# MATHEMATICAL ENGINEERING TECHNICAL REPORTS

# Optimality of Thompson Sampling for Gaussian Bandits Depends on Priors

Junya HONDA and Akimichi TAKEMURA

DEPARTMENT OF MATHEMATICAL INFORMATICS
GRADUATE SCHOOL OF INFORMATION SCIENCE AND TECHNOLOGY
THE UNIVERSITY OF TOKYO
BUNKYO-KU, TOKYO 113-8656, JAPAN

# Optimality of Thompson Sampling for Gaussian Bandits Depends on Priors

**Junya Honda**

**Akimichi Takemura**

The University of Tokyo, Japan.
{honda,takemura}@stat.t.u-tokyo.ac.jp

## Abstract

In stochastic bandit problems, a Bayesian policy called Thompson sampling (TS) has recently attracted much attention for its excellent empirical performance. However, the theoretical analysis of this policy is difficult and its asymptotic optimality is only proved for one-parameter models. In this paper we discuss the optimality of TS for the model of normal distributions with unknown means and variances as one of the most fundamental example of multiparameter models. First we prove that the expected regret of TS with the uniform prior achieves the theoretical bound, which is the first result to show that the asymptotic bound is achievable for the normal distribution model. Next we prove that TS with Jeffreys prior and reference prior cannot achieve the theoretical bound. Therefore the choice of priors is important for TS and non-informative priors are sometimes risky in cases of multiparameter models.

## 1 INTRODUCTION

In reinforcement learning a tradeoff between exploration and exploitation of knowledge is considered. The multiarmed bandit problem is one formulation of the reinforcement learning and is a model of a gambler playing a slot machine with multiple arms. A dilemma for the gambler is that he cannot know whether the expectation of an arm is high or not without pulling it many times but he suffers a loss if he pulls suboptimal (i.e., not optimal) arms many times.

This problem was formulated by Robbins (1952) and its theoretical bound was derived by Lai and Robbins (1985) for single parametric models, which was extended to multiparameter models by Burnetas and Katehakis (1996). These theoretical bounds show that any suboptimal arm has to be pulled at least logarithmic number of rounds and its coefficient is determined by the distributions of suboptimal arms and the expectation of the optimal arm.

Along with the asymptotic bound for this problem, achievability of the bound has also been considered in many models. Lai and Robbins (1985) proved the asymptotic optimality of a policy based on the notion of *upper confidence bound* (UCB) for Laplace distributions (which do not belong to exponential families) and some exponential families including normal distributions with *known* variances. The achievability of the bound was later extended to a subclass of one-parameter exponential families (Garivier & Cappé, 2011).

On the other hand in multiparameter or nonparametric models, Burnetas and Katehakis (1996) and Honda and Takemura (2010) proved the achievability for finite-support distributions and bounded-support distributions, respectively. However, the above two models are both compact and achievability of the bound is not known for non-compact multiparameter models, which include normal distributions with unknown means and variances. Since the normal distribution model is one of the most basic settings of stochastic bandits, many researches have been conducted for this model (Burnetas & Katehakis, 1996; Auer et al., 2002; Kaufmann et al., 2012a). However, to the authors' knowledge, only the UCB-normal policy (Auer et al., 2002) assures a (non-optimal) logarithmic regret for this model[1].

In this paper we discuss the asymptotic optimality of Thompson sampling (TS) (Thompson, 1933) for this normal distribution model with unknown means and variances. TS is a Bayesian policy which chooses an arm randomly according to the posterior probability with which the arm is the optimal. This policy was

---

[1] The theoretical analysis of UCB-normal contains conjectures verified only numerically and the logarithmic regret is not assured in the strict sense.

recently rediscovered and is researched extensively because of its excellent empirical performance for many models (Chapelle & Li, 2012). The theoretical analysis of TS was first given for Bernoulli model (Agrawal & Goyal, 2012; Kaufmann et al., 2012b) and was later extended to general one-parameter exponential families (Korda et al., 2013).

The asymptotic optimality of TS under uniform prior is proved for Bernoulli model in Kaufmann et al. (2012b), whereas it is proved for a more general model, one-parameter exponential family, under Jeffreys prior in Korda et al. (2013). Therefore, TSs with uniform prior and Jeffreys prior are asymptotically equivalent at least for the Bernoulli model. Nevertheless, we prove for the normal distribution model that TS with uniform prior achieves the asymptotic bound whereas TS with Jeffreys prior and reference prior cannot. Furthermore, TS with Jeffreys prior cannot even achieve a logarithmic regret and suffers a polynomial regret in expectation. This result implies that TS may be more sensitive to the choice of priors than expected and noninformative priors are sometimes risky (in other words, too optimistic) for multiparameter models.

This paper is organized as follows. In Sect. 2, we formulate the bandit problem for the normal distribution model and introduce Thompson sampling. We give the main result on the optimality of TS in Sect. 3. The remaining sections are devoted to the proof of the main result. In Sect. 4, we derive inequalities for probabilities which appear in the normal distribution model. We prove the optimality of TS with conservative priors in Sect. 5 and prove the non-optimality of TS with optimistic priors in Sect. 6.

## 2 Preliminaries

We consider the $K$-armed stochastic bandit problem in the normal distribution model. The gambler pulls an arm $i \in \{1, \cdots, K\}$ at each round and receives a reward independently and identically distributed by $\mathcal{N}(\mu_i, \sigma_i^2)$, where $\mathcal{N}(\mu, \sigma^2)$ denotes the normal distribution with mean $\mu$ and variance $\sigma^2$. The gambler does not know the parameter $(\mu_i, \sigma_i^2) \in \mathbb{R} \times (0, \infty)$. The maximum expectation is denoted by $\mu^* = \max_{i \in \{1,2,\cdots,K\}} \mu_i$. Let $J(t)$ be the arm pulled at the $t$-th round and $N_i(t)$ be the number of times that the arm $i$ is pulled before the $t$-th round. Then the regret of the gambler at the $T$-th round is given for $\Delta_i = \mu^* - \mu_i$ by

$$\mathrm{Regret}(T) = \sum_{t=1}^{T} \Delta_{J(t)} = \sum_i \Delta_i N_i(T+1) \,.$$

Let $X_{i,n}$ be the $n$-th reward from the arm $i$. We define

$$\bar{x}_{i,n} = \frac{1}{n} \sum_{m=1}^{n} X_{i,m} \,,$$

$$S_{i,n} = \sum_{m=1}^{n} (X_{i,m} - \bar{x}_{i,n})^2 = \sum_{m=1}^{n} X_{i,m}^2 - n\bar{x}_{i,n}^2 \,,$$

that is, $\bar{x}_{i,n}$ and $S_{i,n}$ denote the sample mean and the sum of squares from $n$ samples from the arm $i$, respectively. We denote the sample mean and the sum of squares before the $t$-th round by $\bar{x}_i(t) = \bar{x}_{i,N_i(t)}$ and $S_i(t) = S_{i,N_i(t)}$. It is well known that

$$\bar{x}_{i,n} \sim \mathcal{N}(\mu_i, \sigma_i^2/n) \,, \qquad \frac{S_{i,n}}{\sigma_i^2} \sim \chi_{n-1}^2 \,, \qquad (1)$$

where the chi-squared distribution $\chi_{n-1}^2$ with degree of freedom $n-1$ has the density

$$\chi_{n-1}^2(s) = \frac{s^{\frac{n-3}{2}} \mathrm{e}^{-\frac{s}{2}}}{2^{\frac{n-1}{2}} \Gamma\left(\frac{n-1}{2}\right)} \,.$$

