants, products, vacation destinations, etc.

We study *incentivized exploration*: the problem faced by the platform in the scenario described above. The platform can recommend actions, but cannot force agents to follow these recommendations. However, the platform controls the flow of information, and can choose what each agent observes about the past. Revealing full information to each agent works badly: agents fail to explore in a broad range of problem instances (*e.g.*, see Ch. 11, Slivkins, 2019). Information asymmetry, when the platform reveals less than it knows, can incentivize the agents to explore.

A common model for incentivized exploration from (Kremer et al., 2014; Mansour et al., 2015) and the subsequent work is as follows. The population of agents faces a *multi-armed bandit* problem, a basic model of exploration-exploitation tradeoff. A bandit algorithm iteratively recommends actions, a.k.a. *arms*. In each round, a new agent arrives, observes a recommendation, chooses an action, and collects a reward for this action. This reward lies in the interval  $[0, 1]$ , and comes from a fixed but unknown action-specific reward distribution. The reward is observed by the algorithm but not by the other agents. The algorithm does not reveal any information other than the recommended action itself; this is w.l.o.g. under standard economic assumptions of Bayesian rationality. In particular, the arms' mean rewards are drawn from a Bayesian prior which is known to everyone. The algorithm needs to be Bayesian incentive-compatible (*BIC*), *i.e.*, incentivize the agents to follow recommendations. The goal is to design a BIC bandit algorithm so as to optimize its learning performance.

Prior work on BIC bandit algorithms compares their learning performance to that of optimal bandit algorithms, BIC or not. In particular, Mansour et al. (2015, 2020) obtain Bayesian regret  $C_{K,\mathcal{P}} \cdot \sqrt{T}$ , where  $T$  is the time horizon and  $C_{K,\mathcal{P}}$  is determined by the number of arms  $K$  and the Bayesian prior  $\mathcal{P}$ ; this dependence on  $T$  is optimal in the worst case.<sup>1</sup> However,  $C_{K,\mathcal{P}}$  can be arbitrarily large depending on the prior, even for  $K = 2$ , and the dependence on  $K$  is exponential in paradigmatic special cases. For contrast, non-BIC bandit algorithms achieve  $O(\sqrt{KT})$  regret rate uniformly over all priors. Similar issues arise for a more basic variant of incentivized exploration, where one only needs to choose each arm at least once. This variant requires  $C_{K,\mathcal{P}}$  rounds in Mansour et al. (2015, 2020), without any non-trivial lower bounds on the number of rounds, or any way to relate upper and lower bounds to one another.

**Our scope.** We focus on the *price of incentivizing exploration* ( $\text{P}\circ\text{I}\text{E}$ ): the penalty in performance incurred for the sake of the BIC property, such as the  $C_{K,\mathcal{P}}$  factor mentioned above. While several refinements of incentivized exploration have been studied, a more fundamental question of **characterizing the optimal  $\text{P}\circ\text{I}\text{E}$**  is largely open. This question is a unifying framing for our results.

---

<sup>1</sup>Bayesian regret is a standard performance measure for Bayesian bandits (*i.e.*, multi-armed bandits with a Bayesian prior). It is defined as the difference in cumulative reward between the algorithm and the best arm, in expectation over the prior.

While intuitive on a high level, the concept of  $\text{P}\circ\text{I}\text{E}$  is subtle to pin down formally. This is because the “penalty in performance” can be expressed via different performance measures. Among these, we are particularly interested in Bayesian regret and *sample complexity*: essentially, how many rounds are needed to choose each arm. Moreover, the increase in Bayesian regret could be multiplicative and/or additive, and is best measured relative to a particular bandit algorithm.<sup>2</sup>

The question of characterizing the optimal  $\text{P}\circ\text{I}\text{E}$  comes in several flavors. First, what is the optimal dependence on  $K$ , the number of arms? For instance, when is this dependence polynomial as opposed to exponential? Second, what is the optimal dependence on the Bayesian prior? It is unclear what are the right parameters to summarize this dependence, and which properties of the prior make the problem difficult. In fact, it is not even clear if the dependence on the prior is needed. Third, while BIC algorithms in prior work suffered from a multiplicative increase in Bayesian regret, it is desirable to make it additive.

We shed light on these issues, focusing on the canonical case of independent priors. That is, the mean reward of each arm  $i$  is drawn independently from the respective Bayesian prior  $\mathcal{P}_i$ .

**Our results: Thompson Sampling is BIC.** We consider Thompson Sampling (Thompson, 1933), a well-known bandit algorithm. We prove that Thompson Sampling is BIC given a *warm-start*: a known number of samples of each arm, denoted  $N_{\text{TS}}$ . More specifically,  $N_{\text{TS}}$  depends only on the Bayesian prior and the number of arms  $K$ , but not on the time horizon  $T$  or the arms’ mean rewards. Further, Thompson Sampling is BIC *as is* when all arms have the same prior mean reward. Thompson Sampling is the first “natural” bandit algorithm found to be BIC, whereas all BIC bandit algorithms from prior work are custom-designed.

This result has far-reaching implications. Thompson Sampling is widely recognized as a state-of-art algorithm for multi-armed bandits (or very close thereto), in terms of provable guarantees as well as empirical performance. In particular, it achieves the optimal  $O(\sqrt{KT})$  Bayesian regret starting from any prior (or any warm-start). We view it as a “gold standard” for bandits, as far as incentivized exploration is concerned. Therefore, the  $\text{P}\circ\text{I}\text{E}$  reduces to the performance loss due to collecting samples to warm-start Thompson Sampling. In particular, the increase in Bayesian regret is additive rather than multiplicative.

The warm-start size  $N_{\text{TS}}$  is an interesting measure of  $\text{P}\circ\text{I}\text{E}$  in its own right, as the initial samples may be collected exogenously, *e.g.*, purchased at a fixed price per sample. We prove that  $N_{\text{TS}}$  is linear in  $K$  under mild assumptions. Moreover, it can be as low as  $O(\log K)$  for the natural example of Beta priors with bounded parameters. This is a huge improvement over (Mansour et al., 2015, 2020), where some (custom-designed) low-regret algorithms are proved to be BIC given some amount of initial data, but the necessary amount is not upper-bounded in terms of  $K$  and can be (at least) exponential in  $K$  for some examples. The  $O(\log K)$  scaling is particularly appealing if each arm is contributed to the platform by a self-interested party, *e.g.*, it represents a restaurant that wishes to be advertised. Then each arm can be asked to pay an entry fee to subsidise the initial samples, and this fee only needs to scale as  $O(\log K)$ .

Lastly, our analysis of Thompson Sampling implies an important *monotonicity* property: its Bayesian-expected per-round reward is non-decreasing over time. Consequently, its Bayesian simple regret<sup>3</sup> at each round  $t$  is at most  $O(\sqrt{K/t})$ . These results appear new, and may be of independent interest.

**Our results: sample complexity.** We turn to collecting initial samples of each arm, arguably the most basic variant of incentivized exploration. The samples can be used to warm-start Thompson Sampling (or some other bandit algorithm with a similar BIC guarantee), and to estimate the expected rewards. More formally, we consider the following problem, called *BIC  $n$ -sampling*: collect  $n$  samples of each arm by a

---

<sup>2</sup>Formally, fix a near-optimal bandit algorithm  $\mathcal{A}^*$  and let  $\mathbb{E}[R(\cdot)]$  denote Bayesian regret. Given a BIC algorithm  $\mathcal{A}$ , write  $\mathbb{E}[R(\mathcal{A})] = \alpha \cdot \mathbb{E}[R(\mathcal{A}^*)] + \beta$ . Then  $\alpha, \beta$  are, resp., multiplicative and additive increase in Bayesian regret.

<sup>3</sup>Bayesian simple regret is another standard performance measure, defined as the difference in reward at round  $t$  between the algorithm and the best arm, in expectation over the prior.

BIC bandit algorithm, in some number of rounds determined by the prior.<sup>4</sup> We are interested in minimizing this number of rounds; we call it the *n-sample complexity*.

We provide “polynomially matching” upper and lower bounds on the optimal *n*-sample complexity. In particular, we obtain *the first non-trivial lower bound* specific to incentivized exploration, for any variant thereof (as opposed to lower bounds on regret from multi-armed bandits). The matching upper bound, *i.e.*, an algorithm and its analysis, is the most technical part of the paper. We use these bounds to resolve exponential vs. polynomial dependence on the number of arms ( $K$ ) and the strength of beliefs, as expressed by one over the smallest variance  $\sigma^2 = \min_i \text{Var}(\mathcal{P}_i)$ . The common case is that the dependence on  $K$  is polynomial, and the dependence on  $1/\sigma$  is exponential. These are also *upper bounds* for Bayesian regret of BIC *n*-sampling, which are new compared to prior work.<sup>5</sup> Thus, we characterize the additive  $\text{POIE}$  of Thompson Sampling: it is polynomial in  $K$  and exponential in  $1/\sigma$  in terms of the number of rounds, and at most that much in terms of Bayesian Regret.

We emphasize that the *n*-sample complexity is an important performance measure on its own. This is because the platform may have various objectives instead of (or in addition to) Bayesian regret, and the *n*-sample complexity is meaningful for most/all of them. **(i)** The platform may be interested in “frequentist” performance guarantees (ones that hold for each realization of the prior), *e.g.*, as in Mansour et al. (2020). In particular, *n*-sample complexity upper-bounds frequentist regret of BIC *n*-sampling, and may plausibly be a good proxy for it. **(ii)** The platform may be interested in “pure exploration”: predicting the best arm after a given number of rounds. In particular, the optimal 1-sample complexity lower-bounds the number of rounds needed for any non-trivial frequentist guarantee on the prediction quality. **(iii)** The platform’s utilities for the actions may be different from the agents’, *e.g.*, the former may be more forgiving for negative outcomes, and/or incorporate platform’s revenue. Also, the platform may treat the prior as (merely) a *belief* shared by the agents, and optimize in expectation over a different belief. In fact, the platform may wish to optimize with respect to multiple versions of utilities and/or beliefs. Yet, the *n*-sample complexity upper-bounds regret of BIC *n*-sampling with respect to any of them. **(iv)** For two arms, the optimal 1-sample complexity is the smallest number of rounds that guarantees *any* non-trivial exploration almost surely.

**Sample complexity in more detail.** We design a BIC bandit algorithm for collecting  $n$  samples of each arm, called `ExponentialExploration`. We prove that it runs for  $T_{\text{UB}}(n)$  rounds for a given  $n$ , where  $T_{\text{UB}}(n)$  is expressed in terms of the prior. We also provide a lower bound  $T_{\text{LB}}$  on 1-sample complexity: the number of rounds needed to choose each arm even once.<sup>6</sup> This lower bound is polynomially related to  $T_{\text{UB}}(N_{\text{TS}})$  if each arm’s prior has at least a constant variance:  $T_{\text{UB}}(N_{\text{TS}}) \leq T_{\text{LB}}^{O(1)}$ ; recall that  $N_{\text{TS}}$  is from Thompson Sampling. Thus, we characterize the optimal *n*-sample complexity for any  $n \in [1, N_{\text{TS}}]$ , *i.e.*, both for warm-starting Thompson Sampling and for choosing each arm once. Moreover,  $T_{\text{UB}}(N_{\text{TS}})$  upper-bounds optimal additive  $\text{POIE}$  in terms of Bayesian regret.

We study how the optimal *n*-sample complexity,  $n \in [1, N_{\text{TS}}]$ , depends on the number of arms  $K$  and the smallest variance  $\sigma^2 = \min_i \text{Var}(\mathcal{P}_i)$ . To isolate the dependence on  $K$ , we stipulate that the priors come from a fixed collection  $\mathcal{C}$ , and study the worst-case dependence on  $K$  over all such problem instances. We find a curious dichotomy: the dependence on  $K$  is either always polynomial or can be exponential,

---

<sup>4</sup>To appreciate why the number of rounds should be determined by the prior, consider a BIC algorithm that collects  $n$  samples of each arm by some round  $T_0$  that depends on the data. Suppose one runs this algorithm for  $T_0$  rounds and then switches to Thompson Sampling. If  $T_0$  depends on the data, then the combined algorithm is not necessarily BIC, as the timing of the switch could potentially leak information to the agents and alter their incentives.

<sup>5</sup>The exponential dependence on  $1/\sigma$  was previously known as an upper bound, but only for 2 arms (Mansour et al., 2020).

<sup>6</sup>Our lower bound requires Bernoulli rewards. If arbitrary reward values are allowed, even a single random sample could reveal a huge amount of information to the algorithm. For example, its binary expansion could encode the mean reward.

depending on  $\mathcal{C}$ . An improved algorithm for the “easy” case of this dichotomy achieves *linear* dependence on  $K$ , which is the best possible for a fixed  $n$ . Next, we focus on truncated Gaussian priors and Beta priors, two paradigmatic examples for Bayesian inference. We find that the optimal  $n$ -sample complexity is polynomial in  $K$  and exponential in  $\text{poly}(1/\sigma)$ . We conclude that the dependence on the priors cannot be avoided, and that *strong beliefs*, as expressed by low-variance priors, is a key factor.

Finally, we zoom in on the important special case when one arm  $j$  represents a well-known alternative and all other arms are new:  $\text{Var}(\mathcal{P}_j) \ll \text{Var}(\mathcal{P}_i)$  for all arms  $i \neq j$ . Focusing on Beta priors, we prove that the exponential dependence on  $\text{poly}(1/\text{Var}(\mathcal{P}_j))$  can be excluded from  $n$ -sample complexity. We replace it with a similar dependence on the second-smallest variance, and only a polynomial dependence on  $1/\text{Var}(\mathcal{P}_j)$ . This result is particularly clean for  $K = 2$  arms, whereby the sample complexity is driven by the larger variance  $\max(\text{Var}(\mathcal{P}_1), \text{Var}(\mathcal{P}_2))$ . For  $K > 2$  arms, this result is the best possible: we prove that *two* arms with small variance cannot be excluded in a similar fashion.

**Explorability characterization.** An important aspect of the  $\text{POIE}$  is whether all arms are *explorable*: can be sampled at least once by a BIC bandit algorithm. It is easy to construct examples when this is not the case. For instance, if there are two arms and  $\mathbb{E}[\mu_2] < \mu_1$  almost surely, where  $\mu_i$  is the mean reward of arm  $i$ , then arm 2 cannot be explored.<sup>7</sup> To ensure that all arms are explorable, we posit that  $\mathbb{E}[\mu_i] > \mu_j$  with positive probability for all arms  $i \neq j$ . This condition, called *pairwise non-dominance*, suffices for our results. Moreover, we prove that this condition is necessary for exploring all arms, under mild non-degeneracy assumptions. In fact, an arm  $i$  is explorable if and only if it satisfies this condition, and all our results can be restricted to explorable arms. Thus, we provide a full characterization for which arms are explorable. This result complements several partial results from (Mansour et al., 2015, 2022, 2020).<sup>8</sup>

**Our techniques.** Algorithm `ExponentialExploration` extends and amplifies the “hidden exploration” approach from Mansour et al. (2015, 2020), whereby one hides low-probability exploration amidst high-probability exploitation. We prove that exploration has a compounding effect: exploration in the present gives the algorithm more leverage to explore in the future, which allows the exploration probability to increase *exponentially* over time. A simple algorithm design which branches out into “pure exploration” and “pure exploitation” in each round is no longer sufficient to realize these improvements. We introduce a new branch which combines exploration and exploitation so as to guarantee a stronger BIC property for one particular arm that is being explored. The latter property allows the algorithm to offset additional exploration of this arm (and does so more efficiently than “pure exploitation”). The policy for this new branch is defined indirectly, as a maximin solution of a certain zero-sum game. The three branches are interleaved in a somewhat intricate way, to achieve a BIC algorithm with the above-mentioned exponential growth.

Our analysis of Thompson Sampling relies on martingale techniques, the FKG inequality (a correlation inequality from statistical mechanics), and a Bayesian version of Chernoff Bounds which appears new. When all prior means are the same, our analysis zooms in on the covariance between the posterior mean reward of one arm and the event that another arm is posterior-best.

**Further discussion.** We focus on a fundamental model of incentivized exploration that combines standard economic assumptions and a basic model of multi-armed bandits. Conceptually, this is the simplest model in which one can study the  $\text{POIE}$ . Reality can be more complex in a variety of ways, both on the economics side and on the machine learning side (see related work for examples). However, our lower bounds immediately apply to models in which incentivized exploration is more difficult.

<sup>7</sup>In this case, Bayesian regret is  $\Omega(T)$  provided that  $\Pr[\mu_2 - \mu_1 > 0] > 0$ .

<sup>8</sup>Specifically, a full characterization for  $K = 2$  arms, some sufficient conditions for  $K > 2$ , and an algorithm that explores all “explorable” arms but does not yield any explicit conditions. Unlike ours, these results extend to correlated priors.



The standard economic assumptions mentioned above include common priors, agents’ rationality and platform’s commitment power. They are shared by all prior work on incentivized exploration (with a notable exception of Immorlica et al. (2020)). Likewise, we assume that rewards are observed by the algorithm. Eliciting informative signals from the agents (*e.g.*, via reviews on an online platform such as Yelp or Amazon) is an important problem that is beyond our scope.

While we do not optimize absolute constants in the performance guarantees, several aspects of our results have practical appeal. We justify the usage of Thompson Sampling in incentivized exploration, reduce the problem to collecting initial data, and calibrate expectations for how much data is needed. An informal takeaway is that Thompson Sampling with a moderately-sized batch of initial data should be BIC. Our results on  $n$ -sample complexity feature exponential improvements in the dependence on the number of arms and (for the scenario with “one known arm”) on the strength of beliefs.

**Open questions.** The most immediate questions concern the dependence on the strength of agents’ beliefs, as expressed by the smallest variance  $\sigma^2 = \min_i \text{Var}(\mathcal{P}_i)$ . The first question is about the warm-start size  $N_{\text{TS}}$  for Thompson Sampling. Can it be made polynomial in  $1/\sigma$ ? While it scales exponentially in  $\text{poly}(1/\sigma)$  in our result, we do not have any lower bounds. The second question is about Bayesian regret for collecting 1 sample of each arm. Can it be made polynomial in  $1/\sigma$ ? We upper-bound it by 1-sample complexity, which in turn is lower-bounded by  $\exp(\text{poly}(1/\sigma))$ . However, Bayesian regret of `ExponentialExploration` is unclear. Our lower bound on 1-sample complexity does not appear to have any bearing on Bayesian regret, so even a *constant* dependence on  $\sigma$  is not ruled out.

Thompson Sampling as a technique applies far beyond the basic version of multi-armed bandits, and our results suggest it as a promising approach for more general models of incentivized exploration. One may hope to handle correlated priors and large-but-tractable bandit problems such as linear bandits. Likewise, one would like to extend our sample-complexity results to such problems.

Going back to independent priors, it is interesting whether other “natural” bandit algorithms can be proved BIC given enough initial data. Our proof techniques are heavily tailored to Thompson Sampling. However, proving that such a result is impossible for a particular algorithm appears quite challenging, too.

**Map of the paper.** First, we analyze BIC properties of Thompson Sampling (Section 4). Then we present and analyze `ExponentialExploration`, our algorithm for collecting initial samples (Section 5). Next, we investigate the sample complexity of incentivized exploration (Section 6): we derive a general lower bound, and mine the upper and lower bounds for the corollaries discussed above. Section 7 contains improved guarantees via fine-tuned versions of `ExponentialExploration`. Explorability characterization can be found in Section 8. Various details are deferred to the appendices.

## 2 Related work

Our model of incentivized exploration was introduced in Kremer et al. (2014), who obtain an optimal policy for the special case of two arms and deterministic rewards.<sup>9</sup> Mansour et al. (2015, 2020) consider the general case of stochastic multi-armed bandits and design BIC bandit algorithms with near-optimal regret rates, up to multiplicative factors that depend on the prior and the number of arms.<sup>10</sup> Further, they obtain a general reduction from bandit algorithms to incentive-compatible ones, and derive several extensions: to correlated priors, partially known priors, contextual bandits, and bandits with auxiliary feedback. They also suggest a

---

<sup>9</sup>The study of incentivized exploration, broadly construed, was initiated in Kremer et al. (2014); Che and Hörner (2018).

<sup>10</sup>In fact, Mansour et al. (2015, 2020) obtain several results of this form: both for Bayesian regret and standard (frequentist) notion of regret, and both in the worst case and for each problem instance.

connection to participation incentives in medical trials. Mansour et al. (2022) focus on exploring all arms than can possibly be explored, and allow for misaligned incentives when the algorithm’s reward is different from the agents’ utility. Several other extensions were considered, under various simplifying assumptions: to repeated games (Mansour et al., 2022), heterogenous agents (Immorlica et al., 2019), social networks (Bahar et al., 2016, 2019), and relaxed economic assumptions (Immorlica et al., 2020). Several related but technically different models have been studied: with time-discounted utilities (Bimpikis et al., 2018), monetary incentives (Frazier et al., 2014; Chen et al., 2018), and continuous information flow (Che and Hörner, 2018). A textbook-like introduction to this line of work can be found in Slivkins (Ch. 11, 2019).

Incentivized exploration is closely related to two important recent topics in theoretical economics. Bayesian Persuasion (*e.g.*, Bergemann and Morris, 2019; Kamenica, 2019) studies, essentially, a single round of our model, where the “principal” uses information asymmetry to persuade the agent to take particular actions. Social learning (*e.g.*, Hörner and Skrzypacz, 2017; Golub and Sadler, 2016) studies various scenarios in which multiple strategic agents interact and learn over time in a shared environment.

Exploration-exploitation problems with incentives issues arise in several other domains, such as dynamic pricing, auction design, and human computation. These problems substantially differ from one another (and from incentivized exploration), depending on who are the self-interested agents and what they control. A review of this literature can be found in Ch. 11.7 of Slivkins (2019).

Exploration-exploitation tradeoff and multi-armed bandits received a huge amount of attention over the past few decades. The diverse and evolving body of research has been summarized in several books: Cesa-Bianchi and Lugosi (2006), Bubeck and Cesa-Bianchi (2012), Gittins et al. (2011), Slivkins (2019), and Lattimore and Szepesvári (2020). Stochastic  $K$ -armed bandits (Lai and Robbins, 1985; Auer et al., 2002) is a canonical “basic” version of the problem, by now it is very well understood.

Thompson Sampling (Thompson, 1933) is a well-known bandit algorithm with much recent progress, see Russo et al. (2018) for background. Most relevantly, it enjoys Bayesian regret bounds which are optimal in the worst case (Russo and Van Roy, 2014; Bubeck and Liu, 2013) and improve for some “nice” priors (Russo and Van Roy, 2014). Also, it attains optimal “frequentist” regret bounds if initialized with some simple priors (Agrawal and Goyal, 2012; Kaufmann et al., 2012; Agrawal and Goyal, 2013).

### 3 Preliminaries

**Problem formulation: incentivized exploration.** There are  $T$  rounds and  $K$  actions, a.k.a. *arms*. In each round  $t \in [T]$ , an algorithm (a.k.a. the *planner*) interacts with a new agent according to the following protocol. The algorithm recommends an arm  $A_t$ , the agent observes the recommendation (and nothing else) and chooses an arm  $A'_t$  (not necessarily the same). The agent collects reward  $r_t \in [0, 1]$  for the chosen action, which is observed by the algorithm, but not by the other agents. The reward of each arm  $i$  is drawn independently from some fixed distribution with mean  $\mu_i \in [0, 1]$ . The reward distributions are initially not known to anybody. If agents always follow recommendations, *i.e.*, if  $A'_t = A_t$  in all rounds  $t$ , then the problem reduces to (Bayesian, stochastic) multi-armed bandits.

We posit Bernoulli rewards, *i.e.*,  $r_t \in \{0, 1\}$  for all rounds  $t$ . This assumption is without loss of generality for all algorithmic results (*i.e.*, all results except the lower bounds). Essentially, this is because one can replace a reward  $r \in [0, 1]$  by a randomized Bernoulli reward with the same expectation.<sup>11</sup>

Let us specify Bayesian priors and incentives. For each arm  $i$ , the mean reward  $\mu_i$  is independently drawn from prior  $\mathcal{P}_i$ . (The joint prior is therefore  $\mathcal{P}_1 \times \dots \times \mathcal{P}_K$ .) The priors are known to all agents and

---

<sup>11</sup>The same trick applies to rewards of larger magnitude after re-scaling them to lie in  $[0, 1]$ .



the algorithm. We require the algorithm to be *Bayesian-incentive compatible (BIC)*: following recommendations is in the agents' best interest. Formally, we condition on the event that recommendations have been followed in the past,  $\mathcal{E}_{t-1} = \{A_s = A'_s : s \in [t-1]\}$ . The BIC condition is as follows: for all rounds  $t$ ,

$$\mathbb{E}[\mu_i - \mu_j \mid A_t = i, \mathcal{E}_{t-1}] \geq 0 \quad \text{all arms } i, j \text{ such that } \Pr[A_t = i] > 0. \quad (3.1)$$

If an algorithm is BIC, we assume that the agents actually follow recommendations.

As a stepstone towards BIC, we use a more restricted condition: a fixed subset  $S$  of rounds is called BIC if (3.1) holds for all rounds  $t \in S$ . In this definition, we still require that recommendations at rounds  $t \notin S$  are followed, as per event  $\mathcal{E}_{t-1}$ , even if these rounds are not necessarily BIC.

