

# Thompson Sampling for Gaussian Entropic Risk Bandits

**Ang Ming Liang**

ANGMINGLIANG@U.NUS.EDU

*Faculty of Science, Computational Biology and Mathematics  
National University of Singapore  
21 Lower Kent Ridge Rd, Singapore 119077*

**Eloise Y. Y. Lim**

LIMELOISEYY@U.NUS.EDU

*Faculty of Engineering, Industrial System Engineering  
National University of Singapore  
21 Lower Kent Ridge Rd, Singapore 119077*

**Joel Q. L. Chang**

JOEL.CHANG@U.NUS.EDU

*Faculty of Science, Mathematics  
National University of Singapore  
21 Lower Kent Ridge Rd, Singapore 119077*

## Abstract

The multi-armed bandit (MAB) problem is a ubiquitous decision-making problem that exemplifies the exploration-exploitation tradeoff. Standard formulations exclude risk in decision making. Risk notably complicates the basic reward-maximising objectives, in part because there is no universally agreed definition of it. In this paper, we consider an entropic risk (ER) measure and explore the performance of a Thompson sampling-based algorithm ERTS under this risk measure by providing regret bounds for ERTS and corresponding instance dependent lower bounds.

**Keywords:** Thompson sampling, entropic risk, multi-armed bandits

## 1. Introduction

The multi-armed bandit (MAB) problem is a classic reinforcement learning problem that analyses sequential decision making, in which the learner has access to partial feedback from her decisions. The problem has been garnering interest in recent years. It informs many critical theoretical questions about the role of exploration vs exploitation in reinforcement learning and applies to both theoretical problems and various real-world applications, such as dynamic pricing, clinical trials, and portfolio optimisation.

In the well-known stochastic MAB setting, a player chooses among  $K$  arms, each characterised by an independent reward distribution. During each period, the player plays one arm and observes a random reward from that arm. She then incorporates the information she receives from pulling that arm in choosing the next arm she selects. The player repeats the process for a horizon of  $n$  periods. In each period, the player faces a dilemma between exploring other arms' potential value or exploiting the arm that the player believes offers the highest estimated reward.

In the usual setting, the risk of pulling an arm is not being taken into account. However, in many practical settings, such as financial portfolio optimisation, the risk is often the clients' main concern.

In this regard, the MAB problem can be tweaked to model such risk-aversion. This paper uses entropic risk measure as the risk measure to minimise due to the simple exponential relationship it has with risk-aversion and utility, and devises a Thompson sampling-based learning algorithm that minimises entropic risk.

### 1.1 Related Work

A variety of analyses on MABs involving risk measures have been carried out. [Sani et al. \(2012\)](#) considered the mean-variance as their risk measure. Each arm  $i$  followed a Gaussian distribution with mean  $\mu_i \in [0, 1]$  and variance  $\sigma_i^2 \in [0, 1]$ . The authors provided an LCB-based algorithm with accompanying regret analyses. [Galichet et al. \(2013\)](#) proposed the Multi-Armed Risk-Aware Bandit (MARAB) algorithm with the goal of minimising the number of pulls of risky arms, using the risk measure CVaR. [Vakili and Zhao \(2016\)](#) demonstrated that the instance-dependent and instance-independent regrets in terms of the mean-variance of the reward process over a horizon  $n$  are lower bounded by  $\Omega(\log n)$  and  $\Omega(n^{2/3})$  respectively. [Sun et al. \(2017\)](#) analysed contextual bandits with risk constraints, and developed a meta algorithm which makes use of the online mirror descent algorithm that achieves near-optimal regret with respect to minimising the total cost. [Zhu and Tan \(2020\)](#) designed the first Thompson sampling algorithm for risk measures, particularly the mean-variance risk measure for Gaussian bandits, and proved near-optimal regret bounds under specific regimes. [Chang et al. \(2021\)](#) designed a Thompson sampling algorithm factoring a user’s “risk tolerance” level, either minimising mean rewards under some “maximum risk” criterion, or simply minimising the risk measure. [Baudry et al. \(2020\)](#) designed and analysed Thompson sampling-based algorithms  $\alpha$ -NPTS for bounded rewards and  $\alpha$ -Multinomial-TS for discrete multinomial distributions.

The papers most related to our work is that by [Zhu and Tan \(2020\)](#) and [Chang et al. \(2021\)](#). [Zhu and Tan \(2020\)](#) considered arms with the *highest* mean-variance to be optimal, and their definitions and methods can be analogously defined for *minimising* the mean-variance. [Chang et al. \(2021\)](#) defined arms with the minimum CVaR as optimal in their “infeasible instance”, which produced theoretical analogues for “feasible instances”. This hints that the heavy duty analysis happens in trying to choose arms with the risk measure minimised. Our paper seeks to explore the efficacy of Thompson sampling in the analogous risk-minimising problem setting proposed by [Zhu and Tan \(2020\)](#), but instead considering the *entropic risk* measure. We demonstrate and prove the asymptotic optimality of ERTS, whose asymptotic upper bound matches the theoretical lower bound for consistent algorithms that solve the entropic risk MAB for Gaussian bandits.

### 1.2 Contributions

- **ERTS Algorithm:** We design ERTS, an algorithm that is similar to the structure of CVaR-TS in [Chang et al. \(2021\)](#) but using entropic risk instead of CVaR as the risk measure. This algorithm uses Thompson sampling ([Thompson, 1933](#)) as explored for mean-variance bandits in [Zhu and Tan \(2020\)](#) and CVaR bandits in [Chang et al. \(2021\)](#).
- **Comprehensive regret bounds:** We provide theoretical analysis of the ERTS algorithm for Gaussian bandits with bounded variances. We state and prove both upper and lower bounds, showing that ERTS is the first asymptotically optimal algorithm that solves the entropic risk

multi-armed bandit problem. Our proof techniques solidify the novel  $\xi$ -trick in [Chang et al. \(2021\)](#), and affirm future analysis on MABs involving generalised risk measures.

