



# Risk-Aversion in Multi-armed Bandits

Amir Sani, Alessandro Lazaric, Rémi Munos

## ► To cite this version:

Amir Sani, Alessandro Lazaric, Rémi Munos. Risk-Aversion in Multi-armed Bandits. [Research Report] 2012. hal-00750298

**HAL Id: hal-00750298**

**<https://inria.hal.science/hal-00750298>**

Submitted on 9 Jan 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# Risk–Aversion in Multi–armed Bandits

---

Amir Sani

Alessandro Lazaric

Rémi Munos

INRIA Lille - Nord Europe, Team SequeL

{amir.sani, alessandro.lazaric, remi.munos}@inria.fr

## Abstract

Stochastic multi–armed bandits solve the Exploration–Exploitation dilemma and ultimately maximize the expected reward. Nonetheless, in many practical problems, maximizing the expected reward is not the most desirable objective. In this paper, we introduce a novel setting based on the principle of risk–aversion where the objective is to compete against the arm with the best risk–return trade–off. This setting proves to be intrinsically more difficult than the standard multi–arm bandit setting due in part to an exploration risk which introduces a regret associated to the variability of an algorithm. Using variance as a measure of risk, we introduce two new algorithms, investigate their theoretical guarantees, and report preliminary empirical results.

## 1 Introduction

The multi–armed bandit [13] elegantly formalizes the problem of on–line learning with partial feedback, which encompasses a large number of real–world applications, such as clinical trials, online advertisements, adaptive routing, and cognitive radio. In the stochastic multi–armed bandit model, a learner chooses among several arms (e.g., different treatments), each characterized by an independent reward distribution (e.g., the treatment effectiveness). At each point in time, the learner selects one arm and receives a noisy reward observation from that arm (e.g., the effect of the treatment on one patient). Given a finite number of  $n$  rounds (e.g., patients involved in the clinical trial), the learner faces a dilemma between repeatedly exploring all arms and collecting reward information versus exploiting current reward estimates by selecting the arm with the highest estimated reward. Roughly speaking, the learning objective is to solve this exploration–exploitation dilemma and accumulate as much reward as possible over  $n$  rounds. In particular, multi–arm bandit literature typically focuses on the problem of finding a learning algorithm capable of maximizing the expected cumulative reward (i.e., the reward collected over  $n$  rounds averaged over all possible observation realizations), thus implying that the best arm returns the highest expected reward. Nonetheless, in many practical problems, maximizing the expected reward is not always the most desirable objective. For instance, in clinical trials, the treatment which works best *on average* might also have considerable *variability*; resulting in adverse side effects for some patients. In this case, a treatment which is less effective on average but consistently effective on different patients may be preferable w.r.t. an effective but risky treatment. More generally, some application objectives require an effective trade–off between risk and reward.

There is no agreed upon definition for risk. A variety of behaviours result in an uncertainty which might be deemed unfavourable for a specific application and referred to as a risk. For example, a solution with guarantees over multiple runs of an algorithm may not satisfy the desire for a solution with low variability over a single implementation of an algorithm. Two foundational risk modeling paradigms are Expected Utility theory [12] and the historically popular and accessible Mean-Variance paradigm [10]. A large part of decision–making theory focuses on defining and managing risk (see e.g., [9] for an introduction to risk from an expected utility theory perspective).

Risk has mostly been studied in on-line learning within the so-called expert advice setting (i.e., adversarial full-information on-line learning). In particular, [8] showed that in general, although it is possible to achieve a small regret w.r.t. to the expert with the best average performance, it is not possible to compete against the expert which best trades off between average return and risk. On the other hand, it is possible to define no-regret algorithms for simplified measures of risk-return. [15] studied the case of pure risk minimization (notably variance minimization) in an on-line setting where at each step the learner is given a covariance matrix and must choose a weight vector that minimizes the variance. The regret is then computed over horizon  $n$  and compared to the fixed weights minimizing the variance in hindsight. In the multi-arm bandit domain, the most interesting results are by [5] and [14]. [5] introduced an analysis of the expected regret and its distribution, revealing that an anytime version of *UCB* [6] and *UCB-V* might have large regret with some non-negligible probability.<sup>1</sup> This analysis is further extended by [14] who derived negative results which show no anytime algorithm can achieve a regret with both a small expected regret and exponential tails. Although these results represent an important step towards the analysis of risk within bandit algorithms, they are limited to the case where an algorithm's cumulative reward is compared to the reward obtained by pulling the arm with the highest expectation.

In this paper, we focus on the problem of competing against the arm with the best risk-return trade-off. In particular, we refer to the first and most popular measure of risk-return, the mean-variance model introduced by [10]. In Section 2 we introduce notation and define the mean-variance bandit problem. In Section 3 we introduce a confidence-bound algorithm and study its theoretical properties. In Section 5 we report a set of numerical simulations aiming at validating the theoretical results. Finally, in Section 7 we conclude with a discussion on possible extensions. The proofs and additional experiments are reported in the appendix.

## 2 Mean-Variance Multi-arm Bandit

In this section we introduce the main notation used throughout the paper and define the mean-variance multi-arm bandit problem.

We consider the standard multi-arm bandit setting with  $K$  arms, each characterized by a distribution  $\nu_i$  bounded in the interval  $[0, 1]$ . Each distribution has a mean  $\mu_i$  and a variance  $\sigma_i^2$ . The bandit problem is defined over a finite horizon of  $n$  rounds. We denote by  $X_{i,s} \sim \nu_i$  the  $s$ -th random sample drawn from the distribution of arm  $i$ . All arms and samples are independent. In the multi-arm bandit protocol, at each round  $t$ , an algorithm selects arm  $I_t$  and observes sample  $X_{I_t, T_{i,t}}$ , where  $T_{i,t}$  is the number of samples observed from arm  $i$  up to time  $t$  (i.e.,  $T_{i,t} = \sum_{s=1}^t \mathbb{I}\{I_s = i\}$ ).

While in the standard literature on multi-armed bandits the objective is to select the arm leading to the highest reward in *expectation* (the arm with the largest expected value  $\mu_i$ ), here we focus on the problem of finding the arm which effectively trades off between its expected reward (i.e., the *return*) and its variability (i.e., the *risk*). Although a large number of models for risk-return trade-off have been proposed, here we focus on the most historically popular and simple model: the mean-variance model proposed by [10],<sup>2</sup> where the return of an arm is measured by the expected reward and its risk by its variance.

**Definition 1.** *The mean-variance of an arm  $i$  with mean  $\mu_i$ , variance  $\sigma_i^2$  and coefficient of absolute risk tolerance  $\rho$  is defined as<sup>3</sup>  $MV_i = \sigma_i^2 - \rho\mu_i$ .*

Thus it easily follows that the best arm minimizes the mean-variance, that is  $i^* = \arg \min_{i=1, \dots, K} MV_i$ . We notice that we can obtain two extreme settings depending on the value of risk tolerance  $\rho$ . As  $\rho \rightarrow \infty$ , the mean-variance of arm  $i$  tends to the opposite of its expected value  $\mu_i$  and the problem reduces to the standard expected reward maximization traditionally considered in multi-arm bandit problems. With  $\rho = 0$ , the mean-variance reduces to minimizing the variance  $\sigma_i^2$  and the objective becomes variance minimization.

<sup>1</sup>Although the analysis is mostly directed to the pseudo-regret, as commented in Remark 2 at page 23 of [5], it can be extended to the true regret.

<sup>2</sup>We discuss the limitations of this model and possible extensions to other models of risk in Section 7.

<sup>3</sup>The coefficient of risk tolerance is the inverse of the more popular coefficient of risk aversion  $A = 1/\rho$ .

Given  $\{X_{i,s}\}_{s=1}^t$  i.i.d. samples from the distribution  $\nu_i$ , we define the empirical mean–variance of an arm  $i$  with  $t$  samples as  $\widehat{\text{MV}}_{i,t} = \hat{\sigma}_{i,t}^2 - \rho \hat{\mu}_{i,t}$ , where

$$\hat{\mu}_{i,t} = \frac{1}{t} \sum_{s=1}^t X_{i,s}, \quad \hat{\sigma}_{i,t}^2 = \frac{1}{t} \sum_{s=1}^t (X_{i,s} - \hat{\mu}_{i,t})^2. \quad (1)$$

We now consider a learning algorithm  $\mathcal{A}$  and its corresponding performance over  $n$  rounds. Similar to a single arm  $i$  we define its empirical mean–variance as

$$\widehat{\text{MV}}_n(\mathcal{A}) = \hat{\sigma}_n^2(\mathcal{A}) - \rho \hat{\mu}_n(\mathcal{A}), \quad (2)$$

where

$$\hat{\mu}_n(\mathcal{A}) = \frac{1}{n} \sum_{t=1}^n Z_t, \quad \hat{\sigma}_n^2(\mathcal{A}) = \frac{1}{n} \sum_{t=1}^n (Z_t - \hat{\mu}_n(\mathcal{A}))^2, \quad (3)$$

with  $Z_t = X_{I_t, T_{i,t}}$ , that is the reward collected by the algorithm at time  $t$ . This leads to a natural definition of the (random) regret at each single run of the algorithm as the difference in the mean–variance performance of the algorithm compared to the best arm.

**Definition 2.** *The regret for a learning algorithm  $\mathcal{A}$  over  $n$  rounds is defined as*

$$\mathcal{R}_n(\mathcal{A}) = \widehat{\text{MV}}_n(\mathcal{A}) - \widehat{\text{MV}}_{i^*, n}. \quad (4)$$

Given this definition, the objective is to design an algorithm whose regret decreases as the number of rounds increases (in high probability or in expectation).