One could consider a more general version of the problem, in which the algorithm can reveal an arbitrary message  $\sigma_t$  in each round  $t$  and does not need to be BIC; then each agent  $t$  chooses an arm  $i$  which maximizes  $\mathbb{E}[\mu_i \mid \sigma_t]$ . However, it is easy to show that restricting to BIC algorithms that (only) recommend arms is w.l.o.g. (Kremer et al., 2014), by a suitable version of the *revelation principle*. We also remark that agents in realistic situations are likely to not know exactly which round they arrive in. However, our BIC condition easily extends if agents instead have beliefs over their arrival times.

We posit a condition called *pairwise non-dominance*: for each arm  $i$ ,

$$\Pr[\mu_j < \mathbb{E}[\mu_i]] > 0 \quad \text{for all arms } j \neq i. \quad (3.2)$$

This condition is w.l.o.g.: essentially, each arm  $i$  is explorable if and only if it satisfies (3.2), see Section 8.

**Recommendation policies.** A *recommendation policy*  $\pi$  is a function that inputs a random signal  $\mathcal{S}$  and outputs an arm. More formally, let *signal*  $\mathcal{S}$  be a random variable (taking values in some abstract set) that is jointly distributed with the mean rewards, in the sense that the tuple  $(\mathcal{S}; \mu_1, \dots, \mu_K)$  comes from some joint distribution. A recommendation policy  $\pi$  given signal  $\mathcal{S}$  is a mapping from  $\text{support}(\mathcal{S})$  to arms. In particular, one round of a bandit algorithm can be interpreted a recommendation policy, with signal  $\mathcal{S}$  being the algorithm's history up to this round. Likewise, if an algorithm invokes a recommendation policy, then (unless specified otherwise) the policy receives the algorithm's current history as a signal.

A natural version of the BIC property considers random variable  $\pi(\mathcal{S})$  and posits that

$$\mathbb{E}[\mu_i - \mu_j \mid \pi(\mathcal{S}) = i] \geq 0 \quad \text{all arms } i, j \text{ such that } \Pr[\pi(\mathcal{S}) = i] > 0. \quad (3.3)$$

If (3.3) holds, we say that policy  $\pi$  is *BIC given signal*  $\mathcal{S}$ .

**Conventions.** We index arms by  $i, j, k \in [K]$ . We refer to them as “arm  $i$ ” or “arm  $a_i$ ” interchangeably.

Let  $\mu_i^0 = \mathbb{E}[\mu_i]$  be the prior mean reward of each arm  $i$ . W.l.o.g., we order arms by their prior mean rewards:  $\mu_1^0 \geq \mu_2^0 \geq \dots \geq \mu_K^0$ . Let  $A^* = \min(\text{argmax}_j \mu_j)$  be the best arm, with this specific tie-breaking.  $\mathcal{F}_t$  denotes the filtration generated by the chosen actions and the realized rewards up to (and not including) a given round  $t$ . We set  $\mathbb{E}^t[\cdot] = \mathbb{E}[\cdot \mid \mathcal{F}_t]$  and  $\Pr^t[\cdot] = \Pr[\cdot \mid \mathcal{F}_t]$ , as a shorthand. Sometimes we condition only on the first  $N_j$  samples of each arm  $j$ , for some fixed  $N_j$ . Such a  $\sigma$ -algebra is called *static* and denoted by  $\mathcal{G}_{(N_1, \dots, N_K)}$ . In the special case when we condition on the first  $N$  samples of arms  $1, \dots, j$ , the  $\sigma$ -algebra is denoted by  $\mathcal{G}_{N,j}$ .

The set of all distributions over arms  $1, \dots, k$  is denoted by  $\Delta_k$ . If  $q = (q_1, \dots, q_K)$  is a distribution over arms, the corresponding mean reward is  $\mu_q := \sum_i \mu_i q_i$ .

**Thompson Sampling.** The core concept in Thompson Sampling is sampling from a Bayesian posterior. Given a random quantity  $X$  determined by the mean rewards  $(\mu_1, \dots, \mu_K)$ , the Bayesian posterior at round  $t$  is the conditional distribution of  $X$  given  $\mathcal{F}_t$ . A *posterior sample* of  $X$  at round  $t$  (equivalently:

given  $\mathcal{F}_t$ ) is an independent random draw from this distribution. Thompson Sampling is a very simple bandit algorithm: in each round  $t$ , the chosen arm  $A_t$  is a posterior sample for the best arm  $A^*$ . In particular,

$$\Pr[A_t = i \mid \mathcal{F}_t] = \Pr[A^* = i \mid \mathcal{F}_t] \quad \text{for each arm } i. \quad (3.4)$$

The algorithm is computationally efficient in various special cases, *e.g.*, for independent Beta priors and Bernoulli rewards, and for independent Gaussian priors (truncated or not) and Gaussian rewards.

**Tools.** We make use of FKG inequality, a correlation inequality which says that increasing functions of independent random variables are non-negatively correlated. We state it in Appendix A.1.

We also use a Bayesian concentration inequality. For a given arm  $i$ , it relates the posterior mean reward and an independent draw from the posterior distribution on  $\mu_i$ . We prove that both quantities are within  $1/\sqrt{\#\text{samples}}$  of  $\mu_i$ . While reminiscent of Chernoff Bounds, which compare  $\mu_i$  to the sample average, this result appears new. We state it below, and prove it in Appendix A.2.

**Lemma 3.1** (Bayesian Chernoff Bound). *Fix round  $t$  and parameters  $\varepsilon, r > 0$ . Suppose  $\mathcal{F}_t$  almost surely contains at least  $\varepsilon^{-2}$  samples of a given arm  $i$ . Let  $\hat{\mu}_i$  be a posterior sample for the mean reward  $\mu_i$ . Then for some universal absolute constant  $C$  we have:*

$$\Pr[|\hat{\mu}_i - \mu_i| \geq r\varepsilon] \leq C e^{-r^2/C}; \quad (3.5)$$

$$\Pr[|\mathbb{E}[\mu_i \mid \mathcal{F}_t] - \mu_i| \geq r\varepsilon] \leq C e^{-r^2/C}. \quad (3.6)$$

More generally, let  $q = (q_1, \dots, q_K)$  be a distribution over arms. Let  $\mu_q = \sum_i \mu_i q_i$  and  $\hat{\mu}_q = \sum_i \hat{\mu}_i q_i$  be the corresponding mean reward and posterior sample. Suppose  $\mathcal{F}_t$  almost surely contains at least  $\varepsilon^{-2}$  samples of each arm  $i$  with  $q_i \neq 0$ . Then (3.5) holds with  $i$  replaced with  $q$ .

## 4 Incentivized Exploration via Thompson Sampling

We prove that Thompson Sampling is BIC if initialized with enough samples of each arm.

**Theorem 4.1.** *Let  $\text{ALG}$  be a BIC bandit algorithm such that by some fixed time  $T_0$  it almost surely collects at least  $N_{\text{TS}} = C_{\text{TS}} \varepsilon_{\text{TS}}^{-2} \log \delta_{\text{TS}}^{-1}$  samples from each arm, for a large enough absolute constant  $C_{\text{TS}}$ , where*

$$\varepsilon_{\text{TS}} = \min_{i,j \in [K]} \mathbb{E}[(\mu_i - \mu_j)_+] \quad \text{and} \quad \delta_{\text{TS}} = \min_{i \in [K]} \Pr[A^* = i]. \quad (4.1)$$

*Then running  $\text{ALG}$  for  $T_0$  rounds followed by Thompson sampling is BIC.*

**Remark 4.2.** It is essential for the BIC property that the “switching time”  $T_0$  in Theorem 4.1 is fixed in advance. In particular, switching to Thompson sampling as soon as  $\text{ALG}$  collects enough samples could leak information and destroy the BIC property. For example, suppose  $\text{ALG}$  has the following property: if arm 2 is good, then  $\text{ALG}$  w.h.p. spends a long time exploring this arm, and then does not play arm 2 during some fixed time interval. Then if arm 2 is recommended by Thompson sampling during the latter time interval, the agent will recognize that arm 2 must be bad, and refuse to play it.

**Remark 4.3.** The BIC property in Theorem 4.1 can be made to hold with a quantifiable margin: specifically, we can ensure that the right-hand side in Eq. (3.1) is  $\varepsilon_{\text{TS}}/2$ .

**Remark 4.4.** Our analysis of Thompson Sampling is oblivious to where the warm-up data is coming from. In particular, the data can be collected by a non-BIC algorithm, and agents' participation may be secured via other means, *e.g.*, monetary payments. However, one needs to ensure that Bayesian update on the warm-up data does not depend on the algorithm used to collect it. One could achieve this by reporting the full history of the data-collection algorithm, or, *e.g.*, only reporting the first  $N_{\text{TS}}$  samples of each arm.

Let us investigate how  $N_{\text{TS}}$ , the sample count from Theorem 4.1, depends on  $K$ , the number of arms. We find that  $N_{\text{TS}} = O_{\mathcal{C}}(K)$  if all priors belong to a finite family  $\mathcal{C}$ , and  $N_{\text{TS}} = O_{\mathcal{C}}(\log K)$  if  $\mathcal{C}$  consists of all Beta priors of bounded strength of beliefs.

**Corollary 4.5.** *Suppose all priors  $\mathcal{P}_i$ ,  $i \in [K]$  come from some fixed, finite collection  $\mathcal{C}$  of priors which satisfy the pairwise non-dominance condition (3.2). Then  $N_{\text{TS}} = O_{\mathcal{C}}(K)$ .*

*Proof.* In fact we have  $N_{\text{TS}} = O(K \varepsilon_{\text{TS}}^{-2} \log \varepsilon_{\text{TS}}^{-1})$ . Indeed for each arm  $i$ ,

$$\Pr[A^* = i] \geq \prod_j \Pr[\mu_i \geq \mu_j] \geq \varepsilon_{\text{TS}}^K.$$

The latter inequality holds simply because  $\Pr[\mu_i \geq \mu_j] \geq \mathbb{E}[(\mu_i - \mu_j)_+]$  for all  $i, j$ . Let  $\varepsilon_{\mathcal{C}}$  be the version of  $\varepsilon_{\text{TS}}$  where the min is over all ordered pairs of (not necessary distinct) priors in  $\mathcal{C}$ . Note that  $\varepsilon_{\text{TS}} \geq \varepsilon_{\mathcal{C}}$ . Since  $\mathcal{C}$  is finite and satisfies pairwise non-dominance,  $\varepsilon_{\mathcal{C}}$  is strictly positive. It remains to show the first part of the inequality above.

We proceed via the FKG inequality. Define the indicator functions  $I_j = 1_{\mu_i \geq \mu_j}$  for each  $j \in [K]$ ; we interpret them as functions of  $(\mu_1, \dots, \mu_K)$ . The functions  $I_j$  are each increasing in  $\mu_i$  and decreasing in  $\mu_\ell$  for all  $\ell \neq i$  (including  $\ell = j$  when  $j \neq i$ ). As the values  $\{\mu_\ell\}_{\ell \in [K]}$  are independent by assumption, the mixed-monotonicity FKG inequality (see Remark A.2) implies that the indicator functions  $I_j$  are non-negatively correlated. In fact, each product  $\prod_{j=1}^{\ell} I_j$  satisfies the same monotonicity properties, so repeated application of mixed-monotonicity FKG implies

$$\begin{aligned} \Pr[A^* = i] &= \mathbb{E} \left[ \prod_{j=1}^K I_j \right] \\ &\geq \mathbb{E} \left[ \prod_{j=1}^{K-1} I_j \right] \cdot \mathbb{E}[I_K] \geq \mathbb{E} \left[ \prod_{j=1}^{K-2} I_j \right] \cdot \mathbb{E}[I_{K-1}] \cdot \mathbb{E}[I_K] \geq \dots \\ &\geq \prod_{j=1}^K \mathbb{E}[I_j] = \prod_j \Pr[\mu_i \geq \mu_j]. \quad \square \end{aligned}$$

**Corollary 4.6.** *Suppose each prior  $\mathcal{P}_i$ ,  $i \in [K]$  is a  $\text{Beta}(\alpha_i, \beta_i)$  distribution with parameters  $\alpha_i, \beta_i \in [1, M]$ , for some fixed  $M$ . Then  $N_{\text{TS}} = O_M(\log K)$ .*

*Proof.* Note that  $\varepsilon_{\text{TS}} = \Omega(9^{-M})$ . This immediately follows from definition of  $\varepsilon_{\text{TS}}$  in Eq. (4.1), because  $\Pr[\mu_i > 2/3] \geq 3^{-M}$  and  $\Pr[\mu_j < 1/3] \geq 3^{-M}$  for all arms  $i, j$ .

To handle  $\delta_{\text{TS}}$ , let  $Q_p[\mathcal{P}_i]$  be the top  $p$ -th quantile of distribution  $\mathcal{P}_i$ . Suppose for some  $\eta > 0$

$$Q_\eta[\mathcal{P}_i] \geq Q_{1/K}[\mathcal{P}_j] \quad \text{for all arms } i, j. \quad (4.2)$$

Then  $\delta_{\text{TS}} \geq \Pr[A^* = a_i] \geq \eta (1 - 1/K)^{K-1} \geq \Omega(\eta)$  for all arms  $i$ , so  $N_{\text{TS}} = O_M(\log 1/\eta)$ .

To complete the proof, we claim that (4.2) holds with  $\eta = (MK)^{-M}$ . This is because for each arm  $i$  we have  $(MK)^{-M} \leq \Pr \left[ \mu_i > 1 - \frac{1}{MK} \right] \leq \frac{1}{K}$ . To verify the last statement, it suffices to focus on the extremal cases  $\text{Beta}(1, M)$  and  $\text{Beta}(M, 1)$ .  $\square$

Moreover, we prove that Thompson sampling is BIC *as is* if all prior mean rewards are the same.

**Theorem 4.7.** *If  $\mu_1^0 = \mu_2^0 = \dots = \mu_K^0$  then Thompson sampling is BIC.*

## 4.1 Proofs

We prove Theorems 4.1 and 4.7. First, we note that for any algorithm and any arms  $i, j$  it holds that

$$\begin{aligned} \mathbb{E}[\mu_i \mid A_t = j] \cdot \Pr[A_t = j] &= \mathbb{E}[\mu_i \cdot \mathbf{1}_{\{A_t=j\}}] \\ &= \mathbb{E}[\mathbb{E}^t[\mu_i \cdot \mathbf{1}_{\{A_t=j\}}]] = \mathbb{E}[\mathbb{E}^t[\mu_i] \cdot \mathbb{E}^t[\mathbf{1}_{\{A_t=j\}}]] \\ &= \mathbb{E}[\mathbb{E}^t[\mu_i] \cdot \Pr^t[A_t = j]]. \end{aligned} \quad (4.3)$$

$$\mathbb{E}[\mu_i - \mu_j \mid A_t = j] \cdot \Pr[A_t = j] = \mathbb{E}[\mathbb{E}^t[\mu_i - \mu_j] \cdot \Pr^t[A_t = j]]. \quad (4.4)$$

(Eq. (4.4) follows by taking a version (4.3) with  $i = j$ , and subtracting it from (4.3).)

Next, we analyze the object inside the expectation in (4.3).

**Lemma 4.8.** *Fix arms  $i, j$ . Let  $H_t := \mathbb{E}^t[\mu_i] \cdot \Pr^t[A^* = j]$ . For any algorithm, the sequence  $(H_1, \dots, H_T)$  is a supermartingale if  $i \neq j$  and a submartingale if  $i = j$ .*

*Proof.* Note that  $(\mathbb{E}^t[\mu_i] : t \in [T])$  and  $(\Pr^t[A^* = j] : t \in [T])$  are martingales by definition.

We consider two cases, depending whether  $A_t = i$ . Recall that an expression such as  $\mathbb{E}^{t+1}[\mu_i]$  is a random variable (with randomness coming from  $\mathcal{F}_{t+1}$ ), and event such as  $A_t = i$  restricts this random variable. First, suppose  $A_t \neq i$ . Then  $\mathbb{E}^{t+1}[\mu_i] = \mathbb{E}^t[\mu_i]$  almost surely, therefore  $H_t$  has expected change 0 on this step since  $\Pr^t[A^* = j]$  is a martingale.

Next, suppose  $A_t = i$ . The crucial claim is that  $\mathbb{E}^{t+1}[\mu_i]$  and  $\Pr^{t+1}[A^* = i]$  are increasing in the time- $t$  reward while  $\Pr^{t+1}[A^* = j]$  is decreasing in the time- $t$  reward. Indeed, Corollary A.9 and Lemma A.7 in the Appendix imply that the conditional distribution of  $\mu_i$  is stochastically increasing in the time- $t$  reward. Observing that the event  $\{A^* = j\}$  is decreasing in the value of  $\mu_i$  now implies the claim. Note that this argument crucially uses both the Bernoulli reward assumption and the fact that we have bandit feedback and independent arms.

Next, recall the FKG inequality (see Appendix A): if  $f, g$  are increasing functions of the same variable then they are positively correlated, i.e.  $\mathbb{E}[fg] \geq \mathbb{E}[f]\mathbb{E}[g]$ . We set  $f = \mathbb{E}^{t+1}[\mu_i]$  and  $g = \Pr^{t+1}[A^* = i]$  and apply FKG conditionally on  $\mathcal{F}_t$  and  $A_t$ , interpreting both  $f$  and  $g$  as functions of the time- $t$  reward. Then

$$\begin{aligned} \mathbb{E}^t[H_{t+1}] &= \mathbb{E}^t[\mathbb{E}^{t+1}[\mu_i] \cdot \Pr^{t+1}[A^* = i]] \\ &\geq \mathbb{E}^{t+1}[\mu_i] \cdot \Pr^{t+1}[A^* = i] \\ &= H_t. \end{aligned}$$

We have just shown that  $H_t$  is a submartingale when  $i = j$ . Similarly, when  $i \neq j$  we apply the FKG inequality to  $f = \mathbb{E}^{t+1}[\mu_i]$  and  $g = \Pr^{t+1}[A^* = j]$ , again conditionally on  $\mathcal{F}_t$  and  $A_t$ . In this case  $g$  is a decreasing function of the time- $t$  reward and so the FKG inequality goes in the opposite direction, stating that  $\mathbb{E}[fg] \leq \mathbb{E}[f]\mathbb{E}[g]$ . We hence obtain

$$\begin{aligned} \mathbb{E}^t[H_{t+1}] &= \mathbb{E}^t[\mathbb{E}^{t+1}[\mu_i] \cdot \Pr^{t+1}[A^* = j]] \\ &\leq \mathbb{E}^{t+1}[\mu_i] \cdot \Pr^{t+1}[A^* = j] \\ &= H_t. \end{aligned} \quad \square$$

The following lemma is essentially an inductive step. It implies Theorem 4.7 by induction on  $t$ , because the premise in the lemma holds trivially when  $t = 0$  and all prior mean rewards are the same.

**Lemma 4.9.** *Let ALG be any bandit algorithm. Fix round  $t$ . Suppose that running ALG for  $t - 1$  steps, followed by Thompson sampling at time  $t$ , is BIC at time  $t$ . Then running ALG for  $t$  steps, followed by Thompson sampling at time  $t + 1$ , is BIC at time  $t + 1$ .*

*Proof.* Thompson sampling is BIC at time  $t$  if and only if  $\mathbb{E}[\mu_i - \mu_j \mid A_t = i] \geq 0$  for all  $(i, j)$ .

$$\begin{aligned} \mathbb{E}[\mu_i - \mu_j \mid A_t = i] &= \frac{\mathbb{E}[\mathbb{E}^t[\mu_i - \mu_j] \cdot \Pr^t[A_t = i]]}{\Pr[A_t = i]} && \text{(by Eq. (4.4))} \\ &= \frac{\mathbb{E}[\mathbb{E}^t[\mu_i - \mu_j] \cdot \Pr^t[A^* = i]]}{\Pr[A^* = i]} && \text{(by Eq. (3.4)).} \end{aligned}$$

In the numerator,  $\mathbb{E}^t[\mu_i - \mu_j] \cdot \Pr^t[A^* = i]$  is a submartingale by Lemma 4.8. In particular its expectation is non-decreasing in  $t$ . On the other hand the denominator  $\Pr[A^* = i]$  is a constant independent of  $t$ .  $\square$

*Proof of Theorem 4.1.* Fix arms  $i, j$  and set  $\delta_i := \Pr[A^* = i] \geq \delta_{\text{TS}}$ . Then

$$\begin{aligned} \mathbb{E}[\mu_i - \mu_j \mid A_{T_0} = i] \cdot \Pr[A_{T_0} = i] &= \mathbb{E}[\mathbb{E}^{T_0}[\mu_i - \mu_j] \cdot \Pr^{T_0}[A_{T_0} = i]] && \text{(by Eq. (4.4))} \\ &= \mathbb{E}[\mathbb{E}^{T_0}[\mu_i - \mu_j] \cdot \Pr^{T_0}[A^* = i]] \\ &= \mathbb{E}[\mathbb{E}^{T_0}[\mathbb{E}^{T_0}[\mu_i - \mu_j] \cdot \mathbf{1}_{\{A^* = i\}}]] \\ &= \mathbb{E}[\mathbb{E}^{T_0}[\mu_i - \mu_j] \cdot \mathbf{1}_{\{A^* = i\}}]. \end{aligned} \tag{4.5}$$

To establish that Thompson Sampling is BIC we prove  $\mathbb{E}[\mathbb{E}^{T_0}[\mu_i - \mu_j] \cdot \mathbf{1}_{\{A^* = i\}}] \geq 0$ . Since the functions  $(\mu_i - \mu_j)_+, \mathbf{1}_{\{A^* = i\}}$  are increasing in  $\mu_i$  and decreasing in  $\mu_k$  for each  $k \neq i$ , the FKG inequality implies

$$\mathbb{E}[(\mu_i - \mu_j) \cdot \mathbf{1}_{\{A^* = i\}}] = \mathbb{E}[(\mu_i - \mu_j)_+ \cdot \mathbf{1}_{\{A^* = i\}}] \geq \varepsilon_{\text{TS}} \delta_i,$$

see Remark A.2. If our estimates  $\mathbb{E}^{T_0}[\mu_i], \mathbb{E}^{T_0}[\mu_j]$  of  $\mu_i, \mu_j$  were exactly correct then we could immediately conclude. Inspired by this, we show the expected absolute error in estimating  $\mathbb{E}^{T_0}[\mu_i - \mu_j] \cdot \mathbf{1}_{\{A^* = i\}}$  by  $(\mu_i - \mu_j) \cdot \mathbf{1}_{\{A^* = i\}}$  is upper bounded by  $\varepsilon_{\text{TS}} \delta_i$ . In other words we aim to show:

$$\mathbb{E}[|\mathbb{E}^{T_0}[\mu_i - \mu_j] \cdot \mathbf{1}_{\{A^* = i\}} - (\mu_i - \mu_j) \cdot \mathbf{1}_{\{A^* = i\}}|] \leq \varepsilon_{\text{TS}} \delta_i. \tag{4.6}$$

By the triangle inequality, establishing Equation (4.6) will complete the proof. By regrouping and using again the triangle inequality as well as  $|x \cdot \mathbf{1}_{\{A^* = i\}}| = |x| \cdot \mathbf{1}_{\{A^* = i\}}$  for any  $x \in \mathbb{R}$ , the left-hand side of (4.6) is upper bounded by

$$\mathbb{E}[|\mathbb{E}^{T_0}[\mu_i] - \mu_i| \cdot \mathbf{1}_{\{A^* = i\}}] + \mathbb{E}[|\mathbb{E}^{T_0}[\mu_j] - \mu_j| \cdot \mathbf{1}_{\{A^* = i\}}]. \tag{4.7}$$

By Lemma 3.1, the values  $|\mathbb{E}^{T_0}[\mu_i] - \mu_i|$  and  $|\mathbb{E}^{T_0}[\mu_j] - \mu_j|$  are  $N_{\text{TS}}^{-1/2}$  times  $O(1)$ -sub-Gaussian random variables. Applying Lemma A.13, we obtain that both terms in Eq. (4.7) are at most  $O\left(\delta_i \sqrt{\log(1/\delta_i)/N_{\text{TS}}}\right)$ . Using  $\delta_{\text{TS}} \leq \delta_i$  and our choice of  $N_{\text{TS}}$  we now conclude.  $\square$

*Proof of Remark 4.3.* Increasing the value of  $N_{\text{TS}}$  by a factor 4, compared to Theorem 4.1, ensures that

$$\mathbb{E}[|\mathbb{E}^{T_0}[\mu_i - \mu_j] \cdot \mathbf{1}_{\{A^* = i\}} - (\mu_i - \mu_j) \cdot \mathbf{1}_{\{A^* = i\}}|] \leq \varepsilon_{\text{TS}} \delta_i/2 \tag{4.8}$$

Revisiting the above proof, this directly implies

$$\mathbb{E}[\mathbb{E}^{T_0}[\mu_i - \mu_j] \cdot \mathbf{1}_{\{A^* = i\}}] \geq \varepsilon_{\text{TS}} \delta_i/2.$$



Recall (4.5) and the fact that  $\delta_i = \Pr[A^* = i] = \Pr[A_{T_0} = i]$  (the latter equality holds by the definition of Thompson sampling and the martingale property of  $\Pr^t[A^* = i]$ ). Combining implies

$$\mathbb{E}[\mu_i - \mu_j \mid A_{T_0} = i] \geq \varepsilon_{\text{TS}}/2. \quad \square$$

## 4.2 Monotonicity of Thompson Sampling

Our analysis of Thompson Sampling implies that its expected reward grows monotonically with time, which in turn allows us to upper-bound its Bayesian simple regret  $\mathbb{E}_{\text{prior}}[\max_{i \in [K]} \mu_i - \mu_{A_t}]$ . These results are new in the literature on Thompson Sampling, to the best of our knowledge.