This paper is structured as follows. We first introduce the formulation of the entropic risk MAB problem in Section 2. In Section 3, we present ERTS algorithm. We present our regret bounds and prove that the upper bound we derived is asymptotically optimal in Section 4. Section 5 provides the proof outlines of the regret bounds in Section 4. We conclude our discussion in Section 6 summarizing our work and suggesting avenues for further research. For brevity, we defer detailed proofs of the theorems to the supplementary material.

## 2. Problem formulation

In this section we define the entropic risk MAB problem. For the rest of the paper, denote  $[k] = \{1, \dots, k\}$  for any  $k \in \mathbb{N}$  and  $(t)^+ = \max\{0, t\}$  for  $t \in \mathbb{R}$ .

**Definition 1.** For any random variable  $X$ , given a risk parameter  $\gamma$ , the *entropic risk* ([Lee et al., 2020](#); [Howard and Matheson., 1972](#)) of  $X$  is defined by

$$\text{ER}_\gamma(X) := \frac{1}{\gamma} \log \mathbb{E}[\exp(-\gamma X)].$$

In this paper, we work with Gaussian random variables  $X \sim \mathcal{N}(\mu, \sigma^2)$ . Direct computations then yield  $\text{ER}_\gamma(X) = -\mu + (\gamma/2)\sigma^2$ , which is consistent with the computation in [Chang et al. \(2021, Section 5\)](#). Setting  $\gamma \rightarrow 0^+$  (resp.  $\gamma \rightarrow +\infty$ ) yields the risk-neutral (resp. risk-averse) setting, since  $\mu$  (resp.  $\sigma^2$ ) dominates in the former (resp. latter) case.

Consider a  $K$ -armed MAB  $\nu = (\nu(i))_{i \in [K]}$  played over a horizon of length  $n$ . Letting  $\text{ER}_\gamma(X)$  denote the entropic risk, our objective is to select the least risky arm, that is, the arm with the lowest entropic risk. Thus, we define an arm  $i$  to be optimal precisely when  $i \in \arg \min_{k \in [K]} \text{ER}_\gamma(\nu(k))$ . Suppose arm 1 is optimal (uniquely, without loss of generality). We can then define  $\Delta_{\text{ER}}(i, \gamma) := \text{ER}_\gamma(i) - \text{ER}_\gamma(1) > 0$  and the *regret* of a policy  $\pi$  by

$$\mathcal{R}_n(\pi) := \sum_{i \in [K] \setminus \{1\}} \mathbb{E}[T_{i,n}] \Delta_{\text{ER}}(i, \gamma),$$

where  $T_{i,n}$  denotes the number of times arm  $i$  was pulled in the first  $n$  rounds. This is a natural definition based on regret decomposition ([Lattimore and Szepesvári, 2020, Chapter 4.5](#)), and in fact corresponds to the regret decomposition in the case  $\gamma \rightarrow 0^+$  (i.e. the risk-neutral setting). In the following, we design and analyse ERTS, which aims to minimise  $\mathcal{R}_n(\pi)$ , and also attain an instance-dependent lower bound, which establishes asymptotic optimality.

## 3. The ERTS Algorithm

In this section, we introduce the Entropic Risk Thompson Sampling (ERTS) algorithm for Gaussian bandits with bounded variances, i.e.,  $\nu \in \mathcal{E}_N^K(\sigma_{\max}^2) := \{\nu = (\nu_1, \dots, \nu_K) : \nu_i \sim \mathcal{N}(\mu_i, \sigma_i^2), \sigma_i^2 \leq \sigma_{\max}^2 \forall i \in [K]\}$  for some  $\sigma_{\max}^2 > 1$ . Similar to [Zhu and Tan \(2020\)](#) and [Chang et al. \(2021\)](#), the algorithm samples from the posteriors of each arm, then chooses the arm according to a multi-criterion procedure.

Denote the mean and precision of the Gaussian by  $\mu$  and  $\psi$  respectively. If  $(\mu, \psi)$  follows the distribution Normal-Gamma( $\mu, T, \alpha, \beta$ ), then  $\psi \sim \text{Gamma}(\alpha, \beta)$ , and  $\mu|\psi \sim \mathcal{N}(\mu, 1/(\psi T))$ . Since the conjugate prior for the Gaussian with unknown mean and variance is the Normal-Gamma distribution, we use Algorithm 1 to update  $(\mu, \psi)$  via Bayes' theorem.

We present the ERTS algorithm. In each round  $t$ , for each arm  $i$ , the player samples the parameters  $(\theta_{it}, \kappa_{it})$  from the posterior distribution of arm  $i$ , then chooses arm  $j = \arg \min_{i \in [k]} \widehat{\text{ER}}_\gamma(i, t)$ , where  $\widehat{\text{ER}}_\gamma(i, t) := -\theta_{i,t} + \gamma/(2\kappa_{i,t})$ , i.e. least risky arm available.

---

**Algorithm 1** Update( $\hat{\mu}_{i,t-1}, T_{i,t-1}, \alpha_{i,t-1}, \beta_{i,t-1}$ )

---

- 1: **Input:** Prior parameters  $(\hat{\mu}_{i,t-1}, T_{i,t-1}, \alpha_{i,t-1}, \beta_{i,t-1})$  and new sample  $X_{i,t}$
  - 2: Update the mean:  $\hat{\mu}_{i,t} = \frac{T_{i,t-1}}{T_{i,t-1}+1} \hat{\mu}_{i,t-1} + \frac{1}{T_{i,t-1}+1} X_{i,t}$
  - 3: Update the number of samples, the shape parameter, and the rate parameter:  $T_{i,t} = T_{i,t-1} + 1$ ,  
 $\alpha_{i,t} = \alpha_{i,t-1} + \frac{1}{2}$ ,  $\beta_{i,t} = \beta_{i,t-1} + \frac{T_{i,t-1}}{T_{i,t-1}+1} \cdot \frac{(X_{i,t} - \hat{\mu}_{i,t-1})^2}{2}$
- 