We notice that the previous definition actually depends on *unobserved* samples. In fact,  $\widehat{\text{MV}}_{i^*, n}$  is computed on  $n$  samples  $i^*$  which are not actually observed when running  $\mathcal{A}$ . This matches the definition of *true* regret in standard bandits (see e.g., [5]). Thus, in order to clarify the main components characterizing the regret, we introduce additional notation. Let

$$Y_{i,t} = \begin{cases} X_{i^*, t} & \text{if } i = i^* \\ X_{i^*, t'} \text{ with } t' = T_{i^*, n} + \sum_{j < i, j \neq i^*} T_{j,n} + t & \text{otherwise} \end{cases}$$

be a renaming of the samples from the optimal arm, such that while the algorithm was pulling arm  $i$  for the  $t$ -th time,  $Y_{i,t}$  is the unobserved sample from  $i^*$ . Then we define the corresponding mean and variance as

$$\tilde{\mu}_{i, T_{i,n}} = \frac{1}{T_{i,n}} \sum_{t=1}^{T_{i,n}} Y_{i,t}, \quad \tilde{\sigma}_{i, T_{i,n}}^2 = \frac{1}{T_{i,n}} \sum_{t=1}^{T_{i,n}} (Y_{i,t} - \tilde{\mu}_{i, T_{i,n}})^2. \quad (5)$$

Given these additional definitions, we can rewrite the regret as (see Appendix A.1)

$$\begin{aligned} \mathcal{R}_n(\mathcal{A}) &= \frac{1}{n} \sum_{i \neq i^*} T_{i,n} \left[ (\hat{\sigma}_{i, T_{i,n}}^2 - \rho \hat{\mu}_{i, T_{i,n}}) - (\tilde{\sigma}_{i, T_{i,n}}^2 - \rho \tilde{\mu}_{i, T_{i,n}}) \right] \\ &\quad + \frac{1}{n} \sum_{i=1}^K T_{i,n} (\hat{\mu}_{i, T_{i,n}} - \hat{\mu}_n(\mathcal{A}))^2 - \frac{1}{n} \sum_{i=1}^K T_{i,n} (\tilde{\mu}_{i, T_{i,n}} - \hat{\mu}_{i^*, n})^2. \end{aligned} \quad (6)$$

Since the last term is always negative and small <sup>4</sup>, our analysis focuses on the first two terms which reveal two interesting characteristics of  $\mathcal{A}$ . First, an algorithm  $\mathcal{A}$  suffers a regret whenever it chooses a suboptimal arm  $i \neq i^*$  and the regret corresponds to the difference in the empirical mean–variance of  $i$  w.r.t. the optimal arm  $i^*$ . Such a definition has a strong similarity to the standard definition of regret, where  $i^*$  is the arm with highest expected value and the regret depends on the number of times suboptimal arms are pulled and their respective gaps w.r.t. the optimal arm  $i^*$ . In contrast to the standard formulation of regret,  $\mathcal{A}$  also suffers an additional regret from the variance  $\hat{\sigma}_n^2(\mathcal{A})$ , which depends on the variability of pulls  $T_{i,n}$  over different arms. Recalling the definition of the mean

<sup>4</sup>More precisely, it can be shown that this term decreases with rate  $O(K \log(1/\delta)/n)$  with probability  $1 - \delta$ .

$\hat{\mu}_n(\mathcal{A})$  as the weighted mean of the empirical means  $\hat{\mu}_{i,T_{i,n}}$  with weights  $T_{i,n}/n$  (see eq. 3), we notice that this second term is a weighted variance of the means and illustrates the exploration risk of the algorithm. In fact, if an algorithm simply selects and pulls a single arm from the beginning, it would not suffer any exploration risk (secondary regret) since  $\hat{\mu}_n(\mathcal{A})$  would coincide with  $\hat{\mu}_{i,T_{i,n}}$  for the chosen arm and all other components would have zero weight. On the other hand, an algorithm accumulates exploration risk through this second term as the mean  $\hat{\mu}_n(\mathcal{A})$  deviates from any specific arm; where the maximum exploration risk peaks at the mean  $\hat{\mu}_n(\mathcal{A})$  furthest from all arm means.

The previous definition of regret can be further elaborated to obtain the upper bound (see App. A.1)

$$\mathcal{R}_n(\mathcal{A}) \leq \frac{1}{n} \sum_{i \neq i^*} T_{i,n} \hat{\Delta}_i + \frac{1}{n^2} \sum_{i=1}^K \sum_{j \neq i} T_{i,n} T_{j,n} \hat{\Gamma}_{i,j}^2, \quad (7)$$

where  $\hat{\Delta}_i = (\hat{\sigma}_{i,T_{i,n}}^2 - \tilde{\sigma}_{i,T_{i,n}}^2) - \rho(\hat{\mu}_{i,T_{i,n}} - \tilde{\mu}_{i,T_{i,n}})$  and  $\hat{\Gamma}_{i,j}^2 = (\hat{\mu}_{i,T_{i,n}} - \hat{\mu}_{j,T_{j,n}})^2$ . Unlike the definition in eq. 6, this upper bound explicitly illustrates the relationship between the regret and the number of pulls  $T_{i,n}$ , suggesting that a bound on the pulls is sufficient to bound the regret.

Finally, we can also introduce a definition of the pseudo-regret.

**Definition 3.** *The pseudo regret for a learning algorithm  $\mathcal{A}$  over  $n$  rounds is defined as*

$$\tilde{\mathcal{R}}_n(\mathcal{A}) = \frac{1}{n} \sum_{i \neq i^*} T_{i,n} \Delta_i + \frac{2}{n^2} \sum_{i=1}^K \sum_{j \neq i} T_{i,n} T_{j,n} \Gamma_{i,j}^2, \quad (8)$$

where  $\Delta_i = \text{MV}_i - \text{MV}_{i^*}$  and  $\Gamma_{i,j} = \mu_i - \mu_j$ .

In the following, we denote the two components of the pseudo-regret as

$$\tilde{\mathcal{R}}_n^\Delta(\mathcal{A}) = \frac{1}{n} \sum_{i \neq i^*} T_{i,n} \Delta_i, \quad \text{and} \quad \tilde{\mathcal{R}}_n^\Gamma(\mathcal{A}) = \frac{2}{n^2} \sum_{i=1}^K \sum_{j \neq i} T_{i,n} T_{j,n} \Gamma_{i,j}^2. \quad (9)$$

Where  $\tilde{\mathcal{R}}_n^\Delta(\mathcal{A})$  constitutes the standard regret derived from the traditional formulation of the multi-arm bandit problem and  $\tilde{\mathcal{R}}_n^\Gamma(\mathcal{A})$  denotes the exploration risk. This regret can be shown to be close to the true regret up to small terms with high probability.

**Lemma 1.** *Given definitions 2 and 3,*

$$\mathcal{R}_n(\mathcal{A}) \leq \tilde{\mathcal{R}}_n(\mathcal{A}) + (5 + \rho) \sqrt{\frac{2K \log(6nK/\delta)}{n}} + 4\sqrt{2} \frac{K \log(6nK/\delta)}{n},$$

with probability at least  $1 - \delta$ .

The previous lemma shows that any (high-probability) bound on the pseudo-regret immediately translates into a bound on the true regret. Thus, we report most of the theoretical analysis according to  $\tilde{\mathcal{R}}_n(\mathcal{A})$ . Nonetheless, it is interesting to notice the major difference between the true and pseudo-regret when compared to the standard bandit problem. In fact, it is possible to show in the risk-averse case that the pseudo-regret is not an unbiased estimator of the true regret, i.e.,  $\mathbb{E}[\mathcal{R}_n] \neq \mathbb{E}[\tilde{\mathcal{R}}_n]$ . Thus, in order to bound the expectation of  $\mathcal{R}_n$  we build on the high-probability result from Lemma 1.

### 3 The Mean-Variance Lower Confidence Bound Algorithm

In this section we introduce a novel risk-averse bandit algorithm whose objective is to identify the arm which best trades off risk and return. The algorithm is a natural extension of *UCB1* [6] and we report a theoretical performance analysis on how well it balances the exploration needed to identify the best arm versus the risk of pulling arms with different means.

```

Input: Confidence  $\delta$ 
for  $t = 1, \dots, n$  do
  for  $i = 1, \dots, K$  do
    Compute  $B_{i,T_{i,t-1}} = \widehat{MV}_{i,T_{i,t-1}} - (5 + \rho) \sqrt{\frac{\log 1/\delta}{2T_{i,t-1}}}$ 
  end for
  Return  $I_t = \arg \min_{i=1, \dots, K} B_{i,T_{i,t-1}}$ 
  Update  $T_{i,t} = T_{i,t-1} + 1$ 
  Observe  $X_{I_t, T_{i,t}} \sim \nu_{I_t}$ 
  Update  $\widehat{MV}_{i,T_{i,t}}$ 
end for

```

Figure 1: Pseudo-code of the *MV-LCB* algorithm.

### 3.1 The Algorithm

We propose an index-based bandit algorithm which estimates the mean-variance of each arm and selects the optimal arm according to the optimistic confidence-bounds on the current estimates. A sketch of the algorithm is reported in Figure 1. For each arm, the algorithm keeps track of the empirical mean-variance  $\widehat{MV}_{i,s}$  computed according to  $s$  samples. We can build high-probability confidence bounds on empirical mean-variance through an application of the Chernoff-Hoeffding inequality (see e.g., [1] for the bound on the variance) on terms  $\hat{\mu}$  and  $\hat{\sigma}^2$ .

**Lemma 2.** *Let  $\{X_{i,s}\}$  be i.i.d. random variables bounded in  $[0, 1]$  from the distribution  $\nu_i$  with mean  $\mu_i$  and variance  $\sigma_i^2$ , and the empirical mean  $\hat{\mu}_{i,s}$  and variance  $\hat{\sigma}_{i,s}^2$  computed as in Equation 1, then*

$$\mathbb{P} \left[ \exists i = 1, \dots, K, s = 1, \dots, n, |\widehat{MV}_{i,s} - MV_i| \geq (5 + \rho) \sqrt{\frac{\log 1/\delta}{2s}} \right] \leq 6nK\delta,$$

The algorithm in Figure 1 implements the principle of optimism in the face of uncertainty used in many multi-arm bandit algorithms. On the basis of the previous confidence bounds, we define a lower-confidence bound on the mean-variance of arm  $i$  when it has been pulled  $s$  times as

$$B_{i,s} = \widehat{MV}_{i,s} - (5 + \rho) \sqrt{\frac{\log 1/\delta}{2s}}, \quad (10)$$

where  $\delta$  is an input parameter of the algorithm. Given the index of each arm at each round  $t$ , the algorithm simply selects the arm with the smallest mean-variance index, i.e.,  $I_t = \arg \min_i B_{i,T_{i,t-1}}$ . We refer to this algorithm as the mean-variance lower-confidence bound (*MV-LCB*) algorithm.

**Remark 1.** We notice that the algorithm reduces to *UCB1* whenever  $\rho \rightarrow \infty$ . This is coherent with the fact that for  $\rho \rightarrow \infty$  the mean-variance problem reduces to the maximization of the cumulative reward, for which *UCB1* is already known to be nearly-optimal. On the other hand, for  $\rho = 0$ , which leads to the problem of cumulative reward variance minimization, the algorithm plays according to a lower-confidence-bound on the variances.

**Remark 2.** The *MV-LCB* algorithm is parameterized by a parameter  $\delta$  which defines the confidence level of the bounds employed in the definition of the index (10). In Theorem 1 we show how to optimize the parameter when the horizon  $n$  is known in advance. On the other hand, if  $n$  is not known, it is possible to design an anytime version of *MV-LCB* by defining a non-decreasing exploration sequence  $(\varepsilon_t)_t$  instead of the term  $\log 1/\delta$ .

### 3.2 Theoretical Analysis

In this section we report the analysis of the regret  $\mathcal{R}_n(\mathcal{A})$  of *MV-LCB* (Fig. 1). As highlighted in eq. 7, it is enough to analyze the number of pulls for each of the arms to recover a bound on the regret. The proofs (reported in the appendix) are mostly based on similar arguments to the proof of *UCB*.

We derive the following regret bound in high probability and expectation.