**Corollary 4.10.** *For Thompson Sampling (starting with an arbitrary prior),  $\mathbb{E}_{\text{prior}}[\mu_{A_t}]$  is non-decreasing in round  $t$ . Consequently, Bayesian simple regret at each round  $t$  is at most  $O(\sqrt{K/t})$ .*

*Proof.* Denote  $H_{t,i} = \mathbb{E}^t[\mu_i \mid \Pr^t[A_t = i]]$ , a key object in the proof of Theorem 4.1. Then

$$\begin{aligned} \mathbb{E}[\mu_{A_{t+1}}] &= \sum_{\text{arms } i} \mathbb{E}[H_{t+1,i}] \geq \sum_{\text{arms } i} \mathbb{E}[H_{t,i}] && \text{(since } (H_{t,i} : t \in \mathbb{N}) \text{ is a submartingale)} \\ &= \sum_{\text{arms } i} \mathbb{E}[\mu_i \mid A_t = i] \Pr[A_t = i] && \text{(by (4.4))} \\ &= \mathbb{E}[\mu_{A_t}]. \end{aligned}$$

The bound on Bayesian simple regret follows simply because the (cumulative) Bayesian regret at round  $t$  is upper-bounded as  $O(\sqrt{Kt})$  by (Bubeck and Liu, 2013, Theorem 1), and equals the sum of Bayesian simple regret over all rounds  $s \in [t]$ .  $\square$

## 5 Collecting Initial Samples

We turn to a basic version of incentivized exploration: collect  $N$  samples of each arm. We design a BIC algorithm, called `ExponentialExploration`, which completes after a pre-determined number of rounds and collects  $N$  samples of each arm almost surely. (This is a desirable property as per Remark 4.2.) We bound the completion time in terms some parameters of the prior.

We describe the algorithm on a high level, and then fill in the details. The algorithm explores the arms in order of increasing index  $j$ , *i.e.*, in the order of decreasing prior mean reward. A given arm  $j$  is explored as follows. We partition time in phases of  $N$  rounds each, where  $N$  is a parameter. Within a given phase, the algorithm recommends the same arm in all rounds. It uses phases of three types: *exploration phases*, when arm  $j$  is always recommended, *exploitation phases*, when the algorithm chooses an arm with the largest posterior mean reward, and *padded phases*, which combine exploration and exploitation. In the exploitation phase, the algorithm conditions on the first  $N'$  samples of each arm  $i < j$ , where  $N'$  is a given *depth* parameter. The algorithm also conditions on the first  $N'$  samples of arm  $j$  if these samples are available before the phase starts. In the padded phase, it leverages the samples from arm  $j$  to guarantee a stronger BIC property for this arm which is “padded” by some prior-dependent amount  $\lambda > 0$ ; this property offsets more exploration for arm  $j$ . A given phase is assigned to one of these three types in a randomized and somewhat intricate way.

Algorithm 1 presents the algorithm with abstract parameters  $N_0 \leq N$  and  $\lambda > 0$ , and recommendation policies for the padded phase. Let us again focus on exploring a particular arm  $j$ . The first two phases (*bootstrapping*) ensure that exploration phase is invoked with probability  $p_j$  for any given vector of mean

---

**Algorithm 1:** ExponentialExploration

---

```
1 Parameters: phase length  $N$ , bootstrapping parameter  $N_0 \leq N$ , padding  $\lambda > 0$ .
2 Given: recommendation policies  $\pi_2, \dots, \pi_K$  for PADDED PHASE.
3 Initialize: EXPLORATION PHASE for arm 1
4 for each arm  $j = 2, 3, \dots, K$  do
5   // Invariant 1: each arm  $i < j$  has been sampled at least  $N$  times.
6   // Bootstrapping: two phases
7   Event  $\text{ZEROS}_{j,N_0} = \{ \text{the first } N_0 \text{ samples of each arm } i < j \text{ return reward } 0 \}$ .
8    $p_j \leftarrow q/(1+q)$ , where  $q = \lambda \cdot \Pr[\text{ZEROS}_{j,N_0}]$ .
9   EXPLOITATION PHASE with depth  $N_0$ 
10  with probability  $p_j$  do
11    EXPLORATION PHASE for arm  $j$ 
12  else if  $\text{ZEROS}_{j,N_0}$  then
13    | PADDED PHASE: use policy  $\text{transform}(\pi_j)$ 
14  else EXPLOITATION PHASE with depth  $N$ 
15
16  // main loop: exponentially grow the exploration probability
17  while  $p_j < 1$  do
18    // Invariant 2:  $\Pr[\text{exploration phase has happened} \mid \mu_1, \dots, \mu_K] = p_j$ .
19    if exploration phase has happened then
20    | PADDED PHASE: use policy  $\pi_j$ 
21    else with probability  $\min\left(1, \frac{p_j}{1-p_j} \cdot \lambda\right)$  do
22      EXPLOITATION PHASE for arm  $j$ 
23    else EXPLOITATION PHASE with depth  $N$ 
24    Update  $p_j \leftarrow \min(1, p_j(1+\lambda))$ .
```

---

rewards  $\mu_1, \dots, \mu_K$ , where  $p_j$  is given in Line 7. We capture this condition as Invariant 2. Then the algorithm enters the main loop, where it exponentially grows the exploration probability. More precisely, consider the *phase-exploration probability*: the probability that the “pure exploration” phase for arm  $j$  has been invoked. The algorithm increases this probability by a  $1 + \lambda$  factor after each iteration, maintaining Invariant 2. This exponential growth is the key aspect of the algorithm, which side-steps the fact that the initial phase-exploration probability in Line 7 may be very small. Inside the main loop, the algorithm randomizes between “pure exploitation” and “pure exploitation” with predetermined probability (chosen so as to guarantee that the  $1 + \lambda$  increase in phase-exploration probability overall), and permanently switch to the padded phase once exploration phase has been invoked. The padded phase offsets the additional exploration in the same iteration. This process continues until the phase-exploration probability reaches 1.

The bootstrapping phases hide a considerable amount of complexity which may be skipped at a first reading. Conceptually, we would like to implement the “hidden exploration” approach from Mansour et al. (2015, 2020), which randomizes between exploration and exploitation phases with some predetermined probability. While this approach may suffice for some “well-behaved” priors, it appears to be insufficient more generally, in the sense that the phase-exploration probability depends on some additional prior-dependent parameters that are difficult to deal with. Instead, we combine exploration and exploitation in a more sophisticated way, as explained later in the section.

Let’s make some observations that are immediate from the algorithm’s specification.

**Claim 5.1.** *Algorithm 1 samples each arm at least  $N$  times with probability 1, and completes in  $(2K - 1)N + N \sum_{j=2}^K \lceil \ln_{1+\lambda}(p_j) \rceil$  rounds, where  $p_j$  is given in Line 7. Both invariants hold at each iteration of the respective loops. Bootstrapping takes exactly two phases, and each iteration of the `while` loop takes exactly one phase.*

**Recommendation policies**  $\pi_j$  need to satisfy several properties. Fix arm  $j$ . First, we require  $\pi_j$  to input only the first  $N$  samples of each arm  $i \leq j$ , ignoring the order in which the arms were sampled by the algorithm.<sup>12</sup> Such policies are called  $(j, N)$ -informed. Note that the algorithm has sufficient data to compute policy  $\pi_j$  thanks to Invariant 1.

Second, we require policy  $\pi_j$  to be BIC. Formally, we let  $\mathcal{S}_{j,N}$  be a signal that consists of exactly  $N$  independently realized samples of each arm  $i \leq j$ , and we require  $\pi_j$  to be BIC w.r.t this signal.

Third, we require  $\pi_j$  to satisfy a stronger BIC property for arm  $j$ :

$$\mathbb{E} \left[ (\mu_j - \mu_i) \cdot \mathbf{1}_{\{\pi_j(\mathcal{S}_{j,N})=j\}} \right] \geq \lambda \quad \text{for each arm } i < j. \quad (5.1)$$

If (5.1) holds, we say that policy  $\pi_j$  is  $(j, \lambda)$ -padded BIC, where  $\lambda$  is the “padding”. The left-hand side in (5.1) is the expected loss for the  $j \rightarrow i$  swap: the expected loss when one starts with policy  $\pi_j$  and replaces arm  $j$  with arm  $i$  whenever arm  $j$  is recommended. Note that we integrate over the event that arm  $j$  being chosen, rather than condition on this event. We recover the “usual” BIC property for arm  $j$  when  $\lambda = 0$ .

If policy  $\pi_j$  satisfies all three properties, it is called  $(j, \lambda, N)$ -suitable.

While BIC and  $(j, \lambda)$ -padded BIC properties of  $\pi_j$  are defined relative to signal  $\mathcal{S}_{j,N}$ , one could also define them relative to any other signal  $\mathcal{S}$  that almost surely contains at least  $N$  samples of each arm  $i \leq j$ . It is easy to see that these definitions are equivalent: any  $(j, N)$ -informed policy is BIC relative to signal  $\mathcal{S}_{j,N}$  if and only if it is BIC relative to signal  $\mathcal{S}$ ; likewise for the  $(j, \lambda)$ -padded BIC property. This point allows us to analyze  $(j, \lambda, N)$ -suitable policies abstractly, regardless of where exactly their input comes from.

**Padding and the main loop.** The key is to specify what happens in the “padded phase”, and argue about incentives that it creates. We use the properties of policies  $\pi_j$ , as listed above, to guarantee that the main loop is BIC. The main point is that the padded-BIC property compensates for the probability of new exploration. We carefully spell out which properties are needed where. In particular, the exploitation phase is only used to skip the round in a BIC way.<sup>13</sup> In fact, the algorithm would work even if the exploitation phase in the main loop would always choose arm 1. However, using the available data for exploitation only improves the algorithm’s efficiency as well as the incentives.

**Lemma 5.2.** *Consider Algorithm 1 with arbitrary parameters  $N_0 \leq N$  and  $\lambda > 0$ . Fix arm  $j \geq 2$  and focus on the respective iteration of the `for` loop of the algorithm. Assume that policy  $\pi_j$  is  $(j, \lambda, N)$ -suitable. Then the `while` loop is well-defined and BIC.*

*Proof.* The `while` loop is well-defined because policy  $\pi_j$  is  $(j, N)$ -informed, so by Invariant 1 the algorithm has a sufficient amount of data to implement it.

Fix some round  $t$  in the `while` loop. Let us restate the BIC property when arm  $j$  is recommended:

$$\mathbb{E} [\Lambda_{i,j}] \geq 0, \quad \text{where } \Lambda_{i,j} := (\mu_j - \mu_i) \cdot \mathbf{1}_{\{A_t=j\}}, \quad (5.2)$$

<sup>12</sup>Formally,  $\pi_j$  inputs an ordered tuple of arm-reward pairs, and pre-processes it as a  $j \times N$  matrix whose  $(i, n)$ -th entry,  $(i, n) \in [j] \times [N]$ , is the reward from arm  $i$  from the  $n$ -th time it was sampled. The policy is then determined by this matrix.

<sup>13</sup>This is an interesting contrast with “hidden exploration” (Mansour et al., 2015, 2020), where the exploitation phase is used to offset exploration. We have little use for this, because the padded phase enables nearly as much additional exploration as possible.

for all arms  $i \neq j$ . Let  $\mathbf{1}_{\text{padded}}$ ,  $\mathbf{1}_{\text{exploit}}$ ,  $\mathbf{1}_{\text{explore}}$  be the indicator variables for the event that round  $t$  is assigned to, resp., a padded, exploitation, or exploration phase. Due to Invariant 2, these indicator variables are independent of mean rewards  $\mu_1, \dots, \mu_K$ . We write

$$\mathbb{E}[\Lambda_{i,j}] = \mathbb{E}[\Lambda_{i,j} \cdot (\mathbf{1}_{\text{padded}} + \mathbf{1}_{\text{exploit}} + \mathbf{1}_{\text{explore}})]$$

and estimate each expectation separately. Let  $\Lambda_{i,j}^{\text{pad}} := (\mu_j - \mu_i) \cdot \mathbf{1}_{\{\pi_j=j\}}$  and observe that

$$\mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{padded}}] = \mathbb{E}[\Lambda_{i,j}^{\text{pad}} \cdot \mathbf{1}_{\text{padded}}] = \mathbb{E}[\Lambda_{i,j}^{\text{pad}}] \cdot \mathbb{E}[\mathbf{1}_{\text{padded}}] = \mathbb{E}[\Lambda_{i,j}^{\text{pad}}] \cdot p_j. \quad (5.3)$$

The last two equalities hold, resp., by independence of  $\mathbf{1}_{\text{padded}}$  and by Invariant 2.

First consider the case  $i < j$ . Then

$$\begin{aligned} \mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{padded}}] &\geq \lambda p_j && \text{(by (5.3) and the padded-BIC property of } \pi_j), \\ \mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{explore}}] &\geq -\mathbb{E}[\mathbf{1}_{\text{explore}}] \geq -\lambda p_j && \text{(a worst-case bound),} \\ \mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{exploit}}] &\geq 0 && \text{(exploitation is always BIC).} \end{aligned}$$

Summing it up gives Eq. (5.2). For  $i > j$ , we have a similar argument, albeit for different reasons:

$$\begin{aligned} \mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{padded}}] &\geq 0 && \text{(by (5.3) and the BIC property of } \pi_j), \\ \mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{explore}}] &= \mathbb{E}[(\mu_j - \mu_i) \cdot \mathbf{1}_{\text{explore}}] \geq 0 && \text{(by independence of } \mathbf{1}_{\text{explore}}), \end{aligned}$$

and  $\mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{exploit}}] \geq 0$  as before. We note that  $\mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{exploit}}] \geq 0$  would hold even if the exploitation phase would simply choose arm 1; this is by independence of  $\mathbf{1}_{\text{exploit}}$ .

It remains to show the BIC property when the algorithm recommends some arm  $\ell \neq j$ , i.e., that  $\mathbb{E}[\Lambda_{i,\ell}] \geq 0$  for all arms  $i \neq \ell$ . We derive  $\mathbb{E}[\Lambda_{i,\ell} \cdot \mathbf{1}_{\text{padded}}] \geq 0$  and  $\mathbb{E}[\Lambda_{i,\ell} \cdot \mathbf{1}_{\text{exploit}}] \geq 0$  same way as before, and  $\mathbb{E}[\Lambda_{i,\ell} \cdot \mathbf{1}_{\text{explore}}] = 0$  because the exploration phase always chooses arm  $j$ .  $\square$

Now, we show that a suitable policy  $\pi_j$  exists, in terms of the following parameters:

$$\begin{aligned} G_{\text{pad}} &= \min_{j, q \in \Delta_{j-1}} \mathbb{E}[(\mu_j - \mu_q)_+], \quad \text{where } \mu_q := \sum_{i \in [K]} q_i \mu_i. \\ N_{\text{pad}} &= C_{\text{pad}} G_{\text{pad}}^{-2} \log(G_{\text{pad}}^{-1}), \end{aligned} \quad (5.4)$$

for large enough absolute constant  $C_{\text{pad}}$ . In words,  $G_{\text{pad}}$  is the smallest “expected advantage” of any arm  $j$  over a convex combination of arms  $i < j$ . We merely guarantee existence of a policy via the minimax theorem, rather than specify how to compute it.

**Lemma 5.3.** *For each arm  $j \geq 2$  there exists a policy  $\pi_j$  which is  $(j, \lambda, N_{\text{pad}})$ -suitable, for  $\lambda = G_{\text{pad}}/10$ .*

*Proof Sketch.* Fix arm  $j \geq 2$ . Maximizing padding  $\lambda$  in Eq. (5.1) is naturally expressed as a zero-sum game between a *planner* who chooses a  $(j, N)$ -informed policy  $\pi_j$  and wishes to maximize the right-hand side of (5.1), and an *agent* who chooses arm  $i$ ; we call it the  $(j, N)$ -*recommendation game*.<sup>14</sup> A policy  $\pi_j$  is  $(j, \lambda)$ -padded if and only if it is guaranteed payoff at least  $\lambda$  in this game. So, any maximin policy in this game is  $(j, \lambda)$ -padded with  $\lambda = V_{j,N}$ , where  $V_{j,N}$  is the game value.

<sup>14</sup>It is a finite game because there are only finitely many deterministic  $(j, N)$ -informed policies.

We connect the game value with  $G_{\text{pad}}$  via the minimax theorem, proving that  $V_{j,N} \geq G_{\text{pad}}/10$ . Indeed, using the minimax theorem and the linearity of expectation, we can write the game value as

$$V_{j,N} = \min_{q \in \Delta_{j-1}} \max_{(j,N)\text{-informed policies } \pi_j} \mathbb{E} \left[ (\mu_j - \mu_q) \cdot \mathbf{1}_{\{\pi_j=j\}} \right]. \quad (5.5)$$

If the policy  $\pi_j$  knew the mean rewards  $(\mu_1, \dots, \mu_j)$  exactly, it could recommend arm  $j$  if and only if  $\mu_j - \mu_q > 0$ . This would guarantee the minimax value of at least  $G_j := \min_{q \in \Delta_{j-1}} \mathbb{E} \left[ (\mu_j - \mu_q)_+ \right]$ . We show that for a large enough  $N$  we can guarantee the minimax value of at least  $\Omega(G_j)$ . Specifically, for  $N = N_{\text{pad}}$  and any given distribution  $q \in \Delta_{j-1}$  there exists a  $(j, N)$ -informed policy  $\pi = \pi_{j,N}^q$  such that

$$\mathbb{E} \left[ (\mu_j - \mu_q) \cdot \mathbf{1}_{\{\pi=j\}} \right] \geq \frac{1}{10} \cdot \mathbb{E} \left[ (\mu_j - \mu_q)_+ \right]. \quad (5.6)$$

This policy is very simple: given the data (the first  $N$  samples of each arm  $i \leq j$ ), recommend arm  $j$  if and only if its empirical reward on this data is larger than the expected reward of distribution  $q$  on the same data. For every given realization of the mean rewards  $(\mu_1, \dots, \mu_j)$ , we use Chernoff Bounds to compare  $(\mu_j - \mu_q) \cdot \mathbf{1}_{\{\pi=j\}}$  and  $(\mu_j - \mu_q)_+$ , which implies (5.6).

From here on, let us focus on the  $(j, N_{\text{pad}})$ -recommendation game. We have proved that there exists a  $(j, N_{\text{pad}})$ -informed policy  $\pi_j$  that is  $(j, \lambda)$ -padded, for  $\lambda = G_{\text{pad}}/10$ , even if we do not specify how to compute such a policy. However, we are not done yet, as we also need this policy to be BIC. First, we need the BIC property when arm  $j$  is recommended. (As spelled out in Eq. (5.2); we already have this property for  $i < j$ , by the padded-BIC property, but we also need it for  $i > j$ .)

Any finite two-player zero-sum game has a minimax-optimal strategy that is *non-weakly-dominated*: not weakly dominated by any other mixed strategy for the max player. So, let us take such a policy  $\pi^*$ . One can prove that any such policy is BIC when recommending arm  $j$ . The proof of this claim focuses on the probability of recommending arm  $j$  as a function of the posterior mean rewards  $\tilde{\mu}_i := \mathbb{E}[\mu_i | \mathcal{G}_{N,j}]$ .<sup>15</sup> We show that the conditional probability  $\Pr[\pi_j = j | \tilde{\mu}_1, \dots, \tilde{\mu}_j]$  is non-decreasing in  $\tilde{\mu}_j$  and non-increasing in  $\tilde{\mu}_i$  for each arm  $i < j$ . This, in turn, allows us to invoke the FKG inequality and derive the claim.

Finally, we extend  $\pi^*$  to a BIC policy. When  $\pi^*$  does *not* recommend arm  $j$ , we choose an arm  $i$  for exploitation, *i.e.*, to maximize  $\mathbb{E}[\mu_i | \mathcal{G}_{N,j}]$ . The resulting policy is BIC, hence  $(j, \lambda, N_{\text{pad}})$ -suitable.  $\square$

**Bootstrapping, revisited.** The bootstrapping proceeds as follows, as per the pseudocode. We start with an exploitation phase with depth  $N_0$ . We choose  $N_0$  so as to guarantee that we explore arm  $j$  in the  $\text{ZEROS}_{j,N_0}$  event, *i.e.*, if all previous arms return 0 rewards in the first  $N_0$  samples. In the second phase, we explore arm  $j$  with a small probability  $p_j$ ; otherwise we do something else to guarantee incentives. Specifically, under event  $\text{ZEROS}_{j,N_0}$  we invoke a version of the padded phase; else we just exploit. This phase is BIC because the small probability of invoking the padded phase compensates for the exploration.

Given that the padded phase is now invoked under event  $\text{ZEROS}_{j,N_0}$ , we need the padded-BIC property to hold conditional on this event. We transform policy  $\pi_j$  into another policy  $\pi = \text{transform}(\pi_j)$  which replaces the BIC and padded-BIC properties with conditional ones:

$$\mathbb{E} \left[ (\mu_j - \mu_i) \cdot \mathbf{1}_{\{\pi=j\}} \mid \text{ZEROS}_{j,N_0} \right] \geq \lambda \quad \text{for each arm } i < j, \quad (5.7)$$

$$\mathbb{E} \left[ (\mu_\ell - \mu_i) \cdot \mathbf{1}_{\{\pi=\ell\}} \mid \text{ZEROS}_{j,N_0} \right] \geq 0 \quad \text{for each arm } i \neq \ell. \quad (5.8)$$

This transformation is generic, in the sense that it works for any policy  $\pi_i$  and any parameters.

<sup>15</sup>Recall that  $\mathcal{G}_{N,j}$  denotes the  $\sigma$ -algebra generated by the first  $N$  samples of each arm  $i \leq j$ .



**Lemma 5.4.** Fix padding  $\lambda > 0$  and any parameters  $N_0, N$ . Given a  $(j, \lambda, N)$ -suitable policy  $\pi_j$ , there exists a BIC policy  $\pi = \text{transform}(\pi_j)$  which is  $(j, N)$ -informed and satisfies (5.7) and (5.8).

*Proof Sketch.* Intuitively, conditioning on arms  $i < j$  being worse-than-usual should only help in finding a  $(j, \lambda)$ -padded BIC policy. To explicitly reduce to  $\pi_j$ , we consider the “true data”: the first  $N$  samples of each arm  $i < j$  which satisfy  $\text{ZEROS}_{j, N_0}$ . We construct a “fake” data set ( $N$  samples for each arm  $i < j$ ), which relates to the “true” data in a certain way. Specifically, the “fake” posterior mean reward  $\hat{\mu}_i$  for each arm  $i < j$  is coupled with the “true” posterior mean reward, denoted  $\tilde{\mu}_i$ , so that  $\hat{\mu}_i \geq \tilde{\mu}_i$  almost surely and the distribution of  $\hat{\mu}_i$  is the unconditional distribution of  $\tilde{\mu}_i$  (i.e., the distribution of  $\tilde{\mu}_i$  without conditioning on the event  $\text{ZEROS}_{j, N_0}$ ). We define  $\text{transform}(\pi_j)$  by applying  $\pi_j$  on the fake data.

Since  $\pi_j$  is  $(j, \lambda, N)$ -suitable, we obtain the BIC (resp., padded-BIC) property against the fake posterior means  $\hat{\mu}_i, i < j$ . Since  $\hat{\mu}_i \geq \tilde{\mu}_i$ , the corresponding properties hold against the values  $\tilde{\mu}_i$  as well.  $\square$

Now we prove that bootstrapping is BIC given a suitable policy  $\pi_j$ . The proof follows the same strategy as that of Lemma 5.2, but everything is conditioned on  $\text{ZEROS}_{j, N_0}$ .

**Lemma 5.5.** Consider Algorithm 1 with arbitrary parameters  $N_0 \leq N$  and  $\lambda > 0$ . Fix some arm  $j \geq 2$ . The bootstrapping phases are BIC as long as policy  $\pi_j$  is  $(j, \lambda, N)$ -suitable and parameter  $N_0$  satisfies

$$\mathbb{E}[\mu_j - \mu_i \mid \text{ZEROS}_{j, N_0}] \geq 0 \quad \text{for all arms } i < j. \quad (5.9)$$

**Putting the pieces together.** Let us formulate an end-to-end guarantee for the algorithm in terms of the appropriate parameters. Recall that we use  $G_{\text{pad}}, N_{\text{pad}}$  from (5.4) to handle the “padded phase”. Given (5.9), we define two more parameters to handle the “bootstrapping phase”:

$$N_{\text{boot}} = \min \{ N_0 \in \mathbb{N} : \mathbb{E}[\mu_j - \mu_i \mid \text{ZEROS}_{j, N_0}] > 0 \text{ for all arms } i < j \}. \quad (5.10)$$

$$p_{\text{boot}} = \min_{\text{arms } j} \Pr[\text{ZEROS}_{j, N_{\text{boot}}}] = \Pr[\text{ZEROS}_{K, N_{\text{boot}}}] . \quad (5.11)$$

In words,  $N_{\text{boot}}$  is the smallest  $N_0$  such that each arm  $j$  is the best arm conditional on seeing  $N_0$  initial samples from all arms  $i < j$  that are all zeroes. Setting  $N_0 = N_{\text{boot}}$  as we do in the theorem,  $p_{\text{boot}}$  is the smallest phase-exploration probability after the bootstrapping phase.