---

**Algorithm 2** Entropic Risk Thompson Sampling (ERTS)

---

- 1: **Input:** Risk parameter  $\gamma$ ,  $\hat{\mu}_{i,0} = 0$ ,  $T_{i,0} = 0$ ,  $\alpha_{i,0} = \frac{1}{2}$ ,  $\beta_{i,0} = \frac{1}{2}$
  - 2: **for**  $t = 1, 2, \dots, K$  **do**
  - 3:   Play arm  $t$  and update  $\hat{\mu}_{t,t} = X_{t,t}$
  - 4:   Update( $\hat{\mu}_{t,t-1}, T_{t,t-1}, \alpha_{t,t-1}, \beta_{t,t-1}$ )
  - 5: **end for**
  - 6: **for**  $t = K + 1, K + 2, \dots$  **do**
  - 7:   Sample  $\kappa_{i,t}$  from  $\text{Gamma}(\alpha_{i,t-1}, \beta_{i,t-1})$
  - 8:   Sample  $\theta_{i,t}$  from  $\mathcal{N}(\hat{\mu}_{i,t-1}, 1/T_{i,t-1})$
  - 9:   Play arm  $j(t) = \arg \min_{i \in [K]} \widehat{\text{ER}}_\gamma(i, t)$  and observe loss  $X_{j(t),t} \sim \nu(j(t))$
  - 10:   Update( $\hat{\mu}_{j(t),t-1}, T_{j(t),t-1}, \alpha_{j(t),t-1}, \beta_{j(t),t-1}$ )
  - 11: **end for**
- 

## 4. Regret Bound for ERTS and Lower Bounds

We present our regret bounds in the following theorems. These verify the conjecture made in Chang et al. (2021, Section 5) regarding risk measures of Gaussian bandits of the form  $af(\mu) + bg(\sigma^2)$ , where  $(f(x), g(x), a, b) = (x, x, -1, \gamma/2)$ . Furthermore, they establish ERTS as asymptotically optimal in the context of Gaussian entropic risk bandits.

**Theorem 2** (Upper Bound). Fix  $\xi \in (0, 1)$ ,  $\gamma \in (0, \infty)$ . Then the asymptotic regret of ERTS for entropic risk Gaussian MAB bandits satisfies

$$\limsup_{n \rightarrow \infty} \frac{\mathcal{R}_n(\text{ERTS})}{\log n} \leq \sum_{i \in [K] \setminus \{1\}} R_i \Delta_{\text{ER}}(i, \gamma),$$

where

$$R_i := \max \left\{ \frac{2}{\xi_\gamma^2 \Delta_{\text{ER}}^2(i, \gamma)}, \frac{1}{h \left( \frac{\gamma \sigma_i^2}{\gamma \sigma_i^2 - 2(1-\xi) \Delta_{\text{ER}}(i, \gamma)} \right)} \right\}.$$

Furthermore, setting

$$\xi_\gamma = 1 - \frac{\gamma \sigma_i^2}{2 \Delta_{\text{ER}}(i, \gamma)} \left( 1 - \frac{1}{h_+^{-1}(\Delta_{\text{ER}}^2(i, \gamma)/2)} \right),$$

yields

$$\frac{1}{h \left( \frac{\gamma \sigma_i^2}{\gamma \sigma_i^2 - 2(1-\xi_\gamma) \Delta_{\text{ER}}(i, \gamma)} \right)} \leq \frac{2}{\xi_\gamma^2 \Delta_{\text{ER}}^2(i)}$$

and  $\xi_\gamma \rightarrow 1^-$  as  $\gamma \rightarrow 0^+$ , where  $h_+^{-1}(y) = \max \{x : h(x) = y\}$ .

**Remark 3.** The final part of the theorem shows that the upper bound is characterised by the quantity  $2/(\xi_\gamma^2 \Delta_{\text{ER}}^2(i, \gamma))$ . By continuity, we obtain the regret bound  $2/(\Delta_{\text{ER}}^2(i, \gamma))$ . Furthermore, we note that  $\Delta_{\text{ER}}(i, \gamma) \rightarrow -\mu_i - (-\mu_1) = \mu_1 - \mu_i$  as  $\gamma \rightarrow 0^+$ , and thus the upper bound simplifies to  $2/(\mu_1 - \mu_i)^2$ . This agrees with our intuition since  $\text{ER}_\gamma(i) = -\mu_i + (\gamma/2)\sigma_i^2 \rightarrow -\mu_i$  as  $\gamma \rightarrow 0^+$ , implying that we are in the risk-neutral setting. Thus, the results correspond to those derived for mean-variance bandits (Zhu and Tan, 2020) and CVaR bandits (Chang et al., 2021).

Next, we establish an instance-dependent lower bound for the expected pulls of non-optimal arms under consistent algorithms. Consider a class  $\mathcal{C}$  of distributions and define  $\mathcal{S}_i = \{\nu'(i) \in \mathcal{C} : \text{ER}(\nu'(i)) < \text{ER}(1)\}$ . Define for each non-optimal arm  $i$ ,

$$\eta(i, \gamma) = \inf_{\nu'(i) \in \mathcal{S}_i} \{\text{KL}(\nu(i), \nu'(i))\},$$

where  $\text{KL}(\mathbb{P}, \mathbb{P}')$  denotes the *KL-divergence* between two probability measures  $\mathbb{P}, \mathbb{P}'$ .

**Theorem 4 (Lower Bound).** Let  $\pi$  be a policy over the class of distributions  $\mathcal{C}$  satisfying  $\mathcal{R}_n(\pi) = o(n^a)$  for any  $a > 0$ . Then for any non-optimal arm  $i$ , we have

$$\liminf_{n \rightarrow \infty} \frac{\mathbb{E}[T_{i,n}]}{\log n} \geq \frac{1}{\eta(i, \gamma)}.$$

In particular, if  $\mathcal{C} = \mathcal{E}_{\mathcal{N}}^K(\sigma_{\max}^2)$ , then

$$\liminf_{n \rightarrow \infty} \frac{\mathcal{R}_n(\pi)}{\log n} \geq \sum_{i \in [K] \setminus \{1\}} R_i \Delta_{\text{ER}}(i, \gamma).$$

**Remark 5.** This implies that the asymptotic lower bound for the regret matches its asymptotic upper bound in Theorem 2 unconditionally. Hence, for the Gaussian entropic risk MAB problem, ERTS is *asymptotically optimal*.