### 2.1 Asymptotic Bound

It is shown in Burnetas and Katehakis (1996) that under any policy satisfying a mild regularity condition the expected regret satisfies

$$\liminf_{T \to \infty} \frac{\mathrm{E}[\mathrm{Regret}(T)]}{\log T}$$
$$\geq \sum_{i:\Delta_i > 0} \frac{\Delta_i}{\inf_{(\mu,\sigma):\mu > \mu^*} D(\mathcal{N}(\mu_i, \sigma_i^2) \| \mathcal{N}(\mu, \sigma^2))} \,, \quad (2)$$

where $D(\cdot \| \cdot)$ is the KL divergence. Since the KL divergence between normal distributions is

$$D(\mathcal{N}(\mu_a, \sigma_b^2) \| \mathcal{N}(\mu_a, \sigma_b^2))$$
$$= \frac{1}{2} \left( \log \frac{\sigma_b^2}{\sigma_a^2} + \frac{\sigma_a^2 + (\mu_b - \mu_a)^2}{\sigma_b^2} - 1 \right) \,,$$

the infimum in (2) is expressed for $\mu_i < \mu^*$ as

$$\frac{1}{2} \log \left( 1 + \frac{(\mu^* - \mu_i)^2}{\sigma_i^2} \right) \,.$$

Therefore, by letting

$$D_{\mathrm{inf}}(\Delta, \sigma^2) = \frac{1}{2} \log \left( 1 + \frac{\Delta^2}{\sigma^2} \right) \,,$$

we can rewrite the theoretical bound in (2) as

$$\liminf_{T \to \infty} \frac{\mathrm{E}[\mathrm{Regret}(T)]}{\log T} \geq \sum_{i:\Delta_i > 0} \frac{\Delta_i}{D_{\mathrm{inf}}(\Delta_i, \sigma_i^2)} \,. \quad (3)$$

---

**Algorithm 1** Thompson Sampling

---

**Parameter:** $\alpha \in \mathbb{R}$.

**Initialization:** Pull each arm $n_0 = \max\{2, 3 - \lceil 2\alpha \rceil\}$ times.

**Loop:**

1. Sample $\tilde{\mu}_i(t)$ from the posterior $\pi_i(\mu_i|\hat{\theta}_i(t))$ under prior $\pi(\mu_i, \sigma_i^2) \sim (\sigma_i^2)^{-1-\alpha}$ for each arm $i$.

2. Pull an arm $i$ maximizing $\tilde{\mu}_i(t)$.

---

## 2.2 Bayesian Theory and Thompson Sampling

Thompson sampling is a policy based on the Bayesian viewpoint. We mainly consider the prior $\pi(\mu_i, \sigma_i^2) \sim (\sigma_i^2)^{-1-\alpha}$, or equivalently, $\pi(\mu_i, \sigma_i) \sim \sigma_i^{-1-2\alpha}$. Since the density of the inverse gamma distribution is

$$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{-1-\alpha} e^{-\frac{\beta}{x}},$$

the above prior for $\sigma_i^2$ corresponds to this distribution with parameters $(\alpha, \beta) = (\alpha, 0)$. The cases $\alpha = -1, -1/2, 0, 1/2$ correspond to uniform for parameter $\sigma_i^2$, uniform for parameter $\sigma_i$, reference and Jeffreys priors, respectively (see, e.g. Robert (2001) for results on Bayesian theory given in this section). Under this prior, the posterior distribution is

$$\pi(\mu_i|\hat{\theta}_{i,n}) \sim \left(1 + \frac{n(\mu_i - \bar{x}_{i,n})^2}{S_{i,n}}\right)^{-\frac{n}{2}-\alpha},$$

where $\hat{\theta}_{i,n} = (\bar{x}_{i,n}, S_{i,n})$. Since the density of $t$-distribution with degree of freedom $\nu$ is

$$f_\nu(x) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}, \tag{4}$$

we see that

$$\pi\left(\sqrt{\frac{n(n+2\alpha-1)}{S_{i,n}}}(\mu_i - \bar{x})\bigg|\hat{\theta}_{i,n}\right) = f_{n+2\alpha-1}. \tag{5}$$

Thompson sampling is the policy which chooses an arm randomly according to the probability with which the arm is the optimal when each $\mu_i$ is distributed independently by the posterior $\pi(\mu_i|\hat{\theta}_i(t))$ for $\hat{\theta}_i(t) = (\bar{x}_i(t), S_i(t))$. This policy is formulated as Algorithm 1. Note that we require $\max\{2, 2 - \lceil 2\alpha \rceil\}$ initial pulls to avoid improper posteriors. We use $n_0 = \max\{2, 3 - \lceil 2\alpha \rceil\}$ for simplicity of the analysis.

## 3 Regret of Thompson Sampling

In this section we give the main result of this paper. First we show that TS achieves the asymptotic bound

for the prior $(\sigma_i^2)^{-1-\alpha}$ with $\alpha < 0$.

**Theorem 1.** *Let $\epsilon > 0$ be arbitrary and assume that there is a unique optimal arm. Under Thompson sampling with $\alpha < 0$, the expected regret is bounded as*

$$\mathrm{E}[\mathrm{Regret}(T)] \leq \sum_{i:\Delta_i>0} \frac{\Delta_i \log T}{D_{\inf}(\Delta_i, \sigma_i^2)} + \mathrm{O}((\log T)^{4/5}).$$

See Lemma 6 for the specific representation of the reminder term $\mathrm{O}((\log T)^{4/5})$. We see from this theorem that TS with $\alpha < 0$ is asymptotically optimal in view of (3).

Next we show that TS with $\alpha \geq 0$ cannot achieve the asymptotic bound. To simplify the analysis we consider a two-armed setting more advantageous to the gambler in which the full information on the arm 2 is known beforehand, that is, the prior on the arm 2 is the unit point mass measure $\delta_{\{(\mu_2,\sigma_2^2)\}}$ instead of $\pi(\mu_2, \sigma_2^2) \sim (\sigma_2^2)^{-1-\alpha}$.

**Theorem 2.** *Assume that there are $K = 2$ arms such that $\mu_1 > \mu_2$. Then, under Thompson sampling such that $\tilde{\mu}_1(t) \sim \pi(\mu_1|\hat{\theta}_1(t))$ with $\alpha \geq 0$ and $\tilde{\mu}_2(t) = \mu_2$, there exists a constant $\xi > 0$ independent of $\sigma_2$ such that*

$$\liminf_{T\to\infty} \frac{\mathrm{E}[\mathrm{Regret}(T)]}{\log T} \geq \xi. \tag{6}$$

*In particular, if $\alpha > 0$ then there exist $\xi' > 0$ and $\eta > 0$ such that*

$$\liminf_{N\to\infty} \frac{\mathrm{E}[\mathrm{Regret}(T)]}{T^\eta} \geq \xi'. \tag{7}$$

Eq. (7) means that TS with $\alpha > 0$ suffers a polynomial regret in expectation. Also note that the asymptotic bound in (3) approaches zero for sufficiently small $\sigma_2$ in the above two-armed setting since $D_{\inf}(\Delta_i, \sigma_i^2) \to \infty$ as $\sigma_i \to 0$. Nevertheless, the LHS of (6) does not go to zero as $\sigma_2 \to 0$ because $\xi > 0$ is independent of $\sigma_2$. Therefore TS with $\alpha = 0$ also does not achieve the asymptotic bound at least for sufficiently small $\sigma_2$.

Recall that Jeffreys and reference priors correspond to $\alpha = 1/2$ and $\alpha = 0$, respectively. Therefore this theorem means that TS with these non-informative priors does not achieve the asymptotic bound.

**Remark 1.** Probability that the sample mean satisfies $\bar{x}_i < \mu$ for any $\mu < \mu_i$ becomes large when $\sigma_i^2$ is large. Therefore the posterior probability that the true expectation $\mu_i$ is larger than $\mu > \bar{x}_i$ becomes large when the prior has heavy weight at large $\sigma_i^2$, that is, $\alpha$ is small. As a result, as $\alpha$ decreases, TS becomes a "conservative" policy which chooses a seemingly suboptimal arm frequently. Theorems 1 and 2 mean that the prior should be conservative to some extent and non-informative priors are too optimistic.

**Remark 2.** Although TS with non-informative priors does not achieve the asymptotic bound in the sense of expectation, this fact does not necessarily mean that these priors are "bad" ones. As we can see from a close inspection of the proof of Theorem 2, the expected regret of TS with these priors becomes large because an enormously large regret arises with fairly small probability. Therefore this policy performs well except for the case arising with this small probability, and the authors think that TS with these priors also becomes a good policy in the probably approximately correct (PAC) framework. In any case, we should be aware that these non-informative priors are "risky" in the sense of expectation.