**Theorem 5.6.** Suppose algorithm `ExponentialExploration` is run with parameters  $N \geq \max(N_{\text{boot}}, N_{\text{pad}})$ ,  $N_0 = N_{\text{boot}}$  and  $\lambda = G_{\text{pad}}/10$ . Suppose each policy  $\pi_i$  is  $(j, \lambda, N)$ -suitable (such policies exist by Lemma 5.3). Then the algorithm is BIC and collects at least  $N$  samples of each arm almost surely in time

$$T_{\text{UB}}(N) = O \left( KN G_{\text{pad}}^{-1} \log \left( G_{\text{pad}}^{-1} p_{\text{boot}}^{-1} \right) \right).$$

**Remark 5.7.** If all prior mean rewards  $\mathbb{E}[\mu_i], i \in [K]$  are pairwise distinct, then the algorithm in Theorem 5.6 is strictly BIC: Eq. (3.1) is satisfied with a strict inequality, as per the same analysis.

Theorem 5.6 provides an explicit formula for the time horizon of `ExponentialExploration`. This formula is an *upper bound* on the the sample complexity in question. It is polynomially optimal and allows for concrete corollaries, as we discuss in Section 6. Setting  $N \geq N_{\text{TS}}$ , we collect enough data to bootstrap Thompson Sampling, as per Section 4.

How to compute a suitable policy  $\pi_j$  for Lemma 5.3? For  $K = 2$  arms, one can take a very simple policy  $\pi_2$ : given the first  $N$  samples of both arms, recommend an arm with a larger empirical reward. (This idea can be extended to an arbitrary  $K$ , but the padding  $\lambda$  degrades exponentially in  $K$ ; we omit the easy details.) We provide a computationally efficient implementation for Beta priors in Section 7. However, we do not provide a computational implementation for the general case.

## 6 Sample Complexity of Incentivized Exploration

We characterize the sample complexity of incentivized exploration: the minimal number of rounds need to collect  $N$  samples of each arm. More precisely, we are interested in the smallest time horizon  $T$  such that some BIC bandit algorithm collects  $N$  samples of each arm almost surely; denote it  $T_{\text{OPT}}(N)$ . Note that  $T_{\text{OPT}}(N)$  is determined by the joint prior, and does not depend on the mean rewards or observed rewards.

We are particularly interested in  $T_{\text{OPT}}(1)$ , the sample complexity of collecting at least one sample of each arm, and  $T_{\text{OPT}}(N_{\text{TS}})$ , the sample complexity of bootstrapping Thompson Sampling. The upper bound  $T_{\text{UB}}(N_{\text{TS}}) \geq T_{\text{OPT}}(N_{\text{TS}})$  comes from Section 5.<sup>16</sup> In this section, we derive a lower bound on  $T_{\text{OPT}}(1)$  and mine the upper/lower bounds for a number of corollaries. In particular, we prove that upper and lower bounds are “polynomially matching”, investigate how the sample complexity scales with  $K$  and smallest variance  $\sigma^2 = \min_{i \in [K]} \text{Var}(\mathcal{P}_i)$ , and work out canonical special cases. We focus on, and largely resolve, the distinction between polynomial and exponential dependence on  $K$  and  $\sigma^{-1}$ .

### 6.1 Lower bound

We provide a lower bound  $T_{\text{LB}} \leq T_{\text{OPT}}(1)$ , and prove that it matches  $T_{\text{UB}}(N_{\text{TS}})$ , in a specific sense that we explain below. The lower bound is driven by the following parameter:

$$L_{\text{LB}} := \max_{j \in [K], q \in \Delta_K} \frac{\mathbb{E}[(\mu_j - \mu_q)_-]}{\mathbb{E}[(\mu_j - \mu_q)_+]} = \max_{j \in [K], q \in \Delta_K} \frac{\mu_q^0 - \mu_j^0}{\mathbb{E}[(\mu_j - \mu_q)_+]} - 1. \quad (6.1)$$

**Theorem 6.1.**  $T_{\text{OPT}}(1) \geq T_{\text{LB}} := \max(K, N_{\text{boot}}, L_{\text{LB}})$ .

*Proof.* The lower bounds of  $K, N_{\text{boot}}$  are clear so we focus on the last one. Suppose it is possible to sample all arms within  $T$  rounds with some BIC algorithm. Fix an arbitrary arm  $j$  and distribution  $q \in \Delta_K$ . In each rounds  $t$  when arm  $j$  is played, it must appear better than  $q$ , i.e.,  $\mathbb{E}[\mathbf{1}_{\{A_t=j\}} \cdot (\mu_j - \mu_q)] \geq 0$ . Summing this up over all rounds  $t$ , and letting  $n_j$  be the number of times arm  $j$  is played (which is a random variable) we obtain  $\mathbb{E}[n_j \cdot (\mu_j - \mu_q)] \geq 0$ . Now, the expression  $n_j \cdot (\mu_j - \mu_q)$  is minimized by taking  $n_j = 1$  when  $\mu_j < \mu_q$ , and  $n_j = T$  when  $\mu_j \geq \mu_q$ . Consequently we have

$$0 \leq \mathbb{E}[n_j \cdot (\mu_j - \mu_q)] \leq T \mathbb{E}[(\mu_j - \mu_q)_+] - \mathbb{E}[(\mu_j - \mu_q)_-].$$

Rearranging shows that  $T \geq L_{\text{LB}}$  as claimed.  $\square$

**Remark 6.2.** While Theorem 6.1 provides a rather tight estimate for 1-sample complexity, we believe it does not imply any non-trivial lower bound on Bayesian regret of BIC 1-sampling. To illustrate this point, we describe a stylized attempt to prove such lower bound via  $L_{\text{LB}}$ , and explain why it fails.

Recall that BIC 1-sampling must proceed for at least  $L_{\text{LB}}$  rounds by Theorem 6.1. Fix arm  $j$  that maximizes Eq. (6.1). Let us attempt to lower-bound Bayesian regret due to *not* playing this arm, under an (overly pessimistic and unjustified) assumption that it is never played in the first  $L_{\text{LB}}$  rounds. Even then, the lower bound we obtain is at most 1.

Fix  $q \in \Delta_K$  that maximizes (6.1) for arm  $j$ , so that

$$L_{\text{LB}} = \frac{\mathbb{E}[(\mu_j - \mu_q)_-]}{\mathbb{E}[(\mu_j - \mu_q)_+]}. \quad (6.2)$$

<sup>16</sup>Since Theorem 5.6 only provides  $T_{\text{UB}}(N)$  for  $N \geq \tilde{N} := \max(N_{\text{boot}}, N_{\text{pad}})$ , we define  $T_{\text{UB}}(N) = T_{\text{UB}}(\tilde{N})$  for  $N < \tilde{N}$ .

Assume  $q_j = 0$  w.l.o.g. (the quantities in (6.1) are unchanged if we linearly rescale the remaining entries by  $q_i \leftarrow \frac{q_i}{1-q_j}$  and set  $q_j = 0$ ). Each of the  $L_{\text{LB}}$  rounds when arm  $j$  is not played contributes Bayesian regret

$$\mathbb{E} \left[ \max_{i \in [K]} \mu_i - \max_{i \in [K] \setminus \{j\}} \mu_i \right] = \mathbb{E} \left[ (\mu_j - \max_{i \neq j} \mu_i)_+ \right] \leq \mathbb{E} [(\mu_j - \mu_q)_+] \leq 1/L_{\text{LB}}. \quad (6.3)$$

(The first inequality in Eq. (6.3) holds because  $\mu_q \leq \max_{i \neq j} \mu_i$  almost surely, which in turn holds because  $q_j = 0$ . The second inequality in Eq. (6.3) holds by (6.2).) Thus, this stylized argument lower-bounds Bayesian regret by  $L_{\text{LB}}$  times the left-hand side of (6.3), which unfortunately is at most 1.

## 6.2 Polynomially matching upper/lower bounds

We express the upper bound  $T_{\text{UB}}(N_{\text{TS}})$  in terms of the parameters in the lower bound and the smallest variance  $\sigma^2$ . We conclude that the upper bound is polynomially optimal up to  $\sigma^{-O(1)}$  factor.

**Corollary 6.3.** *Suppose the prior  $\mathcal{P}_i$  for each arm  $i \in [K]$  has variance at least  $\sigma^2$ . Then*

$$T_{\text{UB}}(N_{\text{TS}}) \leq \tilde{O} \left( \sigma^{-4} K^{4.5} L_{\text{LB}}^3 N_{\text{boot}} + \sigma^{-2} K^{2.5} L_{\text{LB}} N_{\text{boot}}^2 \right) = \tilde{O} \left( \sigma^{-4} \cdot T_{\text{LB}}^{O(1)} \right). \quad (6.4)$$

This result suffices to resolve polynomial vs. exponential dependence on  $\sigma^{-1}$  or  $K$  or any other parameter. To make this statement explicit, suppose we have a family of problem instances indexed by a single parameter  $\sigma^2 = \min_{i \in [K]} \text{Var}(\mathcal{P}_i)$ , and for this family the lower bound scales as  $T_{\text{LB}} = f(1/\sigma)$ . Then

$$f(1/\sigma) = T_{\text{LB}} \leq T_{\text{OPT}}(N) \leq T_{\text{UB}}(N_{\text{TS}}) \leq (f(1/\sigma))^{O(1)} \quad \forall N \in [N_{\text{TS}}]. \quad (6.5)$$

So, both upper and lower bounds are polynomial (resp., exponential) in  $1/\sigma$  when so is  $f(1/\sigma)$ . A similar statement holds for any family of problem instances indexed by  $K$  (or any other parameter), where  $\sigma$  is lower-bounded by an absolute positive constant. In particular, assuming the lower bound scales as  $T_{\text{LB}} = g(K)$ , we have (6.5) with  $f(1/\sigma)$  replaced by  $g(K)$ .

## 6.3 Dependence on the number of arms

To investigate how  $T_{\text{OPT}}(\cdot)$  scales with  $K$ , we need to separate the dependence on  $K$  from the dependence on the priors  $\mathcal{P}_i$ ,  $i \in [K]$ . Therefore, we posit that all priors come from some (possibly infinite) collection  $\mathcal{C}$ ; such problem instances are called  $\mathcal{C}$ -consistent. Keeping  $\mathcal{C}$  fixed, we study the dependence on  $K$  in the worst case over all  $\mathcal{C}$ -consistent instances.

We find a curious dichotomy: under a mild non-degeneracy assumption, either  $T_{\text{UB}}(N_{\text{TS}}) = O_{\mathcal{C}}(K^3)$  for all  $\mathcal{C}$ -consistent instances (under a fairly reasonable assumption on  $\mathcal{C}$ ), and otherwise  $T_{\text{LB}}(1) > \exp(\Omega_{\mathcal{C}}(K))$  for some  $\mathcal{C}$ -consistent instance. In its simplest form, this dichotomy can be stated for a finite collection  $\mathcal{C}$ . For each prior  $\mathcal{P} \in \mathcal{C}$ , we consider the supremum of its support,  $\text{sup}(\mathcal{P}) := \text{sup}(\text{support}(\mathcal{P}))$ . We relate these quantities to the largest prior mean reward over  $\mathcal{C}$ , denoted  $\Phi_{\mathcal{C}} = \sup_{\mathcal{P} \in \mathcal{C}} \mathbb{E}[\mathcal{P}]$ .

**Theorem 6.4.** *Let  $\mathcal{C}$  be a finite collection of distributions over  $[0, 1]$ . Then*

- (a) *If  $\text{sup}(\mathcal{P}) > \Phi_{\mathcal{C}}$  for all priors  $\mathcal{P} \in \mathcal{C}$ , then  $T_{\text{UB}}(N_{\text{TS}}) = O_{\mathcal{C}}(K^3)$  for all  $\mathcal{C}$ -consistent instances.*
  - (b) *If  $\text{sup}(\mathcal{P}) < \Phi_{\mathcal{C}}$  for some  $\mathcal{P} \in \mathcal{C}$ , then  $T_{\text{LB}}(1) > \exp(\Omega_{\mathcal{C}}(K))$  for some  $\mathcal{C}$ -consistent instance.*
- Either (a) or (b) holds assuming that  $\min_{\mathcal{P} \in \mathcal{C}} \text{sup}(\mathcal{P}) \neq \Phi_{\mathcal{C}}$ .*

The assumption in part (a) is very reasonable. For instance, it holds whenever each prior has a positive density everywhere on the  $[0, 1]$  interval.

A quantitative version of Theorem 6.4, which is the version we actually prove, is more difficult to state. For a given parameter  $\delta > 0$ , the easy vs. hard distinction is as follows:

$$\mathcal{C} \text{ is called } \delta\text{-easy} \text{ if } \inf_{\mathcal{P} \in \mathcal{C}} \mathbb{E}_{\mu \sim \mathcal{P}} [(\mu - \Phi_{\mathcal{C}})_+] > \delta \quad (6.6)$$

$$\mathcal{C} \text{ is called } \delta\text{-hard} \text{ if } \inf_{\mathcal{P} \in \mathcal{C}} \Pr_{\mu \sim \mathcal{P}} [\mu \geq \Phi_{\mathcal{C}} - \delta] = 0. \quad (6.7)$$

Moreover, we need a quantitative version of pairwise non-dominance (3.2):  $\mathcal{C}$  is called  $\delta$ -non-dominant if

$$\mathbb{E} [(\mu_j^0 - \mu_i)_+] \geq \delta \quad \text{for every } \mathcal{C}\text{-consistent instance.} \quad (6.8)$$

In terms of these properties, the characterization extends to infinite collections  $\mathcal{C}$ . Note that if  $\mathcal{C}$  is  $\delta$ -easy (resp.,  $\delta$ -non-dominant) then so is any subset of  $\mathcal{C}$ . Likewise, if  $\mathcal{C}$  is  $\delta$ -hard, then so is any superset of  $\mathcal{C}$ .

**Theorem 6.5.** *Let  $\mathcal{C}$  be a (finite or infinite) collection of distributions over  $[0, 1]$ . Then*

- (a) *If  $\mathcal{C}$  is  $\delta$ -easy and  $\delta$ -non-dominant, then  $T_{\text{UB}}(N_{\text{TS}}) = \tilde{O}(K^3/\delta^4)$  for all  $\mathcal{C}$ -consistent instances.*
- (b) *If  $\mathcal{C}$  is  $\delta$ -hard, then  $T_{\text{LB}}(1) > \exp(\Omega_{\delta}(K))$  for some  $\mathcal{C}$ -consistent instance.*

*If  $\mathcal{C}$  is finite, then it is either  $\delta$ -easy or  $\delta$ -hard for some  $\delta > 0$ , provided that  $\min_{\mathcal{P} \in \mathcal{C}} \sup(\mathcal{P}) \neq \Phi_{\mathcal{C}}$ .*

*Proof.* For part (a) we upper-bound in Lemma C.3 the prior-dependent parameters as follows:  $N_{\text{TS}} = \tilde{O}(K\delta^{-2})$  for Thompson Sampling,  $G_{\text{pad}} \geq \delta$  and  $N_{\text{pad}} = \tilde{O}(\delta^{-2})$  for the padded phase, and  $N_{\text{boot}} = \tilde{O}(\delta^{-1})$  and  $\log(p_{\text{boot}}^{-1}) = \tilde{O}(K\delta^{-1})$  for the bootstrapping phase. Then part (a) follows from Theorem 5.6.

For part (b), assume  $\mathcal{C}$  is  $\delta$ -hard. Then for any  $\eta > 0$  there exist priors  $\mathcal{P}, \mathcal{P}' \in \mathcal{C}$  with

$$\Pr_{\mu \sim \mathcal{P}} [\mu \leq \mathbb{E}[\mathcal{P}'] - \delta/2] \geq 1 - \eta.$$

Consider the problem instance in which  $\mu_1, \dots, \mu_{K-1} \sim \mathcal{P}'$  and  $\mu_K \sim \mathcal{P}$ . Let  $q$  be the uniform distribution on  $[K-1]$ . Using Chernoff Bounds, we have

$$\Pr [\mu_q \leq \mathbb{E}[\mathcal{P}'] - \delta/2] \leq e^{-\Omega(K\delta^2)}.$$

Consequently,  $\mathbb{E}[(\mu_K - \mu_q)_+] \leq \Pr[\mu_K \geq \mu_q] \leq 2\eta$ .

Moreover, it holds that  $\mathbb{E}[\mu_q - \mu_K] \geq \delta/2 - \eta$ . Taking  $\eta \leq \min(\delta/4, e^{-\Omega(K\delta^2)})$ , we conclude that

$$T_{\text{LB}} \geq L_{\text{LB}} \geq \frac{\mathbb{E}[\mu_q - \mu_K]}{\mathbb{E}[(\mu_K - \mu_q)_+]} - 1 \geq \delta/4 \cdot e^{-\Omega(K\delta^2)} - 1. \quad \square$$

## 6.4 Canonical priors

We consider two canonical examples of incentivized exploration: when the priors  $\mathcal{P}_i$  are truncated Gaussians and when they are Beta distributions. We find that the optimal sample complexity  $T_{\text{OPT}}(N_{\text{TS}})$  scales polynomially on  $K$  and exponentially in the ‘‘strength of beliefs’’.

For truncated Gaussian priors, we focus on the case when all Gaussians have the same variance  $\sigma^2$ , and we find that the sample complexity is polynomial in  $K$  and exponential in  $\sigma^{-2}$ . Thus, *strong beliefs*, as expressed by small variance  $\sigma^2$ , is what makes incentivized exploration difficult.

**Corollary 6.6.** Let  $\tilde{N}(\nu, \sigma^2)$  be a Gaussian with mean  $\nu$  and variance  $\sigma^2 \leq 1$ , conditioned to lie in  $[0, 1]$ . Suppose  $\mathcal{P}_i \sim \tilde{N}(\nu_i, \sigma^2)$  for each arm  $i$ , where  $\nu_1, \dots, \nu_K \in [0, 1]$ . Then

- (a)  $T_{\text{UB}}(N_{\text{TS}}) = K^3 \cdot \text{poly}\left(\sigma^{-1}, e^{R^2}\right)$  where  $R = \sigma^{-1} \max_{i,j} |\nu_i - \nu_j|$ .
- (b)  $T_{\text{LB}} \geq e^{\Omega(1/\sigma^2)}$  when  $\max_{i,j} |\nu_i - \nu_j|$  is a positive absolute constant.

Strength of beliefs expressed by a particular truncated Gaussian can be usefully interpreted as the number of samples inherent therein. Indeed, any Gaussian distribution with variance  $\sigma^2$  can be represented as a Bayesian update of a unit-variance Gaussian given  $M = \Theta(\sigma^{-2})$  independent data points. Thus, the sample complexity in Corollary 6.6 is exponential in  $M$ .

We obtain a similar result for Beta priors. We define the *strength* of distribution  $\mathcal{P} = \text{Beta}(a, b)$  as the sum  $a + b$ . When  $a, b \in \mathbb{N}$ , the number of data points needed to obtain  $\mathcal{P}$  as a posterior starting from a uniform prior  $\text{Beta}(1, 1)$  is  $a + b - 2$ ; we interpret it as the strength of beliefs expressed by the prior. Note that  $M = \text{strength}(\mathcal{P})$  approximately captures variance: indeed,  $1/\text{Var}(\mathcal{P}) \leq [M, M^2]$ . A problem instance is called *M-strong* if  $\text{strength}(\mathcal{P}_i) \equiv M$ . We prove that  $T_{\text{UB}}(N_{\text{TS}})$  is polynomial in  $K$  for any fixed  $M$ , and exponential in  $M$  even for  $K = 2$ ; the latter dependence is inevitable.

**Corollary 6.7.** Suppose all priors  $\mathcal{P}_1, \dots, \mathcal{P}_K$  are Beta distributions.

- (a)  $T_{\text{UB}}(N_{\text{TS}}) \leq K^3 \cdot (\min(K, M))^{O(M)}$  if  $\text{strength}(\mathcal{P}_i) \leq M$  for all arms  $i$ .
- (b)  $T_{\text{LB}} \geq (\min(K, M))^{\Omega(M)}$  for **some** *M-strong* problem instance.
- (c)  $T_{\text{LB}} \geq 2^{\Omega(M)}$  for **any** *M-strong* problem instance such that  $\mu_1^0 - \mu_K^0 \geq \Omega(1)$ .  
More generally, this holds whenever arms  $i \neq j$  have strength at least  $M$  and  $|\mu_i^0 - \mu_j^0| \geq \Omega(1)$ .

The lower bound in part (b) holds if all arms  $i < K$  have the “smallest” prior  $\mathcal{P}_i = \text{Beta}(M - 1, 1)$ , and arm  $K$  has the “largest” prior  $\mathcal{P}_K \sim \text{Beta}(1, M - 1)$ .

In fact, (slightly weaker versions of) Corollaries 6.6(a) and 6.7(a) can be derived from Theorem 6.5. This is because the corresponding collections of priors are  $\delta$ -easy and  $\delta$ -non-dominant.

**Lemma 6.8.** The collection of all truncated Gaussians  $\tilde{N}(m, \sigma^2)$ ,  $m \in [0, 1]$  with a fixed variance  $\sigma^2 < 1$  is  $\delta$ -easy and  $\delta$ -non-dominant with  $\delta = e^{-\Omega(1/\sigma^2)}$ . Likewise, the collection of all Beta distributions of strength at most  $M \geq 1$  is  $\delta$ -easy and  $\delta$ -non-dominant with  $\delta = M^{-O(M)}$ .

## 6.5 One well-known arm

Let us consider an important special case when some arm  $\ell$  represents a well-known, “default” alternative, and the other arms are new to the agents. Put differently, agents have strong beliefs on one arm but not on all others. We would like to remove the dependence on the well-known arm as much as possible. We obtain two results of this flavor: a version of Corollary 6.3 on polynomially matching upper/lower bounds, and a version of Corollary 6.7 on Beta priors. We strengthen Corollary 6.3 under a mild non-degeneracy condition which ensures that the priors are not extremely concentrated near 1.

**Corollary 6.9.** Eq. (6.4) holds if  $\text{Var}(\mathcal{P}_i) \geq \sigma^2$  for all but one arm  $i$ , provided that

$$\Pr[\mu_i \leq 1 - \sigma] \geq e^{-1/\sigma} \quad \text{for all arms } i. \quad (6.9)$$

We obtain a stronger result focusing on Beta priors. Then the exponential dependence on  $\text{strength}(\mathcal{P}_i)$  holds only for arms  $i \neq \ell$ , whereas the dependence on  $\text{strength}(\mathcal{P}_\ell)$  is only polynomial. There is also a “penalty term” which scales exponentially in  $\ell$ ; this is mild if  $\ell$  is small.



**Corollary 6.10.** *Suppose all priors  $\mathcal{P}_1, \dots, \mathcal{P}_K$  are Beta distributions. Suppose  $\text{strength}(\mathcal{P}_\ell) = M$  for some arm  $\ell$ , and  $\text{strength}(\mathcal{P}_i) \leq m$  for all other arms  $i$ , where  $M \geq m \geq 2$ . Then*

$$T_{\text{UB}}(N_{\text{TS}}) \leq M^2 \cdot K^{O(1)} \cdot \max \left( m, (1 - \mu_\ell^0)^{-1}, (\mu_\ell^0)^{-(\ell-1)} \right)^{O(m)}.$$

*In particular,  $T_{\text{UB}}(N_{\text{TS}}) \leq M^2 \cdot K^{O(1)} \cdot m^{O(\ell m)}$  if  $\mu_\ell^0 \in [1/4, 3/4]$ .*

While dependency on *one* well-known arm can be mitigated, the lower bound in Lemma 6.8(c) rules out a similar improvement if  $\text{strength}(\mathcal{P}_i) \geq M$  for two or more arms  $i$ .

We obtain a particularly clean characterization for  $K = 2$  arms, which is worth stating explicitly.

**Corollary 6.11.** *Assume  $K = 2$  arms and Beta priors  $\mathcal{P}_i \sim \text{Beta}(a_i, b_i)$ . Let  $M = \max_i(a_i + b_i)$  and  $m = \min_i(a_i + b_i) \geq 2$ . Suppose  $\mu_i \in [\varepsilon, 1 - \varepsilon]$  for some  $\varepsilon > 0$  and both arms  $i$ . Then*

$$T_{\text{UB}}(N_{\text{TS}}) \leq M^2 \cdot \max(m, 1/\varepsilon)^{O(m)}.$$

*Moreover,  $L_{\text{LB}} \geq 2^{\Omega(m)}$  provided that  $\mu_1^0 - \mu_2^0 \geq \Omega(1)$ .*

## 7 Extensions via Improved Algorithms

**IMPROVED ALGORITHM FOR “EASY” PROBLEM INSTANCES.** We design a new algorithm for collecting  $N$  samples of each arm, for any given  $N$  (Algorithm 3 in Appendix F). This algorithm is tailored to the “easy” case of the polynomial vs. exponential dichotomy in Section 6.3. Formally, a problem instance is called  $\delta$ -easy,  $\delta > 0$  if the collection of priors  $(\mathcal{P}_1, \dots, \mathcal{P}_K)$  is  $\delta$ -easy and  $\delta$ -non-dominant. We obtain linear dependence on  $K$ , which is obviously optimal for a fixed  $N$ , along with computational efficiency.