## 5. Proof Outlines for Theorem 2 and 4

**Theorem 2:** Denote the sample entropic risk as  $\hat{\text{ER}}_\gamma(i, t) = -\theta_{i,t} + \gamma/(2\kappa_{i,t})$ . Fix  $\varepsilon > 0$  and define  $E_i(t) := \left\{ \hat{\text{ER}}_\gamma(i, t) > \text{ER}_\gamma(1) + \varepsilon \right\}$ , that is, the event that the Thompson sample mean of arm  $i$  is  $\varepsilon$ -riskier than a certain threshold or, more precisely,  $\varepsilon$ -higher than the optimal arm (which has the lowest entropic risk). Intuitively, event  $E_i(t)$  occurs with high probability when the algorithm has explored sufficiently. However, the algorithm does not choose arm  $i$  when  $E_i^c(t)$  occurs with small probability under Thompson sampling, which contributes directly to the regret bound. Therefore, it suffices to bound the number of times  $E_i^c(t)$  occurs.

In order to bound  $\mathbb{E}[T_{i,n}]$ , we can split  $\mathbb{E}[T_{i,n}]$  into two parts using a key lemma by [Lattimore and Szepesvári \(2020\)](#) to yield  $\mathbb{E}[T_{i,n}] \leq \Lambda_1 + \Lambda_2 + 1$ , where  $\Lambda_1 = \mathbb{E}\left[\sum_{s=0}^{n-1} \left(\frac{1}{G_{1,s}} - 1\right)\right]$  and  $\Lambda_2 = \sum_{s=0}^{n-1} \mathbb{P}\left(G_{1,s} > \frac{1}{n}\right)$ . It remains to upper bound  $\Lambda_1$  and  $\Lambda_2$ . The techniques to upper bound  $\Lambda_1$  are similar to those from [Zhu and Tan \(2020, Section 4.6\)](#) and [Chang et al. \(2021, Section 5\)](#). To upper bound  $\Lambda_2$ , we split the event  $E_i^c(t) = \left(\hat{\text{ER}}_\gamma(i, t) \leq \text{ER}_\gamma(1) + \varepsilon\right)$  into

$$\begin{aligned} \Psi_1(\xi) &= \{-\theta_{i,t} + \mu_i \leq -\xi(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)\}, \\ \Psi_2(\xi) &= \left\{ \frac{\gamma}{2} \left( \frac{1}{\kappa_{i,t}} - \sigma_i^2 \right) \leq (-1 + \xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon) \right\} \end{aligned}$$

That is,  $E_i^c(t) \subseteq \Psi_1(\xi) \cup \Psi_2(\xi)$ . We then use the union bound which yields  $\mathbb{P}(E_i^c(t)) \leq \mathbb{P}(\Psi_1(\xi)) + \mathbb{P}(\Psi_2(\xi))$ , which we can upper bound by known concentration bounds. Following the strategy employed by [Chang et al. \(2021\)](#), a judicious selection of the free parameter  $\xi \in (0, 1)$  allows us to allocate "weights" on the bounds of  $\mathbb{P}(\Psi_1(\xi))$  and  $\mathbb{P}(\Psi_2(\xi))$  which then yield  $2/(\xi^2 \Delta_{\text{ER}}^2(i, \gamma))$  and  $\left(h\left(\frac{\gamma \sigma_i^2}{\gamma \sigma_i^2 - 2(1-\xi)\Delta_{\text{ER}}(i, \gamma)}\right)\right)^{-1}$  without incurring further residual terms.

**Theorem 4:** The proof of the lower bound follows immediately from [Kagrecha et al. \(2020, Theorem 4\)](#) by replacing the criterion  $c_\alpha(\nu'(k)) \leq c_\alpha^*$  by  $\text{ER}(\nu'(k)) \leq \text{ER}(1)$ . We then particularize the lower bounds therein by decisively setting the distribution of  $\nu'(i)$  to have a Gaussian distribution with mean  $\mu_i + \sigma_i \sqrt{2/R_i} + \varepsilon$  and variance  $\sigma_i^2$ , which then returns the desired lower bound.

## 6. Conclusion

This paper applies Thompson sampling ([Thompson, 1933](#)) to provide the first solution for entropic risk MAB problems which have not been previously considered before to the best of our knowledge. We proposed a new algorithm ERTS to solve this problem and proved that this proposed algorithm is asymptotically optimal for the ER MAB problem. Further work includes analysing Thompson sampling of Gaussian MABs under general risk measures and exploring Thompson sampling's performance for Entropic-Risk sub-Gaussian bandits. We may also potentially design a general framework for proving the efficacy of Thompson sampling over the state-of-the-art L/UCB-based techniques for generalised risk-averse MABs and a wider class of bandits (under reasonable assumptions, such as the crucial properties of the risk-measures, existence of conjugate prior estimates, as well as relevant concentration bounds).

**Appendix A.**

**Proof** [Proof of Theorem 2] We first state without proof a crucial lemma from [Lattimore and Szepesvári \(2020\)](#) which we will use in our analysis.

**Lemma 6** ([Lattimore and Szepesvári \(2020\)](#)). Let  $\mathbb{P}_t(\cdot) = \mathbb{P}(\cdot | A_1, X_1, \dots, A_{t-1}, X_{t-1})$  be the probability measure conditioned on the history up to time  $t - 1$  and  $G_{is} = \mathbb{P}_t(E_i^c(t) | T_{i,t} = s)$ , where  $E_i(t)$  is any specified event for arm  $i$  at time  $t$ . Then

$$\mathbb{E}[T_{i,n}] \leq \sum_{s=0}^{n-1} \mathbb{E} \left[ \frac{1}{G_{1s}} - 1 \right] + \sum_{s=0}^{n-1} \mathbb{P} \left( G_{is} > \frac{1}{n} \right) + 1.$$