## 4 Inequalities for normal distributions and t-distributions

In this section we derive fundamental inequalities for distributions appearing in Thompson sampling for the normal distribution model. We prove them in Appendix.

First we give a simple inequality to evaluate the ratio of gamma functions which appears in the densities of normal, chi-squared and $t$-distributions.

**Lemma 3.** *For $z \geq 1/2$*

$$\mathrm{e}^{-2/3} \leq \frac{\Gamma(z + \frac{1}{2})}{\Gamma(z)} \leq \mathrm{e}^{1/6}\sqrt{z}\,.$$

Next we give large deviation probabilities (see, e.g., Dembo and Zeitouni (1998)) for empirical means and variances.

**Lemma 4.** *For any $\mu > \mu_i$*

$$\Pr[\bar{x}_{i,n} \geq \mu] \leq \mathrm{e}^{-\frac{n(\mu - \mu_i)^2}{2\sigma_i^2}} \tag{8}$$

*and for any $\sigma^2 > \sigma_i^2$*

$$\Pr[S_{i,n} \geq n\sigma^2] \leq \mathrm{e}^{-nh\left(\frac{\sigma^2}{\sigma_i^2}\right)} \tag{9}$$

*where $h(x) = (x - 1 - \log x)/2 \geq 0$.*

**Remark 3.** It is well known that Mill's ratio (Kendall & Stuart, 1977, Chap. 5) gives a tighter bound for the tail probability of normal distributions, and similar technique can also be applied to the tail weight of $\chi^2$ distributions. However, we use bounds in Lemma 4 based on the large deviation principle because they are simpler and convenient for our analysis.

Finally we evaluate the posterior distribution of the mean for Thompson sampling. Probability that the sample from the posterior is larger than or equal to $\mu$, which is formally defined as

$$p_n(\mu|\hat{\theta}_{i,n}) = \int_\mu^\infty \pi(x|\hat{\theta}_{i,n})\mathrm{d}x\,,$$

is bounded as follows.

**Lemma 5.** *If $\mu > \bar{x}_{i,n}$ and $n \geq n_0$ then*

$$p_n(\mu|\hat{\theta}_{i,n}) \geq A_{n,\alpha}\left(1 + \frac{n(\mu - \bar{x}_{i,n})^2}{S_{i,n}}\right)^{-\frac{n-1}{2} - \alpha} \tag{10}$$

*and*

$$p_n(\mu|\hat{\theta}_{i,n}) \leq \frac{\sqrt{S_{i,n}}}{\mu - \bar{x}_{i,n}}\left(1 + \frac{n(\mu - \bar{x}_{i,n})^2}{S_{i,n}}\right)^{-\frac{n}{2} - \alpha + 1}, \tag{11}$$

*where*

$$A_{n,\alpha} = \frac{1}{2\mathrm{e}^{1/6}\sqrt{\pi(\frac{n}{2} + \alpha)}}\,. \tag{12}$$

## 5 Analysis for Conservative Priors

In this section we show that Thompson sampling achieves the asymptotic bound if $\alpha < 0$. The main result of this section is given as follows.

**Lemma 6.** *Fix any $\alpha < 0$ and assume that $(\mu_1, \sigma_1^2) = (0, 1)$ and the arm 1 is the unique optimal arm. Then, for any $\epsilon < \min_{i:\Delta_i > 0} \Delta_i/2$,*

$\mathrm{E}[\mathrm{Regret}(T)]$

$$\leq \sum_{i:\Delta_i > 0} \Delta_i\left(\frac{\log T}{D_{\inf}(\Delta_i - 2\epsilon, \sigma_i^2 + \epsilon)} + 2 - 2\alpha\right.$$

$$\left. + \frac{\sqrt{\sigma_i^2 + \epsilon}}{\Delta_i - 2\epsilon} + \frac{1}{1 - \mathrm{e}^{-\frac{\epsilon}{2\sigma_i^2}}} + \frac{1}{1 - \mathrm{e}^{-h(1 + \frac{\epsilon}{\sigma_i^2})}}\right)$$

$$+ \Delta_{\max}\left(\frac{1}{1 - \mathrm{e}^{-\frac{\epsilon^2}{2}}} + \frac{1}{1 - \mathrm{e}^{-h(2)}} + \frac{\mathrm{B}(1/2, -\alpha)}{(1 - \mathrm{e}^{-\frac{\epsilon^2}{2}})^2}\right.$$

$$\left. + \frac{2\sqrt{2}}{\epsilon}\frac{(1 + \epsilon^2/8)^{1-\alpha}}{1 - (1 + \epsilon^2/8)^{-1/2}}\right)$$

$$= \sum_{i:\Delta_i > 0} \frac{\Delta_i \log T}{D_{\inf}(\Delta_i - \epsilon, \sigma_i^2 + \epsilon)} + \mathrm{O}(\epsilon^{-4})\,,$$

*where $\Delta_{\max} = \max_i \Delta_i$ and $\mathrm{B}(\cdot, \cdot)$ is the beta function.*

**Corollary 7.** *Under the same assumption as Lemma 6,*

$$\mathrm{E}[\mathrm{Regret}(T)] \leq \sum_{i:\Delta_i > 0} \frac{\Delta_i \log T}{D_{\inf}(\Delta_i, \sigma_i^2)} + \mathrm{O}((\log T)^{4/5})\,.$$

This corollary is straightforward from Lemma 6 with $\epsilon := \mathrm{O}((\log T)^{-1/5})$.

Note that $D_{\inf}(\mu^* - \mu_i, \sigma_2^2)$ is invariant under the location and scale transformation, that is,

$$D_{\inf}(\mu^* - \mu_i, \sigma_i^2) = D_{\inf}\left(\frac{\mu^* - a}{b} - \frac{\mu_i - a}{b}, \frac{\sigma_i^2}{b^2}\right).$$

Thus Theorem 1 easily follows from Corollary 7 by the transformation $((\mu_1, \sigma_1^2), (\mu_2, \sigma_2^2), \cdots, (\mu_K, \sigma_K^2)) \mapsto ((0,1), ((\mu_2 - \mu_1)/\sigma_1, \sigma_2^2/\sigma_1^2), \cdots, ((\mu_K - \mu_1)/\sigma_1, \sigma_K^2/\sigma_1^2))$.

*Proof of Lemma 6.* Define events

$$\mathcal{A}(t) = \{\tilde{\mu}^*(t) \geq -\epsilon\},$$
$$\mathcal{B}_i(t) = \{\bar{x}_i(t) \leq \mu_i + \delta, S_i(t) \leq n(\sigma_i^2 + \epsilon)\},$$

where $\tilde{\mu}^*(t) = \max_i \tilde{\mu}_i(t)$. Then the regret at the round $T$ is bounded as

$\mathrm{Regret}(T)$

$$= \sum_{t=1}^{T} \Delta_{J(t)}$$

$$\leq \Delta_{\max} \sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) \neq 1, \mathcal{A}^c(t)]$$

$$+ \sum_{i=2}^{K} \Delta_i \left(n_0 + \sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) = i, \mathcal{A}(t)]\right)$$

$$\leq \Delta_{\max} \sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) \neq 1, \mathcal{A}^c(t)]$$

$$+ \sum_{i=2}^{K} \Delta_i \left(\sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) = i, \mathcal{A}(t), \mathcal{B}_i(t)]\right.$$

$$\left. + \sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) = i, \mathcal{B}_i^c(l)] + n_0\right), \quad (13)$$

where $\mathbb{1}[\cdot]$ is the indicator function and the superscript "$c$" denotes the complementary set. In the following Lemmas 8–10 we bound the expectation of the above three terms and the proof is completed. □