**Theorem 7.1.** *Given a  $\delta$ -easy problem instance with  $K$  arms, Algorithm 3 is BIC and collects at least  $N$  samples of each arm almost surely in  $\tilde{O} \left( \frac{KN}{\delta} + \frac{K}{\delta^4} \right)$  rounds, for any desired  $N \in \mathbb{N}$ . The running time for each round is  $O(1)$  plus one call to “exploitation” given up to  $\tilde{O}(\delta^{-2})$  samples per arm.*

Thus, going back to the setup in Section 6.3, the dependence on  $K$  admits a crisp *linear vs. exponential dichotomy* for an arbitrary collection  $\mathcal{C}$ . Corollaries for Beta and truncated Gaussian priors, also with linear dependence on  $K$ , follow from Lemma 6.8. In particular for strength- $M$  Beta priors, using Corollary 4.6 shows that we efficiently achieve  $T_{\text{OPT}}(N_{\text{TS}}) \leq \tilde{O}(K \cdot M^{O(M)})$ .

The main insight behind Algorithm 3 is that for  $\delta$ -easy instances the “suitable” policy  $\pi_j$  does not need to depend on the samples from arms  $i < j$ . We choose an arm  $j_0$  uniformly at random and explore it almost surely like in `ExponentialExploration`, but using policy  $\pi_j$  as above. Afterwards we randomize between a padded phase for arm  $j_0$  or an exploration phase for a randomly chosen arm  $i \neq j_0$ . This algorithm is BIC despite only going through the exponential growth process for a single arm  $j_0$ . It is efficient because we are able to let  $\pi_j$  be a form of exploitation.

**FINE-TUNING THE MAIN ALGORITHM.** Algorithm 1 is somewhat wasteful when many samples are desired, i.e.  $N \gg \max(N_{\text{pad}}, N_{\text{boot}})$ . Indeed, only  $N_{\text{boot}}$  samples are needed for the bootstrapping phase to be BIC, and only  $N_{\text{pad}}$  samples are needed for each iteration of the `while` loop to collect  $N_{\text{pad}}$  samples of arm  $j$  almost surely. After that, the remaining samples for arm  $j$  can be collected more efficiently: indeed,

one can directly randomize between an exploration phase and a padded phase. As we show in Appendix E, these modifications reduce the number of rounds for collecting  $N$  samples of each arm to

$$T_{\text{UB2}}(N) = O\left(K G_{\text{pad}}^{-1} \left(N_{\text{pad}} \log(G_{\text{pad}}^{-1} p_{\text{boot}}^{-1}) + N_{\text{boot}} + N\right)\right). \quad (7.1)$$

**EFFICIENT COMPUTATION FOR BETA PRIORS.** We also present a computationally efficient version of the main algorithm. We focus on the special case of Beta priors of strength at most  $M$ , and recover the statistical guarantee in Corollary 6.7(a). The bottleneck is to compute a  $(j, \lambda, N)$ -suitable policy  $\pi_j$ , as all other steps are computationally efficient.<sup>17</sup> Such policy can then be plugged into Algorithm 1. The details can be found in Appendix G.

The key idea is that if we stochastically increase the priors for  $\mu_1, \dots, \mu_{j-1}$  and keep the padding  $\lambda$  fixed, this only makes our task more difficult. (Formalizing this involves a coupling argument between “true” and “artificial” data.) Thus, we reduce the problem to one in which all priors are  $\text{Beta}(1, M)$ . Now that all arms  $i \in [j-1]$  are i.i.d., a symmetry argument shows that the only convex combination  $q$  of these arms that we need to consider in the recommendation game (see Eq. (5.5) in the proof of Lemma 5.3) is the unweighted average. Consequently, the problem of finding a maximin policy for the  $j$ -recommendation game reduces to competing against this unweighted average, which can be done via a simple comparison.

**Theorem 7.2.** *Fix arm  $j$ , the number of arms  $K$ , parameters  $N, M \in \mathbb{N}$  and padding  $\lambda > 0$ . Suppose there exists a  $(j, N)$ -informed policy  $\pi_j$  which is BIC and  $(j, \lambda)$ -padded BIC for all problem instances with  $K$  arms and Beta priors of strength at most  $M$ . Then there is policy  $\pi_j^{\text{eff}}$  with these properties which can be computed efficiently, namely in time  $\text{poly}(K, M, N)$ . In particular, one can take  $\lambda = (\min(K, M))^{O(M)}$ . Plugging these  $\pi_j^{\text{eff}}$  and  $\lambda$  into Algorithm 1 yields the guarantee in Corollary 6.7(a).*

## 8 Explorability Characterization

We prove that the pairwise non-dominance condition (3.2) is in fact necessary, and use this condition to characterize which arms can be explored. Our result is modulo a minor non-degeneracy assumption: an arm is called *support-degenerate* if its true mean reward (according to the prior) is always in the set  $\{x, 1\}$  for some  $x \in [0, 1)$ . The proof is relatively simple: all the “heavy lifting” is done in the algorithmic results.

**Theorem 8.1.** *Suppose all arms are not support-degenerate, in the sense defined above. Let  $S$  be the set of all arms  $i$  which satisfy the pairwise non-dominance condition (3.2).*

- (a) *Only arms in  $S$  can be explored. More formally: for each arm  $i$ , if there exists a BIC algorithm which explores this arm with a positive probability, then  $i \in S$ .*
- (b) *All our algorithms can be restricted to  $S$ . More formally: if an algorithm is guaranteed to be BIC under (3.2), then this algorithm remains BIC if it is restricted to the arms in  $S$ .*
- (c) *All arms in  $S$  can be explored. More formally: there is a BIC algorithm which explores all arms in  $S$  with probability 1 within some finite time  $t^*$  depending on the prior.*

*Proof.* First we show that any arm  $i \notin S$  is not explorable. Let arm  $j$  be a “reason” why  $i \notin S$ , i.e., suppose  $\mu_i^0 \leq m_j$  for some arm  $j \neq i$ , where  $m_j$  is the minimum value in the support of  $\mu_j$ . For sake of contradiction,

<sup>17</sup>The construction of  $\text{transform}(\pi_j)$  in Lemma 5.4 is by computationally efficient reduction to  $\pi_j$ , so this causes no issues.

let  $t$  be the first time at which arm  $i$  is explored with positive probability by a BIC algorithm. Then  $\mathbb{E}^t[\mu_i] = \mu_i^0 \leq m_j < \mathbb{E}^t[\mu_j]$ , contradiction.

The last inequality is strict by non-degeneracy: since arm  $j$  is not support-degenerate, its posterior reward distribution cannot be a point mass at  $m_j$ . Indeed, the only values of  $\mu_j$  that can be ruled out with probability 1 in finite time are 0 and 1 (because of Bernoulli rewards). Thus, if the posterior on  $\mu_j$  is a point mass on  $m_j$ , then  $\mu_j$  must have support in  $\{0, m_j, 1\}$ . If the support includes 0, then  $m_j = 0$  by definition, and the support in  $\{m_j, 1\}$  is ruled out by non-degeneracy. This proves part (a).

Parts (b) and (c) easily follow. For part (b), consider the algorithm restricted to the arms in  $S$ . If the BIC condition is ever violated, it can only be because some arm  $j \notin S$  would be preferred by an agent. Thus, we obtain a BIC algorithm which explores arm  $j$  with positive probability (by recommending arm  $j$  for the same history), which contradicts part (a). Part (c) follows by applying part (b) to Algorithm 1.  $\square$

## References

- Shipra Agrawal and Navin Goyal. Analysis of Thompson Sampling for the multi-armed bandit problem. In *25th Conf. on Learning Theory (COLT)*, 2012.
- Shipra Agrawal and Navin Goyal. Further optimal regret bounds for thompson sampling. In *16th Intl. Conf. on Artificial Intelligence and Statistics (AISTATS)*, pages 99–107, 2013.
- Noga Alon and Joel Spencer. *The Probabilistic Method*. Wiley-Interscience Series in Discrete Mathematics and Optimization. John Wiley & Sons, New York, 4th edition, 2016.
- Krzysztof R. Apt. A primer on strategic games., 2011.
- Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256, 2002.
- Gal Bahar, Rann Smorodinsky, and Moshe Tennenholtz. Economic recommendation systems. In *16th ACM Conf. on Electronic Commerce (ACM-EC)*, 2016.
- Gal Bahar, Rann Smorodinsky, and Moshe Tennenholtz. Social learning and the innkeeper’s challenge. In *ACM Conf. on Economics and Computation (ACM-EC)*, pages 153–170, 2019.
- Dirk Bergemann and Stephen Morris. Information design: A unified perspective. *Journal of Economic Literature*, 57(1):44–95, March 2019.
- Kostas Bimpikis, Yiangos Papanastasiou, and Nicos Savva. Crowdsourcing exploration. *Management Science*, 64(4):1477–1973, 2018.
- Sébastien Bubeck and Nicolo Cesa-Bianchi. Regret Analysis of Stochastic and Nonstochastic Multi-armed Bandit Problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012. Published with *Now Publishers* (Boston, MA, USA). Also available at <https://arxiv.org/abs/1204.5721>.
- Sébastien Bubeck and Che-Yu Liu. Prior-free and prior-dependent regret bounds for thompson sampling. In *26th Advances in Neural Information Processing Systems (NIPS)*, pages 638–646, 2013.
- Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge University Press, Cambridge, UK, 2006.

- Yeon-Koo Che and Johannes Hörner. Recommender systems as mechanisms for social learning. *Quarterly Journal of Economics*, 133(2):871–925, 2018. Working paper since 2013, titled ‘Optimal design for social learning’.
- Bangrui Chen, Peter I. Frazier, and David Kempe. Incentivizing exploration by heterogeneous users. In *Conf. on Learning Theory (COLT)*, pages 798–818, 2018.
- James Allen Fill and Motoya Machida. Stochastic monotonicity and realizable monotonicity. *Annals of probability*, pages 938–978, 2001.
- Peter Frazier, David Kempe, Jon M. Kleinberg, and Robert Kleinberg. Incentivizing exploration. In *ACM Conf. on Economics and Computation (ACM-EC)*, 2014.
- John Gittins, Kevin Glazebrook, and Richard Weber. *Multi-Armed Bandit Allocation Indices*. John Wiley & Sons, Hoboken, NJ, USA, 2nd edition, 2011.
- Benjamin Golub and Evan D. Sadler. Learning in social networks. In Yann Bramoullé, Andrea Galeotti, and Brian Rogers, editors, *The Oxford Handbook of the Economics of Networks*. Oxford University Press, 2016.
- Johannes Hörner and Andrzej Skrzypacz. Learning, experimentation, and information design. In Bo Honoré, Ariel Pakes, Monika Piazzesi, and Larry Samuelson, editors, *Advances in Economics and Econometrics: 11th World Congress*, volume 1, page 63–98. Cambridge University Press, 2017.
- Nicole Immorlica, Jieming Mao, Aleksandrs Slivkins, and Steven Wu. Bayesian exploration with heterogeneous agents. In *The Web Conference (formerly known as WWW)*, pages 751–761, 2019.
- Nicole Immorlica, Jieming Mao, Aleksandrs Slivkins, and Steven Wu. Incentivizing exploration with selective data disclosure. In *ACM Conf. on Economics and Computation (ACM-EC)*, 2020. Working paper available at <https://arxiv.org/abs/1811.06026>.
- Emir Kamenica. Bayesian persuasion and information design. *Annual Review of Economics*, 11(1):249–272, 2019.
- Emilie Kaufmann, Nathaniel Korda, and Rémi Munos. Thompson sampling: An asymptotically optimal finite-time analysis. In *23rd Intl. Conf. on Algorithmic Learning Theory (ALT)*, pages 199–213, 2012.
- Ilan Kremer, Yishay Mansour, and Motty Perry. Implementing the “wisdom of the crowd”. *J. of Political Economy*, 122(5):988–1012, 2014. Preliminary version in *ACM EC 2013*.
- Tze Leung Lai and Herbert Robbins. Asymptotically efficient Adaptive Allocation Rules. *Advances in Applied Mathematics*, 6:4–22, 1985.
- Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, Cambridge, UK, 2020. Versions available at <https://banditalgs.com/> since 2018.
- Yishay Mansour, Aleksandrs Slivkins, and Vasilis Syrgkanis. Bayesian incentive-compatible bandit exploration. In *16th ACM Conf. on Economics and Computation (ACM-EC)*, pages 565–582, 2015.
- Yishay Mansour, Aleksandrs Slivkins, and Vasilis Syrgkanis. Bayesian incentive-compatible bandit exploration. *Operations Research*, 68(4):1132–1161, 2020. Preliminary version in *ACM EC 2015*.

- Yishay Mansour, Aleksandrs Slivkins, Vasilis Syrgkanis, and Steven Wu. Bayesian exploration: Incentivizing exploration in Bayesian games. *Operations Research*, 70(2), 2022. Preliminary version in *ACM EC 2016*.
- Jean-Christophe Mourrat. Free energy upper bound for mean-field vector spin glasses. *arXiv preprint arXiv:2010.09114*, 2020.
- Daniel Russo and Benjamin Van Roy. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243, 2014.
- Daniel Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, and Zheng Wen. A tutorial on thompson sampling. *Foundations and Trends in Machine Learning*, 11(1):1–96, 2018. Published with *Now Publishers* (Boston, MA, USA). Also available at <https://arxiv.org/abs/1707.02038>.
- Aleksandrs Slivkins. Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning*, 12(1-2):1–286, November 2019. Published with *Now Publishers* (Boston, MA, USA). Also available at <https://arxiv.org/abs/1904.07272>. Latest online revision: Jan 2022.
- William R. Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3-4):285–294, 1933.

## A Tools from Probability

### A.1 Fortuin-Kasteleyn-Ginibre (FKG) inequality for correlation

**Lemma A.1.** (Alon and Spencer, 2016, Theorem 6.2.1)[FKG Inequality] Consider measures  $\nu_1, \dots, \nu_n$  on  $\mathbb{R}$ , and let  $\nu = \prod_{i \in [n]} \nu_i$  be the product measure on  $[0, 1]^n$ . Suppose  $f, g : \mathbb{R}^n \rightarrow [0, 1]$  are functions which are increasing in each coordinate. Then  $\mathbb{E}[f \cdot g] \geq \mathbb{E}[f] \cdot \mathbb{E}[g]$ , where  $\mathbb{E}$  denotes expectation relative to  $\nu$ . Likewise, if  $f$  is coordinate-wise increasing and  $g$  is coordinate-wise decreasing then  $\mathbb{E}[f \cdot g] \leq \mathbb{E}[f] \cdot \mathbb{E}[g]$ .

**Remark A.2.** In fact if  $f, g$  are both increasing or both decreasing in each coordinate, then the conclusion above still holds as we can simply negate some coordinates in the parametrization, i.e. view  $f, g$  as increasing functions of  $-x_i$  instead of decreasing functions of  $x_i$ . We will refer to this as the *mixed-monotonicity* FKG inequality to highlight the slight subtlety in its application.

**Corollary A.3.** Suppose  $f, g : \mathbb{R} \rightarrow [0, 1]$  are increasing, and  $X$  is a random variable. Then  $\mathbb{E}[f(X)g(X)] \geq \mathbb{E}[f(X)] \cdot \mathbb{E}[g(X)]$ . If  $f$  is increasing and  $g$  is decreasing then  $\mathbb{E}[f(X)g(X)] \leq \mathbb{E}[f(X)] \cdot \mathbb{E}[g(X)]$ .

The above corollary is sometimes known as the Chebyshev inequality. Not to conflate it with the better known Chebyshev inequality from probability theory, we will also call it the FKG inequality in this work.

**Corollary A.4.** Let  $\mathcal{G}$  be a static  $\sigma$ -algebra and  $F$  a non-negative  $\mathcal{G}$ -measurable function which is increasing in the posterior mean of each arm  $a_s$  for  $s \in S \subseteq [K]$ . Then  $F$  is positively correlated with any coordinate-wise increasing function of the true values  $(\mu_s)_{s \in S}$ .

*Proof.* If  $\mu_s$  increases, this stochastically increases the (binomially distributed) empirical mean of arm  $a_s$ , hence increases the expectation of  $F$ . Hence  $\mathbb{E}[F | (\mu_s)_{s \in S}]$  is a coordinate-wise increasing function, and therefore by FKG is positively correlated with any other coordinate-wise increasing function of  $(\mu_s)_{s \in S}$ .  $\square$

### A.2 Bayesian concentration: proof of Lemma 3.1

**Lemma A.5.** Suppose there exists an estimator  $\theta = \theta(\gamma)$  for a parameter  $\mu$  given an observation signal  $\gamma$  which satisfies a concentration inequality  $\Pr[|\theta - \mu| \geq \varepsilon] \leq \delta$ . Assume we start with a prior  $P$  over  $\mu$  before observing  $\gamma$ . Then if  $\mu$  is in fact chosen according to  $P$  and  $\hat{\mu}$  is a sample from the posterior distribution for  $\mu$  conditional on  $\gamma$ , we have  $\Pr[|\hat{\mu} - \mu| \geq 2\varepsilon] \leq 2\delta$ .

*Proof.* By assumption,  $\Pr[|\theta - \mu| > \varepsilon] \leq \delta$ . Also,  $(\mu, \theta)$  and  $(\hat{\mu}, \theta)$  are identically distributed; choosing  $\hat{\mu}$  amounts to resampling  $\mu$  from the joint law of  $(\mu, \theta)$ . Therefore,  $\Pr[|\theta - \hat{\mu}| > \varepsilon] \leq \delta$ . Combining these two inequalities gives the Lemma by triangle inequality.  $\square$

*Proof of Lemma 3.1.* For simplicity show the claimed inequalities for  $\mu_i$ , explaining the generalization to weighted averages  $\mu_q$  at the end. First we apply Lemma A.5 with  $\theta_i$  being the empirical mean from the first  $\varepsilon^{-2}$  samples of arm  $i$  - this is  $\mathcal{F}_t$  measurable by assumption. With this choice of estimator  $\theta_i$ , for any true mean  $\mu_i$  the ordinary Chernoff bound implies  $\Pr[|\theta_i - \mu_i| \geq r\varepsilon] \leq C_0 e^{-r^2/C_0}$  for some absolute constant  $C_0$ . Then (3.5) follows from Lemma A.5.

We now deduce the second claim (3.6) from the first. Indeed (3.5) combined with the tail-decay definition of sub-Gaussianity implies that  $\frac{\hat{\mu}_i - \mu_i}{\varepsilon}$  is  $O(1)$  sub-Gaussian (averaged over all randomness in the problem). Moreover it manifestly has mean 0. Also,

$$\mathbb{E}[\mu_i | \mathcal{F}_t] - \mu_i = \mathbb{E}[\hat{\mu}_i - \mu_i | \sigma(\mathcal{F}_t, \mu_i)].$$



Using the exponential moment definition of sub-Gaussianity, the fact that  $\frac{\mathbb{E}[\mu_i|\mathcal{F}_i] - \mu_i}{\varepsilon}$  is the conditional expectation of an  $O(1)$  sub-Gaussian variable implies it is also  $O(1)$  sub-Gaussian. (Indeed, exponential moments always decrease under conditional expectation by Jensen's inequality.) Using again the tail-decay definition of sub-Gaussianity completes the proof of (3.6).

The extension to  $\mu_q$  poses no additional challenge because we may apply Lemma A.5 to  $\mu_q = \sum_i q_i \mu_i$ , and the estimator  $\theta_q = \sum_i q_i \theta_i$  obeys exactly the same Chernoff bound.  $\square$

### A.3 Stochastic and MLR Domination

We use a characterization of multivariate (first-order) stochastic domination (Fill and Machida, 2001).

**Lemma A.6.** *Given two probability distributions  $\nu, \nu'$  on  $\mathbb{R}^n$ , the following are equivalent:*

1. *For any coordinate-wise increasing function  $f$  we have  $\mathbb{E}^\nu[f] \geq \mathbb{E}^{\nu'}[f]$ .*
2. *There exists a distribution over pairs  $(X, X') \in \mathbb{R}^n \times \mathbb{R}^n$  such that  $X \geq X'$  coordinate-wise almost surely, and  $X \sim \nu, X' \sim \nu'$ .*

*In both cases we will say  $\nu$  stochastically dominates  $\nu'$ . We may also say that  $\nu$  is stochastically larger than  $\nu'$ , call a change from  $\nu$  to  $\nu'$  a stochastic decrease, etc.*

In the 1-dimensional case, there is a *canonical monotone coupling* which is easy to compute when the CDF functions of  $\nu, \nu'$  are known. It can be realized explicitly as the random pair  $(X, X') = (F_1^{-1}(u), F_2^{-1}(u))$  where  $F_1$  is the CDF for  $\nu$ ,  $F_2$  is the CDF for  $\nu'$ , and  $u \in [0, 1]$  is uniformly random. In the case that  $\nu, \nu'$  have finite support, it is easy to sample  $X'$  conditionally on  $X$  from this coupling by taking

$$u \in \left[ \lim_{\varepsilon \downarrow 0} F_1(X - \varepsilon), F_1(X) \right) \quad (\text{A.1})$$

uniformly at random and setting  $X' = F_2^{-1}(u)$ . This is because (A.1) holds if and only if  $F^{-1}(u) = X$  (when we take CDFs and inverse-CDFs to be right-continuous). When  $\nu$  stochastically dominates  $\nu'$ , this coupling satisfies  $\Pr[X \geq X'] = 1$ . For general random variables, these couplings always have the property that  $X, X'$  are synchronized, i.e. if  $(X_1, X'_1), (X_2, X'_2)$  are two samples from the coupled distribution, then  $X_1 \geq X_2$  implies  $X'_1 \geq X'_2$  with probability 1. See e.g. Mourrat (2020) for more discussion on monotone couplings.

We give several related lemmas in the 1-dimensional case. These results will allow us to construct couplings between various posterior mean distributions, for example in the construction of `transform`( $\pi_j$ ) in Algorithm 1.

**Definition 1.** Given mutually absolutely continuous probability distributions  $\nu, \nu'$  on  $\mathbb{R}^1$ , we say  $\nu$  MLR-dominates  $\nu'$  if the Radon-Nikodym derivative function  $f$  with  $\nu(dx) = g(x)d\nu'(x)$  is increasing.

**Lemma A.7.** *If  $\nu$  MLR-dominates  $\nu'$  then  $\nu$  also stochastically dominates  $\nu'$ .*

*Proof.* Let  $f$  be an increasing function and  $g = \frac{d\nu}{d\nu'}$  the Radon-Nikodym derivative. We have:

$$\frac{\int_{\mathbb{R}} f(x)d\nu(x)}{\int_{\mathbb{R}} f(x)d\nu'(x)} = \frac{\int_{\mathbb{R}} f(x)g(x)d\nu'(x)}{\int_{\mathbb{R}} f(x)d\nu'(x) \cdot \int_{\mathbb{R}} g(x)d\nu'(x)}.$$

We used the fact  $\int_{\mathbb{R}} g(x)d\nu'(x)$  since  $\nu, \nu'$  are probability measures. The result follows by FKG.  $\square$

We now show some comparison results on posterior mean distributions in 1 dimension. In the below,  $\mu, \mu' \in [0, 1]$  always denote the mean reward of bandit arms  $a, a'$  which outputs Bernoulli rewards.

