

THE MULTI-ARMED BANDIT PROBLEM: AN EFFICIENT NON-PARAMETRIC SOLUTION

Hock Peng Chan

stachp@nus.edu.sg

Department of Statistics and Applied Probability

National University of Singapore

Abstract

Lai and Robbins (1985) and Lai (1987) provided efficient parametric solutions to the multi-armed bandit problem, showing that arm allocation via upper confidence bounds (UCB) achieves minimum regret. These bounds are constructed from the Kullback-Leibler information of the reward distributions, estimated from within a specified parametric family. In recent years there has been renewed interest in the multi-armed bandit problem due to new applications in machine learning algorithms and data analytics. Non-parametric arm allocation procedures like ϵ -greedy and Boltzmann exploration were studied, and modified versions of the UCB procedure were also analyzed under a non-parametric setting. However unlike UCB these non-parametric procedures are not efficient under parametric settings. In this paper we propose a subsample comparison procedure that is non-parametric, but still efficient under parametric settings.

1 Introduction

Lai and Robbins (1985) provided an asymptotic lower bound for the regret in the multi-armed bandit problem, and proposed a play-the-leader strategy that is *efficient*, that is it achieves this bound. Lai (1987) showed that allocation to the arm having the highest upper confidence bound (UCB), constructed from the Kullback-Leibler (KL) information between the estimated reward distributions of the arms, is efficient when the distributions belong to a specified exponential family. Agrawal (1995) modified UCB-Lai and showed that efficiency can still be achieved without having to know in advance the total sample size.

Burnetas and Kalehakis (1996) extended the UCB to multi-parameter families, almost showing efficiency in the natural setting of normal rewards with unequal variances. Yakowitz and Lowe (1991) proposed non-parametric procedures that do not make use of KL-information, suggesting logarithmic and polynomial rates of regret under finite exponential and moment conditions respectively.

Auer, Cesa-Bianchi and Fischer (2002) simplified UCB-Agrawal to UCB1, and showed that logarithmic regret is achieved when the reward distributions are supported on $[0,1]$. They also studied the ϵ -greedy algorithm of Sutton and Barto (1998), providing finite-time upper bounds of its regret. Both UCB1 and ϵ -greedy are non-parametric in their applications and, unlike UCB-Lai or -Agrawal, are not expected to be efficient under a general exponential family setting. Other non-parametric methods that have been proposed include reinforcement comparison, Boltzmann exploration (Sutton and Barto, 1998) and pursuit (Thathacher and Sastry, 1985). Kuleshov and Precup (2014) provided numerical comparisons between UCB and these methods. For a description of applications to recommender systems and clinical trials, see Shivaswamy and Joachims (2012). The reader is also strongly encouraged to go over Burtini, Loepky and Lawrence (2015) for a comprehensive survey of the methods, results and applications of the multi-armed bandit problem, developed over the past thirty years.

A strong competitor to UCB under the parametric setting is the use of the Bayesian method, see for example Fabius and van Zwet (1970), Berry (1972) and Kaufmann, Cappé and Garivier (2012). There is also a well-developed literature on optimization under an infinite-time discounted window setting, in which allocation is to the arm maximizing a dynamic allocation (or Gittins) index, see the seminal papers by Gittins (1979) and Gittins and Jones (1979), and also Berry and Fristedt (1985), Chang and Lai (1987), Brezzi and Lai (2002) and Kim and Lim (2016) for more recent advances. Another related problem is the study of the multi-armed bandit with irreversible constraints, initiated by Hu and Wei (1989).

In this paper we propose an arm allocation procedure that though non-parametric, is nevertheless efficient when the reward distributions are from an *unspecified* exponential family. It achieves this by comparing subsample means of the leading arm with the sample means of its competitors. It is empirical in its approach, using more informative subsample means rather than full-sample means alone, for better decision-making. An earlier version of the subsampling strategy, known as best empirical sampled average (BESA), appeared in Baransi, Maillard and Mannor (2014). However there are key differences in their implementation of subsampling from ours, as will be elaborated in Section 2.2.

The layout of the paper is as follows. In Section 2 we describe the subsample comparison strategy for allocating arms. In Section 3 we show that the strategy is efficient for exponential families, including the setting of normal rewards with unequal variances. To the best of our knowledge, this is the first instance that efficiency has been demonstrated under this two-

parameter setting. In Section 4 we show logarithmic regret under the more general setting of Markovian rewards. In Section 5 we provide numerical comparisons against existing methods. In Section 6 we prove the results of Sections 3 and 4.

2 Subsample comparisons

Let Y_{k1}, Y_{k2}, \dots , $1 \leq k \leq K$, be the observations (or rewards) from a statistical population Π_k . We assume here and in Section 3 that the rewards are independent and identically distributed (i.i.d.) within each arm. We extend to Markovian rewards in Section 4. Let $\mu_k = EY_{kt}$ and $\mu_* = \max_{1 \leq k \leq K} \mu_k$. Let $\lfloor \cdot \rfloor$ and $\lceil \cdot \rceil$ denote the greatest and least integer function respectively.

Consider a sequential procedure for selecting the population to be sampled at each time-stage. We refer to it as an arm allocation procedure in accordance to this being the multi-armed bandit problem. Let N_k be the number of observations from Π_k after N stages of sampling, hence $N = \sum_{k=1}^K N_k$. The objective is to minimize the *regret*

$$R_N := \sum_{k=1}^K (\mu_* - \mu_k) E N_k.$$

The Kullback-Leibler information number between two densities f and g , with respect to a common (σ -finite) measure, is

$$D(f|g) = E_f \left[\log \frac{f(Y)}{g(Y)} \right], \quad (2.1)$$

where E_f denotes expectation with respect to $Y \sim f$. An arm allocation procedure is said to converge uniformly fast if

$$R_N = o(N^\epsilon) \text{ for all } \epsilon > 0, \quad (2.2)$$

uniformly over all reward distributions lying within a specified parametric family.

Let f_k be the density of Π_k and let $f_* = f_k$ for k such that $\mu_k = \mu_*$ (assuming f_* is unique). The celebrated result of Lai and Robbins (1985) is that under (2.2) and additional regularity conditions,

$$\liminf_{N \rightarrow \infty} \frac{R_N}{\log N} \geq \sum_{k: \mu_k < \mu_*} \frac{\mu_* - \mu_k}{D(f_k|f_*)}. \quad (2.3)$$

Lai and Robbins (1985) and Lai (1987) went on to propose arm allocation procedures that have regrets achieving the lower bound, and are hence efficient.

2.1 Review of existing methods

In the setting of normal rewards with unit variances, UCB-Lai can be described as the selection of the population Π_k maximizing

$$\bar{Y}_{kn_k} + \sqrt{\frac{2 \log(N/n)}{n}}, \quad (2.4)$$

where $\bar{Y}_{kt} = \frac{1}{t} \sum_{u=1}^t Y_{ku}$, n is the current number of observations from the K populations, and n_k is the current number of observations from Π_k . Agrawal (1995) proposed a modified version of UCB-Lai that does not involve the total sample size N , with the selection instead of the population Π_k maximizing

$$\bar{Y}_{kn_k} + \sqrt{\frac{2(\log n + \log \log n + b_n)}{n_k}}, \quad (2.5)$$

with $b_n \rightarrow \infty$ and $b_n = o(\log n)$. Efficiency holds for (2.4) and (2.5), and there are corresponding versions of (2.4) and (2.5) that are efficient for other one-parameter exponential families.