**Theorem 1.** Let the optimal arm  $i^*$  be unique and  $b = 2(5 + \rho)$ , the *MV-LCB* algorithm achieves a pseudo-regret bounded as

$$\tilde{\mathcal{R}}_n(\mathcal{A}) \leq \frac{b^2 \log 1/\delta}{n} \left( \sum_{i \neq i^*} \frac{1}{\Delta_i} + 4 \sum_{i \neq i^*} \frac{\Gamma_{i^*,i}^2}{\Delta_i^2} + \frac{2b^2 \log 1/\delta}{n} \sum_{i \neq i^*} \sum_{\substack{j \neq i \\ j \neq i^*}} \frac{\Gamma_{i,j}^2}{\Delta_i^2 \Delta_j^2} \right) + \frac{5K}{n},$$

with probability at least  $1 - 6nK\delta$ . Similarly, if *MV-LCB* is run with  $\delta = 1/n^2$  then

$$\mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})] \leq \frac{2b^2 \log n}{n} \left( \sum_{i \neq i^*} \frac{1}{\Delta_i} + 4 \sum_{i \neq i^*} \frac{\Gamma_{i^*,i}^2}{\Delta_i^2} + \frac{4b^2 \log n}{n} \sum_{i \neq i^*} \sum_{\substack{j \neq i \\ j \neq i^*}} \frac{\Gamma_{i,j}^2}{\Delta_i^2 \Delta_j^2} \right) + (17 + 6\rho) \frac{K}{n}.$$

**Remark 1 (the bound).** Let  $\Delta_{\min} = \min_{i \neq i^*} \Delta_i$  and  $\Gamma_{\max} = \max_i |\Gamma_i|$ , then a rough simplification of the previous bound leads to

$$\mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})] \leq O\left(\frac{K}{\Delta_{\min}} \frac{\log n}{n} + K^2 \frac{\Gamma_{\max}^2}{\Delta_{\min}^4} \frac{\log^2 n}{n}\right).$$

First we notice that the regret decreases as  $O(\log^2 n/n)$ , implying that *MV-LCB* is a consistent algorithm. As already highlighted in Definition 2, the regret is mainly composed by two terms. The first term is due to the difference in the mean-variance of the best arm and the arms pulled by the algorithm, while the second term denotes the additional variance introduced by the exploration risk of pulling arms with different means. In particular, it is interesting to note that this additional term depends on the squared difference in the means of the arms  $\Gamma_{i,j}^2$ . Thus, if all the arms have the same mean, this term would be zero.

**Remark 2 (worst-case analysis).** We can further study the result of Theorem 1 by considering the worst-case performance of *MV-LCB*, that is the performance when the distributions of the arms are chosen so as to maximize the regret. In order to illustrate our argument we consider the simple case of  $K = 2$  arms,  $\rho = 0$  (variance minimization),  $\mu_1 \neq \mu_2$ , and  $\sigma_1^2 = \sigma_2^2 = 0$  (deterministic arms).<sup>5</sup> In this case we have a variance gap  $\Delta = 0$  and  $\Gamma^2 > 0$ . According to the definition of *MV-LCB*, the index  $B_{i,s}$  would simply reduce to  $B_{i,s} = \sqrt{\log(1/\delta)/s}$ , thus forcing the algorithm to pull both arms uniformly (i.e.,  $T_{1,n} = T_{2,n} = n/2$  up to rounding effects). Since the arms have the same variance, there is no direct regret in pulling either one or the other. Nonetheless, the algorithm has an additional variance due to the difference in the samples drawn from distributions with different means. In this case, the algorithm suffers a constant (true) regret

$$\mathcal{R}_n(\text{MV-LCB}) = 0 + \frac{T_{1,n}T_{2,n}}{n^2} \Gamma^2 = \frac{1}{4} \Gamma^2,$$

independent from the number of rounds  $n$ . This argument can be generalized to multiple arms and  $\rho \neq 0$ , since it is always possible to design an environment (i.e., a set of distributions) such that  $\Delta_{\min} = 0$  and  $\Gamma_{\max} \neq 0$ .<sup>6</sup> This result is not surprising. In fact, two arms with the same mean-variance are likely to produce similar observations, thus leading *MV-LCB* to pull the two arms repeatedly over time, since the algorithm is designed to try to discriminate between similar arms. Although this behavior does not suffer from any regret in pulling the “suboptimal” arm (the two arms are equivalent), it does introduce an additional variance, due to the difference in the means of the arms ( $\Gamma \neq 0$ ), which finally leads to a regret the algorithm is not “aware” of. This argument suggests that, for any  $n$ , it is always possible to design an environment for which *MV-LCB* has a constant regret. This is particularly interesting since it reveals a huge gap between the mean-variance problem and the standard expected regret minimization problem and will be further investigated in the numerical simulations presented in Section 5. In fact, in the latter case, *UCB* is known to have a worst-case regret per round of  $\Omega(1/\sqrt{n})$  [3], while in the worst case, *MV-LCB* suffers a constant regret. In the next section we introduce a simple algorithm able to deal with this problem and achieve a vanishing worst-case regret.

<sup>5</sup>Note that in this case (i.e.,  $\Delta = 0$ ), Theorem 1 does not hold, since the optimal arm is not unique.

<sup>6</sup>Notice that this is always possible for a large majority of distributions for which the mean and variance are independent or mildly correlated.

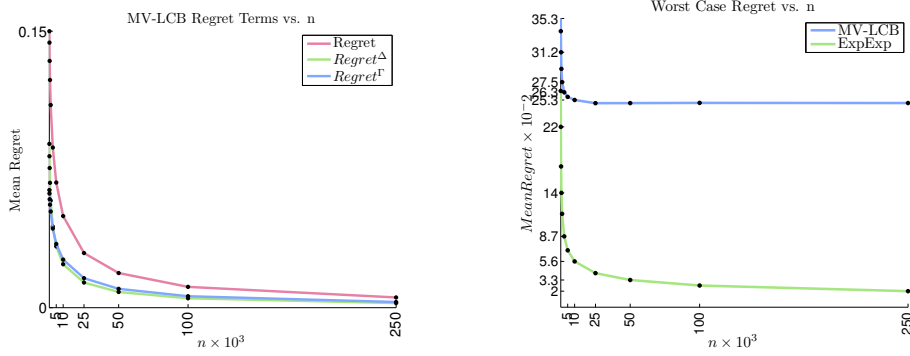


Figure 2: Regret of *MV-LCB* and *ExpExp* in different scenarios.

## 4 The Exploration–Exploitation Algorithm

The *ExpExp* algorithm divides the time horizon  $n$  into two distinct phases of length  $\tau$  and  $n - \tau$  respectively. During the first phase all the arms are explored uniformly, thus collecting  $\tau/K$  samples each<sup>7</sup>. Once the exploration phase is over, the mean–variance of each arm is computed and the arm with the smallest estimated mean–variance  $MV_{i,\tau/K}$  is repeatedly pulled until the end.

The *MV-LCB* is specifically designed to minimize the probability of pulling the wrong arms, so whenever there are two equivalent arms (i.e., arms with the same mean–variance), the algorithm tends to pull them the same number of times, at the cost of potentially introducing an additional variance which might result in a constant regret. On the other hand, *ExpExp* stops exploring the arms after  $\tau$  rounds and then elicits one arm as the best and keeps pulling it for the remaining  $n - \tau$  rounds. Intuitively, the parameter  $\tau$  should be tuned so as to meet different requirements. The first part of the regret (i.e., the regret coming from pulling the suboptimal arms) suggests that the exploration phase  $\tau$  should be long enough for the algorithm to select the empirically best arm  $\hat{i}^*$  at  $\tau$  equivalent to the actual optimal arm  $i^*$  with high probability; and at the same time, as short as possible to reduce the number of times the suboptimal arms are explored. On the other hand, the second part of the regret (i.e., the variance of pulling arms with different means) is minimized by taking  $\tau$  as small as possible (e.g.,  $\tau = 0$  would guarantee a zero regret). The following theorem illustrates the optimal trade-off between these contrasting needs.

**Theorem 2.** *Let  $\text{ExpExp}$  be run with  $\tau = K(n/14)^{2/3}$ , then for any choice of distributions  $\{\nu_i\}$  the expected regret is  $\mathbb{E}[\hat{\mathcal{R}}_n(\mathcal{A})] \leq 2 \frac{K}{n^{1/3}}$ .*

**Remark 1 (the bound).** We first notice that this bound suggests that *ExpExp* performs worse than *MV-LCB* on easy problems. In fact, Theorem 1 demonstrates that *MV-LCB* has a regret decreasing as  $O(K \log(n)/n)$  whenever the gaps  $\Delta$  are not small compared to  $n$ , while in the remarks of Theorem 1 we highlighted the fact that for any value of  $n$ , it is always possible to design an environment which leads *MV-LCB* to suffer a constant regret. On the other hand, the previous bound for *ExpExp* is distribution independent and indicates the regret is still a decreasing function of  $n$  even in the worst case. This opens the question whether it is possible to design an algorithm which works as well as *MV-LCB* on easy problems and as robustly as *ExpExp* on difficult problems.

**Remark 2 (exploration phase).** The previous result can be improved by changing the exploration strategy used in the first  $\tau$  rounds. Instead of a pure uniform exploration of all the arms, we could adopt a best–arm identification algorithms such as *Successive Reject* or *UCB-E*, which maximize the probability of returning the best arm given a fixed budget of rounds  $\tau$  (see e.g., [4]).

## 5 Numerical Simulations

In this section we report numerical simulations aimed at validating the main theoretical findings reported in the previous sections. In the following graphs we study the true regret  $\mathcal{R}_n(\mathcal{A})$  averaged over 500 runs. We first consider the variance minimization problem ( $\rho = 0$ ) with  $K = 2$  Gaussian

<sup>7</sup>In the definition and in the following analysis we ignore rounding effects.



arms set to  $\mu_1 = 1.0$ ,  $\mu_2 = 0.5$ ,  $\sigma_1^2 = 0.05$ , and  $\sigma_2^2 = 0.25$  and run *MV-LCB*<sup>8</sup>. In Figure 2 we report the true regret  $\mathcal{R}_n$  (as in the original definition in eq. 4) and its two components  $\mathcal{R}_n^{\hat{\Delta}}$  and  $\mathcal{R}_n^{\hat{\Gamma}}$  (these two values are defined as in eq. 9 with  $\hat{\Delta}$  and  $\hat{\Gamma}$  replacing  $\Delta$  and  $\Gamma$ ). As expected (see e.g., Theorem 1), the regret is characterized by the regret realized from pulling suboptimal arms and arms with different means (Exploration Risk) and tends to zero as  $n$  increases. Indeed, if we considered two distributions with equal means ( $\mu_1 = \mu_2$ ), the average regret coincides with  $\mathcal{R}_n^{\hat{\Delta}}$ . Furthermore, as shown in Theorem 1 the two regret terms decrease with the same rate  $O(\log n/n)$ .

A detailed analysis of the impact of  $\Delta$  and  $\Gamma$  on the performance of *MV-LCB* is reported in Appendix D. Here we only compare the worst-case performance of *MV-LCB* to *ExpExp* (see Figure 2). In order to have a fair comparison, for any value of  $n$  and for each of the two algorithms, we select the pair  $\Delta_w, \Gamma_w$  which corresponds to the largest regret (we search in a grid of values with  $\mu_1 = 1.5$ ,  $\mu_2 \in [0.4; 1.5]$ ,  $\sigma_1^2 \in [0.0; 0.25]$ , and  $\sigma_2^2 = 0.25$ , so that  $\Delta \in [0.0; 0.25]$  and  $\Gamma \in [0.0; 1.1]$ ). As discussed in Section 4, while the worst-case regret of *ExpExp* keeps decreasing over  $n$ , it is always possible to find a problem for which regret of *MV-LCB* stabilizes to a constant. For numerical results with multiple values of  $\rho$  and 15 arms, please see Appendix D.