Denote the sample entropic risk at  $\gamma$  by  $\widehat{\text{ER}}(i, t) = -\theta_{i,t} + \gamma/(2\kappa_{i,t})$ . Fix  $\varepsilon > 0$ , and define

$$E_i(t) := \left\{ \widehat{\text{ER}}(i, t) > \text{ER}_\gamma(1) + \varepsilon \right\},$$

the event that the Thompson sample entropic risk of arm  $i$  is  $\varepsilon$ -higher than the optimal arm (which has the lowest entropic risk). Intuitively, event  $E_i(t)$  is highly likely to occur when the algorithm has explored sufficiently. However, the algorithm does not choose arm  $i$  when  $E_i^c(t)$ , an event with small probability under Thompson sampling, occurs. By Lemma 6 and the linearity of expectation, we can divide  $\mathbb{E}[T_{i,n}]$  into two parts as

$$\mathbb{E}[T_{i,n}] \leq \sum_{s=0}^{n-1} \mathbb{E} \left[ \frac{1}{G_{1s}} - 1 \right] + \sum_{s=0}^{n-1} \mathbb{P} \left( G_{is} > \frac{1}{n} \right) + 1. \quad (1)$$

By Lemmas 8 and 11 by that which follows, we have

$$\begin{aligned} \sum_{s=1}^n \mathbb{E} \left[ \frac{1}{G_{1s}} - 1 \right] &\leq \frac{C_1}{\varepsilon^3} + \frac{C_2}{\varepsilon^2} + \frac{C_3}{\varepsilon} + C_4, \text{ and} \\ \sum_{s=1}^n \mathbb{P}_t \left( G_{is} > \frac{1}{n} \right) &\leq 1 + \max \left\{ \frac{2 \log(2n)}{\xi^2 (\Delta_{\text{ER}}(i, \gamma) - \varepsilon)^2}, \frac{\log(2n)}{h \left( \frac{\gamma \sigma^2}{\gamma \sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)} \right)} \right\} + \frac{C_5}{\varepsilon^4} + \frac{C_6}{\varepsilon^2}. \end{aligned}$$

Plugging the two displays into (1), we have

$$\mathbb{E}[T_{i,n}] \leq 1 + \max \left\{ \frac{2 \log(2n)}{\xi^2 (\Delta_{\text{ER}}(i, \gamma) - \varepsilon)^2}, \frac{\log(2n)}{h \left( \frac{\gamma \sigma^2}{\gamma \sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)} \right)} \right\} + \frac{C'_1}{\varepsilon^4} + \frac{C'_2}{\varepsilon^3} + \frac{C'_3}{\varepsilon^2} + \frac{C'_4}{\varepsilon} + C'_5, \quad (2)$$

where  $C'_1, C'_2, C'_3, C'_4, C'_5$  are constants. Setting  $\varepsilon = (\log n)^{-\frac{1}{8}}$  into (2), we get

$$\limsup_{n \rightarrow \infty} \frac{\mathcal{R}_n(\text{ERTS})}{\log n} \leq \sum_{i \in [K] \setminus \{1\}} \left( \max \left\{ \frac{2}{\xi^2 \Delta_{\text{ER}}^2(i)}, \frac{1}{h \left( \frac{\gamma \sigma^2}{\gamma \sigma_i^2 - 2(1-\xi) \Delta_{\text{ER}}(i, \gamma)} \right)} \right\} \right) \Delta_{\text{ER}}(i, \gamma). \quad \blacksquare$$

**Lemma 7.** We can lower bound

$$\mathbb{P}_t(E_1^c(t) \mid T_{1,t} = s, \hat{\mu}_{1,s} = \mu, \hat{\sigma}_{1,s} = \sigma) = \mathbb{P}_t\left(\widehat{\text{ER}}_1 \leq \text{ER}_1 + \varepsilon \mid T_{1,t} = s, \hat{\mu}_{1,s} = \mu, \hat{\sigma}_{1,s} = \sigma\right)$$

by

$$\begin{aligned} & \mathbb{P}_t\left(\widehat{\text{ER}}_i \leq \text{ER}_1 + \varepsilon \mid T_{1,t} = s, \hat{\mu}_{1,s} = \mu, \hat{\sigma}_{1,s} = \sigma\right) \\ & \geq \begin{cases} \mathbb{P}_t\left(\theta_{1,t} - \mu_1 \geq -\frac{\varepsilon}{2}\right) \cdot \mathbb{P}_t\left(\frac{1}{\kappa_{1,t}} - \sigma_1 \leq \frac{\varepsilon}{\gamma}\right) & \text{if } \mu \leq \mu_1, \sigma \geq \sigma_1, \\ \frac{1}{2}\mathbb{P}_t\left(\frac{1}{\kappa_{1,t}} - \sigma_1 \leq \frac{\varepsilon}{\gamma}\right) & \text{if } \mu > \mu_1, \sigma \geq \sigma_1, \\ \frac{1}{2}\mathbb{P}_t\left(\theta_{1,t} - \mu_1 \geq -\frac{\varepsilon}{2}\right) & \text{if } \mu \leq \mu_1, \sigma < \sigma_1, \\ \frac{1}{4} & \text{if } \mu > \mu_1, \sigma < \sigma_1. \end{cases} \end{aligned} \quad (3)$$