Auer, Cesa-Bianchi and Fischer (2002) simplified UCB-Agrawal to UCB1, proposing instead that the population Π_k maximizing

$$\bar{Y}_{kn_k} + \sqrt{\frac{2 \log n}{n_k}} \quad (2.6)$$

be selected. They showed that under UCB1, $R_N = O(\log N)$ when the reward distributions are supported on $[0,1]$. In the setting of normal rewards with unequal (and unknown) variances, Auer et al. suggested applying a variant of UCB1 which they called UCB1-Normal, and showed that $R_N = O(\log N)$. Under UCB1-Normal, an observation is taken from any population Π_k with $n_k < \lceil 8 \log n \rceil$. If such a population does not exist, then an observation is taken from the population Π_k maximizing

$$\bar{Y}_{kn_k} + 4\hat{\sigma}_{kn_k} \sqrt{\frac{\log n}{n_k}},$$

where $\hat{\sigma}_{kt}^2 = \frac{1}{t-1} \sum_{u=1}^t (Y_{ku} - \bar{Y}_{kt})^2$.

Auer et al. provided an excellent study of various non-parametric arm allocation procedures, for example the ϵ -greedy procedure proposed by Sutton and Barto (1998), in which an observation is taken from the arm with the largest sample mean with probability $1 - \epsilon$, and randomly with probability ϵ . Auer et al. suggested replacing the fixed ϵ at every stage by a stage-dependent

$$\epsilon_n = \min(1, \frac{cK}{d^2 n}),$$

with c user-specified and $0 < d \leq \min_{k:\mu_k < \mu^*} (\mu^* - \mu_k)$. They showed that if $c > 5$, then logarithmic regret is achieved for reward distributions supported on $[0, 1]$. In practice d is unlikely to be known, in which case the user is effectively selecting a single tuning parameter $\frac{c}{d^2}$. A more recent numerical study by Kuleshov and Precup (2014) considered additional non-parametric procedures, for example Boltzmann exploration in which an observation is taken from a population Π_k with probability proportional to $e^{\bar{Y}_{kn_k}/\tau}$, for some $\tau > 0$.

2.2 Subsample-mean comparisons

A common characteristic of the procedures described in Section 2.1 is that a decision is made based solely on a comparison of the sample means \bar{Y}_{kn_k} , with the exception of UCB1-Normal in which $\hat{\sigma}_{kn_k}$ is also utilized. As we shall illustrate after describing the subsample-mean comparison procedure below, we can utilize subsample-mean information from the leading arm to estimate the same critical value for selecting from inferior arms as UCB-Agrawal and UCB1, and this leads to efficiency despite not specifying the underlying exponential family.

In subsample comparison, we apply the play-the-leader strategy (similar to that) of Lai and Robbins (1985). Let $\bar{Y}_{k,t:u} = \frac{1}{u-t+1} \sum_{v=t}^u Y_{kv}$ and $\bar{Y}_{kt} = \bar{Y}_{1,1:t}$. Let r denote the round number of the challenges.

Subsample-mean comparison

1. $r = 1$. Sample each Π_k exactly once.
2. $r = 2, 3, \dots$
 - (a) Let the leader $\zeta [= \zeta(n)]$ be the population with the most observations, with ties resolved by the larger sample mean \bar{Y}_{kn_k} , followed by randomization.
 - (b) Set up a challenge between Π_ζ and each Π_k for $k \neq \zeta$ in the following manner.
 - i. If $n_k = n_\zeta$, then Π_k loses the challenge automatically.
 - ii. If $n_k < \sqrt{\log n}$, then Π_k wins the challenge automatically.
 - iii. If $\sqrt{\log n} < n_k < n_\zeta$, then Π_k wins the challenge when

$$\bar{Y}_{kn_k} \geq \bar{Y}_{\zeta,t:(t+n_k-1)} \text{ for some } 1 \leq t \leq n_\zeta - n_k + 1. \quad (2.7)$$

- (c) For $k \neq \zeta$, sample from Π_k if Π_k wins its challenge against Π_ζ . Sample from Π_ζ if Π_ζ wins all its challenges.

Note that if Π_ζ wins all its challenges, then ζ and $(n_k : k \neq \zeta)$ are unchanged, and in the next round it suffices to perform the comparison in (2.7) at the largest t only instead of at every t . The computational burden is thus $O(1)$. The computational burden is $O(n)$ in the next round if at least one $k \neq \zeta$ wins its challenge in the current round. Hence when subsample-mean comparison achieves logarithmic regret, the total computational cost is only $O(N \log N)$.

To understand why subsample-mean comparison achieves efficiency, we consider the simple setting of unit variance normal with $K = 2$. Let $z(p)$ be such that

$$P(Z > z(p)) = p \text{ for } Z \sim N(0, 1).$$

Consider unbalanced sample sizes of say $n_2 = O(\log n)$. Since $z(p) \sim \sqrt{2|\log p|}$ for p small,

$$\begin{aligned} \min_{1 \leq t \leq n_1 - n_2 + 1} \bar{Y}_{1,t:(t+n_2-1)} &= \mu_1 - [1 + o_p(1)]z\left(\frac{1}{n_1 - n_2 + 1}\right)\sqrt{\frac{1}{n_2}} \\ &= \mu_1 - [1 + o_p(1)]\sqrt{\frac{2 \log n}{n_2}}. \end{aligned}$$

Hence arm 2 wins the challenge if

$$\bar{Y}_{2n_2} \geq \mu_1 - [1 + o_p(1)]\sqrt{\frac{2 \log n}{n_2}}. \quad (2.8)$$

By (2.5) and (2.6), UCB-Agrawal and UCB1 also select arm 2 when (2.8) holds, since $\bar{Y}_{1n_1} + \sqrt{\frac{2 \log n}{n_1}} = \mu_1 + o_p(1)$. Hence what subsample comparison does is to estimate the critical value of $\mu_1 - [1 + o_p(1)]\sqrt{\frac{2 \log n}{n_2}}$, empirically by using the minimum of the running averages $\bar{Y}_{1,t:(t+n_2-1)}$. The same critical value is similarly estimated by UCB1-Agrawal and subsample-mean for other exponential families. In the case of n_1, n_2 both large compared to $\log n$, $\sqrt{\frac{2 \log n}{n_2}} + \sqrt{\frac{2 \log n}{n_2}} \rightarrow 0$, and subsample-mean comparison, UCB-Agrawal and UCB1 essentially select the population with the larger sample mean.

Baransi, Maillard and Mannor (2014) proposed a subsampling strategy BESA which in the case $K = 2$, involves step 2(b)iii. alone with a single comparison between \bar{Y}_{kn_k} , and an empirical average of a random sample without replacement of size n_k from $\{Y_{\zeta t}\}_{1 \leq t \leq n_\zeta}$. They were able to show logarithmic regret of BESA for rewards supported on $[0,1]$. In contrast our subsample-mean strategy involves considerably more comparisons favoring the ‘‘inferior’’ arms. The additional experimentation is critical to the efficiency of subsample-mean comparison.

2.3 Subsample- t comparisons

For efficiency outside one-parameter exponential families, we need to work with test statistics beyond sample means. For example to achieve efficiency for normal reward distributions with unknown variances, the analogue of mean comparisons is t -statistic comparisons

$$\frac{\bar{Y}_{kn_k} - \mu_\zeta}{\hat{\sigma}_{kn_k}} \geq \frac{\bar{Y}_{\zeta,t:(t+n_k-1)} - \mu_\zeta}{\hat{\sigma}_{\zeta,t:(t+n_k-1)}},$$

where $\hat{\sigma}_{k,t;u}^2 = \frac{1}{u-t} \sum_{v=t}^u (Y_{kv} - \bar{Y}_{k,t;u})^2$ and $\hat{\sigma}_{kt} = \hat{\sigma}_{k,1:t}$. Since μ_ζ is unknown, we estimate it by $\bar{Y}_{\zeta n_\zeta}$.