## 6 Discussion

In this paper we evaluate the *risk* of an algorithm in terms of the variability of the sequences of samples that it actually generates. Although this notion might resemble other analyses of UCB-based algorithms (see e.g., the high-probability analysis in [5]), it captures different features of the learning algorithm. Whenever a bandit algorithm is run over  $n$  rounds, its behavior, combined with the arms' distributions, generates a probability distribution over sequences of  $n$  rewards. While the *quality* of this sequence is usually defined by its cumulative sum (or average), here we say that a sequence of rewards is *good* if it displays a good trade-off between its (empirical) mean and variance. It is important to notice that this notion of risk-return tradeoff does not coincide with the variance of the algorithm over multiple runs.

Let us consider a simple case with two arms that deterministically generate 0s and 1s respectively, and two different algorithms. Algorithm  $\mathcal{A}_1$  pulls the arms in a fixed sequence at each run (e.g., arm 1, arm 2, arm 1, arm 2, and so on), so that each arm is always pulled  $n/2$  times. Algorithm  $\mathcal{A}_2$  chooses one arm uniformly at random at the beginning of the run and repeatedly pulls this arm for  $n$  rounds. Algorithm  $\mathcal{A}_1$  generates sequences such as 010101... which have high variability within each run, incurs a high regret (e.g., if  $\rho = 0$ ), but has no variance over multiple runs because it always generates the same sequence. On the other hand,  $\mathcal{A}_2$  has no variability in each run, since it generates sequences with only 0s or only 1s, suffers no regret in the case of variance minimization, but has high variance over multiple runs since the two completely different sequences are generated with equal probability. This simple example demonstrates that an algorithm with a very small standard regret w.r.t. the cumulative reward (e.g.,  $\mathcal{A}_1$ ), might result in a very high variability in a single run of the algorithm, while an algorithm with small mean-variance regret (e.g.,  $\mathcal{A}_2$ ) could have a high variance over multiple runs.

## 7 Conclusions

The majority of multi-armed bandit literature focuses on the problem of minimizing the regret w.r.t. the arm with the highest return in expectation. We study the notion of risk associated to the variance over multiple runs and risk of variability associated to a single run of an algorithm. The later case highlights an interesting effect on the regret due to the need to estimate variability within a single sequence of finite random samples before making a risk-averse decision. Further, controlling the variance risk over multiple runs does not necessarily control the risk of variability over a single run. In this paper, we introduced a novel multi-armed bandit setting where the objective is to perform as well as the arm with the best risk-return trade-off. In particular, we relied on the mean-variance model introduced in [10] to measure the performance of the arms and define the regret of a learning algorithm. We proposed two novel algorithms to solve the mean-variance bandit problem and we reported their corresponding theoretical analysis. While *MV-LCB* shows a small regret of order  $O(\log n/n)$  on “easy” problems (i.e., where the mean-variance gaps  $\Delta$  are big w.r.t.  $n$ ), we showed that it has a constant worst-case regret. On the other hand, we proved that *ExpExp* has a vanishing

<sup>8</sup>Notice that although in the paper we assumed the distributions to be bounded in  $[0, 1]$  all the results can be extended to sub-Gaussian distributions.

worst-case regret at the cost of worse performance on “easy” problems. To the best of our knowledge this is the first work introducing risk-aversion in the multi-armed bandit setting and it opens a series of interesting questions.

**Lower bound.** In this paper we introduced two algorithms, *MV-LCB* and *ExpExp*. As discussed in the remarks of Theorem 1 and Theorem 2, *MV-LCB* has a regret of order  $O(\sqrt{K/n})$  on easy problems and  $O(1)$  on difficult problems, while *ExpExp* achieves the same regret  $O(K/n^{1/3})$  over all problems. The primary open question is whether  $O(K/n^{1/3})$  is actually the best possible achievable rate (in the worst-case) for this problem or a better rate is possible. This question is of particular interest since the standard reward expectation maximization problem has a known lower-bound of  $\Omega(\sqrt{1/n})$ , and a minimax rate of  $\Omega(1/n^{1/3})$  for the mean-variance problem would imply that the risk-averse bandit problem is intrinsically more difficult than standard bandit problems.

**Different measures of return-risk.** Considering alternative notions of risk is a straightforward extension to the previous setting. In fact, over the years the mean-variance model has often been criticized. From a point of view of the expected utility theory, the mean-variance model is only justified under a Gaussianity assumption on the arm distributions. It also violates the monotonicity condition due to the different orders of the mean and variance and is not a coherent measure of risk [2]. Furthermore, the variance is a symmetric measure of risk, while it is often the case that only one-sided deviations from the mean are undesirable (e.g., in finance only losses w.r.t. to the expected return are considered as a risk, while any positive deviation is not considered as a real risk). A popular replacement for the mean-variance is to use the  $\alpha$  value-at-risk (i.e., the quantile) to measure the risk of a random variable. The main challenge in this case is the estimation of the value-at-risk for each arm. In fact, while the cumulative distribution of a random variable can be reliably estimated (see e.g., [11]), estimating the quantile might be more difficult.

In [2] axiomatic rules are listed to define coherent measures of risk. Though  $\alpha$  value-at-risk violates these rules, Conditional Value at Risk (otherwise known as average value at risk, tail value at risk, expected shortfall and lower tail risk) passes these rules as a coherent measure of risk. One can easily imagine a lower confidence bound algorithm based on [7] in the same composition as *MV-LCB* which replaces the variance by the conditional value at risk.

The notion of optimality in the risk sensitive setting also depends on the selection of a single-period or multi-period risk evaluation. While the single-period risk of an arm is simply the risk of its distribution, in a multi-period evaluation we consider the risk of the sum of rewards obtained by repeatedly pulling the same arm over  $n$  rounds. Unlike the variance, for which the variance of a sum of  $n$  independent realizations of the same random variable is simply  $n$  times its variance, for other measures of risk (e.g.,  $\alpha$  value-at-risk) this is not necessarily the case. As a result, an arm with the smallest single-period risk might not be the optimal choice over an horizon of  $n$  rounds. Therefore, the performance of a learning algorithm should be compared to the smallest risk that can be achieved by any sequence of arms over  $n$  rounds, thus requiring a new definition of regret.

**Linear bandits.** In linear bandits, each arm is characterized by a marginal distribution with expected value  $\mu_i$  and a covariance matrix  $C$ . At each step the learner chooses a combination of arms and observes the corresponding combined reward. In this case, the best combination is obtained by solving the mean-variance quadratic program  $\min_{\mathbf{x}} (\mathbf{x}^\top C \mathbf{x} - \rho \mathbf{x}^\top \boldsymbol{\mu})$  where  $\mathbf{x}$  is usually a point in the  $K$ -dimensional simplex (e.g., in finance  $\mathbf{x}$  is in the simplex when no short-selling is allowed). Similar to the multi-arm case, the objective is to define an algorithm able to achieve a mean-variance as small as the best point in the simplex over  $n$  rounds.

**Simple regret.** Finally, an interesting related problem is the simple regret setting where the learner is allowed to explore over  $n$  rounds and it only suffers a regret defined on the solution returned at the end. It is known that it is possible to design algorithm able to effectively estimate the mean of the arms and finally return the best arm with high probability. In the risk-return setting, the objective would be to return the arm with the best risk-return tradeoff.

**Acknowledgments** This work was supported by Ministry of Higher Education and Research, Nord-Pas de Calais Regional Council and FEDER through the “contrat de projets état region 2007–2013”, French National Research Agency (ANR) under project LAMPADA  $n^\circ$  ANR-09-EMER-007, European Community’s Seventh Framework Programme (FP7/2007-2013) under grant agreement  $n^\circ$  270327, and PASCAL2 European Network of Excellence.

## References

- [1] András Antos, Varun Grover, and Csaba Szepesvári. Active learning in heteroscedastic noise. *Theoretical Computer Science*, 411:2712–2728, June 2010.
- [2] P Artzner, F Delbaen, JM Eber, and D Heath. Coherent measures of risk. *Mathematical finance*, (June 1996):1–24, 1999.
- [3] Jean-Yves Audibert and Sébastien Bubeck. Regret bounds and minimax policies under partial monitoring. *Journal of Machine Learning Research*, 11:2785–2836, 2010.
- [4] Jean-Yves Audibert, Sébastien Bubeck, and Rémi Munos. Best arm identification in multi-armed bandits. In *Proceedings of the Twenty-third Conference on Learning Theory (COLT’10)*, 2010.
- [5] Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Exploration-exploitation trade-off using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410:1876–1902, 2009.
- [6] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multi-armed bandit problem. *Machine Learning*, 47:235–256, 2002.
- [7] David B. Brown. Large deviations bounds for estimating conditional value-at-risk. *Operations Research Letters*, 35:722–730, 2007.
- [8] Eyal Even-Dar, Michael Kearns, and Jennifer Wortman. Risk-sensitive online learning. In *Proceedings of the 17th international conference on Algorithmic Learning Theory (ALT’06)*, pages 199–213, 2006.
- [9] Christian Gollier. *The Economics of Risk and Time*. The MIT Press, 2001.
- [10] Harry Markowitz. Portfolio selection. *The Journal of Finance*, 7(1):77–91, 1952.
- [11] Pascal Massart. The tight constant in the dvoretzky-kiefer-wolfowitz inequality. *The Annals of Probability*, 18(3):pp. 1269–1283, 1990.
- [12] J Neumann and O Morgenstern. Theory of games and economic behavior. *Princeton University, Princeton*, 1947.
- [13] Herbert Robbins. Some aspects of the sequential design of experiments. *Bulletin of the AMS*, 58:527–535, 1952.
- [14] Antoine Salomon and Jean-Yves Audibert. Deviations of stochastic bandit regret. In *Proceedings of the 22nd international conference on Algorithmic learning theory (ALT’11)*, pages 159–173, 2011.
- [15] Manfred K. Warmuth and Dima Kuzmin. Online variance minimization. In *Proceedings of the 19th Annual Conference on Learning Theory (COLT’06)*, pages 514–528, 2006.