**Proof** Given  $T_{1,t} = s, \hat{\mu}_{1,s} = \mu, \hat{\sigma}_{1,s} = \sigma$ , a direct calculation gives us,

$$\begin{aligned} & \mathbb{P}_t\left(\widehat{\text{ER}}_1 \leq \text{ER}_1 + \varepsilon \mid T_{1,t} = s, \hat{\mu}_{1,s} = \mu, \hat{\sigma}_{1,s} = \sigma\right) \\ & = \mathbb{P}_t\left(-\theta_{1,t} + \frac{\gamma}{2\kappa_{1,t}} - (-\mu_1 + (\gamma/2)\sigma_1^2) \leq \varepsilon \mid T_{1,t} = s, \hat{\mu}_{1,s} = \mu, \hat{\sigma}_{1,s} = \sigma\right) \\ & = \mathbb{P}_t\left(-(\theta_{1,t} - \mu_1) + \frac{\gamma}{2}\left(\frac{1}{\kappa_{1,t}} - \sigma_1^2\right) \leq \varepsilon \mid T_{1,t} = s, \hat{\mu}_{1,s} = \mu, \hat{\sigma}_{1,s} = \sigma\right) \\ & \geq \mathbb{P}_t\left(-(\theta_{1,t} - \mu_1) \leq \varepsilon/2 \mid T_{1,t} = s, \hat{\mu}_{1,s} = \mu, \hat{\sigma}_{1,s} = \sigma\right) \\ & \quad \mathbb{P}_t\left(\frac{\gamma}{2}\left(\frac{1}{\kappa_{1,t}} - \sigma_1^2\right) \leq \varepsilon/2 \mid T_{1,t} = s, \hat{\mu}_{1,s} = \mu, \hat{\sigma}_{1,s} = \sigma\right) \\ & \geq \begin{cases} \mathbb{P}_t\left(\theta_{1,t} - \mu_1 \geq -\frac{\varepsilon}{2}\right) \cdot \mathbb{P}_t\left(\frac{1}{\kappa_{1,t}} - \sigma_1^2 \leq \frac{\varepsilon}{\gamma}\right) & \text{if } \mu \leq \mu_1, \sigma^2 \geq \sigma_1^2, \\ \frac{1}{2}\mathbb{P}_t\left(\frac{1}{\kappa_{1,t}} - \sigma_1^2 \leq \frac{\varepsilon}{\gamma}\right) & \text{if } \mu > \mu_1, \sigma^2 \geq \sigma_1^2, \\ \frac{1}{2}\mathbb{P}_t\left(\theta_{1,t} - \mu_1 \geq -\frac{\varepsilon}{2}\right) & \text{if } \mu \leq \mu_1, \sigma^2 < \sigma_1^2, \\ \frac{1}{4} & \text{if } \mu > \mu_1, \sigma^2 < \sigma_1^2. \end{cases} \end{aligned}$$

Then the lemma holds since  $\mathbb{P}_t(\theta_{1,t} - \mu_1 \geq -\varepsilon/2) > 1/2$  if  $\mu > \mu_1$ , and  $\mathbb{P}_t\left(\frac{1}{\kappa_{1,t}} - \sigma_1^2 \leq \frac{\varepsilon}{\gamma}\right) \geq 1/2$  if  $\sigma < \sigma_1^2$ , by using properties of the median of the Gaussian and Gamma distributions respectively.  $\blacksquare$

**Lemma 8** (Upper bounding the first term of (1)). We have

$$\sum_{s=1}^n \mathbb{E}\left[\frac{1}{G_{1s}} - 1\right] \leq \frac{C_1}{\varepsilon^2} + \frac{C_2}{\varepsilon} + C_3,$$

where  $C_1, C_2, C_3$ .

**Proof** The proof follows immediately from Lemma 7 and Zhu and Tan (2020, S-3.3) by scaling  $\varepsilon > 0$ .  $\blacksquare$

**Lemma 9.** For  $\xi \in (0, 1)$ , we have

$$\begin{aligned} & \mathbb{P}\left(\widehat{\text{ER}}_i \leq \text{ER}_1 + \varepsilon \mid T_{i,t} = s, \hat{\mu}_{i,t} = \mu, \hat{\sigma}_{i,t}^2 = \sigma^2\right) \\ & \leq \exp\left(-\frac{s}{2}(\mu_i - \mu + \xi(\Delta_{\text{ER}}(i, \gamma) - \varepsilon))^2\right) + \exp\left(-sh\left(\frac{\gamma\sigma^2}{\gamma\sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)}\right)\right), \end{aligned}$$

where  $h(x) = \frac{1}{2}(x - 1 - \log x)$ .

**Proof** For  $\xi \in (0, 1)$ , we have

$$\begin{aligned} & \mathbb{P}\left(\widehat{\text{ER}}_i \leq \text{ER}_1 + \varepsilon \mid T_{i,t} = s, \hat{\mu}_{i,t} = \mu, \hat{\sigma}_{i,t}^2 = \sigma^2\right) \\ & = \mathbb{P}_t\left(-\theta_{i,t} + \mu_i + \frac{\gamma}{2}\left(\frac{1}{\kappa_{i,t}} - \sigma_i^2\right) \leq -\Delta_{\text{ER}}(i, \gamma) + \varepsilon \mid T_{i,t} = s, \hat{\mu}_{i,s} = \mu, \hat{\sigma}_{i,s} = \sigma\right) \\ & \leq \mathbb{P}_t\left(-\theta_{i,t} + \mu_i \leq -\xi(\Delta_{\text{ER}}(i, \gamma) - \varepsilon) \mid T_{i,t} = s, \hat{\mu}_{i,s} = \mu, \hat{\sigma}_{i,s} = \sigma\right) + \\ & \quad \mathbb{P}_t\left(\frac{\gamma}{2}\left(\frac{1}{\kappa_{i,t}} - \sigma_i^2\right) \leq -(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon) \mid T_{i,t} = s, \hat{\mu}_{i,s} = \mu, \hat{\sigma}_{i,s} = \sigma\right) \\ & = \mathbb{P}_t\left(\theta_{i,t} - \mu \geq (\mu_i - \mu) + \xi(\Delta_{\text{ER}}(i, \gamma) - \varepsilon) \mid T_{i,t} = s, \hat{\mu}_{i,s} = \mu, \hat{\sigma}_{i,s} = \sigma\right) + \\ & \quad \mathbb{P}_t\left(\kappa_{i,t} \geq \frac{\gamma}{\gamma\sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)} \mid T_{i,t} = s, \hat{\mu}_{i,s} = \mu, \hat{\sigma}_{i,s} = \sigma\right) \\ & \leq \exp\left(-\frac{s}{2}(\mu_i - \mu + \xi(\Delta_{\text{ER}}(i, \gamma) - \varepsilon))^2\right) + \exp\left(-sh\left(\frac{\gamma\sigma^2}{\gamma\sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)}\right)\right), \end{aligned}$$

where  $h(x) = \frac{1}{2}(x - 1 - \log x)$ .

The lemma holds by the Chernoff upper bound for  $\mathbb{P}_t(\theta_{i,t} \geq \cdot)$  and Lemma 10 below to upper-bound  $\mathbb{P}_t(\kappa_{i,t} \geq \cdot)$ .