Subsample- t comparison

Proceed as in subsample-mean comparison, with step 2(b)iii.' below replacing step 2(b)iii.

iii.' If $\sqrt{\log n} < n_k < n_\zeta$, then Π_k wins the challenge when either $\bar{Y}_{kn_k} \geq \bar{Y}_{\zeta n_\zeta}$ or

$$\frac{\bar{Y}_{kn_k} - \bar{Y}_{\zeta n_\zeta}}{\hat{\sigma}_{kn_k}} \geq \frac{\bar{Y}_{\zeta,t:(t+n_k-1)} - \bar{Y}_{\zeta n_\zeta}}{\hat{\sigma}_{\zeta,t:(t+n_k-1)}} \text{ for some } 1 \leq t \leq n_\zeta - n_k + 1. \quad (2.9)$$

Note that as in subsample-mean comparison, only $O(N \log N)$ computations are needed when the regret is $O(\log N)$. This is because it suffices to record the range of $\bar{Y}_{\zeta n_\zeta}$ that satisfies (2.9) for each $k \neq \zeta$, and the actual value of $\bar{Y}_{\zeta n_\zeta}$. The updating of these requires $O(1)$ computations when both ζ and $(n_k : k \neq \zeta)$ are unchanged.

3 Efficiency

Consider firstly an exponential family of density functions

$$f(x; \theta) = e^{\theta x - \psi(\theta)} f(x; 0), \quad \theta \in \Theta, \quad (3.1)$$

with respect to some measure ν , where $\psi(\theta) = \log[\int e^{\theta x} f(x; 0) \nu(dx)]$ is the log moment generating function and $\Theta = \{\theta : \psi(\theta) < \infty\}$. Let $f_k = f(\cdot; \theta_k)$ for some $\theta_k \in \Theta$, $1 \leq k \leq K$. Let $\theta_* = \max_{1 \leq k \leq K} \theta_k$ and $f_* = f(\cdot; \theta_*)$. Note that by (2.1) and (3.1), the KL-information in (2.3),

$$\begin{aligned} D(f_k | f_*) &= \int \{(\theta_k - \theta_*)x - [\psi(\theta_k) - \psi(\theta_*)]\} f(x; \theta_k) \nu(dx) \\ &= (\theta_k - \theta_*) \mu_k - [\psi(\theta_k) - \psi(\theta_*)] = I_*(\mu_k), \end{aligned}$$

where I_* is the large deviations rate function of f_* .

Theorem 1. Under (3.1), subsample-mean comparison has regret R_N satisfying

$$\limsup_{N \rightarrow \infty} \frac{R_N}{\log N} \leq \sum_{k: \mu_k < \mu_*} \frac{\mu_* - \mu_k}{D(f_k | f_*)},$$

and is thus efficient.

We next consider normal rewards with unequal variances, that is with densities

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad (3.2)$$

with respect to Lebesgue measure. Let $M(g) = \frac{1}{2} \log(1 + g^2)$. Burnetas and Katehakis (1996) showed that if $f_k = f(\cdot; \mu_k, \sigma_k^2)$, then under uniformly fast convergence and additional regularity conditions, an arm allocation procedure must have regret R_N satisfying

$$\liminf_{N \rightarrow \infty} \frac{R_N}{\log N} \geq \sum_{k: \mu_k < \mu_*} \frac{\mu_* - \mu_k}{M\left(\frac{\mu_* - \mu_k}{\sigma_k}\right)}.$$

They proposed an extension of UCB-Lai but needed the verification of a technical condition to show efficiency. In the case of UCB1-Normal, logarithmic regret also depended on tail bounds of the χ^2 - and t -distributions that were only shown to hold numerically by Auer et al. (2002). In Theorem 2 we show that subsample- t comparison attains the goal of efficiency.

Theorem 2. Under (3.2), subsample- t comparison has regret R_N satisfying

$$\limsup_{N \rightarrow \infty} \frac{R_N}{\log N} \leq \sum_{k: \mu_k < \mu_*} \frac{\mu_* - \mu_k}{M\left(\frac{\mu_* - \mu_k}{\sigma_k}\right)},$$

and is thus efficient

4 Logarithmic regret

We show here that logarithmic regret can be achieved by subsample-mean comparison under Markovian assumptions. This is possible because in subsample comparison we compare blocks of observations that retain the Markovian structure.

For $1 \leq k \leq K$, let X_{k1}, X_{k2}, \dots be a \mathcal{X} -valued Markov chain, with σ -field \mathcal{A} and transition kernel

$$P_k(x, A) = P(X_{kt} \in A | X_{k,t-1} = x), \quad x \in \mathcal{X}, A \in \mathcal{A}.$$

Let Y_{k1}, Y_{k2}, \dots be real-valued and conditionally independent given $(X_{kt})_{t \geq 1}$, which we shall assume for convenience to be stationary, and having conditional densities $\{f_k(\cdot|x) : 1 \leq k \leq K, x \in \mathcal{X}\}$, with respect to some measure ν , such that

$$P(Y_{kt} \in B | X_{k1} = x_1, X_{k2} = x_2, \dots) = \int_B f_k(y|x_t) \nu(dy).$$

We assume that the K Markov chains are independent of each other and that the following Doeblin-type condition holds.

(C1) For $1 \leq k \leq K$, there exists a non-trivial measure λ_k on $(\mathcal{X}, \mathcal{A})$ such that

$$P_k(x, A) \geq \lambda_k(A), \quad x \in \mathcal{X}, A \in \mathcal{A}.$$

As before let $\mu_k = EY_{kt}$, $\mu_* = \max_{1 \leq k \leq K} \mu_k$ and the regret

$$R_N = \sum_{k: \mu_k < \mu_*} (\mu_* - \mu_k) EN_k.$$

In addition to (C1) we assume the following sample mean large deviations.

(C2) For $1 \leq k \leq K$ and $\epsilon > 0$, there exists $\theta(= \theta_{k\epsilon}) > 0$ such that

$$P(|\bar{Y}_{kt} - \mu_k| \geq \epsilon) = O(e^{-t\theta}) \text{ as } t \rightarrow \infty.$$

(C3) For k such that $\mu_k < \mu_*$ and j such that $\mu_j = \mu_*$, there exists $b_\omega > 0$ for all $\omega < \mu_k$, such that

$$P(\bar{Y}_{jt} \leq \omega) = O(e^{-tb_\omega} P(\bar{Y}_{kt} \leq \omega)) \text{ as } t \rightarrow \infty.$$

Theorem 3. *Under (C1)–(C3), subsample-mean comparison has regret $R_N = O(\log N)$.*

EXAMPLE 1. Consider the setting of $Y_{kt} \stackrel{\text{i.i.d.}}{\sim} f_k$ within each arm k , with f_k positive on the real line. Let $F_k(x) = \int_{-\infty}^x f_k(y) dy$. We check that (C1) holds with $\lambda_k \equiv f_k$. If I_k , the large deviations rate function of f_k , is positive at $\omega \neq \mu_k$, then (C2) holds for $0 < \theta < \min(I_k(\mu_k - \epsilon), I_k(\mu_k + \epsilon))$. If in addition $F_k(x) \leq F_j(x)$ for all x and k, j such that $\mu_k < \mu_j = \mu_*$, then at $\omega < \mu_k$, $I_k(\omega) < I_j(\omega)$ and so (C3) holds for $0 < b_\omega < I_j(\omega) - I_k(\omega)$.