## A The Regret

### A.1 The True Regret

We recall the definition of the (empirical) regret as

$$\mathcal{R}_n(\mathcal{A}) = \widehat{\mathbf{M}\mathbf{V}}_n(\mathcal{A}) - \widehat{\mathbf{M}\mathbf{V}}_{i^*,n}.$$

Given the definitions reported in the main paper, we first elaborate on the two mean terms in the regret as

$$\hat{\mu}_{i^*,n} = \frac{1}{n} \sum_{i=1}^K \sum_{t=1}^{T_{i,n}} Y_{i,t} = \frac{1}{n} \sum_{i=1}^K T_{i,n} \tilde{\mu}_{i,T_{i,n}},$$

and

$$\hat{\mu}_n(\mathcal{A}) = \frac{1}{n} \sum_{i=1}^K \sum_{t=1}^{T_{i,n}} X_{i,t} = \frac{1}{n} \sum_{i=1}^K T_{i,n} \hat{\mu}_{i,T_{i,n}}.$$

Similarly, the two variance terms can be written as

$$\begin{aligned} \hat{\sigma}_n^2(\mathcal{A}) &= \frac{1}{n} \sum_{i=1}^K \sum_{t=1}^{T_{i,n}} (X_{i,t} - \hat{\mu}_n(\mathcal{A}))^2 \\ &= \frac{1}{n} \sum_{i=1}^K \sum_{t=1}^{T_{i,n}} (X_{i,t} - \hat{\mu}_{i,T_{i,n}})^2 + \frac{1}{n} \sum_{i=1}^K \sum_{t=1}^{T_{i,n}} (\hat{\mu}_{i,T_{i,n}} - \hat{\mu}_n(\mathcal{A}))^2 + \frac{2}{n} \sum_{i=1}^K \sum_{t=1}^{T_{i,n}} (X_{i,t} - \hat{\mu}_{i,T_{i,n}})(\hat{\mu}_{i,T_{i,n}} - \hat{\mu}_n(\mathcal{A})) \\ &= \frac{1}{n} \sum_{i=1}^K T_{i,n} \hat{\sigma}_{i,T_{i,n}}^2 + \frac{1}{n} \sum_{i=1}^K T_{i,n} (\hat{\mu}_{i,T_{i,n}} - \hat{\mu}_n(\mathcal{A}))^2 + 0, \end{aligned}$$

and

$$\begin{aligned} \sigma_{i^*,n}^2 &= \frac{1}{n} \sum_{i=1}^K \sum_{t=1}^{T_{i,n}} (Y_{i,t} - \hat{\mu}_{i^*,n})^2 \\ &= \frac{1}{n} \sum_{i=1}^K \sum_{t=1}^{T_{i,n}} (Y_{i,t} - \tilde{\mu}_{i,T_{i,n}})^2 + \frac{1}{n} \sum_{i=1}^K \sum_{t=1}^{T_{i,n}} (\tilde{\mu}_{i,T_{i,n}} - \hat{\mu}_{i^*,n})^2 + \frac{2}{n} \sum_{i=1}^K \sum_{t=1}^{T_{i,n}} (Y_{i,t} - \tilde{\mu}_{i,T_{i,n}})(\tilde{\mu}_{i,T_{i,n}} - \hat{\mu}_{i^*,n}) \\ &= \frac{1}{n} \sum_{i=1}^K T_{i,n} \tilde{\sigma}_{i,T_{i,n}}^2 + \frac{1}{n} \sum_{i=1}^K T_{i,n} (\tilde{\mu}_{i,T_{i,n}} - \hat{\mu}_{i^*,n})^2 + 0. \end{aligned}$$

Putting together these terms, we obtain the regret (see eq. 4)

$$\begin{aligned} \mathcal{R}_n(\mathcal{A}) &= \frac{1}{n} \sum_{i \neq i^*} T_{i,n} \left[ (\hat{\sigma}_{i,T_{i,n}}^2 - \tilde{\sigma}_{i,T_{i,n}}^2) - \rho(\hat{\mu}_{i,T_{i,n}} - \tilde{\mu}_{i,T_{i,n}}) \right] \\ &\quad + \frac{1}{n} \sum_{i=1}^K T_{i,n} (\hat{\mu}_{i,T_{i,n}} - \hat{\mu}_n(\mathcal{A}))^2 - \frac{1}{n} \sum_{i=1}^K T_{i,n} (\tilde{\mu}_{i,T_{i,n}} - \hat{\mu}_{i^*,n})^2 \end{aligned}$$

If we further elaborate the second term, we obtain

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^K T_{i,n} (\hat{\mu}_{i,T_{i,n}} - \hat{\mu}_n(\mathcal{A}))^2 &= \frac{1}{n} \sum_{i=1}^K T_{i,n} \left( \hat{\mu}_{i,T_{i,n}} - \frac{1}{n} \sum_{j=1}^K T_{j,n} \hat{\mu}_{j,T_{j,n}} \right)^2 \\
&= \frac{1}{n} \sum_{i=1}^K T_{i,n} \left( \sum_{j=1}^K \frac{T_{j,n}}{n} (\hat{\mu}_{i,T_{i,n}} - \hat{\mu}_{j,T_{j,n}}) \right)^2 \\
&\leq \frac{1}{n} \sum_{i=1}^K T_{i,n} \sum_{j=1}^K \frac{T_{j,n}}{n} (\hat{\mu}_{i,T_{i,n}} - \hat{\mu}_{j,T_{j,n}})^2 \\
&= \frac{1}{n^2} \sum_{i=1}^K \sum_{j \neq i}^K T_{i,n} T_{j,n} (\hat{\mu}_{i,T_{i,n}} - \hat{\mu}_{j,T_{j,n}})^2.
\end{aligned}$$

Using the definitions  $\hat{\Delta}_i = (\hat{\sigma}_{i,T_{i,n}}^2 - \tilde{\sigma}_{i,T_{i,n}}^2) - \rho(\hat{\mu}_{i,T_{i,n}} - \tilde{\mu}_{i,T_{i,n}})$  and  $\hat{\Gamma}_{i,j}^2 = (\hat{\mu}_{i,T_{i,n}} - \hat{\mu}_{j,T_{j,n}})^2$  we finally obtain an upper-bound on the regret of the form

$$\mathcal{R}_n(\mathcal{A}) \leq \frac{1}{n} \sum_{i \neq i^*} T_{i,n} \hat{\Delta}_i + \frac{1}{n^2} \sum_{i=1}^K \sum_{j \neq i}^K T_{i,n} T_{j,n} \hat{\Gamma}_{i,j}^2.$$

In the following we refer to the two terms as  $\mathcal{R}_n^{\hat{\Delta}}$  and  $\mathcal{R}_n^{\hat{\Gamma}}$ .

## A.2 The Pseudo-Regret

Similar to what is done in the standard bandit problem, we can introduce a different notion of regret. Starting from the last equation in the previous section, we define the pseudo-regret

$$\tilde{\mathcal{R}}_n(\mathcal{A}) = \frac{1}{n} \sum_{i \neq i^*} T_{i,n} \Delta_i + \frac{2}{n^2} \sum_{i=1}^K \sum_{j \neq i}^K T_{i,n} T_{j,n} \Gamma_{i,j}^2,$$

where the empirical values  $\hat{\Delta}_i$  and  $\hat{\Gamma}_{i,j}$  are substituted by their corresponding exact values<sup>9</sup>. In the following we show that the true and pseudo regrets differ for values that tend to zero with high probability.

*Proof.* (Lemma 1)

We define a high-probability event in which the empirical values and the true values only differ for small quantities

$$\mathcal{E} = \left\{ \forall i = 1, \dots, K, \forall s = 1, \dots, n, |\hat{\mu}_{i,s} - \mu_i| \leq \sqrt{\frac{\log 1/\delta}{2s}} \text{ and } |\hat{\sigma}_{i,s}^2 - \sigma_i^2| \leq 5\sqrt{\frac{\log 1/\delta}{2s}} \right\}.$$

Using Chernoff–Hoeffding inequality and a union bound over arms and rounds, we have that  $\mathbb{P}[\mathcal{E}^C] \leq 6nK\delta$ . Under this event we rewrite the empirical  $\hat{\Delta}_i$  as

$$\begin{aligned}
\hat{\Delta}_i &= \Delta_i - (\sigma_i^2 - \sigma_{i^*}^2) + \rho(\mu_i - \mu_{i^*}) + (\hat{\sigma}_{i,T_{i,n}}^2 - \tilde{\sigma}_{i,T_{i,n}}^2) - \rho(\hat{\mu}_{i,T_{i,n}} - \tilde{\mu}_{i,T_{i,n}}) \\
&\leq \Delta_i + 2(5 + \rho) \sqrt{\frac{\log 1/\delta}{2T_{i,n}}}.
\end{aligned}$$

Similarly,  $\hat{\Gamma}_{i,j}$  is upper-bounded as

$$\begin{aligned}
|\hat{\Gamma}_{i,j}| &= |\Gamma_{i,j} - \mu_i + \mu_j + \hat{\mu}_{i,T_{i,n}} - \hat{\mu}_{j,T_{j,n}}| \\
&\leq |\Gamma_{i,j}| + \sqrt{\frac{\log 1/\delta}{2T_{i,n}}} + \sqrt{\frac{\log 1/\delta}{2T_{j,n}}}.
\end{aligned}$$

<sup>9</sup>Notice that the factor 2 in front of the second term is due to a rough upper bounding used in the proof of Lemma 1.

Thus the regret can be written as

$$\begin{aligned}
\mathcal{R}_n(\mathcal{A}) &\leq \frac{1}{n} \sum_{i \neq i^*} T_{i,n} \left( \Delta_i + 2(5 + \rho) \sqrt{\frac{\log 1/\delta}{2T_{i,n}}} \right) + \frac{1}{n^2} \sum_{i=1}^K \sum_{j \neq i} T_{i,n} T_{j,n} \left( |\Gamma_{i,j}| + \sqrt{\frac{\log 1/\delta}{2T_{i,n}}} + \sqrt{\frac{\log 1/\delta}{2T_{j,n}}} \right)^2 \\
&\leq \frac{1}{n} \sum_{i \neq i^*} T_{i,n} \Delta_i + \frac{5 + \rho}{n} \sum_{i \neq i^*} \sqrt{2T_{i,n} \log 1/\delta} + \frac{2}{n^2} \sum_{i=1}^K \sum_{j \neq i} T_{i,n} T_{j,n} \Gamma_{i,j}^2 \\
&\quad + \frac{2\sqrt{2}}{n^2} \sum_{i=1}^K \sum_{j \neq i} T_{j,n} \log 1/\delta + \frac{2\sqrt{2}}{n^2} \sum_{i=1}^K \sum_{j \neq i} T_{i,n} \log 1/\delta \\
&\leq \frac{1}{n} \sum_{i \neq i^*} T_{i,n} \Delta_i + \frac{2}{n^2} \sum_{i=1}^K \sum_{j \neq i} T_{i,n} T_{j,n} \Gamma_{i,j}^2 + (5 + \rho) \sqrt{\frac{2K \log 1/\delta}{n}} + 4\sqrt{2} \frac{K \log 1/\delta}{n}.
\end{aligned}$$

where in the next to last passage we used Jensen's inequality for concave functions and rough upper bounds on other terms ( $K - 1 < K$ ,  $\sum_{i \neq i^*} T_{i,n} \leq n$ ). By recalling the definition of  $\tilde{\mathcal{R}}_n(\mathcal{A})$  we finally obtain

$$\mathcal{R}_n(\mathcal{A}) \leq \tilde{\mathcal{R}}_n(\mathcal{A}) + (5 + \rho) \sqrt{\frac{2K \log 1/\delta}{n}} + 4\sqrt{2} \frac{K \log 1/\delta}{n},$$

with probability  $1 - 6nK\delta$ . Thus we can conclude that any upper bound on the pseudo-regret  $\tilde{\mathcal{R}}_n(\mathcal{A})$  is a valid upper bound for the true regret  $\mathcal{R}_n(\mathcal{A})$ , up to a decreasing term of order  $O(\sqrt{K/n})$ .  $\square$

## B MV-LCB Theoretical Analysis

In order to simplify the notation in the following we use  $b = 2(5 + \rho)$ .