**Lemma 8.** *If $\alpha < 0$ then*

$$\mathrm{E}\left[\sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) \neq 1 \cup \mathcal{A}^c(t)]\right]$$

$$\leq \frac{1}{1 - \mathrm{e}^{-\frac{\epsilon^2}{8}}} + \frac{1}{1 - \mathrm{e}^{-h(2)}} + \frac{2\sqrt{2}}{\epsilon} \frac{\left(1 + \frac{\epsilon^2}{8}\right)^{1-\alpha}}{1 - \left(1 + \frac{\epsilon^2}{8}\right)^{-1/2}}$$

$$+ \frac{\mathrm{B}(1/2, -\alpha)}{\left(1 - \mathrm{e}^{-\frac{\epsilon^2}{2}}\right)^2}$$

$$= \mathrm{O}(\epsilon^{-4}).$$

**Lemma 9.** *For any $i \neq 1$,*

$$\mathrm{E}\left[\sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) = i, \mathcal{A}(t), \mathcal{B}_i(t)]\right]$$

$$\leq \frac{\log N}{D_{\inf}(\Delta_i - 2\epsilon, \sigma_i^2 + \epsilon)} + 2 - 2\alpha + \frac{\sqrt{\sigma_i^2 + \epsilon}}{\Delta_i - 2\epsilon}$$

$$= \frac{\log N}{D_{\inf}(\Delta_i - 2\epsilon, \sigma_i^2 + \epsilon)} + \mathrm{O}(1).$$

**Lemma 10.** *For any $i \neq 1$,*

$$\mathrm{E}\left[\sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) = i, \mathcal{B}_i^c(t)]\right]$$

$$\leq \frac{1}{1 - \mathrm{e}^{-\frac{\epsilon^2}{2\sigma_i^2}}} + \frac{1}{1 - \mathrm{e}^{-h\left(1 + \frac{\epsilon}{\sigma_i^2}\right)}} = \mathrm{O}(\epsilon^{-2}).$$

We prove Lemma 10 in Appendix and prove Lemmas 8 and 9 in this section.

Whereas the second term of (13) becomes the main term of the regret, the evaluation of the first term is the most difficult point of the proof, which corresponds to Lemma 8. In fact, it is reported in Burnetas and Katehakis (1996) that they were not able to prove the asymptotic optimality of a policy for the normal distribution model because of difficulty of the evaluation corresponding to this term. Also note that this is the term which does not become a constant in the case $\alpha \geq 0$ and is considered in the proof of Theorem 2.

In this paper we evaluate this term by first bounding this term for a fixed statistic $\hat{\theta}_{i,n} = (\bar{x}_{i,n}, S_{i,n})$ and finally taking its expectation, whereas a probability on this statistic is first evaluated in Burnetas and Katehakis (1996). By leaving the evaluation on the distribution of $\hat{\theta}_{i,n}$ to the latter part, we can significantly simplify the integral by variable transformation.

*Proof of Lemma 8.* First we bound the summation as

$$\sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) \neq 1, \mathcal{A}^c(t)]$$

$$= \sum_{n=n_0}^{T} \sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) \neq 1, \mathcal{A}^c(t), N_1(t) = n]$$

$$= \sum_{n=n_0}^{T} \sum_{m=1}^{T}$$

$$\mathbb{1}\left[m \leq \sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) \neq 1, \mathcal{A}^c(t), N_1(t) = n]\right].$$

Note that

$$m \leq \sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) \neq 1, \mathcal{A}^c(t), N_1(t) = n]$$

implies that $\tilde{\mu}_1(t) \leq \tilde{\mu}^*(t) \leq -\epsilon$ occurred for the first $m$ elements of $\{t : \mathcal{A}^c(t), N_1(t) = n\}$. Therefore,

$$
\Pr\left[ m \leq \sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) \neq 1, \ \mathcal{A}^c(t), \ N_1(t) = n] \right]
$$
$$
\leq (1 - p_n(-\epsilon|\hat{\theta}_{1,n}))^m
$$

and we have

$$
\mathrm{E}\left[ \sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) \neq 1 \cup \mathcal{A}^c(t)] \right]
$$
$$
\leq \mathrm{E}\left[ \sum_{n=n_0}^{T} \sum_{m=1}^{T} (1 - p_n(-\epsilon|\hat{\theta}_{1,n}))^m \right]
$$
$$
\leq \sum_{n=n_0}^{T} \mathrm{E}\left[ \frac{1 - p_n(-\epsilon|\hat{\theta}_{1,n})}{p_n(-\epsilon|\hat{\theta}_{1,n})} \right]. \tag{14}
$$

Since $p_n(-\epsilon|\hat{\theta}_{i,n}) \geq 1/2$ for $-\epsilon \leq \bar{x}_{i,n}$ from the symmetry of $t$-distribution, this expectation is partitioned into

$$
\mathrm{E}\left[ \frac{1 - p_n(-\epsilon|\hat{\theta}_{1,n})}{p_n(-\epsilon|\hat{\theta}_{1,n})} \right]
$$
$$
\leq 2\mathrm{E}\left[ \mathbb{1}[-\epsilon \leq \bar{x}_{1,n}] \left(1 - p_n(-\epsilon|\hat{\theta}_{1,n})\right) \right]
$$
$$
+ \mathrm{E}\left[ \frac{\mathbb{1}[\bar{x}_{1,n} \leq -\epsilon]}{p_n(-\epsilon|\hat{\theta}_{1,n})} \right]
$$
$$
\leq \Pr\left[ -\epsilon < \bar{x}_{1,n} \leq -\epsilon/2 \right]
$$
$$
+ \Pr\left[ -\epsilon/2 < \bar{x}_{1,n}, \ S_{1,n} \geq 2n \right]
$$
$$
+ 2\mathrm{E}\left[ \mathbb{1}[-\epsilon/2 < \bar{x}_{i,n}, \ S_{1,n} \leq 2n] \left(1 - p_n(-\epsilon|\hat{\theta}_{i,n})\right) \right]
$$
$$
+ \mathrm{E}\left[ \frac{\mathbb{1}[\bar{x}_{1,n} \leq -\epsilon]}{p_n(-\epsilon|\hat{\theta}_{1,n})} \right]. \tag{15}
$$

From Lemma 4, the first and second terms of (15) are bounded as

$$
\Pr\left[ -\epsilon < \bar{x}_{i,n} \leq -\epsilon/2 \right] \leq \mathrm{e}^{-\frac{n\epsilon^2}{8}},
$$
$$
\Pr\left[ -\epsilon/2 < \bar{x}_{1,n}, \ S_{1,n} \geq 2n \right] \leq \mathrm{e}^{-nh(2)}, \tag{16}
$$

respectively. Next, recall that $\hat{\theta}_{1,n} = (\bar{x}_{i,n}, S_{i,n})$. Then, from the symmetry of $t$-distribution

$$
1 - p_n(-\epsilon|\bar{x}_{i,n}, S_{i,n}) = 1 - p_n(-\bar{x}_{1,n} - \epsilon|0, S_{1,n})
$$
$$
= p_n(\bar{x}_{1,n} + \epsilon|0, S_{1,n})
$$
$$
= p_n(2\bar{x}_{1,n} + \epsilon|\bar{x}_{1,n}, S_{1,n})
$$

and the third term of (15) is bounded from (11) as

$$
\mathrm{E}\left[ \mathbb{1}[-\epsilon/2 < \bar{x}_{i,n}, \ S_{1,n} \leq 2] \left(1 - p_n(-\epsilon|\hat{\theta}_{i,n})\right) \right]
$$
$$
= \mathrm{E}\left[ \mathbb{1}[-\epsilon/2 < \bar{x}_{i,n}, \ S_{1,n} \leq 2] \, p_n(2\bar{x}_{i,n} - \epsilon|\hat{\theta}_{i,n}) \right]
$$
$$
\leq \frac{2\sqrt{2}}{\epsilon} \left(1 + \frac{\epsilon^2}{8}\right)^{-\frac{n}{2} - \alpha + 1}. \tag{17}
$$