Procedure	c	Regret	
		$N = 1000$	$N = 10000$
Subsample-mean		11.2	19.2
UCB1		37.6	91.6
ϵ -greedy	0.1	35.3	297.5
	0.2	26.9	254.5
	0.5	21.3	178.0
	1	18.3	97.6
	2	17.2	65.9
	5	23.8	50.5
	10	36.9	64.1
	20	59.4	102.0

Table 1: Comparison of the regrets of subsample-mean, UCB1 and ϵ -greedy, for Bernoulli reward distributions.

5 Numerical studies

We compare here subsample-mean and $-t$ against the state-of-the-art procedures described in Section 2.1. In Examples 2–5 we consider exponential families and the comparisons are chiefly against procedures in which either efficiency or logarithmic regret has been established. In Example 6 we consider a non-exponential family and there the comparisons are against procedures that have been shown to perform well numerically. In each study, 1000 datasets are generated for each N , and the regret of a procedure is estimated by averaging over $\sum_{k:\mu_k < \mu_*} (\mu_* - \mu_k)N_k$. The same datasets are used for all the procedures compared within a study.

Two procedures are considered to be comparable if their regrets differ by no more than 10% (of the larger regret). A procedure is considered to be significantly better than another if its regret is less than two-thirds of the other one. The general conclusion from the numerical studies is that UCB1-tuned (to be introduced in Example 5) is the best performer at $N = 1000$, whereas subsample-mean and $-t$ are the best performers as $N = 10000$. A properly-tuned ϵ -greedy does well when the noise levels are high.

EXAMPLE 2. Consider Y_{kt} i.i.d. Bernoulli(μ_k) for $1 \leq k \leq 3$. We compare subsample-mean against UCB1 and ϵ -greedy. In each dataset we generate $\mu_k \sim \text{Uniform}(0,1)$. For ϵ -greedy we consider

$$\epsilon_n = \min(1, \frac{3c}{n}), \quad (5.1)$$

μ_2	μ_3	Regret	
		Subsample-mean	UCB-Lai
0.475	0.450	2.26	2.20
0.450	0.401	4.05	3.60
0.450	0.269	5.09	3.95
0.450	0.119	5.24	3.65
0.401	0.269	6.38	5.20
0.354	0.119	6.92	5.05
0.269	0.119	7.61	4.80
0.119	0.047	7.15	3.90
0.018	0.002	6.86	3.40
0.0003	0.0003	6.81	3.30

Table 2: Comparison of the regrets of subsample-mean and UCB-Lai for $K = 3$, $\mu_1 = 0.5$ and $N = 100$.

and experiment with various values of c .

Table 1 shows that ϵ -greedy, with $c = 2$ or 5 , outperforms UCB1. This is largely consistent with what has been reported in Auer, Cesa-Bianchi and Fischer (2002). However subsample-mean is significantly better than ϵ -greedy, uniformly over c .

EXAMPLE 3. We consider Bernoulli rewards here, but unlike in Example 2 we follow Lai (1987) with fixed values of $\mu_1 = 0.5$ and $\mu_3 \leq \mu_2 < 0.5$. Tables 2 and 3 summarize the comparison between subsample-mean and UCB-Lai. The regrets for UCB-Lai, taken from Tables 2 and 3 of Lai (1987), are smaller than those of subsample mean, significantly so in the settings of very small μ_2 and μ_3 . This should not be surprising given that UCB-Lai is specific to a given exponential family of reward distributions, whereas subsample-mean is non-parametric. In addition UCB-Lai is targeted towards a specific sample size N whereas subsample-mean operates independently of N . In Lai (1987) it is assumed that μ_k is known to lie between 0.01 and 0.99, though the last two lines of Table 2 indicate that UCB-Lai performs well even when this assumption is violated.

EXAMPLE 4. Consider $Y_{kt} \sim N(\mu_k, 1)$, $1 \leq k \leq 10$. In Table 4 we see that subsample-mean improves upon UCB1 at $N = 1000$, and outperforms UCB-Agrawal [setting $b_n = \log \log \log n$ in (2.5)] at both $N = 1000$ and 10000. Here we generate $\mu_k \sim N(0,1)$ in each dataset.

EXAMPLE 5. Consider $Y_{kt} \sim N(\mu_k, \sigma_k^2)$, $1 \leq k \leq 10$. We compare

μ_2	μ_3	Regret	
		Subsample-mean	UCB-Lai
0.495	0.490	11.19	10.75
0.490	0.480	19.65	17.50
0.490	0.450	23.86	19.75
0.490	0.401	23.14	16.25
0.480	0.450	30.32	23.75
0.470	0.401	32.50	25.50
0.450	0.401	33.71	25.75
0.401	0.354	27.20	16.75
0.310	0.231	20.10	11.00
0.168	0.168	15.60	8.50

Table 3: Comparison of the regrets of subsample-mean and UCB-Lai for $K = 3$, $\mu_1 = 0.5$ and $N = 2500$.

	Regret	
	$N = 1000$	$N = 10000$
Subsample-mean	89	137
UCB1	91	154
UCB-Agrawal	113	195

Table 4: The regrets of subsample-mean, UCB1 and UCB-Agrawal for $K = 10$ populations. The rewards have normal distributions with unit variances.

subsample- t against UCB1-tuned and UCB1-Normal. UCB1-tuned was suggested by Auer et al. and shown to perform well numerically. Under UCB1-tuned the population Π_k maximizing

$$\bar{Y}_{kn_k} + \sqrt{\frac{\log n}{n_k} \min(\frac{1}{4}, V_{kn})},$$

where $V_{kn} = \hat{\sigma}_{kn_k}^2 + \sqrt{\frac{2 \log n}{n_k}}$, is selected. In Table 5 we see that UCB-tuned is significantly better at $N = 1000$ whereas subsample- t is better at $N = 10000$. UCB1-Normal performs quite badly. Here we generate $\mu_k \sim N(0, 1)$ and $\sigma_k^{-2} \sim \text{Exp}(1)$ in each dataset.

EXAMPLE 6. Consider double exponential rewards $Y_{kt} \sim f_k$, with densities

$$f_k(x) = \frac{1}{2\lambda} e^{-|x-\mu_k|/\lambda}, \quad 1 \leq k \leq 10,$$

	Regret	
	$N = 1000$	$N = 10000$
Subsample- t	251	570
UCB-tuned	126	761
UCB1-Normal	1549	5091

Table 5: The regrets of subsample- t , UCB-tuned and UCB1-Normal, for $K = 10$ populations. The rewards have normal distributions with unequal variances.

with respect to Lebesgue measure. We compare subsample-mean against UCB1-tuned, Boltzmann exploration and ϵ -greedy [see (5.1)]. We generate $\mu_k \sim N(0,1)$ in each dataset. Table 6 shows that UCB1-tuned has the best performances at $N = 1000$, whereas subsample-mean has the best performances at $N = 10000$. A properly-tuned ϵ -greedy or Boltzmann exploration does well at $N = 1000$ and $\lambda = 2$, and a properly-tuned ϵ -greedy also does well at $\lambda = 5$ for both $N = 1000$ and 10000.

6 Proofs of Theorems 1–3

Since subsample comparison is index-blind, we may assume without loss of generality that $\mu_1 = \mu_*$. We provide here the statements and proofs of supporting Lemmas 1 and 2, and follow up with the proofs of Theorems 1–3 in Sections 6.1–6.3.

Recall that the leading arm refers to the population that has been sampled the most times. We label an arm k as optimal if $\mu_k = \mu_*$, otherwise we label it as inferior. Let A_n be the event that at stage n the leading arm is inferior.

Let B_m be the event that at stage m the leading arm is optimal, and it loses a challenge to an inferior arm that has been sampled at least $\frac{m}{\log m}$ times. Let C_m be the event that at stage m the leading arm is inferior, and it wins a challenge against an optimal arm. In Lemma 1 we show how bounds on $P(B_m)$ and $P(C_m)$ for $1 \leq m \leq n$ lead to an important bound on $P(A_n)$. Let $n_k(m_k)$ be the number of times Π_k has been sampled at the start of stage $n(m)$. Hence $\sum_{k=1}^K n_k = n - 1$ and $\sum_{k=1}^K m_k = m - 1$.