**Lemma A.8.** *Suppose  $\mu \sim \mathcal{P}$  and  $\mu' \sim \mathcal{P}'$  where distribution  $\mathcal{P}$  MLR-dominates distribution  $\mathcal{P}'$ . Then if we observe 1 sample from each arm and the reward of arm  $a$  is at least that of arm  $a'$ , the posterior distribution of  $\mu$  continues to MLR-dominate that of  $\mu'$ .*

*Proof.* Receiving a reward corresponds to changing  $\mathcal{P}(dx) \rightarrow cx\mathcal{P}(dx)$  for some appropriate constant  $c$ . Receiving no rewards corresponds to  $\mathcal{P}(x) \rightarrow c'(1-x)\mathcal{P}(dx)$ . If  $a$  receives at least the reward of  $a'$ , the MLR property is preserved so the result follows.  $\square$

**Corollary A.9.** *Let  $\mu \sim \mathcal{P}$ , and let  $\mathcal{P}_0, \mathcal{P}_1$  be the posterior distributions for  $\mu$  after observing 0, 1 respectively reward from 1 sample. Then  $\mathcal{P}$  MLR-dominates  $\mathcal{P}$  which MLR-dominates  $\mathcal{P}_0$ .*

*Proof.* The result follows from first principles, since observing a reward again corresponds to changing  $\mathcal{P}(dx) \rightarrow cx\mathcal{P}(dx)$  for some appropriate constant  $c$ . Alternatively, since  $\mathcal{P}$  is a mixture of  $\mathcal{P}_1, \mathcal{P}_0$  we can reason that  $\mathcal{P}_1$  stochastically dominates  $\mathcal{P}_0$  using Lemma A.8 and conclude.  $\square$

**Lemma A.10.** *Suppose  $\mu \sim \mathcal{P}$  and  $\mu' \sim \mathcal{P}'$  where  $\mathcal{P}$  MLR-dominates  $\mathcal{P}'$ . Let  $\mathcal{P}_N, \mathcal{P}'_N$  be the distributions of their posterior means after observing  $N$  samples each arm. Then  $\mathcal{P}_N$  stochastically dominates  $\mathcal{P}'_N$ .*

*Proof.* We reveal the reward of the samples for each arm simultaneously in  $N$  steps. We claim there exists a Markovian coupling between the two reward processes so that the posterior distribution of  $\mu$  always MLR-dominates that of  $\mu'$ . This holds by induction: as long as MLR-dominance holds at time  $t$ , the posterior mean of  $\mu$  is higher than that of  $\mu'$ , so we can couple the rewards so  $\mu$  gets at least as much reward. Applying Lemma A.8, this preserves the MLR property. At the end, the final posterior of  $\mu$  MLR-dominates that of  $\mu'$ . Therefore we have coupled the rewards for  $\mu, \mu'$  so that  $\mu$  always has higher posterior mean, implying the claimed stochastic domination.  $\square$

**Corollary A.11.** *Let  $\mu' \sim \text{Beta}(a, b)$  with  $1 \leq a, b \leq M$ , and  $\mu \sim \text{Beta}(1, M)$ . Then for any  $N$ , the distribution of the posterior mean of  $\mu$  after  $N$  samples stochastically dominates that of  $\mu'$ .*

*Proof.* The MLR-domination of  $\mu$  over  $\mu'$  implies the result.  $\square$

**Lemma A.12.** *Suppose  $\mu \sim \mathcal{P}$ . Let  $\mathcal{L}_N$  be the distribution for the posterior mean of  $\mu$  after  $N$  samples, and  $\mathcal{L}_{N, N_0}$  the distribution for the posterior mean of  $\mu$  after  $N$  samples, conditioned on the first  $N_0$  samples having 0 reward. Then  $\mathcal{L}_N$  stochastically dominates  $\mathcal{L}_{N, N_0}$ .*

*Proof.* If  $N_0 \geq N$  the result is clear so we assume  $N_0 \leq N$ . Suppose  $\mu'$  is an independent copy of  $\mu$  corresponding to an arm which observed  $N_0$  samples with no reward at the start. After observing the first  $N_0$  samples of  $\mu$ , we obtain a posterior distribution for  $\mu$  which MLR-dominates that of  $\mu'$ . Applying from this point Lemma A.10, we have established the result conditionally on any fixed values for the first  $N_0$  samples for  $\mu$ . Averaging over the random choice of first  $N_0$  samples now implies the result.  $\square$

#### A.4 Tails of sub-Gaussian distributions

**Lemma A.13.** *If random variable  $X$  is  $O(1)$ -sub-Gaussian and event  $E$  has probability  $\Pr[E] \leq p$ , then  $\mathbb{E}[|X \cdot 1_E|] \leq O(p\sqrt{\log(p^{-1})})$ .*

*Proof.* Let  $Y := X \cdot 1_E$ . Then we estimate:

$$\Pr[|Y| \geq t] \leq \min(p, \Pr[|X| \geq t]) \leq \min(p, e^{-\Omega(t^2)}).$$

Consequently,

$$\mathbb{E}[|Y|] \leq \int_{t=0}^{\infty} \Pr[|Y| \geq t] dt \leq \int_{t=0}^{\infty} \min(p, \Pr[|X| \geq t]) dt.$$

The integrand is  $t$  for  $t = O(\sqrt{\log(p^{-1})})$  and  $e^{-\Omega(t^2)}$  elsewhere. Combining gives the claimed bound.  $\square$

## B Initial Sampling: proofs for Section 5

### B.1 Existence of a Suitable Policy: proof of Lemma 5.3

We first formally define the  $j$ -recommendation game. We proceed more generally than in the main body, defining it relative to any  $K$ -tuple  $(N_1, \dots, N_K)$ . We recall the definition of the static  $\sigma$ -algebra  $\mathcal{G}_{N_1, \dots, N_K}$  which is generated by  $N_i$  samples of each arm  $i$ . When considering arm  $j$ , we will always have  $N_k = 0$  for all  $k > j$ . If  $N_i = N$  for all  $i$  we recover the  $(j, N)$ -recommendation game as defined in Section 5. The definition of  $(j, N)$ -informed generalizes readily to  $(j, \mathcal{G})$ -informed.

**Definition 2.** The  $j$ -recommendation game is a two-player zero sum game played between a *planner* and an *agent*. The players share a common independent prior over the true mean rewards  $(\mu_i)_{i \leq j}$ , and the planner gains access to the static  $\sigma$ -algebra  $\mathcal{G} = \mathcal{G}_{(N_1, \dots, N_j)}$ . The planner then either recommends arm  $j$  or does nothing. The agent picks a mixed strategy on  $[j-1]$ , resulting in the choice of an arm  $i \leq j-1$ . Suppose the arms chosen are  $a_i$  for the planner,  $a_k$  for the agent (where  $k = j$  or  $k$  is undefined). The payoff for the planner is  $(\mu_k - \mu_i) \cdot 1_{k=j}$ .

We remark we can without loss of generality view all  $j$ -recommendation game strategies as depending only on the posterior means of each arm condition on  $\mathcal{G}$ . In the below, we set  $\tilde{\mu}_i = \mathbb{E}[\mu_i | \mathcal{G}]$  for the relevant static  $\sigma$ -algebra  $\mathcal{G}$ . Given a planner strategy in the  $j$ -recommendation game, we naturally obtain a corresponding  $(j, \mathcal{G})$ -informed policy for our original problem in which we recommend arm  $j$  when the planner would, and recommend the  $\mathcal{G}$ -conditional-expectation-maximizing arm otherwise. The key point of this game is as follows.

**Lemma B.1.** A strategy for the planner in the  $j$ -recommendation game corresponds to a  $(j, \lambda)$ -padded BIC policy if and only if it has minimax value at least  $\lambda$ .

*Proof.* Because this policy only plays an arm  $i \neq j$  when  $\mathbb{E}[\mu_i | \mathcal{G}]$  is maximal among all arms, it suffices to show that recommending arm  $k$  is  $(j, \lambda)$ -padded BIC. This follows directly from the definition.  $\square$

**Lemma B.2.** The minimax value of the  $j$ -recommendation game with  $\sigma$  algebra  $\mathcal{G}$  is equal to

$$\inf_{q \in \Delta_{j-1}} \mathbb{E} [(\mathbb{E}[\mu_j - \mu_q | \mathcal{G}])_+] = \inf_{q \in \Delta_{j-1}} \mathbb{E} [(\tilde{\mu}_j - \tilde{\mu}_q)_+].$$

*Proof.* This expression is the value of the best response when the agent plays mixed strategy  $q$ . By the minimax theorem, the infimum over all  $q \in \Delta_{j-1}$  is exactly the value of the game.  $\square$

**Lemma B.3.** (Apt, 2011, Corollary 5.5) In any finite two-player zero-sum game there is a Nash equilibrium of mixed strategies which are not weakly dominated by any other mixed strategy.

The  $j$ -recommendation game can be viewed as a finite game because there are only finitely many datasets due to the Bernoulli reward assumption. Hence the above lemma applies.

**Definition 3.** We say a policy  $\pi$  (for the original bandit problem) is  $j$ -correlated if  $\mathbb{P}[\pi = j | (\tilde{\mu}_i)_{i \in [K]}]$  is increasing in  $\tilde{\mu}_j$ .

**Lemma B.4.** *If a strategy  $S$  for the  $j$ -recommendation game is not weakly dominated by any other policy, then the resulting policy for the original bandit problem is  $j$ -correlated.*

*Proof.* If  $\mathbb{P}[S = j | (\tilde{\mu}_i)_{i \in [K]}]$  is not increasing in  $\mathbb{E}[\mu_j | \mathcal{G}]$  then there exist two  $j$ -tuples

$$(e_1, \dots, e_j), (e_1, \dots, e_{j-1}, e'_j)$$

of values for  $(\tilde{\mu}_i)_{i \leq j}$  which agree except on  $e_j > e'_j$  and yet

$$\Pr[S = j | (\tilde{\mu}_i)_{i \in [K]} = (e_1, \dots, e_j)] < \Pr[S' \text{ recommends } a_j | (\tilde{\mu}_i)_{i \in [K]} = (e_1, \dots, e'_j)].$$

In this case we may decrease the probability to play  $a_j$  on former event that  $(\tilde{\mu}_i)_{i \in [K]} = (e_1, \dots, e_j)$ , and proportionally increase this probability given  $(\tilde{\mu}_i)_{i \in [K]} = (e_1, \dots, e'_j)$ , so that the total probability to play arm  $j$  is preserved. This strategy weakly dominates  $S'$ , yielding a contradiction.  $\square$

**Lemma B.5.** *If a strategy for the  $j$ -recommendation game is  $j$ -correlated and  $(j, 0)$ -padded, then the resulting  $(j, \mathcal{G})$ -informed policy is BIC.*

*Proof.* The BIC property against arms  $i < j$  is equivalent to the  $(j, 0)$ -padded property. By the first property of  $j$ -correlation it follows that the event of playing arm  $j$  is increasing in  $\mathbb{E}[\mu_j | \mathcal{G}]$ , hence the probability to play arm  $j$  is stochastically increasing in  $\mu_j$ . Therefore FKG implies that

$$\mathbb{P}[\mu_j | A_t = j] \geq \mu_j^0 \geq \mu_i^0$$

for any  $i > j$ . It remains to show the BIC property for playing arms other than arm  $j$ , and this follows from the fact that we simply exploit conditionally on  $\mathcal{G}$  whenever this happens.  $\square$

**Lemma B.6.** *Given a static  $\sigma$ -algebra  $\mathcal{G}$ , suppose there exists a  $\mathcal{G}$ -measurable  $(j, \lambda)$ -padded BIC strategy  $\hat{\pi}_j$  for  $\lambda \geq 0$ . Then there exists a  $\mathcal{G}$ -measurable  $(j, \lambda)$ -suitable BIC strategy  $\pi_j$ .*

*Proof.* The assumption implies there exists a  $j$ -recommendation game strategy  $S$ , and any such strategy yields a  $(j, \lambda)$ -padded BIC policy. By Lemma B.3 we may assume  $S$  is not weakly dominated. In such case, by Lemmas B.4 and B.5 the resulting bandit strategy  $\pi_j$  is BIC as well.  $\square$

Combining the previous lemmas, we obtain the main guarantee for  $(j, \lambda)$ -suitable strategies.

**Lemma B.7.** *Fix  $j \in [K]$  and  $\lambda > 0$ . Suppose there exists a static  $\sigma$ -algebra  $\mathcal{G} = \mathcal{G}_{N_1, N_2, \dots, N_k, 0, \dots, 0}$  which is independent of  $(\mu_{j+1}, \dots, \mu_K)$  and satisfies*

$$\mathbb{E} [(\mathbb{E}[\mu_j - \mu_q | \mathcal{G}])_+] \geq \lambda \quad \text{for all } q \in \Delta_{j-1}.$$

*Then there exists a  $\mathcal{G}$ -measurable  $(j, \lambda)$ -suitable strategy for the planner.*

*Proof.* Follows from combining Lemmas B.1, B.2, B.6.  $\square$

The next lemma upper bounds the number of samples we must include in  $\mathcal{G}$  to ensure  $\lambda \geq \Omega(G_{\text{pad}})$ .

**Lemma B.8.** Fix arm  $j$  and time  $T_0$ , and a sufficiently large universal constant  $C_{\text{pad}} \geq 5$ . Suppose for some  $q \in \Delta_{j-1}$  and  $\varepsilon_{\text{pad}}, \delta_{\text{pad}} > 0$ ,

$$\Pr [(\mu_j - \mu_q)_+ \geq \varepsilon_{\text{pad}}] \geq \delta_{\text{pad}}.$$

Let ALG be a BIC algorithm which by time  $T_0$  collects at least  $N_{\text{pad}}$  samples of each arm  $i \in [j]$  a.s., where

$$N_{\text{pad}} \geq C_{\text{pad}} \varepsilon_{\text{pad}}^{-2} (1 + \log(\varepsilon_{\text{pad}}^{-1} \cdot \delta_{\text{pad}}^{-1}))$$

Let  $\mathcal{G}_{N_{\text{pad}},k}$  be the  $\sigma$ -algebra generated by the first  $N_{\text{pad}}$  samples of each of these  $k$  arms. Then

$$\mathbb{E} \left[ \left( \mathbb{E} [\mu_j - \mu_q \mid \mathcal{G}_{N_{\text{pad}},k}] \right)_+ \right] \geq \frac{\mathbb{E}[(\mu_j - \mu_q)_+] - \delta_{\text{pad}} - \varepsilon_{\text{pad}}}{2}.$$

*Proof.* Applying Lemma 3.1 and using that  $C_{\text{pad}} \geq 5$  is sufficiently large, for any  $r \geq 1$  we have

$$\Pr \left[ \left| \mathbb{E}[\mu_q \mid \mathcal{G}_{N_{\text{pad}},j}] - \mu_q \right| \geq \frac{r\varepsilon_{\text{pad}}}{10} \right] \leq C_{\text{pad}} e^{-100C_{\text{pad}}^2 r^2 \log(\varepsilon_{\text{pad}}\delta_{\text{pad}})} \leq C_{\text{pad}} (\varepsilon_{\text{pad}}\delta_{\text{pad}})^{100C_{\text{pad}}r}.$$

This easily implies by integration that

$$\mathbb{E} \left[ \left( \left| \mathbb{E}[\mu_q \mid \mathcal{G}_{N_{\text{pad}},j}] - \mu_q \right| - \frac{\varepsilon_{\text{pad}}}{10} \right)_+ \right] \leq \frac{\varepsilon_{\text{pad}}\delta_{\text{pad}}}{10}.$$

Similarly,

$$\mathbb{E} \left[ \left( \left| \mathbb{E}[\mu_j \mid \mathcal{G}_{N_{\text{pad}},j}] - \mu_j \right| - \frac{\varepsilon_{\text{pad}}}{10} \right)_+ \right] \leq \frac{\varepsilon_{\text{pad}}\delta_{\text{pad}}}{10}.$$

However we have

$$\begin{aligned} \left( \mathbb{E} [\mu_j - \mu_q \mid \mathcal{G}_{N_{\text{pad}},j}] \right)_+ &\geq \left( \mu_j - \mu_q - \frac{\varepsilon_{\text{pad}}}{5} \right) \cdot 1_{\mu_j - \mu_q \geq \varepsilon_{\text{pad}}} - \left( \left| \mathbb{E}[\mu_j \mid \mathcal{G}_{N_{\text{pad}},j}] - \mu_j \right| - \frac{\varepsilon_{\text{pad}}}{10} \right)_+ \\ &\quad - \left( \left| \mathbb{E}[\mu_q \mid \mathcal{G}_{N_{\text{pad}},j}] - \mu_q \right| - \frac{\varepsilon_{\text{pad}}}{10} \right)_+. \end{aligned}$$

Note that  $(\mu_j - \mu_q - \frac{\varepsilon_{\text{pad}}}{5}) \cdot 1_{\mu_j - \mu_q \geq \varepsilon_{\text{pad}}} \geq \frac{(\mu_j - \mu_q)_+}{2} \cdot 1_{\mu_j - \mu_q \geq \varepsilon_{\text{pad}}}$ . Substituting and taking expectations, the conclusion of the lemma follows.  $\square$

We now obtain Lemma 5.3, guaranteeing the existence of the policies  $\pi_j$  needed in Algorithm 1.

*Proof of Lemma 5.3.* Letting  $\varepsilon_{\text{pad}} = \delta_{\text{pad}} = G_{\text{pad}}/3$  in Lemma B.8, the result follows immediately. To see that the assumption of Lemma B.8 holds note that  $(\mu_j - \mu_q)_+ \in [0, 1]$  almost surely. Therefore

$$\Pr[(\mu_j - \mu_q)_+ \geq G_{\text{pad}}/3] \geq G_{\text{pad}}/3$$

as long as  $\mathbb{E}[(\mu_j - \mu_q)_+] \geq G_{\text{pad}}$ .  $\square$

## B.2 A suitable policy for bootstrapping: proof of Lemma 5.4

We consider  $\text{transform}(\pi_j)$ , a transformed version of policy  $\pi_j$  which is, essentially,  $(j, \lambda, N)$ -suitable conditioned on  $\text{ZEROS}_{j, N_0}$ .

Existence of such a policy is straightforward, as the conditioning on  $\text{ZEROS}_{j, N_0}$  stochastically decreases the values  $(\tilde{\mu}_i)_{i < j}$ , hence increasing the minimax value of the  $j$ -recommendation game. However we give a reduction to  $\pi_j$ , so that to make Algorithm 1 efficient it essentially suffices to compute policies  $\pi_j$ .

We use the following fake data technique. Let  $\tilde{\mu}_i$  be the posterior mean  $\mathbb{E}[\mu_i | \mathcal{G}_{j, N}]$ , and  $\mathcal{L}_i$  be the distribution of  $\tilde{\mu}_i$ . Let  $\mathcal{L}'_i$  be the distribution of  $\tilde{\mu}_i$  conditioned on  $\text{ZEROS}_{j, N_0}$ . It follows from Lemma A.12 that  $\mathcal{L}_i$  is stochastically dominated by  $\mathcal{L}'_i$ . Hence by Lemma A.6 there exists a coupling  $(\tilde{\mu}_i, \hat{\mu}_i)$  with  $\tilde{\mu}_i \leq \hat{\mu}_i$  almost surely, with  $\tilde{\mu}_i \sim \mathcal{L}_i, \hat{\mu}_i \sim \mathcal{L}'_i$ .

We now define  $\text{transform}(\pi_j)$ . We view  $\pi_j = \pi_j(\tilde{\mu}_1, \dots, \tilde{\mu}_j)$  as a policy of posterior mean rewards; this makes no difference because for fixed prior on some  $\mu_i$  and a fixed number of samples, specifying the posterior mean  $\tilde{\mu}_i$  is equivalent to specifying the empirical average reward. We set  $\text{transform}(\pi_j)(\tilde{\mu}_1, \dots, \tilde{\mu}_j) = j$  when  $\pi_j(\hat{\mu}_1, \dots, \hat{\mu}_{j-1}, \tilde{\mu}_j) = j$ . That is, we use the couplings  $\mathcal{X}_i$  to generate overly optimistic posterior mean rewards for all arms except arm  $j$ , and then apply  $\pi_j$ . When  $\pi_j(\hat{\mu}_1, \dots, \hat{\mu}_{j-1}, \tilde{\mu}_j) \neq j$  we simply exploit based on the first  $N$  samples of each arm  $1, \dots, j$  (which may result in sampling arm  $j$  anyway).

We proceed to show  $\text{transform}(\pi_j)$  is  $(j, \lambda)$ -suitable. It suffices to show arm  $j$  recommendations are BIC since otherwise  $\text{transform}(\pi_j)$  exploits. Fix  $i < j$ . Below we let  $\pi_j$  denote the value  $\pi_j(\hat{\mu}_1, \dots, \hat{\mu}_{j-1}, \tilde{\mu}_j)$ . The key point is that the distribution of  $(\hat{\mu}_1, \dots, \hat{\mu}_{j-1}, \tilde{\mu}_j)$  conditioned on  $\text{ZEROS}_{j, N_0}$  is, by construction, the unconditioned distribution for  $(\tilde{\mu}_1, \dots, \tilde{\mu}_j)$ . Therefore, by the  $(j, \lambda)$ -padded BIC property of  $\pi_j$ :

$$\mathbb{E}[(\tilde{\mu}_j - \hat{\mu}_i) \cdot \mathbf{1}_{\pi_j=j} | \text{ZEROS}_{j, N_0}] \geq \lambda.$$

As  $\tilde{\mu}_i \geq \hat{\mu}_i$  almost surely we obtain

$$\mathbb{E}[(\tilde{\mu}_j - \tilde{\mu}_i) \cdot \mathbf{1}_{\pi_j=j} | \text{ZEROS}_{j, N_0}] \geq \lambda.$$

for each  $i < j$ . On the event  $\pi_j \neq j$  we always exploit, implying that

$$\mathbb{E}[(\tilde{\mu}_j - \tilde{\mu}_i) \cdot \mathbf{1}_{\pi_j \neq j} | \text{ZEROS}_{j, N_0}] \geq 0.$$

Adding, we obtain the padded BIC property conditional on  $\text{ZEROS}_{j, N_0}$ . Compared against arms  $i > j$ , the conditional BIC property for  $\text{transform}(\pi_j)$  follows directly from the unconditioned property of  $\pi_j$ . Indeed, the couplings restore the distribution of the first  $j - 1$  arm samples back to normal, and do not affect the remaining  $N + 1 - j$  arms. This concludes the proof.

## B.3 BIC property for bootstrapping: proof of Lemma 5.5

**Lemma 5.5.** *Consider Algorithm 1 with arbitrary parameters  $N_0 \leq N$  and  $\lambda > 0$ . Fix some arm  $j \geq 2$ . The bootstrapping phases are BIC as long as policy  $\pi_j$  is  $(j, \lambda, N)$ -suitable and parameter  $N_0$  satisfies*

$$\mathbb{E}[\mu_j - \mu_i | \text{ZEROS}_{j, N_0}] \geq 0 \quad \text{for all arms } i < j. \quad (5.9)$$

*Proof of Lemma 5.5.* The proof follows the same strategy as that of Lemma 5.2. Line 8 is BIC since it is just exploitation. For the next phase, recommending any arm  $i \neq j$  is BIC since it can only happen via exploitation, or via the BIC policy  $\text{transform}(\pi_j)$ . Hence we fix an arm  $i \neq j$  and show that arm  $j$  is BIC against arm  $i$  in this phase. We define  $\mathbf{1}_{\text{padded}} + \mathbf{1}_{\text{exploit}} + \mathbf{1}_{\text{explore}} = 1$  and  $\Lambda_{i,j}$  as in the proof of



Lemma 5.2. Note that we have probability  $\Pr[\text{ZEROS}_{j,N_0}] \cdot (1 - p_j)$  to reach line 12. Applying the guarantee of Lemma 5.4 conditional on this event, we obtain:

$$\begin{aligned}
\mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{padded}}] &= \mathbb{E}[(\mu_j - \mu_i) \cdot \mathbf{1}_{\text{transform}(\pi_j)=j} \cdot \mathbf{1}_{\text{padded}}] \\
&= \mathbb{E}[(\tilde{\mu}_j - \tilde{\mu}_i) \cdot \mathbf{1}_{\text{transform}(\pi_j)=j} \cdot \mathbf{1}_{\text{padded}}] \\
&\geq \mathbb{E}[(\tilde{\mu}_j - \hat{\mu}_i) \cdot \mathbf{1}_{\text{transform}(\pi_j)=j} \cdot \mathbf{1}_{\text{padded}}] \\
&\geq \Pr[\text{ZEROS}_{j,N_0}] \cdot \lambda(1 - p_j) \\
&= p_j.
\end{aligned}$$

The second equality above follows from the fact that both  $\mathbf{1}_{\text{transform}(\pi_j)=j}$  and  $\mathbf{1}_{\text{padded}}$  only depend on the first  $N$  samples of the first  $j$  arms and independent external randomness.

Moreover using a worst-case bound on exploration, and that exploitation is always BIC, we have:

$$\mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{explore}}] \geq -\mathbb{E}[\mathbf{1}_{\text{explore}}] = -p_j \quad \text{and} \quad \mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_{\text{exploit}}] \geq 0.$$

Adding, we obtain Equation 5.2 for  $i < j$ . For  $i > j$ , similar reasoning again leads to the same three inequalities with  $p_j$  replaced by 0, i.e.  $\mathbb{E}[\Lambda_{i,j} \cdot \mathbf{1}_E] \geq 0$  for  $\mathbf{1}_E \in \{\mathbf{1}_{\text{padded}}, \mathbf{1}_{\text{exploit}}, \mathbf{1}_{\text{explore}}\}$ . Altogether we have shown that playing arm  $j$  is BIC against any arm  $i$ , establishing the BIC property for the bootstrapping phases.  $\square$

## C Sample Complexity for Arbitrary Priors (for Section 6)

The purpose of this Appendix is to estimate the parameters  $G_{\text{pad}}, N_{\text{TS}}, N_{\text{pad}}, N_{\text{boot}}, \log(\Pr[\text{ZEROS}_{j,N_{\text{boot}}}]^{-1})$  appearing in the sample complexity upper bound of Theorem 5.6. In Theorem C.1 we achieve upper bounds based on the lower bounds  $K, N_{\text{boot}}, L_{\text{LB}}$  for  $T_{\text{UB}}(1)$  as well as a standard deviation lower bound  $\sigma$ . Then in Lemma C.3 we show upper bounds depending on  $\delta$  for  $\delta$ -easy and  $\delta$ -non-dominant problem instances.

In Theorem C.1 we use the following definition, a type of anti-concentration assumption on the priors. It is implied (up to constant factors) by all  $K$  arms having standard deviation at least  $\sigma$ . However it also allows for a single highly concentrated arm, as long as this arm has some chance to be away from 1.

**Definition 4.** The priors are called  $\sigma$ -non-degenerate if (i)  $\text{StdDev}(\mu_i) \geq \sigma$  for all but at most one arm  $i$ , and (ii)  $\Pr[\mu_i \in [0, 1 - \sigma]] \geq e^{-1/\sigma}$  for all arms  $i$ .

The next result is a more refined statement of Corollary 6.3, incorporating Corollary 6.9 as well via Definition 4.