Finally we evaluate the fourth term of (15). From (1) and (10), we have

$$
\mathrm{E}\left[ \frac{\mathbb{1}[\bar{x}_{1,n} \leq -\epsilon]}{p_n(-\epsilon|\hat{\theta}_{1,n})} \right]
$$
$$
\leq \frac{1}{A_{n,\alpha}} \int_{-\infty}^{-\epsilon} \int_0^{\infty} \left(1 + \frac{n(x+\epsilon)^2}{s}\right)^{\frac{n-1}{2}+\alpha}
$$
$$
\cdot \sqrt{\frac{n}{2\pi}} \mathrm{e}^{-\frac{nx^2}{2}} \frac{s^{\frac{n-3}{2}} \mathrm{e}^{-\frac{s}{2}}}{2^{\frac{n-1}{2}} \Gamma(\frac{n-1}{2})} \mathrm{d}s \mathrm{d}x
$$
$$
\leq \frac{\mathrm{e}^{-n\epsilon^2/2}}{A_{n,\alpha}} \int_{-\infty}^{-\epsilon} \int_0^{\infty} \left(1 + \frac{n(x+\epsilon)^2}{s}\right)^{\frac{n-1}{2}+\alpha}
$$
$$
\cdot \sqrt{\frac{n}{2\pi}} \mathrm{e}^{-\frac{n(x+\epsilon)^2}{2}} \frac{s^{\frac{n-3}{2}} \mathrm{e}^{-\frac{s}{2}}}{2^{\frac{n-1}{2}} \Gamma(\frac{n-1}{2})} \mathrm{d}s \mathrm{d}x \tag{18}
$$

by $x^2 \geq (x+\epsilon)^2 + \epsilon^2$ for $x \leq -\epsilon \leq 0$. By letting

$$
(x, s) = \left(-\epsilon - \sqrt{\frac{2zw}{n}}, \ 2z(1-w)\right),
$$

we have

$$
\mathrm{d}x\mathrm{d}s = \det\begin{pmatrix} -\sqrt{\frac{w}{2nz}} & 2(1-w) \\ -\sqrt{\frac{z}{2nw}} & -2z \end{pmatrix} \mathrm{d}z\mathrm{d}w
$$
$$
= \sqrt{\frac{2z}{nw}} \mathrm{d}z\mathrm{d}w
$$

and (18) is rewritten as

$$
\mathrm{E}\left[ \frac{\mathbb{1}[\bar{x}_{i,n} \leq -\epsilon]}{p_n(-\epsilon|\hat{\theta}_{i,n})} \right]
$$
$$
\leq \frac{\mathrm{e}^{-n\epsilon^2/2}}{A_{n,\alpha}} \int_0^1 \int_0^{\infty} \left(1 + \frac{w}{1-w}\right)^{\frac{n-1}{2}+\alpha}
$$
$$
\cdot \sqrt{\frac{n}{2\pi}} \mathrm{e}^{-zw} \frac{(z(1-w))^{\frac{n-3}{2}} \mathrm{e}^{-z(1-w)}}{2\Gamma(\frac{n-1}{2})} \sqrt{\frac{2z}{nw}} \mathrm{d}z\mathrm{d}w
$$
$$
= \frac{\mathrm{e}^{-n\epsilon^2/2}}{2\sqrt{\pi} A_{n,\alpha} \Gamma(\frac{n-1}{2})} \int_0^1 w^{-1/2}(1-w)^{-1-\alpha} \mathrm{d}w
$$
$$
\cdot \int_0^{\infty} \mathrm{e}^{-z} z^{\frac{n}{2}-1} \mathrm{d}z
$$
$$
= \frac{\mathrm{e}^{-n\epsilon^2/2}}{2\sqrt{\pi} A_{n,\alpha} \Gamma(\frac{n-1}{2})} \mathrm{B}(1/2, -\alpha)\Gamma(n/2)
$$
$$
\leq \frac{\mathrm{e}^{1/3} \sqrt{(n+2\alpha)(n-1)}}{2} \mathrm{e}^{-n\epsilon^2/2} \mathrm{B}(1/2, -\alpha)
$$
$$
\qquad \text{(by Lemma 3 and (12))}
$$
$$
\leq n\mathrm{e}^{-n\epsilon^2/2}\mathrm{B}(1/2, -\alpha). \quad \text{(by } \mathrm{e}^{1/3} \leq 2 \text{ and } \alpha < 0) \tag{19}
$$

By combining (14), (15), (16), (17) with (19), we ob-

tain

$$
\mathrm{E}\left[\sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) \neq 1,\ \mathcal{A}^c(t)]\right]
$$

$$
\leq \sum_{n=n_0}^{T} \left( \mathrm{e}^{-\frac{n\epsilon^2}{8}} + \mathrm{e}^{-nh(2)} + \frac{2\sqrt{2}}{\epsilon}\left(1+\frac{\epsilon^2}{8}\right)^{-\frac{n}{2}-\alpha+1} \right.
$$

$$
\left. + n\mathrm{e}^{-n\epsilon^2/2}\mathrm{B}(1/2,-\alpha) \right)
$$

$$
\leq \frac{1}{1-\mathrm{e}^{-\frac{\epsilon^2}{8}}} + \frac{1}{1-\mathrm{e}^{-h(2)}} + \frac{2\sqrt{2}}{\epsilon}\frac{\left(1+\frac{\epsilon^2}{8}\right)^{1-\alpha}}{1-\left(1+\frac{\epsilon^2}{8}\right)^{-1/2}}
$$

$$
+ \frac{\mathrm{B}(1/2,-\alpha)}{\left(1-\mathrm{e}^{-\frac{\epsilon^2}{2}}\right)^2}
$$

$$
= \mathrm{O}(\epsilon^{-2}) + \mathrm{O}(1) + \mathrm{O}(\epsilon^{-3}) + \mathrm{O}(\epsilon^{-4}) = \mathrm{O}(\epsilon^{-4}). \quad \square
$$

*Proof of Lemma 9.* Let $n_i > 0$ be arbitrary. Then

$$
\sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t)=i,\ \mathcal{A}(t),\ \mathcal{B}_i(t)]
$$

$$
\leq \sum_{t=Kn_0+1}^{T} \mathbb{1}[\tilde{\mu}_i(t) \geq -\epsilon,\ B_i(l)]
$$

$$
\leq n_i + \sum_{t=Kn_0+1}^{T} \mathbb{1}[\tilde{\mu}_i(t) \geq -\epsilon,\ \mathcal{B}_i(t),\ N_i(t) \geq n_i]. \tag{20}
$$

Under the condition $\{\mathcal{B}_i(t),\ N_i(t)=n\}$, the probability of the event $\tilde{\mu}_i(t) \geq -\epsilon = \mu^* - \epsilon$ is bounded from Lemma 5 as

$$
p_n(-\epsilon|\hat{\theta}_{i,n}) \leq \frac{\sqrt{\sigma_i^2+\epsilon}}{\Delta_i - 2\epsilon}\left(1+\frac{(\Delta_i-2\epsilon)^2}{\sigma_i^2+\epsilon}\right)^{-\frac{n}{2}-\alpha+1}
$$

$$
= \frac{\sqrt{\sigma_i^2+\epsilon}}{\Delta_i - 2\epsilon}\mathrm{e}^{-(n+2\alpha-2)D_{\inf}(\Delta_i-2\epsilon,\sigma_i^2+\epsilon)}.
$$

Therefore the expectation of (20) is bounded as

$$
\mathrm{E}\left[\sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t)=i,\ \mathcal{A}(t),\ \mathcal{B}_i(t)]\right]
$$

$$
\leq n_i + \sum_{t=Kn_0+1}^{T} \mathrm{Pr}\left[\tilde{\mu}_i(t) \geq -\epsilon,\ \mathcal{B}_i(t),\ N_i(t) \geq n_i\right]
$$

$$
\leq n_i + T\frac{\sqrt{\sigma_i^2+\epsilon}}{\Delta_i - 2\epsilon}\mathrm{e}^{-(n_i+2\alpha-2)D_{\inf}(\Delta_i-2\epsilon,\sigma_i^2+\epsilon)}
$$

and we complete the proof by letting $n_i = (\log T)/D_{\inf}(\Delta_i - 2\epsilon, \sigma_i^2 + \epsilon) + 2 - 2\alpha.$ $\quad\square$

## 6 Analysis for Optimistic Priors

In this section we prove Theorem 2. As mentioned before, the evaluation in the proof corresponds to Lemma 8, in which $\alpha < 0$ is required so that $\mathrm{B}(1/2, -\alpha)$ becomes finite. We show in the following proof that this requirement is actually necessary to achieve the asymptotic bound.