*Proof.* (Theorem 1)

We begin by defining a high-probability event  $\mathcal{E}$  as

$$\mathcal{E} = \left\{ \forall i = 1, \dots, K, \forall s = 1, \dots, n, \quad |\hat{\mu}_{i,s} - \mu_i| \leq \sqrt{\frac{\log 1/\delta}{2s}} \quad \text{and} \quad |\hat{\sigma}_{i,s}^2 - \sigma_i^2| \leq 5\sqrt{\frac{\log 1/\delta}{2s}} \right\}.$$

Using Chernoff-Hoeffding inequality and a union bound over arms and rounds, we have that  $\mathbb{P}[\mathcal{E}^c] \leq 6nK\delta$ .

We now introduce the definition of the algorithm. Consider any time  $t$  when arm  $i \neq i^*$  is pulled (i.e.,  $I_t = i$ ). By definition of the algorithm in Figure 1,  $i$  is selected if its corresponding index  $B_{i,T_{i,t-1}}$  is bigger than for any other arm, notably the best arm  $i^*$ . By recalling the definition of the index and the empirical mean-variance at time  $t$ , we have

$$\begin{aligned}
\hat{\sigma}_{i,T_{i,t-1}}^2 - \rho \hat{\mu}_{i,T_{i,t-1}} - (5 + \rho) \sqrt{\frac{\log 1/\delta}{2T_{i,t-1}}} &= B_{i,T_{i,t-1}} \leq \\
&\leq B_{i^*,T_{i^*,t-1}} = \hat{\sigma}_{i^*,T_{i^*,t-1}}^2 - \rho \hat{\mu}_{i^*,T_{i^*,t-1}} - (5 + \rho) \sqrt{\frac{\log 1/\delta}{2T_{i^*,t-1}}}.
\end{aligned}$$

Over all the possible realizations, we now focus on the realizations in  $\mathcal{E}$ . In this case, we can rewrite the previous condition as

$$\sigma_i^2 - \rho \mu_i - 2(5 + \rho) \sqrt{\frac{\log 1/\delta}{2T_{i,t-1}}} \leq B_{i,T_{i,t-1}} \leq B_{i^*,T_{i^*,t-1}} \leq \sigma_{i^*}^2 - \rho \mu_{i^*}.$$

Let time  $t$  be the last time when arm  $i$  is pulled until the final round  $n$ , then  $T_{i,t-1} = T_{i,n} - 1$  and

$$T_{i,n} \leq \frac{2(5+\rho)^2}{\Delta_i^2} \log \frac{1}{\delta} + 1,$$

which suggests that the suboptimal arms are pulled only few times with high probability. Plugging the bound in the regret in eq. 8 leads to the final statement

$$\tilde{\mathcal{R}}_n(\mathcal{A}) \leq \frac{1}{n} \sum_{i \neq i^*} \frac{b^2 \log 1/\delta}{\Delta_i} + \frac{1}{n} \sum_{i \neq i^*} \frac{4b^2 \log 1/\delta}{\Delta_i^2} \Gamma_{i^*,i}^2 + \frac{1}{n^2} \sum_{i \neq i^*} \sum_{\substack{j \neq i \\ j \neq i^*}} \frac{2b^4 (\log 1/\delta)^2}{\Delta_i^2 \Delta_j^2} \Gamma_{i,j}^2 + \frac{5K}{n},$$

with probability  $1 - 6nK\delta$ .

We now move from the previous high-probability bound to a bound in expectation. The pseudo-regret is (roughly) bounded as  $\tilde{\mathcal{R}}_n(\mathcal{A}) \leq 2 + \rho$  (by bounding  $\Delta_i \leq 1 + \rho$  and  $\Gamma \leq 1$ ), thus

$$\mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})] = \mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})\mathbb{I}\{\mathcal{E}\}] + \mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})\mathbb{I}\{\mathcal{E}^C\}] \leq \mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})\mathbb{I}\{\mathcal{E}\}] + (2 + \rho)\mathbb{P}[\mathcal{E}^C].$$

By using the previous high-probability bound and recalling that  $\mathbb{P}[\mathcal{E}^C] \leq 6nK\delta$ , we have

$$\begin{aligned} \mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})] &\leq \frac{1}{n} \sum_{i \neq i^*} \frac{b^2 \log 1/\delta}{\Delta_i} + \frac{1}{n} \sum_{i \neq i^*} \frac{4b^2 \log 1/\delta}{\Delta_i^2} \Gamma_{i^*,i}^2 + \frac{1}{n^2} \sum_{i \neq i^*} \sum_{\substack{j \neq i \\ j \neq i^*}} \frac{2b^4 (\log 1/\delta)^2}{\Delta_i^2 \Delta_j^2} \Gamma_{i,j}^2 \\ &\quad + \frac{5K}{n} + (2 + \rho)6nK\delta. \end{aligned}$$

The final statement of the lemma follows by tuning the parameter  $\delta = 1/n^2$  so as to have a regret bound decreasing with  $n$ .  $\square$

While a high-probability bound for  $\mathcal{R}_n$  can be immediately obtained from Lemma 1, the expectation of  $\mathcal{R}_n$  is reported in the next corollary.

*Proof.* Since the mean-variance  $-\rho \leq \widehat{\text{MV}} \leq 1/4$ , the regret is bounded by  $-1/4 - \rho \leq \mathcal{R}_n(\mathcal{A}) \leq 1/4 + \rho$ . Thus we have

$$\mathbb{E}[\mathcal{R}_n(\mathcal{A})] \leq u\mathbb{P}[\mathcal{R}_n(\mathcal{A}) \leq u] + \left(\frac{1}{4} + \rho\right)\mathbb{P}[\mathcal{R}_n(\mathcal{A}) > u].$$

By taking  $u$  equal to the previous high-probability bound and recalling that  $\mathbb{P}[\mathcal{E}^C] \leq 6nK\delta$ , we have

$$\begin{aligned} \mathbb{E}[\mathcal{R}_n(\mathcal{A})] &\leq \frac{1}{n} \sum_{i \neq i^*} \frac{b^2 \log 1/\delta}{\Delta_i} + \frac{1}{n} \sum_{i \neq i^*} \frac{4b^2 \log 1/\delta}{\Delta_i^2} \Gamma_{i^*,i}^2 + \frac{1}{n^2} \sum_{i \neq i^*} \sum_{\substack{j \neq i \\ j \neq i^*}} \frac{2b^4 (\log 1/\delta)^2}{\Delta_i^2 \Delta_j^2} \Gamma_{i,j}^2 \\ &\quad + \frac{5K}{n} + b\sqrt{\frac{K \log 1/\delta}{2n}} + 4\sqrt{2}\frac{K \log 1/\delta}{n} + \left(\frac{1}{4} + \rho\right)6nK\delta. \end{aligned}$$

The final statement of the lemma follows by tuning the parameter  $\delta = 1/n^2$  so as to have a regret bound decreasing with  $n$ .  $\square$

## C Exp-Exp Theoretical Analysis

During the exploitation phase the algorithm pulls arm  $\hat{i}^*$  with the smallest empirical variance estimated during the exploration phase of length  $\tau$ . As a result, the number of pulls of each arm is

$$T_{i,n} = \frac{\tau}{K} + (n - \tau)\mathbb{I}\{i = \hat{i}^*\} \quad (11)$$

We analyze the two terms of the regret separately.

$$\tilde{\mathcal{R}}_n^\Delta = \frac{1}{n} \sum_{i \neq i^*} \left( \frac{\tau}{K} + (n - \tau) \mathbb{I}\{i = \hat{i}^*\} \right) \Delta_i = \frac{\tau}{nK} \sum_{i \neq i^*} \Delta_i + \frac{n - \tau}{n} \sum_{i \neq i^*} \underbrace{\Delta_i \mathbb{I}\{i = \hat{i}^*\}}_{(a)}.$$

We notice that the only random variable in this formulation is the best arm  $\hat{i}^*$  at the end of the exploration phase. We thus compute the expected value of  $\tilde{\mathcal{R}}_n^\Delta$ .

$$\begin{aligned} \mathbb{E}[(a)] &= \mathbb{P}[i = \hat{i}^*] \Delta_i = \mathbb{P}[\forall j \neq i, \hat{\sigma}_{i,\tau/K}^2 \leq \hat{\sigma}_{j,\tau/K}^2] \Delta_i \\ &\leq \mathbb{P}[\hat{\sigma}_{i,\tau/K}^2 \leq \hat{\sigma}_{i^*,\tau/K}^2] \Delta_i = \mathbb{P}[(\hat{\sigma}_{i,\tau/K}^2 - \sigma_i^2) + (\sigma_{i^*}^2 - \hat{\sigma}_{i^*,\tau/K}^2) \leq \Delta_i] \Delta_i \\ &\leq 2\Delta_i \exp\left(-\frac{\tau}{K} \Delta_i^2\right) \end{aligned}$$

The second term in the regret can be bounded as follows.

$$\begin{aligned} \tilde{\mathcal{R}}_n^\Gamma &= \frac{1}{n^2} \sum_{i=1}^K \sum_{j \neq i} \left( \frac{\tau}{K} + (n - \tau) \mathbb{I}\{i = \hat{i}^*\} \right) \left( \frac{\tau}{K} + (n - \tau) \mathbb{I}\{j = \hat{i}^*\} \right) \Gamma_{i,j}^2 \\ &= \frac{1}{n^2} \sum_{i=1}^K \sum_{j \neq i} \left( \frac{\tau^2}{K^2} + (n - \tau)^2 \mathbb{I}\{i = \hat{i}^*\} \mathbb{I}\{j = \hat{i}^*\} + \frac{\tau}{K} (n - \tau) \mathbb{I}\{j = \hat{i}^*\} + \frac{\tau}{K} (n - \tau) \mathbb{I}\{i = \hat{i}^*\} \right) \Gamma_{i,j}^2 \\ &= \frac{\tau^2}{n^2 K^2} \sum_{i=1}^K \sum_{j \neq i} \Gamma_{i,j}^2 + 2 \frac{(n - \tau) \tau}{K n^2} \sum_{i=1}^K \sum_{j \neq i} \Gamma_{i,j}^2 \mathbb{I}\{i = \hat{i}^*\} \\ &\leq \frac{\tau}{n^2} + 2 \frac{(n - \tau) \tau}{n^2} \leq 2 \frac{\tau}{n} \end{aligned}$$

Grouping all the terms,  $\text{ExpExp}$  has an expected regret bounded as

$$\mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})] \leq 2 \frac{\tau}{n} + 2 \sum_{i \neq i^*} \Delta_i \exp\left(-\frac{\tau}{K} \Delta_i^2\right)$$