**Lemma 10 (Harremoës (2016)).** For a Gamma r.v.  $X \sim \text{Gamma}(\alpha, \beta)$  with shape  $\alpha \geq 2$  and rate  $\beta > 0$ , we have

$$\mathbb{P}(X \geq x) \leq \exp\left(-2\alpha h\left(\frac{\beta x}{\alpha}\right)\right), \quad x > \frac{\alpha}{\beta},$$

where  $h(x) = \frac{1}{2}(x - 1 - \log x)$ . ■

**Lemma 11** (Upper bounding the second term of (1)). We have

$$\sum_{s=1}^n \mathbb{P}_t\left(G_{is} > \frac{1}{n}\right) \leq 1 + \max\left\{\frac{2 \log(2n)}{\xi^2 (\Delta_{\text{ER}}(i, \gamma) - \varepsilon)^2}, \frac{\log(2n)}{h\left(\frac{\gamma\sigma^2}{\gamma\sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)}\right)}\right\} + \frac{C_4}{\varepsilon^4} + \frac{C_5}{\varepsilon^2},$$

where  $C_4, C_5$  are constants.

**Proof** Following from Lemma 9, we have the following inclusions:

$$\begin{aligned} & \left\{ \hat{\mu}_{i,t} + \sqrt{\frac{2 \log 2n}{s}} \leq \mu_i + \xi(\Delta_{\text{ER}}(i, \gamma) - \varepsilon) \right\} \\ & \subseteq \left\{ \exp\left(-\frac{s}{2}(\mu_i - \mu + \xi(\Delta_{\text{ER}}(i, \gamma) - \varepsilon))^2\right) \leq \frac{1}{2n} \right\} \end{aligned}$$

and

$$\begin{aligned} & \left\{ \frac{\gamma \hat{\sigma}_{i,t}^2}{\gamma \sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)} \leq h_-^{-1}\left(\frac{\log 2n}{s}\right) \right\} \\ & \cup \left\{ \frac{\gamma \hat{\sigma}_{i,t}^2}{\gamma \sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)} \geq h_+^{-1}\left(\frac{\log 2n}{s}\right) \right\} \\ & \subseteq \left\{ \exp\left(-sh\left(\frac{\gamma \sigma^2}{\gamma \sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)}\right)\right) \leq \frac{1}{2n} \right\}, \end{aligned}$$

where  $h_+^{-1}(y) = \max\{x : h(x) = y\}$  and  $h_-^{-1}(y) = \min\{x : h(x) = y\}$ . Hence, for

$$s \geq u = \max \left\{ \frac{2 \log(2n)}{\xi^2 (\Delta_{\text{ER}}(i, \gamma) - \varepsilon)^2}, \frac{\log(2n)}{h\left(\frac{\gamma \sigma^2}{\gamma \sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)}\right)} \right\},$$

by replacing  $\left(\mu_1 - \varepsilon, \frac{\hat{\sigma}_i^2}{\sigma_1^2 + \varepsilon}\right)$  with  $\left(\mu_i + \xi(\Delta_{\text{ER}}(i, \gamma) - \varepsilon), \frac{\gamma \hat{\sigma}_{i,t}^2}{\gamma \sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)}\right)$  in [Zhu and Tan \(2020, S-3.4\)](#), we get

$$\mathbb{P}_t \left( G_{is} > \frac{1}{n} \right) \leq \exp\left(-\frac{s\varepsilon^2}{\sigma_i^2}\right) + \exp\left(-(s-1)\frac{\varepsilon^2}{\sigma_i^4}\right).$$

Summing over  $s$ ,

$$\begin{aligned} \sum_{s=1}^n \mathbb{P}_t \left( G_{is} > \frac{1}{n} \right) & \leq u + \sum_{s=\lceil u \rceil}^n \left[ \exp\left(-\frac{s\varepsilon^2}{\sigma_i^2}\right) + \exp\left(-(s-1)\frac{\varepsilon^2}{\sigma_i^4}\right) \right] \\ & \leq 1 + \max \left\{ \frac{2 \log(2n)}{\xi^2 (\Delta_{\text{ER}}(i, \gamma) - \varepsilon)^2}, \frac{\log(2n)}{h\left(\frac{\gamma \sigma^2}{\gamma \sigma_i^2 - 2(1-\xi)(\Delta_{\text{ER}}(i, \gamma) - \varepsilon)}\right)} \right\} + \frac{C_4}{\varepsilon^4} + \frac{C_5}{\varepsilon^2}. \end{aligned}$$

Finally, set

$$\xi_\gamma = 1 - \frac{\gamma \sigma_i^2}{2\Delta_{\text{ER}}(i, \gamma)} \left( 1 - \frac{1}{h_+^{-1}(\Delta_{\text{ER}}^2(i, \gamma)/2)} \right) \in (0, 1),$$

where  $h_+^{-1}(y) = \max \{x : h(x) = y\}$ . By algebra,

$$\begin{aligned}
 & h\left(\frac{\gamma\sigma_i^2}{\gamma\sigma_i^2 - 2(1 - \xi_\gamma)\Delta_{\text{ER}}(i, \gamma)}\right) \\
 &= h\left(\frac{\gamma\sigma_i^2}{\gamma\sigma_i^2 - 2 \cdot \frac{\gamma\sigma_i^2}{2\Delta_{\text{ER}}(i, \gamma)} \left(1 - \frac{1}{h_+^{-1}(\Delta_{\text{ER}}^2(i, \gamma)/2)}\right) \cdot \Delta_{\text{ER}}(i, \gamma)}\right) \\
 &= h\left(\frac{\gamma\sigma_i^2}{\gamma\sigma_i^2 - \gamma\sigma_i^2 \left(1 - \frac{1}{h_+^{-1}(\Delta_{\text{ER}}^2(i, \gamma)/2)}\right)}\right) = h\left(\frac{\gamma\sigma_i^2}{\left(\frac{\gamma\sigma_i^2}{h_+^{-1}(\Delta_{\text{ER}}^2(i, \gamma)/2)}\right)}\right) \\
 &= h\left(h_+^{-1}(\Delta_{\text{ER}}^2(i, \gamma)/2)\right) = \Delta_{\text{ER}}^2(i, \gamma)/2 \geq \xi_\gamma^2 \Delta_{\text{ER}}^2(i, \gamma)/2
 \end{aligned}$$