*Proof of Theorem 2.* We assume $(\mu_1, \sigma_1^2) = (0, 1)$ without loss of generality. Fix any $n \geq n_0$ and let $T_{1,n} \in \mathbb{N} \cup \{+\infty\}$ be the first round at which the number of samples from the arm 1 is $n$, that is, we define $T_{1,n} = \min\{t : N_1(t) = n\}$. Since $T_{1,n} = t$ implies that the arm 2 is pulled $t - n - 1$ times through the first $t - 1$ rounds, we have

$$
\mathrm{E}[\mathrm{Regret}(T)]
$$

$$
= \Delta_i \sum_{t=1}^{\infty} \mathrm{Pr}[T_{1,n}=t]\mathrm{E}\left[\sum_{m=1}^{2T} \mathbb{1}[J(m)=2]\,\middle|\,T_{1,n}=t\right]
$$

$$
\geq \Delta_i \sum_{t=1}^{T} \mathrm{Pr}[T_{1,n}=t]\mathrm{E}\left[\sum_{m=1}^{2T} \mathbb{1}[J(m)=2]\,\middle|\,T_{1,n}=t\right]
$$

$$
+ \Delta_i \sum_{t=T+1}^{\infty} \mathrm{Pr}[T_{1,n}=t](T-n). \tag{21}
$$

Now we consider the case $t \leq T$. The conditional expectation in (21) is bounded as

$$
\mathrm{E}\left[\sum_{m=1}^{2T} \mathbb{1}[J(m)=2]\,\middle|\,T_{1,n}=t\right]
$$

$$
\geq \mathrm{E}\left[\sum_{m=t}^{T+t-1} \mathbb{1}[J(m)=2,\ N_1(m)=n]\,\middle|\,T_{1,n}=t\right]
$$

Note that if $\{J(m) \neq 2,\ N_1(m) = n\}$ then $N_1(m') > n$ for any $m' > m$. Therefore, for any $m \geq T_{1,n}$,

$$
\{J(m)=2,\ N_1(m)=n\} \Leftrightarrow \bigcup_{k=0}^{m-T_{1,n}} \{J(T_{1,n}+k)=2\}
$$

$$
\Leftarrow \bigcup_{k=0}^{m-T_{1,n}} \{\tilde{\mu}_1(T_{1,n}+k) < \mu_2\}
$$

since $\tilde{\mu}_2(t) = \mu_2$ always holds. By defining $\bar{p}_n(\mu_2|\hat{\theta}_{1,n}) = 1 - p_n(\mu_2|\hat{\theta}_{1,n})$ we have

$$
\mathrm{E}\left[\sum_{m=1}^{T+t-1} \mathbb{1}[J(m)=2]\,\middle|\,T_{1,n}=t\right]
$$

$$
\geq \mathrm{E}\left[\sum_{m=t}^{T+t-1} \mathbb{1}\left[\bigcup_{k=0}^{m-t} \{\tilde{\mu}_1(t+k) < \mu_2\}\right]\,\middle|\,T_{1,n}=t\right]
$$

$$= \mathrm{E}\left[\sum_{m=t}^{T+t-1}(\bar{p}_n(\mu_2|\hat{\theta}_{1,n}))^{m-t+1}\right]$$

$$\geq \mathrm{E}\left[\sum_{m=1}^{T}(\bar{p}_n(\mu_2|\hat{\theta}_{1,n}))^m\right]$$

$$= \mathrm{E}\left[\left(1-(\bar{p}_n(\mu_2|\hat{\theta}_{1,n}))^T\right)\frac{\bar{p}_n(\mu_2|\hat{\theta}_{1,n})}{p_n(\mu_2|\hat{\theta}_{1,n})}\right]$$

$$\geq \frac{1}{2}\mathrm{E}\left[\mathbb{1}\left[(\bar{p}_n(\mu_2|\hat{\theta}_{1,n}))^T \leq 1/2\right]\frac{\bar{p}_n(\mu_2|\hat{\theta}_{1,n})}{p_n(\mu_2|\hat{\theta}_{1,n})}\right]$$

$$\geq \frac{1}{2}\mathrm{E}\left[\frac{\mathbb{1}\left[(\bar{p}_n(\mu_2|\hat{\theta}_{1,n}))^T \leq 1/2\right]}{p_n(\mu_2|\hat{\theta}_{1,n})}\right] - \frac{1}{2}. \qquad (22)$$

$$\text{(by } (1-p)/p = 1/p - 1)$$

Here we obtain from (10) that

$$(\bar{p}_n(\mu_2|\hat{\theta}_{1,n}))^T \leq 1/2$$

$$\Leftrightarrow p_n(\mu_2|\hat{\theta}_{1,n}) \geq 1 - 2^{-\frac{1}{T}}$$

$$\Leftarrow \left(1+\frac{n(\mu_2-\bar{x}_{1,n})^2}{S_{1,n}}\right)^{-\frac{n-1}{2}-\alpha} \geq \frac{1-2^{-\frac{1}{T}}}{A_{n,\alpha}}$$

$$\Leftarrow \left(1+\frac{n(\mu_2-\bar{x}_{1,n})^2}{S_{1,n}}\right)^{-\frac{n-1}{2}-\alpha} \geq \frac{\log 2}{A_{n,\alpha}T}$$

$$\text{(by } 2^x \geq 1 + x\log 2)$$

$$\Leftrightarrow \frac{n(\mu_2-\bar{x}_{1,n})^2}{S_{1,n}} \leq \left(\frac{A_{n,\alpha}T}{\log 2}\right)^{\frac{1}{\frac{n-1}{2}+\alpha}} - 1 =: C_T. \qquad (23)$$

Therefore the expectation in (22) is bounded from (11) as

$$\mathrm{E}\left[\frac{\mathbb{1}\left[(\bar{p}_n(\mu_2|\hat{\theta}_{1,n}))^T \leq 1/2\right]}{p_n(\mu_2|\hat{\theta}_{1,n})}\right]$$

$$\geq \iint\limits_{\substack{x\leq\mu_2,\,s\geq 0,\\ \frac{n(\mu_2-x)^2}{s}\leq C_T}} \frac{\mu_2-x}{\sqrt{s}}\left(1+\frac{n(\mu_2-x)^2}{s}\right)^{\frac{n}{2}+\alpha-1}$$

$$\cdot\sqrt{\frac{n}{2\pi}}e^{-\frac{nx^2}{2}}\frac{s^{\frac{n-3}{2}}e^{-\frac{s}{2}}}{2^{\frac{n-1}{2}}\Gamma(\frac{n-1}{2})}\mathrm{d}x\mathrm{d}s$$

$$= \frac{\sqrt{n}e^{-\mu_2^2/2}}{\sqrt{\pi}2^{\frac{n}{2}}\Gamma(\frac{n-1}{2})}\iint\limits_{\substack{x\leq\mu_2,\,s\geq 0,\\ \frac{n(\mu_2-x)^2}{s}\leq C_T}} e^{-\frac{n(\mu-x)^2}{2}+2\mu_2(\mu_2-x)}$$

$$\cdot(\mu_2-x)\left(1+\frac{n(\mu_2-x)^2}{s}\right)^{\frac{n}{2}+\alpha-1}s^{\frac{n}{2}-2}e^{-\frac{s}{2}}\mathrm{d}x\mathrm{d}s.$$