We can now move to the worst-case analysis of the regret. Let  $f(\Delta_i) = \Delta_i \exp\left(-\frac{\tau}{K} \Delta_i^2\right)$ , the “adversarial” choice of the gap is determined by maximizing the regret, which corresponds to

$$\begin{aligned} f'(\Delta_i) &= \exp\left(-\frac{\tau}{K} \Delta_i^2\right) + \Delta_i \left( -2 \frac{\tau}{K} \Delta_i \exp\left(-\frac{\tau}{K} \Delta_i^2\right) \right) \\ &= \left( 1 - 2 \frac{\tau}{K} \Delta_i^2 \right) \exp\left(-\frac{\tau}{K} \Delta_i^2\right) = 0, \end{aligned}$$

and leads to a worst-case choice for the gap of

$$\Delta_i = \sqrt{\frac{K}{2\tau}}.$$

The worst-case regret is then



$$\mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})] \leq 2\frac{\tau}{n} + (K-1)\sqrt{2K}\frac{1}{\sqrt{\tau}}\exp(-0.5) \leq 2\frac{\tau}{n} + K^{3/2}\frac{1}{\sqrt{\tau}}$$

We can now choose the parameter  $\tau$  minimizing the worst-case regret. Taking the derivative of the regret w.r.t.  $\tau$  we obtain

$$\frac{d\mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})]}{d\tau} = \frac{2}{n} - \frac{1}{2}\left(\frac{K}{\tau}\right)^{3/2} = 0,$$

thus leading to the optimal parameter  $\tau = (n/4)^{2/3}K$ . The final regret is thus bounded as

$$\mathbb{E}[\tilde{\mathcal{R}}_n(\mathcal{A})] \leq 3\frac{K}{n^{1/3}}.$$

## D Additional Simulations

### D.1 Comparison between *MV-LCB* and *ExpExp* with $K = 2$

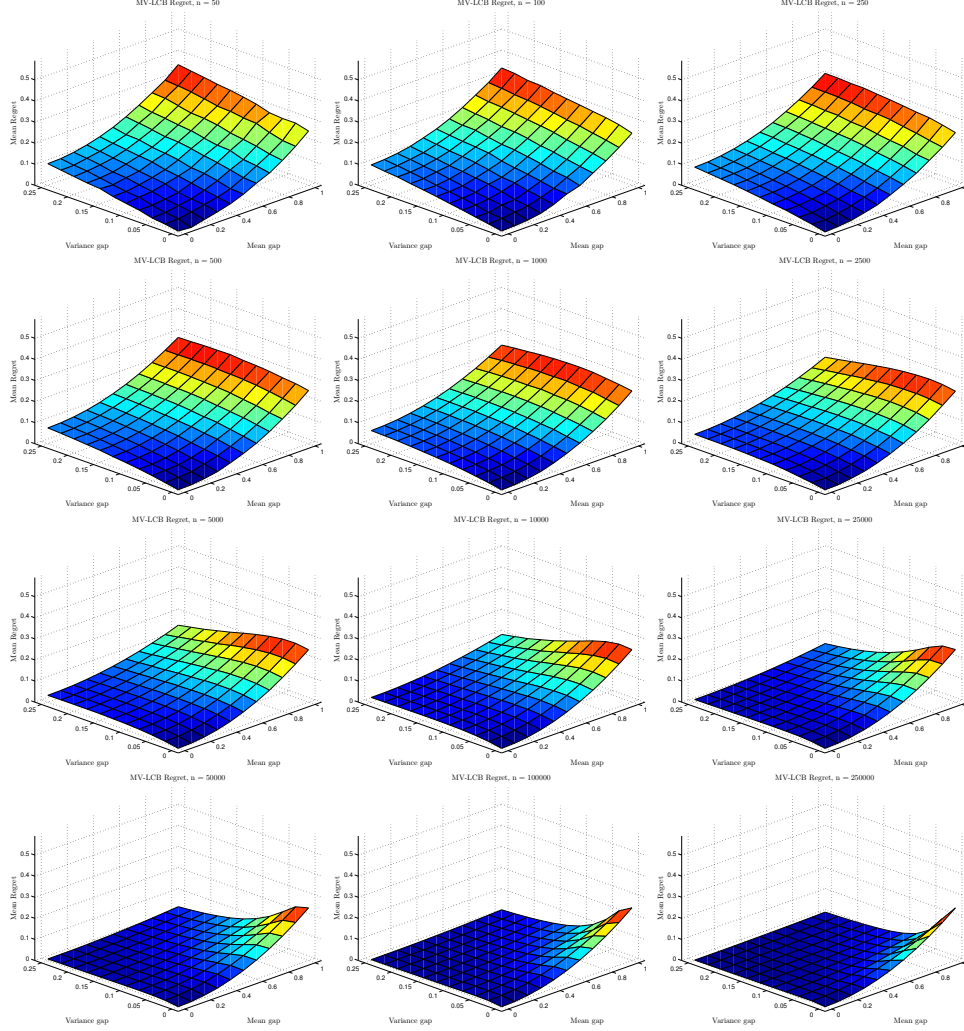


Figure 3: Regret  $\mathcal{R}_n$  of *MV-LCB*.

We consider the variance minimization problem ( $\rho = 0$ ) with  $K = 2$  Gaussian arms with different means and variances. In particular, we consider a grid of values with  $\mu_1 = 1.5$ ,  $\mu_2 \in [0.4; 1.5]$ ,  $\sigma_1^2 \in [0.0; 0.25]$ , and  $\sigma_2^2 = 0.25$ , so that  $\Delta \in [0.0; 0.25]$  and  $\Gamma \in [0.0; 1.1]$  and number of rounds  $n \in [50; 2.5 \times 10^5]$ . Figures 3 and 4 report the mean regret for different values of  $n$ . The colors are renormalized in each plot so that dark blue corresponds to the smallest regret and red to the largest regret. The results confirm the theoretical findings of Theorem 1 and 2. In fact, for simple problems (large gaps  $\Delta$ ) *MV-LCB* converges to a zero-regret faster than *ExpExp*, while for  $\Delta$  close to zero (i.e., equivalent arms), *MV-LCB* has a constant regret which does not decrease with  $n$  and the regret of *ExpExp* slowly decreases to zero.

### D.2 Risk tolerance sensitivity

In Section 5 we report numerical results demonstrating the composition of the regret and performance of algorithms with only 2 arms in the case of variance minimization. Here we report results for a wide range of risk tolerance  $\rho \in [0.0; 10.0]$  and  $K = 15$  arms. We set the mean and variance for each of the 15 arms so that a subset of arms is always dominated (i.e., for any  $\rho$ ,  $MV_i^\rho > MV_{i^*}^\rho$ ) demonstrating the effect of different  $\rho$  values on the position of the optimal arm  $i^*$ .

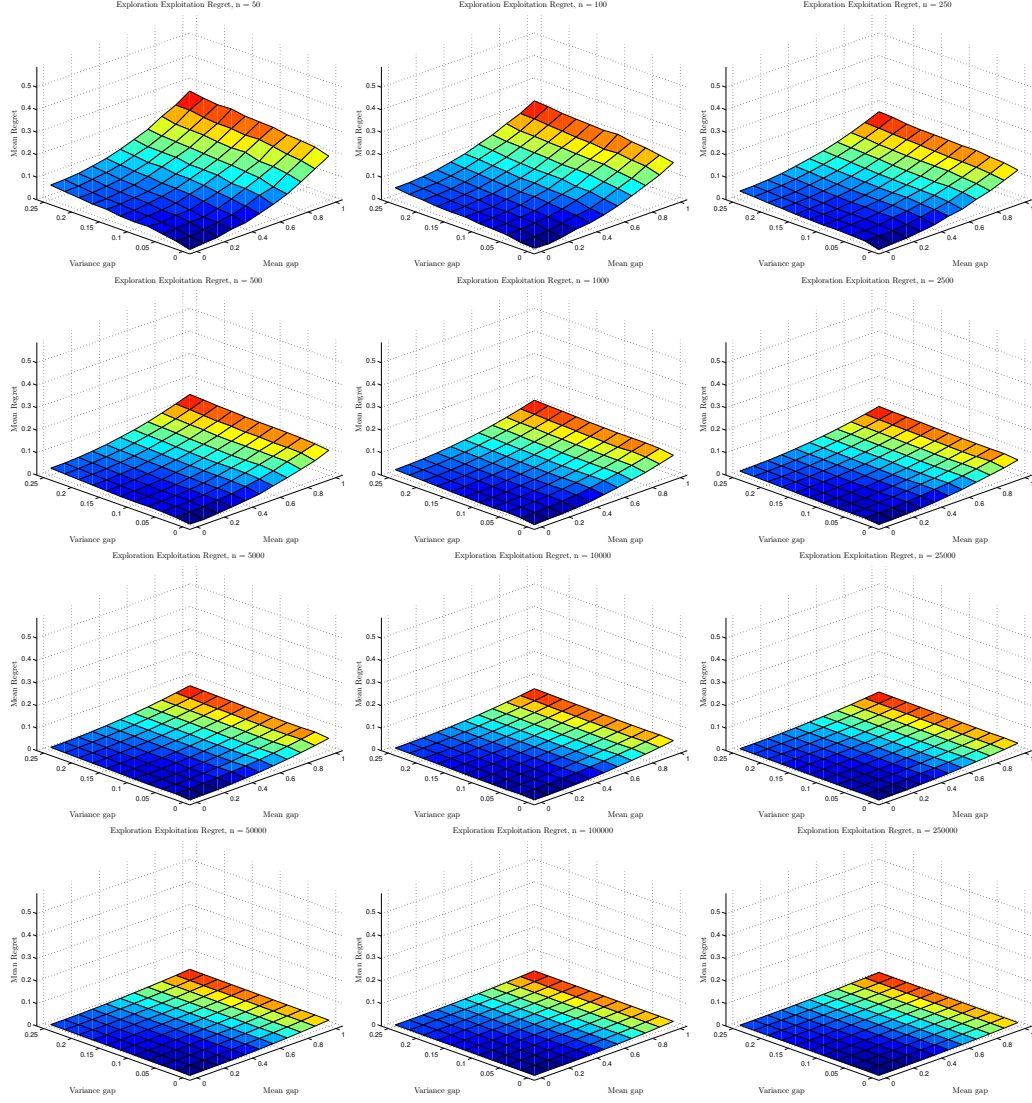


Figure 4: Regret  $\mathcal{R}_n$  of *ExpExp*.