Lemma 1. *If there exists $a > 0$ such that $P(B_m) = O(e^{-\frac{am}{\log m}})$ and $P(C_m) =$*

Method	c/τ	Regret					
		$N = 1000$			$N = 10000$		
		$\lambda = 1$	$\lambda = 2$	$\lambda = 5$	$\lambda = 1$	$\lambda = 2$	$\lambda = 5$
Subsample-mean		144	333	799	239	653	2394
UCB1-tuned		100	261	647	479	1923	4879
Boltzmann	0.01	174	334	722	1460	3195	6173
	0.02	166	330	720	1401	3132	6131
	0.05	147	316	716	1145	2919	5914
	0.1	122	294	700	802	2342	5455
	0.2	125	271	662	802	1705	4433
	0.5	333	393	657	3106	3130	4255
	1	733	737	826	7276	7265	7082
ϵ -greedy	0.1	164	322	696	1242	2799	5747
	0.2	161	318	684	1103	2524	5344
	0.5	152	278	638	775	1842	4487
	1	157	286	600	597	1357	3691
	2	197	285	588	495	1077	3000
	5	332	394	622	649	994	2465
	10	521	567	738	959	1189	2316
	20	805	844	951	1559	1749	2651

Table 6: Regret comparisons for double exponential density rewards.

$O(\frac{1}{m(\log m)^2})$, uniformly over $1 \leq m \leq n$, then

$$P(A_n) = o(\frac{1}{n}). \quad (6.1)$$

PROOF. Let $m_0 = \lfloor (\log \log n)^2 \rfloor$. If at some stage $m \in [m_0, n - 1]$ the leading arm is optimal, then the probability that there will be at least one win by an inferior arm with at least $\frac{m}{\log m}$ observations, between stages m and $n - 1$, is no more than

$$\sum_{q=m}^{n-1} P(B_q) = O(ne^{-\frac{am_0}{\log m_0}}) = o(\frac{1}{n}).$$

It remains for us to show that the probability the leader is inferior, at all stages between m_0 and n , is $o(\frac{1}{n})$. Since

$$\sum_{m=m_0}^{n-1} P(C_m) = O(\frac{1}{\log n}),$$

the probability that an inferior leading arm wins at least $\frac{n}{\sqrt{\log n}}$ times between stages m_0 and $n - 1$, against optimal arms, is $\frac{\sqrt{\log n}}{n} O(\frac{1}{\log n}) = o(\frac{1}{n})$. But it is not possible to have inferior leading arms at all stages between m to n , with them winning no more than $\frac{n}{\sqrt{\log n}}$ times against optimal arms between stages m and $n - 1$. This is because $\max_{k:\mu_k < \mu_*} m_k - \max_{k:\mu_k = \mu_*} m_k$ reduces by 1 after each round in which the leading inferior arm loses to all optimal arms. Note in particular that by step 2(b)i., the leading inferior arm wins against all inferior arms with the same number of observations. With that we conclude (6.1). \square

Let B_{nk} be the event that at stage n the leading arm is optimal, and it loses a challenge to an inferior arm k .

Lemma 2. *Consider the following conditions for an inferior arm k .*

(I) *There exists $\xi_k > 0$ such that for all $\epsilon > 0$, as $N \rightarrow \infty$,*

$$P(B_{nk} \text{ occurs for some } 1 \leq n \leq N \text{ with } n_k = \lfloor (1 + \epsilon)\xi_k \log N \rfloor) \rightarrow 0.$$

(II) *There exists $J_k > 0$ such that as $N \rightarrow \infty$,*

$$P(B_{nk} \text{ occurs for some } 1 \leq n \leq N \text{ with } n_k = \lfloor J_k \log N \rfloor) = O(N^{-1}).$$

If (6.1), (I) and (II) hold, then

$$\limsup_{N \rightarrow \infty} \frac{EN_k}{\log N} \leq \xi_k. \quad (6.2)$$

PROOF. By (6.1), $\sum_{n=1}^N P(A_n) = o(\log N)$, and (6.2) thus follows from (I) and (II). \square

6.1 Proof of Theorem 1

We consider here subsample-mean comparisons. Let I_j be the large deviations rate function of f_j . By Lemmas 1 and 2 it suffices, in Lemmas 3–5 below, to verify the conditions needed to show that (6.2) holds with $\xi_k = 1/I_1(\mu_k)$.

Lemma 3. *Under (3.1), $P(B_m) = O(e^{-\frac{am}{\log m}})$ for some $a > 0$.*

PROOF. Consider the case that at stage m the leading optimal arm is arm 1, and it loses a challenge to an inferior arm k with $m_k \geq \frac{m}{\log m}$. Let $\max_{j:\mu_j < \mu_1} \mu_j < \omega < \mu_1$. It follows from large deviations that

$$P(\bar{Y}_{1,t:(t+m_k-1)} \leq \omega \text{ for some } 1 \leq t \leq m) \leq m e^{-m_k I_1(\omega)},$$

$$P(\bar{Y}_{km_k} \geq \omega) \leq e^{-m_k I_k(\omega)}.$$

Since $m_k \geq \frac{m}{\log m}$, the above inequalities imply that Lemma 3 holds for $0 < a < \min_{1 \leq k \leq K} I_k(\omega)$. \square

Lemma 4. Under (3.1), $P(C_m) = O(\frac{1}{m(\log m)^2})$.

PROOF. Consider the case that k is the leading inferior arm, and it wins a challenge against optimal arm 1 at stage m . By step 2(b)ii., arm k loses automatically when $m_1 < \sqrt{\log m}$, hence we need only consider $m_1 > \sqrt{\log m}$.

Case 1: $m_1 > (\log m)^2$. Let $\mu_k < \omega < \mu_1$. By large deviations,

$$P(\bar{Y}_{1m_1} < \omega \text{ for some } m_1 > (\log m)^2) = O(\frac{1}{m(\log m)^2}), \quad (6.3)$$

$$P(\bar{Y}_{km_1} > \omega \text{ for some } m_1 > (\log m)^2) = O(\frac{1}{m(\log m)^2}). \quad (6.4)$$

Case 2: $\sqrt{\log m} < m_1 < (\log m)^2$. Since $m_k \geq \frac{m-1}{K}$, it suffices to show that there exists ω such that

$$P(\bar{Y}_{1m_1} < \omega) = O(\frac{1}{m(\log m)^4}), \quad (6.5)$$

$$\begin{aligned} P(\bar{Y}_{k,t:(t+m_1-1)} > \omega \text{ for } 1 \leq t \leq \frac{m-1}{K} - m_1 + 1) \\ (\leq [P(\bar{Y}_{km_1} > \omega)]^{[(m-1)/K - m_1 + 1]/m_1}) = O(\frac{1}{m(\log m)^4}). \end{aligned} \quad (6.6)$$

Since $\theta_1 > \theta_k$, if $\sum_{t=1}^{m_1} y_t \leq m_1 \mu_k$, then by (3.1),

$$\begin{aligned} \prod_{t=1}^{m_1} f(y_t; \theta_1) &= e^{(\theta_1 - \theta_k) \sum_{t=1}^{m_1} y_t - m_1 [\psi(\theta_1) - \psi(\theta_k)]} \prod_{t=1}^{m_1} f(y_t; \theta_k) \\ &\leq e^{-m_1 I_1(\mu_k)} \prod_{t=1}^{m_1} f(y_t; \theta_k). \end{aligned}$$

Hence if $\omega \leq \mu_k$, then as $m_1 > \sqrt{\log m}$,

$$P(\bar{Y}_{1m_1} < \omega) \leq e^{-m_1 I_1(\mu_k)} P(\bar{Y}_{km_1} < \omega) = O(\frac{1}{(\log m)^8} P(\bar{Y}_{km_1} < \omega)). \quad (6.7)$$

Let $\omega (\leq \mu_k$ for large m) be such that

$$P(\bar{Y}_{km_1} < \omega) \leq \frac{(\log m)^4}{m} \leq P(\bar{Y}_{km_1} \leq \omega). \quad (6.8)$$

We check that (6.5) follows from (6.7) and the first inequality in (6.8), whereas (6.6) follows from the second inequality in (6.8) and $m_1 \leq (\log m)^2$. Lemma 3 follows from (6.3)–(6.6). \square

Lemma 5. Under (3.1), both (I) and (II), in the statement of Lemma 2, hold for $\xi_k = 1/I_1(\mu_k)$.