By letting

$$(x,s) = \left(\mu_2 - \sqrt{\frac{2zw}{n}},\, 2z(1-w)\right),$$

we obtain in a similar way to (19) that

$$\mathrm{E}\left[\frac{\mathbb{1}\left[(\bar{p}_n(\mu_2|\hat{\theta}_{1,n}))^T \leq 1/2\right]}{p_n(\mu_2|\hat{\theta}_{1,n})}\right]$$

$$\geq \frac{e^{-\mu_2^2/2}}{2\sqrt{\pi n}\Gamma(\frac{n-1}{2})}\int_0^{\frac{1}{1+\frac{1}{C_T}}}\int_0^\infty e^{-z}e^{2\mu_2\sqrt{\frac{2zw}{n}}}$$

$$\cdot z^{\frac{n}{2}-1}(1-w)^{-1-\alpha}\mathrm{d}z\mathrm{d}w$$

$$\geq \frac{e^{-\mu_2^2/2}}{2\sqrt{\pi n}\Gamma(\frac{n-1}{2})}\int_0^\infty e^{2\mu_2\sqrt{\frac{2z}{n}}}e^{-z}z^{\frac{n}{2}-1}\mathrm{d}z$$

$$\cdot\int_0^{\frac{1}{1+\frac{1}{C_T}}}(1-w)^{-1-\alpha}\mathrm{d}w. \qquad \text{(by } \mu_2 < \mu_1 = 0)$$

Here note that

$$\int_0^{\frac{1}{1+\frac{1}{C_T}}}(1-w)^{-1-\alpha}\mathrm{d}w = \begin{cases} \log(1+C_T), & \alpha = 0, \\ \frac{(1+C_T)^\alpha-1}{\alpha}, & \alpha > 0. \end{cases}$$

Then there exists a constant $B_{n,\alpha,\mu_2}$ such that

$$\mathrm{E}\left[\frac{\mathbb{1}\left[(\bar{p}_n(\mu_2|\hat{\theta}_{1,n}))^T \leq 1/2\right]}{p_n(\mu_2|\hat{\theta}_{1,n})}\right]$$

$$\geq \begin{cases} B_{n,\alpha,\mu_2}\log(1+C_T), & \alpha = 0, \\ B_{n,\alpha,\mu_2}((1+C_T)^\alpha-1), & \alpha > 0. \end{cases} \qquad (24)$$

Finally by putting (21), (22) and (24) together we obtain for $\alpha = 0$ that

$$\mathrm{E}[\mathrm{Regret}(2T)]$$
$$\geq \Delta_2 \min\left\{T-n,\, (1/2)B_{n,\alpha,\mu_2}\log(1+C_T) - 1/2\right\}.$$

Eq. (6) follows since $n \geq n_0$ is fixed and $C_T$ defined in (23) is polynomial in $T$. Eq. (7) for $\alpha > 0$ is obtained in the same way. $\qquad\square$

## 7 Conclusion

We considered the stochastic multiarmed bandit problem such that each reward follows a normal distribution with an unknown mean and variance. We proved that Thompson sampling with prior $\pi(\mu_i,\sigma_i^2) \sim (\sigma_i^2)^{-1-\alpha}$ achieves the asymptotic bound if $\alpha < 0$ but cannot if $\alpha \geq 0$, which includes reference prior $\alpha = 0$ and Jeffreys prior $\alpha = 1/2$.

A future work is to examine whether TS with non-informative priors is risky or not for other multiparameter models where TS is used without theoretical analysis (see e.g., Chapelle and Li (2012)). Since the analysis of this paper heavily depends on the specific form of normal distributions, it is currently unknown whether the technique of this paper can be applied to other models and this generalization remains as an important open problem.

# Appendix: Proof of Lemmas

We prove Lemmas 3, 4, 5 and 10 in this appendix.

First we prove Lemma 3 on the ratio of gamma functions.

*Proof of Lemma 3.* Since

$$\sqrt{2\pi}z^{z-1/2}\mathrm{e}^{-z} \leq \Gamma(z) \leq \sqrt{2\pi}\mathrm{e}^{1/6}z^{z-1/2}\mathrm{e}^{-z}$$

for $z \geq 1/2$ from Stirling's formula (Olver et al., 2010, Sect. 5.6(i)), we have

$$\begin{aligned}
\frac{\Gamma(z+\frac{1}{2})}{\Gamma(z)} &\geq \mathrm{e}^{-2/3}\sqrt{z}\left(1+\frac{1}{2z}\right)^z \\
&\geq \mathrm{e}^{-2/3}\sqrt{1/2}\left(1+\frac{1}{2\cdot 1/2}\right)^{1/2} \\
&= \mathrm{e}^{-2/3}\,.
\end{aligned}$$

Similarly we have

$$\begin{aligned}
\frac{\Gamma(z+\frac{1}{2})}{\Gamma(z)} &\leq \mathrm{e}^{-1/3}\sqrt{z}\left(1+\frac{1}{2z}\right)^z \\
&\leq \mathrm{e}^{1/6}\sqrt{z}\,,
\end{aligned}$$

which completes the proof. $\qquad\square$

Next we prove Lemma 4 based on Cramér's theorem (Dembo & Zeitouni, 1998) given below.

**Proposition 11** (Cramér's theorem). *Let $Z_1, Z_2, \cdots$ be i.i.d. random variables on $\mathbb{R}^d$. Then, for $\bar{Z} = n^{-1}\sum_{m=1}^n Z_m \in \mathbb{R}^d$ and any convex set $C \in \mathbb{R}^d$,*

$$\Pr[\bar{Z} \in C] \leq \exp\left(-\inf_{z\in C}\Lambda^*(z)\right),$$

*where*

$$\Lambda^*(z) = \sup_{\lambda\in\mathbb{R}^d}\{\lambda\cdot z - \log\mathrm{E}[\mathrm{e}^{\lambda\cdot Z_1}]\}\,.$$

*Proof of Lemma 4.* Eq. (8) is straightforward from Cramér's theorem with $Z_m := X_{i,m}$ (see also e.g. Dembo & Zeitouni, 1998, Ex. 2.2.23).

Now we show (9). Let $Z_m = (Z_m^{(1)}, Z_m^{(2)}) := (X_{i,m}, X_{i,m}^2) \in \mathbb{R}^2$. Then it is easy to see that the Fenchel-Legendre transform of the cumulant generating function of $Z_i$ is given by

$$\begin{aligned}
&\Lambda^*(z^{(1)}, z^{(2)}) \\
&= \begin{cases} h\left(\frac{z^{(2)}-(z^{(1)})^2}{\sigma_i^2}\right) + \frac{(z^{(1)}-\mu_i)^2}{2\sigma_i^2}, & z^{(2)} > (z^{(1)})^2, \\ +\infty, & z^{(2)} \leq (z^{(1)})^2. \end{cases}
\end{aligned}$$

Eq. (9) follows from

$$\begin{aligned}
&\Pr[S_{i,n} \geq n\sigma^2] \\
&= \Pr[\bar{Z}^{(2)} - (\bar{Z}^{(1)})^2 \geq \sigma^2] \\
&\leq \exp\left(-n \inf_{(z^{(1)},z^{(2)}):\, z^{(2)}-(z^{(1)})^2\geq\sigma^2}\Lambda^*(z^{(1)}, z^{(2)})\right) \\
&\leq \exp\left(-nh\left(\frac{\sigma^2}{\sigma_i^2}\right)\right),
\end{aligned}$$

where the first and the second inequalities follow because $\{(z^{(1)}, z^{(2)}):\, z^{(2)} - (z^{(1)})^2 \geq \sigma^2\}$ is a convex set and $h(x)$ is increasing in $x \geq 1$, respectively. $\qquad\square$

Next we prove Lemma 5 based on Lemma 3.