Arm	$\mu$	$\sigma^2$
$\alpha_1$	0.10	0.05
$\alpha_2$	0.20	0.34
$\alpha_3$	0.23	0.28
$\alpha_4$	0.27	0.09
$\alpha_5$	0.32	0.23
$\alpha_6$	0.32	0.72
$\alpha_7$	0.34	0.19
$\alpha_8$	0.41	0.14
$\alpha_9$	0.43	0.44
$\alpha_{10}$	0.54	0.53
$\alpha_{11}$	0.55	0.24
$\alpha_{12}$	0.56	0.36
$\alpha_{13}$	0.67	0.56
$\alpha_{14}$	0.71	0.49
$\alpha_{15}$	0.79	0.85

Arm	$\mu$	$\sigma^2$
$\alpha_1$	0.1	0.05
$\alpha_2$	0.2	0.0725
$\alpha_3$	0.27	0.09
$\alpha_4$	0.32	0.11
$\alpha_5$	0.41	0.145
$\alpha_6$	0.49	0.19
$\alpha_7$	0.55	0.24
$\alpha_8$	0.59	0.28
$\alpha_9$	0.645	0.36
$\alpha_{10}$	0.678	0.413
$\alpha_{11}$	0.69	0.445
$\alpha_{12}$	0.71	0.498
$\alpha_{13}$	0.72	0.53
$\alpha_{14}$	0.765	0.72
$\alpha_{15}$	0.79	0.854

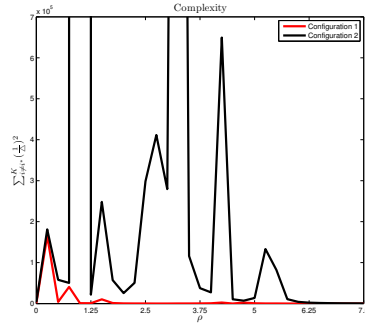


Figure 5: Configuration 1 and Configuration 2 and their corresponding complexity.

In Figure 2 we arranged the true values of each arm along the red frontier and the  $\rho$ -directed performance of the algorithms in a standard deviation–mean plot. The green and blue lines show the standard deviation and mean for the performance of each algorithm for a specific  $\rho$  setting and fi-

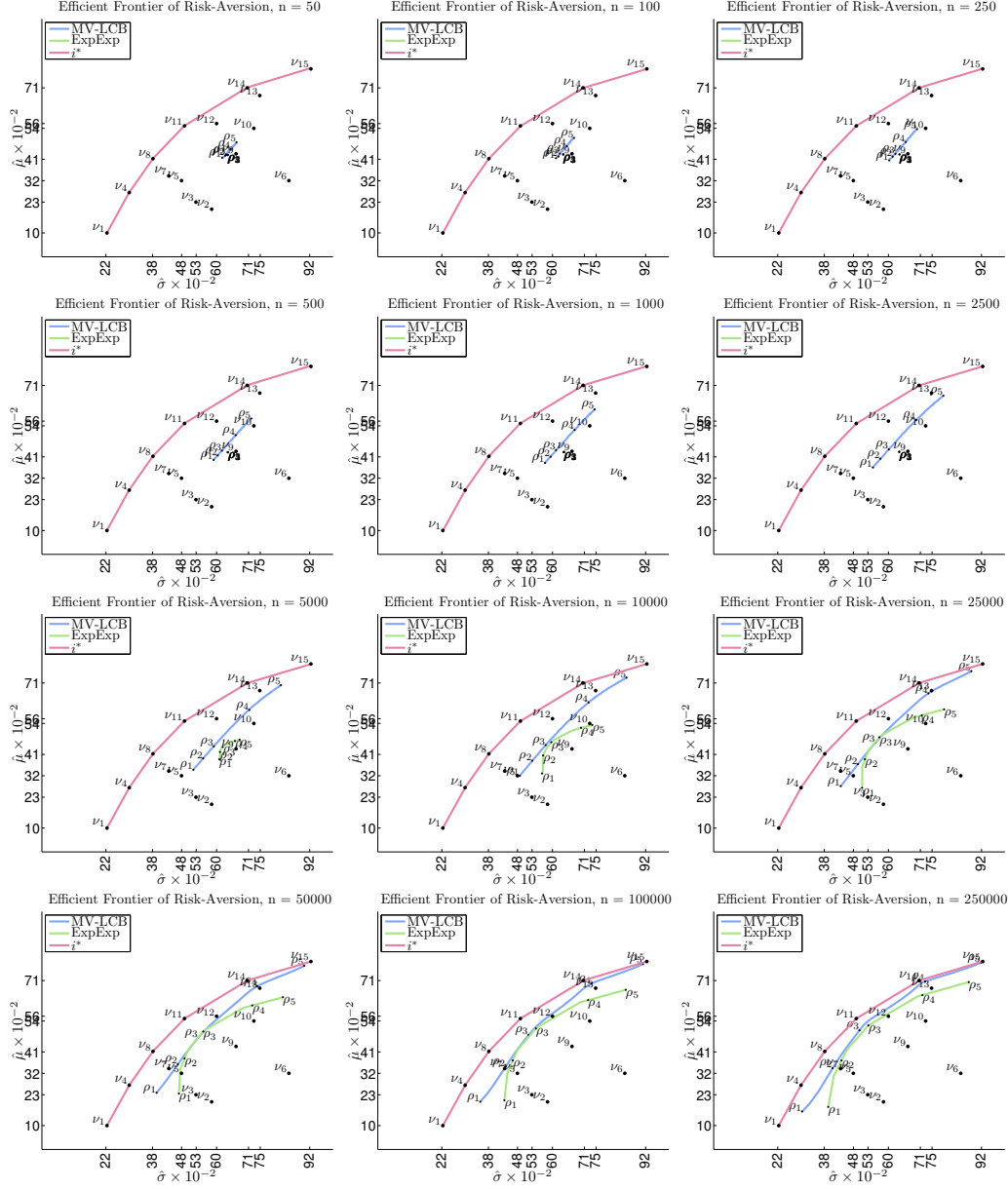


Figure 6: Risk tolerance sensitivity of MV-LCB and ExpExp for Configuration 1.

nite time  $n$ , where each point represents the resulting mean–standard deviation of the sequence of pulls on the arms by the algorithm with that specific value of  $\rho$ . The gap between the  $\rho$  specific performance of the algorithm and the corresponding optimal arm along the red frontier represents the regret for the specific  $\rho$  value. Accordingly, the gap between the algorithm performance curves represents the gap in performance with regard to MV-LCB versus ExpExp. Where a lot of arms have big gaps (e.g., all the dominated arms have a large gap for any value of  $\rho$ ), MV-LCB tends to perform better than ExpExp. The series of plots represent increasing values of  $n$  and demonstrate the relative algorithm performance versus the optimal red frontier. The set of plots represent the two settings reported in Figure 5. We chose the values of the arms so as to have configurations with different complexities. In particular, configuration 1 corresponds to “easy” problems for MV-LCB since the arms all have quite large gaps (for different values of  $\rho$ ) and this should allow it to perform well. On the other hand, the second configuration has much smaller gaps and, thus, higher complexity. According to the bounds for MV-LCB we know that a good proxy for its learning complexity is represented by the term  $\sum_i 1/\Delta_{i,\rho}^2$ . In Figure 5 we report such complexity for different values of  $\rho$  and, as it can be noticed, configuration 2 has always a higher complexity than configuration 1.

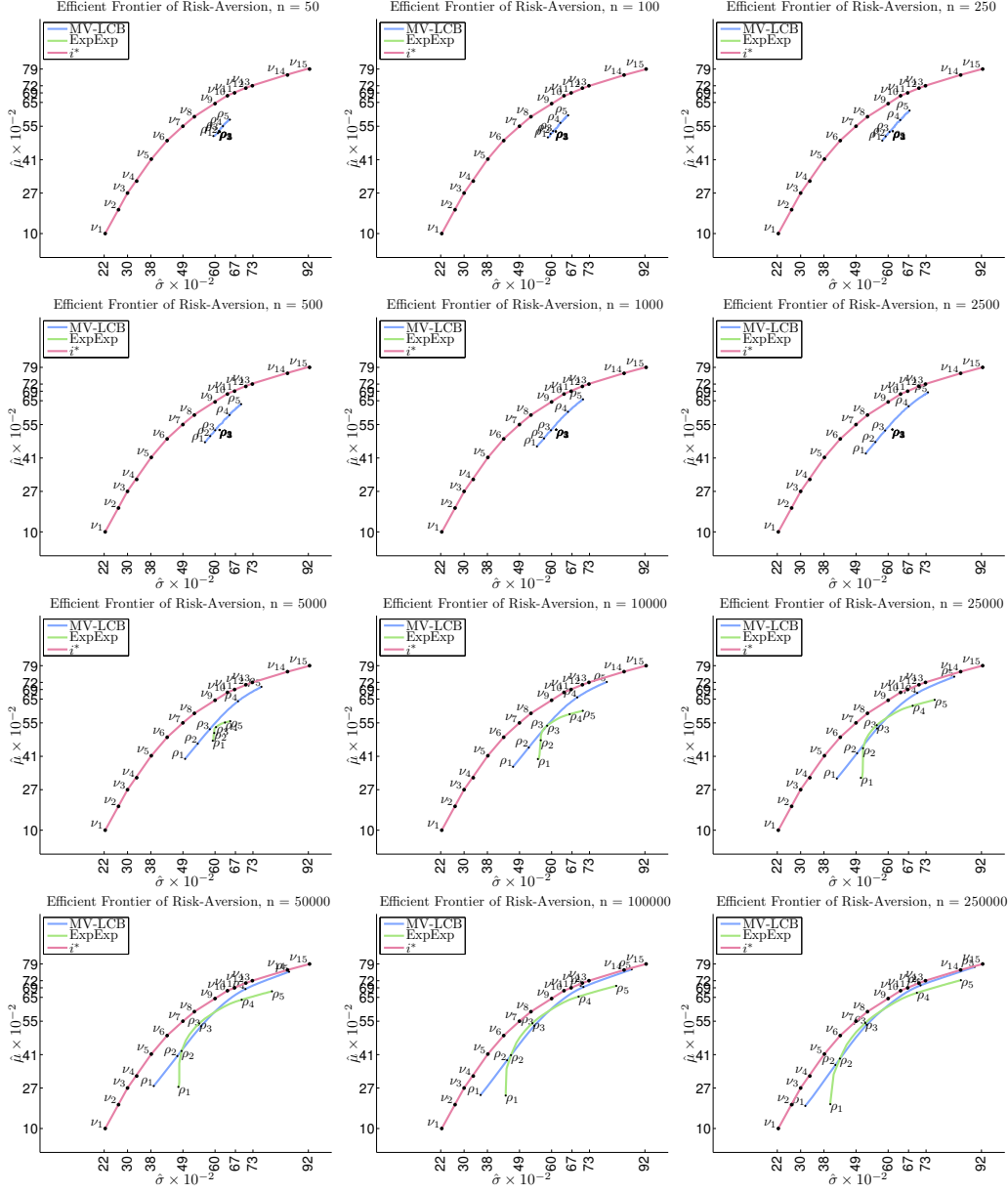


Figure 7: Risk tolerance sensitivity of MV-LCB and ExpExp for Configuration 2.

As we notice, in both configurations the performance of MV-LCB and ExpExp approach one of the optimal arms  $i^*_\rho$  for each specific  $\rho$  as  $n$  increases. Nonetheless, in configuration 1 the large number of suboptimal arms (e.g., arms with large gaps) allows MV-LCB to outperform ExpExp and converge faster to the optimal arm (and thus zero regret). On the other hand, in configuration 2 there are more arms with similar performance and for some values of  $\rho$  ExpExp eventually achieves better performance than MV-LCB.