**Theorem C.1.** *Assume the priors are  $\sigma$ -non-degenerate. Then Algorithm 1 achieves  $N_{\text{TS}}$  samples of each arm within this many rounds:*

$$T_{\text{UB}}(N_{\text{TS}}) = \tilde{O} \left( \frac{K^{9/2} L_{\text{LB}}^3 N_{\text{boot}}}{\sigma^4} + \frac{K^{5/2} L_{\text{LB}} N_{\text{boot}}^2}{\sigma^2} \right).$$

Moreover Algorithm 2 (defined later), achieves  $N_{\text{TS}}$  samples of each arm within this many rounds:

$$T_{\text{UB}}(N_{\text{TS}}) = \tilde{O} \left( \frac{K^{7/2} L_{\text{LB}}^3 N_{\text{boot}}}{\sigma^4} \right).$$

*Proof.* We use the guarantees from Theorem 5.6 and (7.1), and we upper-bound each of the relevant parameters in terms of  $L_{\text{LB}}, N_{\text{boot}}, K, \sigma$ . Recall from Theorem 4.1 that  $N_{\text{TS}} = C_{\text{TS}} \varepsilon_{\text{TS}}^{-2} \log \delta_{\text{TS}}^{-1}$ , where  $\delta_{\text{TS}} \geq \varepsilon_{\text{TS}}^K$  by FKG (see the proof of Lemma 4.5) so that  $N_{\text{TS}} = O\left(\frac{K \log(1/\varepsilon_{\text{TS}})}{\varepsilon_{\text{TS}}^2}\right)$ .

From the condition that  $\Pr[\mu_i \in [0, 1 - \sigma]] \geq e^{-K\sigma^{-1}}$ , it is not hard to obtain

$$\Pr[\text{ZEROS}_{j, N_{\text{boot}}}] \geq e^{-O(KN_{\text{boot}}\sigma^{-1})},$$

which implies

$$\log\left(\frac{1}{\Pr[\text{ZEROS}_{j, N_{\text{boot}}}]}\right) = O(KN_{\text{boot}}\sigma^{-1}).$$

Define

$$\hat{L} := \sup_{q \in \Delta, j \in [K]: q_j=0} \frac{1}{\mathbb{E}[(\mu_j - \mu_q)_+]}.$$

We have  $\frac{1}{G_{\text{pad}}} \leq \hat{L}$  as  $\hat{L}$  minimizes the same objective over a large set. We may take  $\varepsilon_{\text{TS}} = \frac{1}{10\hat{L}}$ , so that  $N_{\text{TS}} = \tilde{O}(K\hat{L}^2)$ . Recall that  $N_{\text{pad}} = O\left(\frac{\log(G_{\text{pad}}^{-1})}{G_{\text{pad}}^2}\right) = \tilde{O}(\hat{L}^2)$ . Combining we obtain the bounds

$$\tilde{O}\left(K^2 N_{\text{boot}} \hat{L} (K\hat{L}^2 + N_{\text{boot}})\sigma^{-1}\right) \quad \text{and} \quad \tilde{O}\left(K^2 \hat{L}^3 N_{\text{boot}}\sigma^{-1}\right)$$

on the number of rounds needed by Algorithms 1 and 2 respectively. Next we estimate:

$$\begin{aligned} \hat{L} &= \sup_{q \in \Delta_K, j \in [K]: q_j=0} \left( \frac{1}{\mathbb{E}[(\mu_j - \mu_q)_+]} \right) \\ &\leq \sup_{q \in \Delta_K, j \in [K]: q_j=0} \frac{\mathbb{E}[(\mu_j - \mu_q)_+ + (\mu_j - \mu_q)_-]}{\mathbb{E}[(\mu_j - \mu_q)_+]}. \quad \sup_{q \in \Delta_K, j \in [K]: q_j=0} \frac{1}{\mathbb{E}[(\mu_j - \mu_q)_+ + (\mu_j - \mu_q)_-]} \\ &\leq \sup_{q \in \Delta_K, j \in [K]: q_j=0} \frac{\mathbb{E}[(\mu_j - \mu_q)_+ + (\mu_j - \mu_q)_-]}{\mathbb{E}[(\mu_j - \mu_q)_+]}. \quad \sup_{q \in \Delta_K, j \in [K]: q_j=0} \frac{1}{\mathbb{E}[|\mu_j - \mu_q|]} \\ &\leq O\left(\frac{L_{\text{LB}}\sqrt{K}}{\sigma}\right). \end{aligned}$$

In the last step, we use the assumption that at most one  $\mu_i$  has standard deviation less than  $\sigma$  to bound the last supremum by  $O\left(\frac{\sqrt{K}}{\sigma}\right)$ . This bound follows because either  $\mu_k$  has standard deviation at least  $\sigma$ , or  $\mu_q$  has standard deviation at least  $\frac{\sigma}{\sqrt{K}}$ . The claimed bounds follow.  $\square$

**Lemma C.2.** *If  $X \in [0, 1]$  almost surely then  $\mathbb{P}[X \leq \mathbb{E}[X] + \varepsilon] \geq \frac{\varepsilon}{1+\varepsilon}$ .*

*Proof.* Let  $p = \mathbb{P}[X \leq \mathbb{E}[X] + \varepsilon]$ . Then  $\mathbb{E}[X - \mathbb{E}[X]] \geq p\varepsilon - (1-p)$ . However, this is false if  $p < \frac{\varepsilon}{1+\varepsilon}$ .  $\square$

In the next lemma we estimate all the parameters appearing in the guarantee of Theorem 5.6 for Algorithm 1 for  $\delta$ -easy and  $\delta$ -non-dominant problem instances.

**Lemma C.3.** *If  $\mathcal{C}$  is  $\delta$ -easy and  $\delta$ -non-dominant, then:  $N_{\text{TS}} = \tilde{O}(K\delta^{-2})$  and  $G_{\text{pad}} \geq \delta$  and  $N_{\text{pad}} = \tilde{O}(\delta^{-2})$  and  $N_{\text{boot}} = \tilde{O}(\delta^{-1})$  and  $\log(p_{\text{boot}}^{-1}) = \tilde{O}(K\delta^{-1})$ .*

*Proof.* We begin with the first assertion. Since  $\mathbb{E}[(\mu - \Phi_{\mathcal{C}}^{\text{sup}})_+] > \delta$  for any  $\mu_j \sim \mathcal{P}_j \in \mathcal{C}$  we see that  $\mathbb{P}[\mu_j \geq \Phi^{\text{sup}} + \frac{\delta}{2}] \geq \frac{\delta}{2}$ . On the other hand by Lemma C.2 with  $\varepsilon = \frac{\delta}{3}$  we have for each  $\mu_i \in \mathcal{P}_i$  that

$$\mathbb{P}[\mu_i \leq \Phi^{\text{sup}} + \frac{\delta}{3}] \geq \mathbb{P}[\mu_i \leq \mathbb{E}[\mu_i] + \frac{\delta}{3}] \geq \Omega(\delta).$$

We conclude that  $\mathbb{P}[A^* = \mu_j] \geq \Omega(\delta)^K$  for any  $j$  and therefore that  $\delta_{\text{TS}}^{-1} = O(K \log(\delta^{-1}))$ . On the other hand by definition:

$$\mathbb{E}[(\mu_j - \mu_i)_+] \geq \mathbb{E}[(\mu_j - \Phi_{\mathcal{C}}^{\text{sup}})_+] \geq \delta$$

which implies  $\varepsilon_{\text{TS}} \geq \delta$  and hence  $N_{\text{TS}} = \tilde{O}(K\delta^{-2})$ .

To see that  $G_{\text{pad}} \geq \delta$  we observe that for any  $i \in [K], q \in \Delta_{i-1}$  we have by convexity:

$$\begin{aligned} \mathbb{E}[(\mu_i - \mu_q)_+] &\geq \mathbb{E}[(\mu_i - \mu_q^0)_+] \\ &\geq \mathbb{E}[(\mu_i - \Phi_{\mathcal{C}}^{\text{sup}})_+] \\ &\geq \delta. \end{aligned}$$

The fact that  $G_{\text{pad}} \geq \delta$  immediately implies  $N_{\text{pad}} = \tilde{O}(\delta^{-2})$ .

Finally, for  $N_{\text{boot}}$  and  $p_{\text{boot}}$  we rely also on  $\delta$ -explorability. To estimate  $N_{\text{boot}}$  we observe that every  $\delta^{-1}$  zero-reward samples of arm  $i$  give a constant factor likelihood ratio advantage of any  $\mu_i \leq \mu_j^0 - \delta/2$  over any  $\mu_i \geq \mu_j^0$ . Since  $\mathbb{P}[\mu_i \leq \mu_j^0 - \delta/2] \geq \frac{\delta}{2}$  it is easy to see that  $N_{\text{boot}} = \tilde{O}(\delta)$  from this. Since each zero-reward event has probability at least  $1 - \mu_i^0 \geq \Omega(\delta)$  given the previous ones,  $\log(p_{\text{boot}}^{-1}) = \tilde{O}(K/\delta)$ .  $\square$

## C.1 Necessity of the Non-Degeneracy Assumption

We previously showed that under a  $\sigma$ -non-degeneracy assumption, Thompson sampling can be made BIC after an amount of time polynomial in  $\sigma^{-1}$  and the time required to sample every action. It is natural to wonder if the dependence on an additional parameter  $\sigma$  is necessary, or if these two quantities are always polynomially related. Here we show that dependence on an additional parameter is unavoidable. In particular we give an  $\varepsilon$ -non-degenerate problem instance in which both arms can be sampled in 2 rounds, but  $\Omega(\varepsilon^{-1})$  rounds are needed to make Thompson sampling BIC. In particular, the polynomial dependence on  $\sigma$  in Corollary 6.3 is necessary in this example.

**Proposition C.4.** *Consider an initial prior on two arms where  $\mu_1 = \frac{1}{2} \pm \varepsilon$  where the sign is chosen uniformly at random, and  $\mu_2 = \frac{1}{2} - \frac{\varepsilon^2}{10}$  almost surely. Then it is possible to sample both arms in time  $O(1)$ . However, if there is a BIC algorithm which almost surely uses Thompson sampling on round  $t$ , then  $t \geq \Omega(\varepsilon^{-1})$ .*

*Proof.* First we explain how to sample both arms in time  $O(1)$ . We first sample arm 1. Note that if the observed reward  $r_1$  is 0, then the posterior mean for  $\mu_1$  is  $\frac{1}{2} - 2\varepsilon^2$  while if the reward is 1 then the posterior mean is  $\frac{1}{2} + 2\varepsilon^2$ . Therefore, if we in the next round sample  $\mu_2$  with probability 1 when  $r_1 = 0$ , and with probability  $\frac{1}{2}$  when  $r_1 = 1$ , this is also BIC. Doing this twice with appropriate coupling allows us to sample both arms almost surely in 3 rounds.

Now we prove the lower bound for Thompson sampling. We recall that Thompson sampling at time  $t$  is BIC if and only if

$$\mathbb{E}^t[\mu_i - \mu_j] \cdot \text{Pr}^t[A_t = i] \geq 0.$$

The proof of Lemma 4.9 shows that the left-hand side of the above equation, which we will call  $X_t$ , is a submartingale. Now, sampling arm 2 gives no information, hence no change in  $X_t$ . Therefore we may

define  $Y_m$  to be the value of  $X_t$  at any time when arm 2 has been sampled  $m$  times. It is not hard to see that  $Y_m$  is also a submartingale, and is independent of the choice of which arm to sample. (We may define  $Y_m$  for all  $m$  by generating infinity many samples of arm 2 “in secret” and thinking of sampling arm 2 as revealing the next sample.) Since  $Y$  is a submartingale, we by definition have  $\mathbb{E}^t[Y_t] \geq \mathbb{E}^t[X_t]$  almost surely for any bandit algorithm. Therefore  $\mathbb{E}[Y^t] \geq \mathbb{E}[X^t]$  for any algorithm. We conclude that if there exists a BIC algorithm which uses Thompson sampling almost surely at time  $t$ , then sampling arm 1 for the first  $t - 1$  rounds also works. After  $t - 1$  samples of  $a_1$ , we gain  $O(t\varepsilon^2)$  bits of information on the value of  $\mu_1$ . Therefore any function of our observations is correlated at most  $O(t\varepsilon^2)$  with  $1_{\mu_1 > \frac{1}{2}}$ , and hence  $O(t\varepsilon^3)$  correlated with the value  $\mu_1$ . This implies that  $t = \Omega(\varepsilon^{-1})$  is required for TS to be BIC.  $\square$

## D Sample Complexity for Truncated Gaussians and Beta priors

In this Appendix we consider concrete problems with truncated Gaussian and Beta priors, determining the sample complexity for Thompson sampling up to polynomial dependence in several situations. Our strategy throughout is to methodically estimate the parameters  $G_{\text{pad}}, N_{\text{boot}}, L_{\text{LB}}$ , and so on. The next lemma allows us to simplify many of these computations using symmetry.

**Lemma D.1.** *Suppose  $\mu_1, \dots, \mu_{j-1}$  are independent and identically distributed. Then*

$$\sup_{q \in \Delta_{j-1}} \frac{\mathbb{E}[\mu_q^0 - \mu_j^0]}{\mathbb{E}[(\mu_j - \mu_q)_+]}$$

is achieved at  $q = \left(\frac{1}{j-1}, \frac{1}{j-1}, \dots, \frac{1}{j-1}\right)$ .

*Proof.* It suffices to show that  $\mathbb{E}[(\mu_j - \mu_q)_+]$  is minimized at the claimed value of  $q$ , since  $\mu_q^0$  is independent of  $q \in \Delta_{j-1}$ . Indeed  $\mathbb{E}[(\mu_j - \mu_q)_+]$  is a convex, symmetric function of  $q \in \Delta_{j-1}$  so it must be minimized when all coordinates are equal. The symmetry is clear while convexity holds for any fixed values of  $(\mu_i)_{1 \leq i \leq j}$ , hence in expectation.  $\square$

The next lemma allows comparison with stochastically dominating problem instances to estimate  $L_{\text{LB}}, G$ .

**Lemma D.2.** *For any  $j$  and  $q \in \Delta_K$  with  $q_j = 0$ , the expressions*

$$\frac{\mathbb{E}[\mu_q^0 - \mu_j^0]}{\mathbb{E}[(\mu_j - \mu_q)_+]}, \quad \frac{1}{\mathbb{E}[(\mu_j - \mu_q)_+]}$$

are stochastically decreasing in the prior  $\mathcal{P}_j$  for  $\mu_j$  and stochastically increasing in all the priors  $\mathcal{P}_i$  for  $\mu_i$  when  $i \neq j$ .

*Proof.* Based on Lemma A.6 it suffices to show the relevant monotonicity of each part of the expressions without expectations in  $\mu_j$  and  $\mu_i$ , which is clear.  $\square$

**Lemma D.3.** *For any  $\mu < \mu'$ , the distribution  $\tilde{N}(\mu, \sigma)$  is stochastically smaller than the distribution  $\tilde{N}(\mu', \sigma)$ .*

*Proof.* In fact MLR domination holds, and hence stochastic domination follows. Observe that the densities  $f(x), g(x)$  for  $\tilde{N}(\mu, \sigma)$  and  $\tilde{N}(\mu', \sigma)$  satisfy  $f(a)/g(a) \leq f(b)/g(b)$  for any  $a \leq b$ . In fact the ratio is proportional to  $e^{(a-b) \cdot (\mu' - \mu)}$ .  $\square$

**Corollary 6.6.** *Let  $\tilde{N}(\nu, \sigma^2)$  be a Gaussian with mean  $\nu$  and variance  $\sigma^2 \leq 1$ , conditioned to lie in  $[0, 1]$ . Suppose  $\mathcal{P}_i \sim \tilde{N}(\nu_i, \sigma^2)$  for each arm  $i$ , where  $\nu_1, \dots, \nu_K \in [0, 1]$ . Then*

- (a)  $T_{\text{UB}}(N_{\text{TS}}) = K^3 \cdot \text{poly}\left(\sigma^{-1}, e^{R^2}\right)$  where  $R = \sigma^{-1} \max_{i,j} |\nu_i - \nu_j|$ .
- (b)  $T_{\text{LB}} \geq e^{\Omega(1/\sigma^2)}$  when  $\max_{i,j} |\nu_i - \nu_j|$  is a positive absolute constant.

*Proof.* We assume without loss of generality that  $\sigma$  is at most a small constant. First, the distributions are clearly  $\Omega(\sigma)$ -nondegenerate so we focus on bounding the values  $\delta$  and  $L_{\text{LB}}$ . The upper bound then follows based on Lemma C.3.

We note that the mean of  $\tilde{N}(\nu_i, \sigma^2)$  is  $\nu_i \pm O(\sigma)$  because (as  $\sigma = O(1)$ ) we are conditioning on a probability  $\Omega(1)$  event in truncating the Gaussian to  $[0, 1]$ . Similarly any such distribution has mean  $\Omega(\sigma)$  and  $1 - \Omega(\sigma)$ .

Let  $\mathcal{C}$  be the set of  $K$  priors in the problem, we claim  $\mathcal{C}$  is  $\delta$ -easy and  $\delta$ -non-dominating for  $\delta^{-1} \leq \text{poly}(\sigma^{-1}, e^{R^2})$ . Since the conditions are symmetric under  $\nu_i \rightarrow 1 - \nu_i$  we focus on  $\delta$ -non-dominance. First it is easy to see that  $\mu_i^0 \geq \Omega(\sigma)$  for each arm  $i$ . The fact that for any  $i, j$  and constants  $C, c$  we have

$$\mu_i^0 \geq \max(c\sigma, \nu_j - (R + C)\sigma),$$

which implies that

$$\mathbb{E}[(\mu_i^0 - \mu_j)_+] \geq \Omega(\sigma \mathbb{P}[\mu_j \leq \max(c\sigma/2, \nu_j - 2(R + C)\sigma)]).$$

It is easy to see that the density of  $\mu_j$  is at least  $e^{-O(R^2)}$  in the interval

$$\max(c\sigma/4, (\nu_j - 2(R + C) - c/4)\sigma).$$

Hence  $\mathbb{P}[\mu_j \leq \max(c\sigma/2, \nu_j - 2(R + C)\sigma)] \geq \Omega(\sigma e^{R^2})$  and so we obtain the claimed lower bound on  $\delta$ . Using the general  $\delta$ -dependent upper bounds now yields the result.

We next turn to estimating  $L_{\text{LB}}$  to achieve a lower bound. Suppose  $m_i + R\sigma = m_j$  with  $R\sigma = \Omega(1)$  of constant order. Then we easily see that

$$\mathbb{P}[\mu_j \geq \mu_i] \leq e^{-\Omega(R^2)}.$$

Indeed this is clear for non-truncated Gaussians, and conditioning  $\mu_i, \mu_j \in [0, 1]$  restricts to a constant probability event, hence can only increase probabilities by a constant factor. Therefore

$$L_{\text{LB}} \geq \frac{\nu_j - \nu_i}{\mathbb{P}[\mu_j \geq \nu_i]} = e^{\Omega(R^2)}(\nu_j - \nu_i).$$

Since we may take  $i, j$  to maximize  $|\nu_i - \nu_j|$  this shows the claimed lower bound (as the assumptions imply  $R = \Omega(\sigma^{-2})$ ).  $\square$

**Corollary 6.7.** *Suppose all priors  $\mathcal{P}_1, \dots, \mathcal{P}_K$  are Beta distributions.*

- (a)  $T_{\text{UB}}(N_{\text{TS}}) \leq K^3 \cdot (\min(K, M))^{O(M)}$  if  $\text{strength}(\mathcal{P}_i) \leq M$  for all arms  $i$ .
- (b)  $T_{\text{LB}} \geq (\min(K, M))^{\Omega(M)}$  for **some**  $M$ -strong problem instance.

(c)  $T_{\text{LB}} \geq 2^{\Omega(M)}$  for **any**  $M$ -strong problem instance such that  $\mu_1^0 - \mu_K^0 \geq \Omega(1)$ .

More generally, this holds whenever arms  $i \neq j$  have strength at least  $M$  and  $|\mu_i^0 - \mu_j^0| \geq \Omega(1)$ .

*Proof.* For convenience we replace the strength condition with the (up to constant factors) equivalent condition that all  $\text{Beta}$  parameters are at most  $M$ . We first estimate  $N_{\text{TS}}, N_{\text{boot}}, \log(p_{\text{boot}}^{-1})$  towards establishing the upper bound. We estimate  $G_{\text{pad}}, L_{\text{LB}}$  together at the end, which implies the upper bound via Theorem E.1 and the lower bound directly. Corollary 4.6 gives  $N_{\text{TS}} = 2^{O(M)}$ . Computing  $N_{\text{boot}}$  may be done exactly in the case when one variable is  $\text{Beta}(1, M)$  and the other is  $\text{Beta}(M, 1)$  which is easily seen to be the worst case. In fact for all the priors this is the worst case, by Lemma D.2 for  $G, L_{\text{LB}}$  and similar arguments for the others. Recall that the mean of a  $\text{Beta}(a, b)$  random variable is exactly  $\frac{a}{a+b}$ . So, we need

$$\frac{M}{M+1+N_{\text{boot}}} < \frac{1}{M+1},$$

and so  $N_{\text{boot}} = M^2$ . Since we need at most  $K \cdot N_{\text{boot}} \leq KM^2$  total zero-rewards to make all arms have 0 reward for their first  $M^2$  samples, and this has probability at least  $p_{\text{boot}} \geq M^{-KM^2}$ , we obtain  $\log(p_{\text{boot}}^{-1}) = \tilde{O}(KM)$ .

To estimate  $G_{\text{pad}}$  we show  $\mathbb{E}[(\mu_K - \mu_q)_+] \geq M^{-O(M)}$ . By Lemma D.1 and the fact that  $\text{Beta}(M, 1)$  is the stochastically largest possible prior, we see that  $G_{\text{pad}}$  is minimized when  $\mu_1, \dots, \mu_{K-1} \sim \text{Beta}(M, 1)$ ,  $\mu_K \sim \text{Beta}(1, M)$ , and  $q = \left(\frac{1}{K-1}, \frac{1}{K-1}, \dots, \frac{1}{K-1}\right)$ . We will show that in this case:

$$\mathbb{E}[(\mu_K - \mu_q)_+] = \min(K, M)^{\Theta(M)}.$$

Because also  $\mathbb{E}[(\mu_K - \mu_q)_+] = \frac{\Omega(1)}{L_{\text{LB}}}$ , this establishes both the upper and lower bounds. We first show  $\mathbb{E}[(\mu_K - \mu_q)_+] \geq K^{-O(M)}$ . Observe that with probability  $K^{-O(M)}$  we have

$$\mu_K > 1 - \frac{1}{K}, \quad \mu_1, \mu_2, \mu_3 < \frac{1}{3}.$$

In this situation,  $\mu_q \leq \frac{K-3}{K-1} < 1 - \frac{2}{K}$  and  $\mu_K \geq 1 - \frac{1}{K}$  and so  $(\mu_K - \mu_q)_+ \geq \frac{1}{K}$ . We conclude that:

$$\mathbb{E}[(\mu_K - \mu_q)_+] \geq K^{-O(M)}.$$

We next show  $\mathbb{E}[(\mu_K - \mu_q)_+] \geq M^{-O(M)}$ . By Jensen's inequality we have:  $\mathbb{E}[(\mu_K - \mu_q)_+] \geq \mathbb{E}[(\mu_K - \mathbb{E}[\mu_q])_+]$ . Now,  $\mathbb{E}[\mu_q] = 1 - \frac{1}{M}$  and there is at least an  $M^{-O(M)}$  chance that  $\mu_K \geq 1 - \frac{1}{2M}$  so this shows the desired lower bound of  $M^{-O(M)}$ .

We now turn to the matching lower bound. Assume first that  $K \leq \frac{M}{10}$ . We estimate the  $M/2$  exponential moment of a  $\text{Beta}(1, M)$  variable denoted by  $Z$ . We note that  $(1-x)e^x \leq 1$  for  $x \in [0, 1]$ . Therefore

$$\mathbb{E}[e^{MZ/2}] = M \int_0^1 (1-x)^{M-1} e^{Mx/2} dx \leq M \int_0^1 (1-x)^{M/2-1} dx = 2.$$

The random variable  $1 - \mu_i$  for  $i \leq K-1$  is distributed as  $\text{Beta}(1, M)$  and so  $\mathbb{E}[e^{M(1-\mu_1)}] \leq 2$ . Moreover since exponential moments multiply under independent sum, we obtain:

$$\mathbb{E}[e^{M(K-1)(1-\mu_q)/2}] \leq 2^{K-1} \leq 2^K.$$



Therefore:

$$\begin{aligned} \mathbb{P}\left[\mu_q \leq 1 - \frac{\log(K)}{K-1}\right] &= \mathbb{P}\left[1 - \mu_q \geq \frac{\log(K)}{K-1}\right] \\ &\leq \mathbb{E}\left[e^{M(K-1)(1-\mu_q)/2}\right] \cdot e^{-\frac{M \log(K)}{2}} \\ &\leq 2^K K^{-M/2}. \end{aligned}$$

Since we assume  $K \leq \frac{M}{10}$  we have

$$\mathbb{P}\left[\mu_q \leq 1 - \frac{\log(K)}{K-1}\right] \leq 2^K K^{-M/2} \leq K^{-\Omega(M)}.$$

We also trivially have  $\Pr[\mu_K \geq 1 - \frac{\log(K)}{K-1}] \leq K^{-\Omega(M)}$ . Therefore if  $K \leq \frac{M}{10}$  we have established the upper bound  $\mathbb{E}[(\mu_K - \mu_q)_+] \leq K^{-\Omega(M)}$ . Now let us denote by  $f(K, M)$  the value of  $\mathbb{E}[(\mu_K - \mu_q)_+]$  in the worst case we are working in. Observe that  $f(K, M)$  is decreasing in  $K$ , since decreasing  $K \rightarrow (K-1)$  is equivalent to replacing  $q = \left(\frac{1}{K-1}, \dots, \frac{1}{K-1}\right) \in \Delta_{K-1}$  with  $q' = \left(\frac{1}{K-2}, \dots, \frac{1}{K-2}, 0\right) \in \Delta_{K-1}$ , which by Lemma D.1 gives a larger value of  $\mathbb{E}[(\mu_K - \mu_q)_+]$ .