PROOF. Consider an inferior arm k and stage n with $n_k = \lfloor [(1 + \epsilon) \log N]/I_1(\mu_k) \rfloor$. Let $\mu_k < \omega < \mu_1$ be such that $(1 + \epsilon)I_1(\omega) > I_1(\mu_k)$. It follows from large deviations that

$$\begin{aligned} P(\bar{Y}_{1,t:(t+n_k-1)} \leq \omega \text{ for some } 1 \leq t \leq N) &\leq N e^{-n_k I_1(\omega)} \rightarrow 0, \\ P(\bar{Y}_{kn_k} \geq \omega) &\leq e^{-n_k I_k(\omega)} \rightarrow 0, \end{aligned}$$

and (I) therefore holds.

Next consider $J_k > \max(\frac{1}{I_k(\omega)}, \frac{2}{I_1(\omega)})$. If $n_k = \lfloor J_k \log N \rfloor$, then

$$\begin{aligned} P(\bar{Y}_{1,t:(t+n_k-1)} \leq \omega \text{ for some } 1 \leq t \leq N) &\leq N e^{-n_k I_1(\omega)} = O(N^{-1}), \\ P(\bar{Y}_{kn_k} \geq \omega) &\leq e^{-n_k I_k(\omega)} = O(N^{-1}), \end{aligned}$$

and (II) holds as well. \square

6.2 Proof of Theorem 2

We consider here subsample- t comparisons. By Lemmas 1 and 2 it suffices, in Lemmas 6–8 below, to verify the conditions needed to show that (6.2) holds with $\xi_k = 1/M(\frac{\mu_k - \mu_1}{\sigma_k})$. Let $\bar{\Phi}(z) = P(Z > z)$ for $Z \sim N(0,1)$.

Lemma 6. Under (3.2), $P(B_m) = O(e^{\frac{am}{\log m}})$ for some $a > 0$.

PROOF. Consider the case that at stage m the leading optimal arm is arm 1, and that it loses a challenge to an inferior arm k with $m_k \geq \frac{m}{\log m}$. Let $\epsilon > 0$ be such that $\omega := \frac{\mu_k - \mu_1 + \epsilon}{2\sigma_k} < 0$, noting that

$$P\left(\frac{\bar{Y}_{km_k} - \bar{Y}_{1m_1}}{\hat{\sigma}_{km_k}} \geq \omega\right) \leq P\left(\frac{\bar{Y}_{km_k} - \bar{Y}_{1m_1}}{2\sigma_k} \geq \omega\right) + P(\hat{\sigma}_{km_k} \geq 2\sigma_k). \quad (6.9)$$

Since $\bar{Y}_{km_k} - \bar{Y}_{1m_1} \sim N(\mu_k - \mu_1, \frac{\sigma_1^2}{m_1} + \frac{\sigma_k^2}{m_k})$ and $m_1 \geq m_k$,

$$\begin{aligned} P\left(\frac{\bar{Y}_{km_k} - \bar{Y}_{1m_k}}{2\sigma_k} \geq \omega\right) &\leq \bar{\Phi}\left(\epsilon \sqrt{\frac{m_k}{\sigma_1^2 + \sigma_k^2}}\right) \\ &= O(m^{-2} e^{-\frac{am}{\log m}}) \text{ for } 0 < a < \frac{\epsilon^2}{2(\sigma_1^2 + \sigma_k^2)}. \end{aligned} \quad (6.10)$$

Since $m_k \geq \frac{m}{\log m}$, by large deviations,

$$P(\hat{\sigma}_{km_k} > 2\sigma_k) = O(m^{-2} e^{-\frac{am}{\log m}}) \text{ for some } a > 0. \quad (6.11)$$

It follows from arguments similar to those in (6.10) and (6.11) that

$$\begin{aligned}
& P\left(\frac{\bar{Y}_{1,t:(t+m_k-1)} - \bar{Y}_{1m_1}}{\hat{\sigma}_{1,t:(t+m_k-1)}} \leq \omega \text{ for some } 1 \leq t \leq m\right) \quad (6.12) \\
& \leq m \left[P\left(\frac{\bar{Y}_{1m_k} - \bar{Y}_{1m_1}}{\sigma_{1/2}} \leq \omega\right) + P(\hat{\sigma}_{1m_k} \leq \frac{\sigma_1}{2}) \right] \\
& = O(m^{-2} e^{-\frac{am}{\log m}}) \text{ for some } a > 0.
\end{aligned}$$

It also follows from large deviations arguments that

$$P(\bar{Y}_{km_k} \geq \bar{Y}_{1m_1}) = O(m^{-2} e^{-\frac{am}{\log m}}) \text{ for some } a > 0,$$

and Lemma 6 thus follows from (6.9)–(6.12). \square

Lemma 7. *Under (3.2), $P(C_m) = O(\frac{1}{m(\log m)^2})$.*

PROOF. Consider the case that k is the leading inferior arm, and it wins against optimal arm 1 at stage m . By step 2(b)ii. of subsample- t comparisons, we need only consider $m_1 > \sqrt{\log m}$. Note that $m_k \geq \frac{m-1}{K}$.

Case 1. $m_1 > (\log m)^2$. Let $\omega = \frac{\mu_1 + \mu_k}{2}$ and check that

$$\begin{aligned}
& P(\bar{Y}_{1m_1} \leq \omega) + P(\bar{Y}_{km_k} \geq \omega) \quad (6.13) \\
& \leq e^{-m_1(\mu_1 - \mu_k)^2 / (8\sigma_1^2)} + e^{-m_k(\mu_1 - \mu_k)^2 / (8\sigma_k^2)} = O(m^{-3}).
\end{aligned}$$

Case 2. $\sqrt{\log m} < m_1 < (\log m)^2$. Let us condition on $\bar{Y}_{km_k} = \mu_k - \gamma$. Since

$$P(|\bar{Y}_{km_k} - \mu_k| \geq m^{-\frac{1}{3}}) \leq e^{-\frac{m_k}{2m^{2/3}}} = O(m^{-3}),$$

it suffices to consider $|\gamma| \leq m^{-\frac{1}{3}}$. Let $\omega (\leq 0 \text{ for } m \text{ large})$ be such that

$$P\left(\frac{\bar{Y}_{km_1} - \mu_k + \gamma}{\hat{\sigma}_{km_1}} \leq \omega\right) = \frac{(\log m)^4}{m}. \quad (6.14)$$