*Proof of Lemma 5.* Letting

$$\begin{aligned}
\tilde{A} &= \frac{\Gamma(\frac{n}{2}+\alpha)}{\sqrt{\pi(n+2\alpha-1)}\Gamma(\frac{n-1}{2}+\alpha)}\,, \\
x_0 &= \sqrt{\frac{n(n+2\alpha-1)}{S_{i,n}}}(\mu - \bar{x}_{i,n})\,,
\end{aligned}$$

we can express $p_n(\mu|\hat{\theta}_{i,n})$ from (4) and (5) as

$$p_n(\mu|\hat{\theta}_{i,n}) = \tilde{A}\int_{x_0}^\infty \left(1+\frac{x^2}{n+2\alpha-1}\right)^{-\frac{n}{2}-\alpha}\mathrm{d}x. \quad (25)$$

This integral is bounded from below by

$$\begin{aligned}
&p_n(\mu|\hat{\theta}_{i,n}) \\
&= \tilde{A}\int_{x_0}^\infty \left(1+\frac{x^2}{n+2\alpha-1}\right)^{\frac{1}{2}} \\
&\qquad\qquad \cdot\left(1+\frac{x^2}{n+2\alpha-1}\right)^{-\frac{n+1}{2}-\alpha}\mathrm{d}x \\
&\geq \tilde{A}\int_{x_0}^\infty \frac{x}{\sqrt{n+2\alpha-1}}\left(1+\frac{x^2}{n+2\alpha-1}\right)^{-\frac{n+1}{2}-\alpha}\mathrm{d}x \\
&\qquad\qquad\qquad\qquad (\text{by } \sqrt{1+u^2} \geq u) \\
&= \frac{\tilde{A}}{\sqrt{n+2\alpha-1}}\left(1+\frac{x_0^2}{n+2\alpha-1}\right)^{-\frac{n-1}{2}-\alpha}.
\end{aligned}$$

From Lemma 3

$$\begin{aligned}
\frac{\tilde{A}}{\sqrt{n+2\alpha-1}} &= \frac{\Gamma(\frac{n}{2}+\alpha)}{2\sqrt{\pi}\Gamma(\frac{n+1}{2}+\alpha)} \\
&\geq \frac{1}{2\mathrm{e}^{1/6}\sqrt{\pi(\frac{n}{2}+\alpha)}}
\end{aligned}$$

and we obtain (10).

On the other hand, the integral (25) is bounded from

above by

$$p_n(\mu|\hat{\theta}_{i,n})$$

$$= \tilde{A} \int_{x_0}^{\infty} \frac{1}{x} \cdot x \left(1 + \frac{x^2}{n + 2\alpha - 1}\right)^{-\frac{n}{2} - \alpha} \mathrm{d}x$$

$$= \tilde{A} \left[\frac{1}{x} \cdot \frac{\frac{n-1}{2} + \alpha}{\frac{n-2}{2} + \alpha} \left(1 + \frac{x^2}{n + 2\alpha - 1}\right)^{-\frac{n}{2} - \alpha + 1}\right]_{\infty}^{x_0}$$

$$\quad - \tilde{A} \int_{x_0}^{\infty} \frac{2}{x^2} \frac{\frac{n-1}{2} + \alpha}{\frac{n-2}{2} + \alpha} \left(1 + \frac{x^2}{n + 2\alpha - 1}\right)^{-\frac{n}{2} - \alpha + 1} \mathrm{d}x$$

$$\leq \frac{\tilde{A}}{x_0} \frac{\frac{n-1}{2} + \alpha}{\frac{n-2}{2} + \alpha} \left(1 + \frac{x_0^2}{n + 2\alpha - 1}\right)^{-\frac{n}{2} - \alpha + 1} .$$

From Lemma 3

$$\frac{\tilde{A}}{x_0} \frac{\frac{n-1}{2} + \alpha}{\frac{n-2}{2} + \alpha} = \frac{\Gamma(\frac{n-2}{2} + \alpha)}{2\sqrt{\pi n}\Gamma(\frac{n-1}{2} + \alpha)} \frac{\sqrt{S_{i,n}}}{\mu - \bar{x}_{i,n}}$$

$$\leq \frac{1}{2\sqrt{\pi n}\mathrm{e}^{-2/3}} \frac{\sqrt{S_{i,n}}}{\mu - \bar{x}_{i,n}}$$

$$\leq \frac{\sqrt{S_{i,n}}}{\mu - \bar{x}_{i,n}}$$

and we complete the proof. $\qquad\square$

Finally we prove Lemma 10 for the proof of Lemma 6.

*Proof of Lemma 10.* First we have

$$\sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) = i, \, \mathcal{B}_i^c(t)]$$

$$= \sum_{n=n_0}^{T} \mathbb{1}\left[\bigcup_{t=Kn_0+1}^{T} \{J(t) = i, \, \mathcal{B}_i^c(t), \, N_i(t) = n\}\right]$$

$$\leq \sum_{n=n_0}^{T} \mathbb{1}\left[\bar{x}_{i,n} \geq \mu_i + \delta \text{ or } S_{i,n} \geq n(\sigma_i^2 + \epsilon)\right] .$$

Therefore, from Lemma 4,

$$\mathrm{E}\left[\sum_{t=Kn_0+1}^{T} \mathbb{1}[J(t) = i, \, \mathcal{B}_i^c(t)]\right]$$

$$\leq \sum_{n=n_0}^{T} \left(\mathrm{e}^{-n\frac{\epsilon^2}{2\sigma_i^2}} + \mathrm{e}^{-nh\left(1 + \frac{\epsilon}{\sigma_i^2}\right)}\right)$$

$$\leq \frac{1}{1 - \mathrm{e}^{-\frac{\epsilon^2}{2\sigma_i^2}}} + \frac{1}{1 - \mathrm{e}^{-h\left(1 + \frac{\epsilon}{\sigma_i^2}\right)}}$$

$$= \mathrm{O}(\epsilon^{-2}) + \mathrm{O}(\epsilon^{-2}) = \mathrm{O}(\epsilon^{-2}) . \qquad\square$$

# References

Agrawal, S., & Goyal, N. (2012). Analysis of thompson sampling for the multi-armed bandit problem. *Proceedings of COLT 2012*, *23*, 39.1–39.26.

Auer, P., Cesa-Bianchi, N., & Fischer, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, *47*, 235–256.

Burnetas, A. N., & Katehakis, M. N. (1996). Optimal adaptive policies for sequential allocation problems. *Advances in Applied Mathematics*, *17*, 122–142.

Chapelle, O., & Li, L. (2012). An empirical evaluation of Thompson sampling. *Proceedings of NIPS 2011* (pp. 1252–1260). Granada, Spain.

Dembo, A., & Zeitouni, O. (1998). *Large deviations techniques and applications*, vol. 38 of *Applications of Mathematics*. New York: Springer-Verlag. Second edition.

Garivier, A., & Cappé, O. (2011). The KL-UCB algorithm for bounded stochastic bandits and beyond. *Proceedings of COLT 2011*. Budapest, Hungary.

Honda, J., & Takemura, A. (2010). An asymptotically optimal bandit algorithm for bounded support models. *Proceedings of COLT 2010* (pp. 67–79). Haifa, Israel.

Kaufmann, E., Cappé, O., & Garivier, A. (2012a). On bayesian upper confidence bounds for bandit problems. *Proceedings of Fifteenth International Conference on Artificial Intelligence and Statistics (AISTATS2012)* (pp. 592–600).

Kaufmann, E., Korda, N., & Munos, R. (2012b). Thompson sampling: an asymptotically optimal finite-time analysis. *Proceedings of the 23rd international conference on Algorithmic Learning Theory (ALT'12)* (pp. 199–213). Berlin, Heidelberg: Springer-Verlag.

Kendall, M. G., & Stuart, A. (1977). *The advanced theory of statistics*, vol. 1: Distribution theory. C. Griffin London. 4th edition.

Korda, N., Kaufmann, E., & Munos, R. (2013). Thompson sampling for 1-dimensional exponential family bandits. *Proceedings of NIPS 2013*. Lake Tahoe, NV, USA.

Lai, T. L., & Robbins, H. (1985). Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, *6*, 4–22.

Olver, F. W., Lozier, D. W., Boisvert, R. F., & Clark, C. W. (2010). *NIST handbook of mathematical functions*. New York, NY, USA: Cambridge University Press. 1st edition.

Robbins, H. (1952). Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical Society*, *58*, 527–35.

Robert, C. P. (2001). *The Bayesian choice.* New York: Springer. 2nd edition.

Thompson, W. R. (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika, 25*, 285–294.