To finish, it is not hard to see that the statement  $f(M, K) = \min(K, M)^{-\Theta(M)}$  follows from combining the three statements below:

- $\min(K, M)^{-O(M)} \leq f(M, K)$
- $f(M, K) \leq K^{-\Omega(M)}$  for  $K \leq \frac{M}{10}$ .
- $f(M, K) \geq f(M, K+1)$ .

Indeed, to complete the upper bound, if  $K \geq \frac{M}{10}$  we have  $f(M, K) \leq f(M, M/10) = M^{-\Omega(M)}$ . This concludes the proof.

Finally in the case that  $\mu_i, \mu_j$  have constant-separated prior means, the lower bound on  $L_{\text{LB}}$  follows by a simple Chernoff bound showing that

$$\mathbb{E}[(\mu_j - \mu_i)_+] \leq \mathbb{P}[\mu_j \geq \mu_i] \leq \mathbb{P}[\mu_j \geq t] + \mathbb{P}[\mu_i \leq t] \leq 2^{-\Omega(M)} \quad \text{for } t = 1/2 (\mu_i^0 + \mu_j^0). \quad \square$$

**Lemma 6.8.** *The collection of all truncated Gaussians  $\tilde{N}(m, \sigma^2)$ ,  $m \in [0, 1]$  with a fixed variance  $\sigma^2 < 1$  is  $\delta$ -easy and  $\delta$ -non-dominant with  $\delta = e^{-\Omega(1/\sigma^2)}$ . Likewise, the collection of all Beta distributions of strength at most  $M \geq 1$  is  $\delta$ -easy and  $\delta$ -non-dominant with  $\delta = M^{-O(M)}$ .*

*Proof.* The Gaussian case was already proved inside the proof of Corollary 6.6. In the Beta case, it suffices to lower bound  $\mathbb{E}\left[\left(\mu - 1 - \frac{1}{M}\right)_+\right]$  for  $\mu \sim \text{Beta}(1, M)$ . This is at least

$$\frac{1}{2M} \mathbb{P}\left[\mu \geq 1 - \frac{1}{2M}\right] \geq \Omega(M^{-O(M)}) \quad \square$$

**Corollary 6.10.** *Suppose all priors  $\mathcal{P}_1, \dots, \mathcal{P}_K$  are Beta distributions. Suppose  $\text{strength}(\mathcal{P}_\ell) = M$  for some arm  $\ell$ , and  $\text{strength}(\mathcal{P}_i) \leq m$  for all other arms  $i$ , where  $M \geq m \geq 2$ . Then*

$$T_{\text{UB}}(N_{\text{TS}}) \leq M^2 \cdot K^{O(1)} \cdot \max\left(m, (1 - \mu_\ell^0)^{-1}, (\mu_\ell^0)^{-(\ell-1)}\right)^{O(m)}.$$

*In particular,  $T_{\text{UB}}(N_{\text{TS}}) \leq M^2 \cdot K^{O(1)} \cdot m^{O(\ell m)}$  if  $\mu_\ell^0 \in [1/4, 3/4]$ .*

*Proof.* Again we assume that both Beta parameters are bounded by  $M$  or  $m$ . It is easy to see that  $N_{\text{boot}} = O(M + m)$  based on the mean value of Beta distributions, and so as before we have  $\log(p_{\text{boot}}^{-1}) = \tilde{O}(K \cdot \max(m, (1 - \mu_j^0)^{-1}))$ .

We next estimate  $G_{\text{pad}}$ . Note that for any  $q \in \Delta_K$  we have

$$\mathbb{E}[\mu_q] \leq \max\left(1 - \frac{1}{2m}, \mu_j^0\right).$$

Then we simply observe that for each  $i \geq 2$  and  $t \in [0, 1]$ , we have

$$\mathbb{P}[\mu_i \geq 1 - t] \geq t^m.$$

In particular letting  $m_0 = \min\left(\frac{1}{2m}, 1 - \mu_j^0\right)$  and taking  $t = m_0/2$  we obtain:

$$\mathbb{E}[(\mu_i - (1 - m_0))_+] \geq \frac{m_0}{2} \mathbb{P}\left[\mu_i \geq \left(1 - \frac{m_0}{2}\right)\right] \geq \left(\frac{m_0}{2}\right)^{m+1}.$$

This handles everything except the case of  $\mathbb{E}[(\mu_j - \mu_q)_+]$  for  $q \in \Delta_{j-1}$  (which is irrelevant when  $j = 1$ ). Note that for any  $i \neq j$  we have:

$$\mathbb{P}\left[\mu_i \leq \frac{\mu_j^0}{2}\right] \geq \left(\frac{\mu_j^0}{2}\right)^m.$$

Therefore we obtain  $\mathbb{E}[(\mu_j - \mu_q)_+] \geq \left(\frac{\mu_j^0}{2}\right)^{m(j-1)+1}$  for  $q \in \Delta_{j-1}$ . Altogether this implies

$$G_{\text{pad}} \geq \max(m_0, (\mu_j^0)^{(j-1)})^{O(m)}.$$

To estimate  $N_{\text{TS}}$ , in the case  $j = 1$ , for the highly informed arm  $\mu_j$  we don't need to do anything since Thompson sampling is already BIC for arm 1 at time 1. For all the other arms  $i \geq 2$ , the computation above shows that

$$\mathbb{P}[\mu_i \geq \max_{j \neq i} \mu_j + \varepsilon_{\text{TS}}] \geq \delta_{\text{TS}}$$

for  $\varepsilon_{\text{TS}} \geq \frac{m_0}{2}$  and  $\delta_{\text{TS}} \geq m_0^m$ . Therefore  $N_{\text{TS}} \leq \tilde{O}(m_0^{-3})$ . When  $j \geq 2$  we also need to estimate  $\varepsilon_{\text{TS}}, \delta_{\text{TS}}$  for arm  $j$ . We have

$$\mathbb{E}[(\mu_j - \mu_i)_+] \geq (\mu_j^0/2) \mathbb{P}[\mu_i \leq \mu_j^0/2] \geq (\mu_j^0/2)^{m+1}.$$

Similarly the chance that arm  $j$  is the best is at least

$$\delta_{\text{TS}} \geq \Omega(1) \cdot (\mu_j^0)^{(j-1)m} \cdot 2^{-O(K)}.$$

Indeed there is a constant chance for beta random variables to be on either side of their mean, and the above therefore lower bounds the chance that  $\mu_j \geq \mu_j^0$  and  $\mu_i \leq \mu_j^0$  for all  $i \neq j$ . Applying Theorem 5.6 or Theorem E.1 concludes the proof.  $\square$

---

**Algorithm 2:** ExponentialExploration with tighter phases

---

```
1 Parameters: Desired number of samples  $N$ , calculated value and phase length  $N_{\text{pad}}$ 
2 Given: recommendation policies  $\pi_1, \dots, \pi_K$  for PADDED PHASE.
3 Initialize: EXPLORATION PHASE for arm 1 of length  $\max(N, N_{\text{boot}}, N_{\text{pad}})$ 
4 for each arm  $j = 2, 3, \dots, K$  do
5   // Invariant 1: each arm  $i < j$  has been sampled at least  $\max(N, N_{\text{boot}}, N_{\text{pad}})$  times.
   // Bootstrapping: two phases
6   Event  $\text{ZEROS}_{j, N_0} = \{\text{the first } N_0 \text{ samples of each arm } i < j \text{ return reward } 0\}$ .
7    $p_j \leftarrow q/(1+q)$ , where  $q = \lambda \cdot \Pr[\text{ZEROS}_{j, N_0}]$ .
8   EXPLOITATION PHASE with depth  $N_0$ 
9   with probability  $p_j$  do
10     EXPLORATION PHASE for arm  $j$ 
11   else if  $\text{ZEROS}_{j, N_0}$  then
12     | PADDED PHASE: use policy  $\text{transform}(\pi_j)$ 
13   else EXPLOITATION PHASE with depth  $N$ 
   // main loop: exponentially grow the exploration probability
14  while  $p_j < 1$  do
15    // Invariant 2:  $\Pr[\text{exploration phase has happened} \mid \mu_1, \dots, \mu_K] = p_j$ .
16    if exploration phase has happened then
17      | PADDED PHASE: use policy  $\pi_j$ 
18    else with probability  $\min\left(1, \frac{p_j}{1-p_j} \cdot \lambda\right)$  do
19      | EXPLORATION PHASE for arm  $j$ 
20    else EXPLOITATION PHASE with depth  $N$ 
21    Update  $p_j \leftarrow \min(1, p_j(1+\lambda))$ .
   // Post-processing: collect remaining samples
22  for each arm  $j = 1, 2, \dots, K$  do
23    Choose phase  $\ell_0$  uniformly at random from  $[1 + \lceil \lambda^{-1} \rceil]$ .
24    for each phase  $\ell = 1, 2, \dots, \lceil \lambda^{-1} \rceil + 1$  do
25      | if  $\ell = \ell_0$  then
26        | Explore arm  $j$  for  $\max(N, N_{\text{boot}}, N_{\text{pad}})$  rounds
27      | else
28        | Use policy  $\pi_j$  for  $\max(N, N_{\text{boot}}, N_{\text{pad}})$  rounds
```

---

## E Extension: A more efficient version of `ExponentialExploration`

We now give Algorithm 2, a version of Algorithm 1 which requires fewer rounds. Algorithm 2 uses the observation that when  $N, N_{\text{boot}} \geq N_{\text{pad}}$  is rather large, only the initial  $N_{\text{pad}}$  samples of each arm  $j$  require the hard work of exponentially growing exploration probability. As a result, we can use the same technique as in Algorithm 1 to obtain the first  $N_{\text{pad}}$  samples of arm  $j$ , and then obtain the remaining samples more quickly. The “post-processing” stage of each loop is easily seen to be BIC by the  $(j, \lambda)$ -padded BIC property of policy  $\pi_j$ . By inspection, Algorithm 2 completes in this many rounds:

$$O\left(K G_{\text{pad}}^{-1} \left(N_{\text{pad}} \log(G_{\text{pad}}^{-1} p_{\text{boot}}^{-1}) + N_{\text{boot}} + N\right)\right). \quad (\text{E.1})$$

**Theorem E.1.** *Given a parameter  $N \geq N_{\text{pad}}$ , Algorithm 2 with  $N_0 = N_{\text{boot}}$  collects  $N$  rounds of each arm almost surely and completes in the number of rounds given by (E.1).*

Note that the phase length in the main part of the algorithm is only  $N_{\text{pad}}$  in Algorithm 2. We remark that we do not need to assume  $N_{\text{boot}} \leq N_{\text{pad}}$ , even though Algorithm 1 required  $N_{\text{boot}} \leq N$ . This is because the post-processing phase ensures that Invariant 1 continues to hold, and Lemma 5.4 does not require  $N_{\text{boot}} \leq N$ .

## F Extension: Improved Algorithm for “Easy” problem Instances

We now explain Algorithm 3, which achieves the guarantee of Theorem 7.1. We fix an  $N \geq 1$  and show it samples each arm  $N$  times in  $\tilde{O}\left(\frac{KN}{\delta} + \frac{K}{\delta^4}\right)$  rounds.

The algorithm’s structure is again similar to Algorithm 1, featuring an initial bootstrap phase followed phase of exponentially growing exploration probability facilitated by a padded phase. The main difference is that we only carry out these steps for a single arm  $j_0$  chosen randomly, so that we manage to sample arm  $j_0$  many times without needing to first “unlock” the previous arms. In general problem instances this may not be possible, but the  $\delta$ -easy assumption ensures that it is.

The algorithm continues with a `for` loop to complete the exploration, balanced by a padded phase. This is reminiscent of Algorithm 2, but in this case we have only thoroughly explored arm  $j_0$ . A key new insight is that having to explore the single arm  $j_0$  allows us to explore the remaining arms without requiring another exponential growth phase. This is achieved by randomizing between exploiting arm  $j_0$  and exploring a random arm  $i$ . Because  $j_0$  is random, the agent seeing the recommendation does not know whether we are exploring or exploiting, so this is BIC when the exploration probability is small. We couple the random choices of exploration arms  $i$  so that there are no repeats via a uniformly random permutation  $\theta : [K] \rightarrow [K]$ .

**Lemma F.1.** *Suppose  $\mathcal{C}$  is  $\delta$ -easy and  $\delta$ -non-dominant, and  $N \geq N_{\text{pad}} = \tilde{O}(\delta^{-2})$  as guaranteed by Lemma C.3. For fixed  $j \in [K]$ , consider exploitation based on  $N_{\text{pad}}$  samples from arm  $j$  and no other information. This policy is  $(j, \frac{\delta}{10})$ -suitable.*

*Proof.* Since we collect no information on the first  $j - 1$  arms, it is equivalent to replace the random values  $\mu_1, \dots, \mu_{j-1}$  with their expectations  $\mathbb{E}[\mu_i]$ . From the fact that  $\mathcal{C}$  is  $\delta$ -easy to explore we see that  $G_{\text{pad}} \geq \delta$  in this case. Since we replace the  $\mu_i$  with their (deterministic) expectations, we may include  $N_{\text{pad}}$  samples of each arm  $a_1, \dots, a_{j-1}$  without making any difference. That a  $(j, \frac{\delta}{10})$ -suitable policy exists now follows from Lemma 5.3. However because we only observe samples of arm  $j$  it is not difficult to see that for any value of  $N$ , exploitation based on  $N$  samples from arm  $j$  yields the optimal strategy in the  $j$ -recommendation game. The equivalence of Lemma B.1 implies the result.  $\square$

---

**Algorithm 3: Collect Samples For Easy Collections**

---

```
1 Parameters: number of target samples  $N \geq N_{\text{boot}}, N_{\text{pad}}$ , padding  $\lambda = \frac{\delta}{10} > 0$ ;  
2 Given: A uniformly random permutation  $\theta : [K] \rightarrow [K]$   
3 Choose a random arm  $j_0$ .  
4 Event  $\text{ZEROS}_{j_0, N_{\text{boot}}} = \{ \text{the first } N_{\text{boot}} \text{ samples of each arm } i < j_0 \text{ return reward } 0 \}$ .  
   // Setup: get  $N$  samples of arm  $j_0$  with positive probability  
5 for  $j = 1, 2, \dots, K$  do  
6   | EXPLOITATION PHASE with depth  $N_{\text{boot}}$  and length  $N_{\text{boot}}$   
7 if  $\text{ZEROS}_{j_0, N_{\text{boot}}}$  then  
8   | EXPLORATION PHASE for arm  $j_0$  with length  $N_{\text{pad}}$   
9 else  
10  | EXPLOITATION PHASE with depth  $N_{\text{boot}}$  and length  $N_{\text{pad}}$   
   // Bootstrapping phase  
11  $p_{j_0} \leftarrow \frac{\lambda \cdot \Pr[\text{ZEROS}_{j_0, N_{\text{boot}}}]}{1 + \lambda \cdot \Pr[\text{ZEROS}_{j_0, N_{\text{boot}}}]}$ .  
12 with probability  $p_{j_0}$  do  
13   EXPLORATION PHASE for arm  $j_0$  with length  $N_{\text{pad}}$   
14 else  
15   | if exploration phase has happened then  
16     | PADDED PHASE of length/depth  $N_{\text{pad}}$ : use  $\hat{\pi}_{j_0}$  to decide whether to play arm  $j$ . If  $\hat{\pi}_{j_0} \neq j$ ,  
17     | exploit based on all available data.  
18   | else  
19     | EXPLOITATION PHASE with depth  $N_{\text{pad}}$  with length  $N_{\text{pad}}$ .  
   // grow the exploration probability  
20 Set  $p_{j_0} = \Pr[\text{ZEROS}_{j_0, N_{\text{boot}}}]$ .  
21 while  $p_{j_0} < 1$  do  
22   // Invariant :  $\Pr[\text{exploration phase has happened} \mid \mu_1, \dots, \mu_K] = p_{j_0}$ .  
23   | if exploration phase has happened then  
24     | PADDED PHASE of length and depth  $N_{\text{pad}}$ : use  $\hat{\pi}_{j_0}$  for  $N_{\text{pad}}$  steps.  
25   | else  
26     | With probability  $\min\left(1, \frac{p_{j_0}}{1-p_{j_0}} \cdot \lambda\right)$ , EXPLORATION PHASE for arm  $j_0$  with length  $N_{\text{pad}}$ ;  
27     | with the remaining probability: EXPLOITATION PHASE with depth  $N_{\text{pad}}$ .  
    $p_{j_0} \leftarrow \min(1, p_{j_0}(1 + \lambda))$   
   // explore the other arms  
28 for each  $j = 1, 2, \dots, K$  do  
29   | Pick a phase  $\ell_0$  uniformly at random from  $[n_\lambda]$ , where  $n_\lambda := 1 + \lceil \lambda^{-1} \rceil$   
30   | for each phase  $\ell = 1, 2, \dots, n_\lambda$  do  
31     | if  $\ell = \ell_0$  then  
32       | Play arm  $\theta(j)$  for  $N$  rounds.  
33     | else  
34     | PADDED PHASE using  $\hat{\pi}_{j_0}$  with depth  $N_{\text{pad}}$  for  $N$  rounds
```

---

Based on the lemma above, we define for each  $j \in [K]$  the policy  $\hat{\pi}_j$  which decides which arm to play by exploitation based on  $N_{\text{pad}}$  samples from arm  $j$ . In particular  $\hat{\pi}_j$  will only play arm 1 or  $j$ .

**Theorem 7.1.** *Given a  $\delta$ -easy problem instance with  $K$  arms, Algorithm 3 is BIC and collects at least  $N$  samples of each arm almost surely in  $\tilde{O}\left(\frac{KN}{\delta} + \frac{K}{\delta^4}\right)$  rounds, for any desired  $N \in \mathbb{N}$ . The running time for each round is  $O(1)$  plus one call to “exploitation” given up to  $\tilde{O}(\delta^{-2})$  samples per arm.*

*Proof.* The algorithm uses  $O\left(KN_{\text{boot}} + \frac{N_{\text{pad}} \log(\lambda^{-1} p_{\text{boot}}^{-1})}{\lambda} + \frac{KN}{\lambda}\right) = \tilde{O}\left(\frac{KN}{\delta} + \frac{K}{\delta^4}\right)$  rounds by construction, according to the general estimates of Lemma C.3. Hence we focus on the BIC property. Lines 6, 8, 10 are clearly BIC. Line 8, if executed, always samples arm  $j_0$  for  $N_{\text{pad}}$  rounds. By our definition of  $\hat{\pi}_j$ , before the second while loop recommending any arm  $i \neq j_0$  is always BIC, hence we focus on the recommendations of  $j_0$ . We first consider the combination of Lines 13, 16, and 18. The proof is essentially identical to that of Lemma 5.5, where the point is that conditioned on reaching line 16, playing from  $\hat{\pi}_{j_0}$  is  $(j, \lambda)$ -padded BIC. That we condition on  $\text{ZEROS}_{j_0, N_{\text{boot}}}$  only helps the  $(j, \lambda)$ -padded property, because this conditioning decreases the mean reward of arm  $i$  for all  $i < j_0$ . This counterbalances the exploration in line 13.

The while loop is BIC for the same reason as in the proof of 5.2. The key point is again that the padded, exploration, and exploitation phases occur independently of the true mean rewards. To show that the final for loop is BIC, we observe:

$$\mathbb{E}[(\mu_j - \mu_i) \cdot 1_{A_t=a_j}] \geq \frac{1}{K} \cdot \left( \lambda \cdot \frac{n_\lambda - 1}{n_\lambda} - \frac{1}{n_\lambda} \right) \geq 0.$$

Here the first term comes from the exploitation phase while the second term comes from the event  $t = s$ . The factor  $\frac{1}{K}$  comes from the randomness in choosing  $j_0$  and  $\theta : [K] \rightarrow [K]$ . This concludes the proof that the algorithm is BIC.  $\square$

## G Extension: Efficient Computation for Beta Priors

**Theorem 7.2.** *Fix arm  $j$ , the number of arms  $K$ , parameters  $N, M \in \mathbb{N}$  and padding  $\lambda > 0$ . Suppose there exists a  $(j, N)$ -informed policy  $\pi_j$  which is BIC and  $(j, \lambda)$ -padded BIC for all problem instances with  $K$  arms and Beta priors of strength at most  $M$ . Then there is policy  $\pi_j^{\text{eff}}$  with these properties which can be computed efficiently, namely in time  $\text{poly}(K, M, N)$ . In particular, one can take  $\lambda = (\min(K, M))^{O(M)}$ . Plugging these  $\pi_j^{\text{eff}}$  and  $\lambda$  into Algorithm 1 yields the guarantee in Corollary 6.7(a).*

*Proof.* First suppose that we are in the worst case  $\mu_1, \dots, \mu_{j-1} \sim \text{Beta}(M, 1)$ ,  $\mu_j \sim \text{Beta}(1, M)$  and aim to explore  $a_j$ . The key point is that by the argument of Lemma D.1, assuming without loss of generality that our strategy is symmetric in  $(\mu_1, \dots, \mu_{j-1})$ , the uniform distribution  $q = q_j := \left(\frac{1}{j-1}, \frac{1}{j-1}, \dots, \frac{1}{j-1}\right)$  over arms  $i < j$  is always a best response in the  $j$ -recommendation game. Therefore the optimal  $j$ -recommendation strategy is to recommend arm  $K$  exactly when  $\tilde{\mu}_j \geq \tilde{\mu}_{q_j}$ . Here as usual we use  $\tilde{\mu}_i$  to denote posterior mean with respect to the relevant data, in this case the first  $N$  samples of arm  $j$ .

To efficiently compute the resulting value of  $\lambda = \mathbb{E}[(\tilde{\mu}_j \geq \tilde{\mu}_{q_j})_+] / 10$  is not difficult. One can simply compute the distribution of  $\mathbb{E}[\mu_i | \mathcal{G}]$  for each  $i < j$  and then compute convolutions to find the distribution of  $\mathbb{E}[\mu_{q_0} | \mathcal{G}]$ . Since Beta distributions with a fixed strength have closed-form probability mass functions supported on an arithmetic progression, this is computationally efficient.



Of course, we might not have  $\mu_1, \dots, \mu_{j-1} \sim \text{Beta}(M, 1)$ ,  $\mu_K \sim \text{Beta}(1, M)$ . The second key point is that we may reduce to this case in a similar manner to the construction of  $\text{transform}(\pi_j)$  previously. Indeed, let  $\mu'_i \sim \text{Beta}(M, 1)$ . Likewise let  $\mu'_j \sim \text{Beta}(1, M)$ , and let  $\mathcal{G}' = \mathcal{G}'_{N,j}$  be a  $\sigma$ -algebra encapsulating  $N$  samples of arms with means  $\mu'_1, \dots, \mu'_j$ .

By Corollary A.11 we know that  $\tilde{\mu}_i = \mathbb{E}[\mu_i | \mathcal{G}]$  is stochastically smaller than  $\mathbb{E}[\mu'_i | \mathcal{G}']$ , for each  $i < j$ , and that the opposite holds for  $\mu_j, \mu'_j$ . Moreover it is computationally easy to compute these distributions exactly.<sup>18</sup> Once the two distributions are computed we then compute the canonical monotone coupling between the two conditional expected values, as in defining the couplings  $\mathcal{X}_i$  in the description of Algorithm 1. Finally, given the values  $\tilde{\mu}_i = \mathbb{E}[\mu_i | \mathcal{G}]$  we sample the value  $\hat{\mu}_i = \mathbb{E}[\mu'_i | \mathcal{G}']$  according to  $\mathcal{X}_i$ , for each  $i \leq j$ . Hence we have  $\tilde{\mu}_i \leq \hat{\mu}_i$  if  $i < j$ , and  $\tilde{\mu}_j \geq \hat{\mu}_j$  otherwise.

$\pi_j^{\text{eff}}$  can now be defined.  $\pi_j^{\text{eff}}$  first decides whether to recommend arm  $j$ , doing so whenever

$$\hat{\mu}_j \geq \hat{\mu}_q = \frac{1}{j-1} \sum_{i \in [j-1]} \hat{\mu}_i.$$

When this does not happen,  $\pi_j^{\text{eff}}$  exploits conditional on  $\mathcal{G}$  as usual.

To see that  $\pi_j^{\text{eff}}$  is  $(j, \lambda)$ -padded BIC we use the same strategy as in previous comparison arguments. Exploitation parts are automatically BIC. For  $i < j$  we have:

$$\begin{aligned} \mathbb{E}[(\mu_j - \mu_i) \cdot 1_{\pi_j^{\text{eff}}=j}] &= \mathbb{E}[(\tilde{\mu}_j - \tilde{\mu}_i) \cdot 1_{\pi_j^{\text{eff}}=j}] \\ &\geq \mathbb{E}[(\hat{\mu}_j - \hat{\mu}_i) \cdot 1_{\pi_j^{\text{eff}}=j}] \\ &\geq \lambda. \end{aligned}$$

The ordinary BIC property against the other arms  $i > j$  holds similarly. □

---

<sup>18</sup>For instance, the sequence of 0/1-reward values with a Beta-prior is a simple Markov chain following Laplace's rule of succession (with prior-dependent initialization), so the probabilities to earn  $k$  reward from  $N$  samples can be computed easily by dynamic programming.