which implies

$$\frac{1}{h\left(\frac{\gamma\sigma_i^2}{\gamma\sigma_i^2 - 2(1 - \xi_\gamma)\Delta_{\text{ER}}(i, \gamma)}\right)} \leq \frac{2}{\xi_\gamma^2 \Delta_{\text{ER}}^2(i)}.$$

and  $\xi_\gamma \rightarrow 1^-$  as  $\gamma \rightarrow 0^+$ . ■

**Proof [Proof of Theorem 4]** We note that for any arm  $i$  with distribution  $\nu(i) \sim \mathcal{N}(\mu_i, \sigma_i^2)$  and  $\nu'(i) \sim \mathcal{N}(\mu'_i, (\sigma'_i)^2)$ , the KL-divergence given by

$$\text{KL}(\nu(i), \nu'(i)) = \log \frac{\sigma'_i}{\sigma_i} + \frac{\sigma_i^2 + (\mu_i - \mu'_i)^2}{2(\sigma'_i)^2} - \frac{1}{2}$$

is well-known. Denote  $\mathcal{S}_i = \{\nu'(i) \in \mathcal{E}_N^K : \text{ER}(\nu'(i)) < \text{ER}(1)\}$ . Denote

$$R_i := \max \left\{ \frac{2}{\xi^2 \Delta_{\text{ER}}^2(i)}, \frac{1}{h\left(\frac{\gamma\sigma^2}{\gamma\sigma_i^2 - 2(1 - \xi)\Delta_{\text{ER}}(i, \gamma)}\right)} \right\} > 0,$$

and fix  $\varepsilon > 0$  and consider the arm with the distribution  $\mathcal{N}(\mu_i + \sigma_i\sqrt{2/R_i} + \varepsilon, \sigma_i^2)$ . Then a direct computation gives

$$\begin{aligned}
 \text{ER}(\nu'(i)) - \text{ER}(\nu(1)) &= -(\mu_i + \sigma_i\sqrt{2/R_i} + \varepsilon) + \frac{\gamma}{2}\sigma_i^2 - \left(-\mu_i + \frac{\gamma}{2}\sigma_i^2\right) \\
 &= -(\sigma_i\sqrt{2/R_i} + \varepsilon) < 0,
 \end{aligned}$$

thus  $\text{ER}(\nu'(i)) < \text{ER}(\nu(1))$  and  $\nu'(i) \in \mathcal{S}_i$ . Furthermore,

$$\begin{aligned}
 \text{KL}(\nu(i), \nu'(i)) &= \log \frac{\sigma_i}{\sigma'_i} + \frac{\sigma_i^2 + \left(\mu_i - \left(\mu_i + \sigma_i\sqrt{2/R_i} + \varepsilon\right)\right)^2}{2\sigma_i'^2} - \frac{1}{2} \\
 &= \frac{1}{R_i} + \frac{(2\sigma_i\sqrt{2/R_i} + \varepsilon)\varepsilon}{2\sigma_i'^2}.
 \end{aligned}$$

By the definition of  $\eta$ ,

$$\eta(i, \gamma) \leq \lim_{\varepsilon \rightarrow 0^+} \left[ \frac{1}{R_i} + \frac{(2\sigma_i \sqrt{2/R_i} + \varepsilon)\varepsilon}{2\sigma_i^2} \right] = \frac{1}{R_i} \implies \frac{1}{\eta(i, \gamma)} \geq R_i.$$

Hence,

$$\liminf_{n \rightarrow \infty} \frac{\mathcal{R}_n(\pi)}{\log n} = \sum_{i \in [K] \setminus \{1\}} \left( \liminf_{n \rightarrow \infty} \frac{\mathbb{E}[T_{i,n}]}{\log n} \right) \Delta_{\text{ER}}(i, \gamma) \geq \sum_{i \in [K] \setminus \{1\}} R_i \Delta_{\text{ER}}(i, \gamma).$$

Thus, we have that ERTS is asymptotically optimal unconditionally. ■

## References

- Dorian Baudry, Romain Gautron, Emilie Kaufmann, and Odalric-Ambryn Maillard. Thompson sampling for CVaR bandits. *arXiv preprint arXiv:2012.05754*, 2020.
- Joel Q. L. Chang, Qiuyu Zhu, and Vincent Y. F. Tan. Risk-constrained thompson sampling for cvar bandits, 2021.
- Nicolas Galichet, Michele Sebag, and Olivier Teytaud. Exploration vs exploitation vs safety: Risk-aware multi-armed bandits. In *Asian Conference on Machine Learning*, pages 245–260, 2013.
- Peter Harremoës. Bounds on tail probabilities for negative binomial distributions. *Kybernetika*, 52(6):943–966, 2016.
- Ronald A. Howard and James E. Matheson. Risk-sensitive Markov decision processes. *Management Science*, 18(7):356–369, 1972.
- Anmol Kagrecha, Jayakrishnan Nair, and Krishna Jagannathan. Constrained regret minimization for multi-criterion multi-armed bandits. *arXiv preprint arXiv:2006.09649*, 2020.
- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Jaeho Lee, Sejun Park, and Jinwoo Shin. Learning bounds for risk-sensitive learning. In *Advances in Neural Information Processing Systems*, 2020.
- Amir Sani, Alessandro Lazaric, and Rémi Munos. Risk-aversion in multi-armed bandits. In *Advances in Neural Information Processing Systems*, pages 3275–3283, 2012.
- Wen Sun, Debadepta Dey, and Ashish Kapoor. Risk-aversion in multi-armed bandits. In *International Conference on Machine Learning*, pages 3280–3288, 2017.
- 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.
- Sattar Vakili and Qing Zhao. Risk-averse multi-armed bandit problems under mean-variance measure. *IEEE Journal of Selected Topics in Signal Processing*, 10(6):1093–1111, 2016.
- Qiuyu Zhu and Vincent YF Tan. Thompson sampling algorithms for mean-variance bandits. In *International Conference on Machine Learning*, pages 2645–2654, 2020.