It follows from (6.14) that

$$\begin{aligned}
& P\left(\frac{\bar{Y}_{k,t:(t+m_1-1)} - \mu_k + \gamma}{\hat{\sigma}_{k,t:(t+m_1-1)}} > \omega \text{ for } 1 \leq t \leq \frac{m-1}{K} - m_1 + 1\right) \quad (6.15) \\
& \leq \left[P\left(\frac{\bar{Y}_{km_1} - \mu_k + \gamma}{\hat{\sigma}_{km_1}} > \omega\right) \right]^{[(m-1)/K - m_1 + 1]/m_1} = O\left(\frac{1}{m(\log m)^4}\right).
\end{aligned}$$

Lemma 7 then follows from (6.13) and showing that

$$P\left(\frac{\bar{Y}_{1m_1} - \mu_k + \gamma}{\hat{\sigma}_{1m_1}} \leq \omega\right) = O\left(\frac{1}{m(\log m)^4}\right). \quad (6.16)$$

To show (6.16) from (6.14), we note that conditioned on $\widehat{\sigma}_{km_1}^2 = \sigma_k^2/\tau$ for some $\tau > 0$,

$$\frac{\bar{Y}_{km_1 - \mu_k + \gamma}}{\widehat{\sigma}_{km_1}} \sim N\left(\frac{\gamma\sqrt{\tau}}{\sigma_k}, \frac{\tau}{m_1}\right), \quad (6.17)$$

whereas conditioned on $\widehat{\sigma}_{1m_1}^2 = \sigma_1^2/\tau$,

$$\frac{\bar{Y}_{1m_1 - \mu_k + \gamma}}{\widehat{\sigma}_{1m_1}} \sim N\left(\frac{(\mu_1 - \mu_k + \gamma)\sqrt{\tau}}{\sigma_1}, \frac{\tau}{m_1}\right). \quad (6.18)$$

By a change-of-measure argument on $Z \sim N(\beta, \frac{\tau}{m_1})$ for $\beta_1 \geq \beta_2$ (and since $\omega \leq 0$),

$$\begin{aligned} P_{\beta_1}(Z \leq \omega) &\leq e^{m_1[(\omega - \beta_1)^2 - (\omega - \beta_2)^2]/(2\tau)} P_{\beta_2}(Z \leq \omega) \\ &\leq e^{m_1(\beta_2^2 - \beta_1^2)/(2\tau)} P_{\beta_2}(Z \leq \omega). \end{aligned}$$

In view of (6.17) and (6.18), we consider the above inequalities with $\beta_1 = c + \frac{\gamma\sqrt{\tau}}{\sigma_1}$ (where $c = \frac{\mu_1 - \mu_k}{\sigma_1}$) and $\beta_2 = \frac{\gamma\sqrt{\tau}}{\sigma_k}$, which leads us to

$$\begin{aligned} P_{\beta_1}(Z \leq \omega) &\leq [1 + o(1)]e^{-m_1 c^2/(2\tau)} P_{\beta_2}(Z \leq \omega) \\ &= O\left(\frac{1}{(\log m)^8} P_{\beta_2}(Z \leq \omega)\right), \end{aligned} \quad (6.19)$$

uniformly over $0 < \tau < \frac{m_1}{(\log \log m)^2} (:= \tau_m)$ and $|\gamma| \leq m^{-\frac{1}{3}}$. Since $\widehat{\sigma}_{1m_1}^2/\sigma_1^2 \stackrel{d}{=} \widehat{\sigma}_{km_1}^2/\sigma_k^2 \stackrel{d}{=} \sum_{t=1}^{m_1-1} Z_t^2$ with $Z_t^2 \stackrel{\text{i.i.d.}}{\sim} N(0,1)$, it follows from a change-of-measure of the distribution of Z_t^2 to $N(0, \tau_m^{-1})$ that

$$P(\widehat{\sigma}_{1m_1}^2/\sigma_1^2 \leq \tau_m^{-1}) \leq e^{-\frac{m_1-1}{2\tau_m}} / (\tau_m^{\frac{m_1-1}{2}} e^{-\frac{m_1-1}{2}}) = O\left(\frac{1}{m(\log m)^4}\right).$$

Therefore by (6.14), (6.19) and the independence between $\widehat{\sigma}_{km_1}$ and $\widehat{\sigma}_{1m_1}$ with \bar{Y}_{km_1} and \bar{Y}_{1m_1} ,

$$P\left(\frac{\bar{Y}_{1m_1 - \mu_k + \gamma}}{\widehat{\sigma}_{1m_1}} \leq \omega\right) = O\left(\frac{1}{(\log m)^8} \frac{(\log m)^4}{m}\right),$$

and (6.16) indeed holds. \square

Lemma 8. *Under (3.2), both (I) and (II), in the statement of Lemma 2, hold for $\xi_k = 1/M(\frac{\mu_* - \mu_k}{\sigma_k})$.*

PROOF. By considering the rewards $\frac{Y_{kt} - \mu_1}{\sigma_1}$, we may assume without loss of generality that $(\mu_* =) \mu_1 = 0$ and $\sigma_1^2 = 1$. Consider an inferior arm k and

stage n with $n_k = \lfloor [(1 + \epsilon) \log N] / M(g_k) \rfloor$ for some $\epsilon > 0$. Let $g_k = \frac{\mu_k}{\sigma_k}$ and let g_ω be such that

$$0 > g_\omega > g_k \text{ and } (1 + \epsilon)M(g_\omega) > M(g_k). \quad (6.20)$$

Select $\tau > 0$ small enough such that

$$P(\widehat{\sigma}_{1,t:(t+n_k-1)} \leq \tau \text{ for some } 1 \leq t \leq N) \rightarrow 0, \quad (6.21)$$

and let $\beta > 1$ and $\gamma_0 > 0$ be such that

$$\frac{\mu_k / \sqrt{\beta} - \gamma}{\sigma_k \sqrt{\beta}} \leq g_\omega - \frac{\gamma}{\tau} \text{ for } |\gamma| \leq \gamma_0. \quad (6.22)$$

By large deviations, there exists $a_k > 0$ such that,

$$P(\bar{Y}_{kn_k} \geq \frac{\mu_k}{\sqrt{\beta}}) + P(\widehat{\sigma}_{kn_k} \geq \sigma_k \sqrt{\beta}) \leq 2e^{-n_k a_k} \rightarrow 0, \quad (6.23)$$

$$P(|\bar{Y}_{1n_1}| \geq \gamma_0 \text{ for some } n_k \leq n_1 \leq N) \rightarrow 0. \quad (6.24)$$

Moreover

$$P\left(\frac{\bar{Y}_{1,t:(t+n_k-1)}}{\widehat{\sigma}_{1,t:(t+n_k-1)}} \leq g_\omega \text{ for some } 1 \leq t \leq N\right) \leq NP\left(\frac{\bar{Y}_{1n_k}}{\widehat{\sigma}_{1n_k}} \leq g_\omega\right) \rightarrow 0, \quad (6.25)$$

since by large deviations and the independence of \bar{Y}_{1n_k} and $\widehat{\sigma}_{1n_k}$,

$$\begin{aligned} & n_k^{-1} |\log P\left(\frac{\bar{Y}_{1n_k}}{\widehat{\sigma}_{1n_k}} \leq g_\omega\right)| \\ &= \inf_{\sigma > 0} \{n_k^{-1} [|\log P(\bar{Y}_{1n_k} \leq g_\omega \sigma)| + |\log P(\widehat{\sigma}_{1n_k} \leq \sigma)|]\} + o(1) \\ &\rightarrow \inf_{\sigma > 0} \left[\frac{(g_\omega \sigma)^2}{2} + \frac{1}{2}(\sigma^2 - 1 - \log \sigma^2)\right] = M(g_\omega), \end{aligned}$$

and $n_k M(g_\omega) > \log N$ (note infimum above attained at $\sigma^2 = \frac{1}{g_\omega^2 + 1}$).

By (6.21)–(6.25),

$$P\left(\frac{\bar{Y}_{kn_k} - \bar{Y}_{1n_1}}{\widehat{\sigma}_{kn_k}} \geq \frac{\bar{Y}_{1,t:(t+n_k-1)} - \bar{Y}_{1n_1}}{\widehat{\sigma}_{1,t:(t+n_k-1)}} \text{ for some } 1 \leq t \leq N\right) \rightarrow 0,$$

and this, with

$$P(\bar{Y}_{kn_k} \geq \bar{Y}_{1n_1} \text{ for some } n_k \leq n_1 \leq N) \rightarrow 0, \quad (6.26)$$

shows (I).

To show (II), we check that for $n_k = \lfloor J_k \log N \rfloor$ with J_k large enough, the relations in (6.21) and (6.23)–(6.26) hold with “ $= O(N^{-1})$ ” replacing “ $\rightarrow 0$ ”. \square

6.3 Proof of Theorem 3

We consider here subsample-mean comparisons. Assume (C1)–(C3) and let $\tilde{\mu} = \max_{k: \mu_k < \mu^*} \mu_k$. By Lemmas 1 and 2 it suffices, in Lemmas 9–11 below, to verify the conditions needed to show that (6.2) holds for some $\xi_k > 0$.

Lemma 9. *Under (C2), $P(B_m) = O(e^{-\frac{am}{\log m}})$ for some $a > 0$.*

PROOF. Let $\epsilon = \frac{1}{2}(\mu_* - \tilde{\mu})$. It follows from (C2) and arguments in the proof of Lemma 3, setting $\omega = \frac{1}{2}(\mu_* + \tilde{\mu})$, that $P(B_m) = O(m^2 e^{-\frac{\theta m}{\log m}})$ for some $\theta > 0$. Hence Lemma 9 holds for $0 < a < \theta$. \square

Lemma 10. *Under (C1)–(C3), $P(C_m) = O(\frac{1}{m(\log m)^2})$.*

PROOF. Consider the case that k is the leading inferior arm, and it wins against optimal arm 1 at stage m . By step 2(b)ii. of subsample-mean comparison, we need only consider $m_1 > \sqrt{\log m}$.

Case 1: $m_1 > (\log m)^2$. Let ω and ϵ be as in the proof of Lemma 9, and check that by (C2), there exists $\theta > 0$ such that

$$P(\bar{Y}_{1m_1} \leq \omega) + P(\bar{Y}_{km_k} \geq \omega) = O(e^{-m_1\theta}) = O(m^{-3}). \quad (6.27)$$

Case 2: $\sqrt{\log m} < m_1 < (\log m)^2$. Select $\omega (\leq \mu_k$ for m large) such that

$$P(\bar{Y}_{km_1} < \omega) \leq \frac{(\log m)^4}{m} \leq P(\bar{Y}_{km_1} \leq \omega). \quad (6.28)$$

Let $p_\omega = P(\bar{Y}_{km_1} > \omega)$ and let $\kappa = \lfloor \frac{(m-1)/K - m_1}{2m_1} \rfloor$. By (C1) and the second inequality of (6.28),

$$\begin{aligned} & P(\bar{Y}_{k,t:(t+m_1-1)} > \omega \text{ for } 1 \leq t \leq \frac{m-1}{K} - m_1 + 1) \quad (6.29) \\ & \leq P(\bar{Y}_{k,t:(t+m_1-1)} > \omega \text{ for } t = 1, 2m_1 + 1, \dots, 2\kappa m_1 + 1) \\ & \leq p_\omega^{\kappa+1} + \kappa[1 - \lambda_k(\mathbf{R})]^{m_1} = O(\frac{1}{m(\log m)^4}). \end{aligned}$$

It follows from (C3) and the first inequality of (6.28) that

$$P(\bar{Y}_{1m_1} < \omega) = O(\frac{1}{m(\log m)^4}),$$

and Lemma 10 thus follows from (6.27) and (6.29). \square

Lemma 11. *Under (C2), statement (II) in Lemma 2 holds.*

PROOF. Let ϵ and ω be as in the proof of Lemma 9, and let $J_k > \frac{2}{\theta}$, where θ is given in (C2). It follows that for $n = \lfloor J_k \log N \rfloor$,

$$\begin{aligned} P(\bar{Y}_{1,t:(t+n_k-1)} \leq \omega \text{ for some } 1 \leq t \leq N) &= O(Ne^{-n_k\theta}) = O(N^{-1}), \\ P(\bar{Y}_{kn_k} \geq \omega) &= O(e^{-n_k\theta}) = O(N^{-1}), \end{aligned}$$

and (II) indeed holds. \square

References

- [1] AGRAWAL, R. (1995). Sample mean based index policies with $O(\log n)$ regret for the multi-armed bandit problem. *Adv. Appl. Probab.* **17** 1054–1078.
- [2] AUER, P., CESA-BIANCHI, N. and FISCHER, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine Learning* **47** 235–256.
- [3] BARANSI, A., MAILLARD, O.A. and MANNOR, S. (2014). Sub-sampling for multi-armed bandits. *Proceedings of the European Conference on Machine Learning* pp.13.
- [4] BERRY, D. and FRISTEDT, B. (1985). *Bandit problems*. Chapman and Hall, London.
- [5] BREZZI, M. and LAI, T.L. (2002). Optimal learning and experimentation in bandit problems. *J. Econ. Dynamics Cont.* **27** 87–108.
- [6] BURNETAS, A. and KATEHAKIS, M. (1996). Optimal adaptive policies for sequential allocation problems. *Adv. Appl. Math.* **17** 122–142.
- [7] BURTINI, G., LOEPPKY, J. and LAWRENCE, R. (2015). A survey of online experiment design with the stochastic multi-armed bandit. arXiv:1510.00757.
- [8] CHANG, F. and LAI, T.L. (1987). Optimal stopping and dynamic allocation. *Adv. Appl. Probab.* **19** 829–853.
- [9] GITTINS, J.C. Bandit processes and dynamic allocation indices. *JRSS‘B’* **41** 148–177.
- [10] GITTINS, J.C. and JONES, D.M. (1979). A dynamic allocation index for the discounted multi-armed bandit problem. *Biometrika* **66** 561–565.
- [11] HU, I. and WEI, C.Z. (1989). Irreversible adaptive allocation rules. *Ann. Statist.* **17** 801–823.
- [12] KAUFMANN, E., CAPPÉ and GARIVIER, A. (2012). On Bayesian upper confidence bounds for bandit problems. *J. Machine Learning Res.*
- [13] KIM, M. and LIM, A. (2016). Robust multiarmed bandit problems. *Management Sci.* **62** 264–285.

- [14] KULESHOV, V. and PRECUP, D. (2014). Algorithms for the multi-armed bandit problem. arXiv:1402.6028.
- [15] LAI, T.L. (1987). Adaptive treatment allocation and the multi-armed bandit problem. *Ann. Statist.* **15** 1091–1114.
- [16] LAI, T.L. and ROBBINS, H. (1985). Asymptotically efficient adaptive allocation rules. *Adv. Appl. Math.* **6** 4–22.
- [17] SHIVASWAMY, P. and JOACHIMS, T. (2012). Multi-armed bandit problems with history. *J. Machine Learning Res.*
- [18] SUTTON, B. and BARTO, A. (1998). *Reinforcement Learning, an Introduction*. MIT Press, Cambridge.
- [19] THATHACHER V. and SASTRY, P.S. (1985). A class of rapidly converging algorithms for learning automata. *IEEE Trans. Systems, Man Cyber.* **16** 168–175.
- [20] YAKOWITZ, S. and LOWE, W. (1991). Nonparametric bandit problems. *Ann. Oper. Res.* **28** 297–312.