
Finite-Time Analysis of Stratified Sampling for Monte Carlo

Alexandra Carpentier
INRIA Lille - Nord Europe
alexandra.carpentier@inria.fr

Rémi Munos
INRIA Lille - Nord Europe
remi.munos@inria.fr

Abstract

We consider the problem of stratified sampling for Monte-Carlo integration. We model this problem in a multi-armed bandit setting, where the arms represent the strata, and the goal is to estimate a weighted average of the mean values of the arms. We propose a strategy that samples the arms according to an upper bound on their standard deviations and compare its estimation quality to an ideal allocation that would know the standard deviations of the strata. We provide two regret analyses: a distribution-dependent bound $\tilde{O}(n^{-3/2})$ that depends on a measure of the disparity of the strata, and a distribution-free bound $\tilde{O}(n^{-4/3})$ that does not.

1 Introduction

Consider a polling institute that has to estimate as accurately as possible the average income of a country, given a finite budget for polls. The institute has call centers in every region in the country, and gives a part of the total sampling budget to each center so that they can call random people in the area and ask about their income. A naive method would allocate a budget proportionally to the number of people in each area. However some regions show a high variability in the income of their inhabitants whereas others are very homogeneous. Now if the polling institute knew the level of variability within each region, it could adjust the budget allocated to each region in a more clever way (allocating more polls to regions with high variability) in order to reduce the final estimation error.

This example is just one of many for which an efficient method of sampling a function with natural strata (i.e., the regions) is of great interest. Note that even in the case that there are no natural strata, it is always a good strategy to design arbitrary strata and allocate a budget to each stratum that is proportional to the size of the stratum, compared to a crude Monte-Carlo. There are many good surveys on the topic of stratified sampling for Monte-Carlo, such as (Rubinstein and Kroese, 2008)[Subsection 5.5] or (Glasserman, 2004).

The main problem for performing an efficient sampling is that the variances within the strata (in the previous example, the income variability per region) are usually unknown. One possibility is to estimate the variances *online* while sampling the strata. There is some interesting research along this direction, such as (Arouna, 2004) and more recently (Etoré and Jourdain, 2010, Kawai, 2010). The work of Etoré and Jourdain (2010) matches exactly our problem of designing an efficient adaptive sampling strategy. In this article they propose to sample according to an empirical estimate of the variance of the strata, whereas Kawai (2010) addresses a computational complexity problem which is slightly different from ours. The recent work of Etoré et al. (2011) describes a strategy that enables to sample *asymptotically* according to the (unknown) standard deviations of the strata and at the same time adapts the shape (and number) of the strata online. This is a very difficult problem, especially in high dimension, that we will not address here, although we think this is a very interesting and promising direction for further researches.

These works provide asymptotic convergence of the variance of the estimate to the targeted stratified variance¹ divided by the sample size. They also prove that the number of pulls within each stratum converges to the desired number of pulls i.e. the optimal allocation if the variances per stratum were known. Like Etoré and Jourdain (2010), we consider a stratified Monte-Carlo setting with fixed strata. Our contribution is to design a sampling strategy for which we can derive a finite-time analysis (where 'time' refers to the number of samples). This enables us to predict the quality of our estimate for any given budget n .

We model this problem using the setting of multi-armed bandits where our goal is to estimate a weighted average of the mean values of the arms. Although our goal is different from a usual bandit problem where the objective is to play the best arm as often as possible, this problem also exhibits an *exploration-exploitation trade-off*. The arms have to be pulled both in order to estimate the initially unknown variability of the arms (exploration) and to allocate correctly the budget according to our current knowledge of the variability (exploitation).

Our setting is close to the one described in (Antos et al., 2010) which aims at estimating *uniformly well* the mean values of all the arms. The authors present an algorithm, called GAFS-MAX, that allocates samples proportionally to the empirical variance of the arms, while imposing that each arm is pulled at least \sqrt{n} times to guarantee a sufficiently good estimation of the true variances.

Note though that in the Master Thesis (Grover, 2009), the author presents an algorithm named GAFS-WL which is similar to GAFS-MAX and has an analysis close to the one of GAFS-MAX. It deals with stratified sampling, i.e. it targets an allocation which is proportional to the standard deviation (and not to the variance) of the strata times their size². Some questions remain open in this work, notably that no distribution independent regret bound is provided for GAFS-WL. We clarify this point in Section 4. Our objective is similar, and we extend the analysis of this setting.

Contributions: In this paper, we introduce a new algorithm based on Upper-Confidence-Bounds (UCB) on the standard deviation. They are computed from the empirical standard deviation and a confidence interval derived from Bernstein's inequalities. We provide a finite-time analysis of its performance. The algorithm, called MC-UCB, samples the arms proportionally to an UCB³ on the standard deviation times the size of the stratum. Note that the idea is similar to the one in (Carpentier et al., 2011). Our contributions are the following:

- We derive a *finite-time analysis* for the stratified sampling for Monte-Carlo setting by using an algorithm based on upper confidence bounds. We show how such a family of algorithm is particularly interesting in this setting.
- We provide two regret analysis: (i) a distribution-dependent bound $\tilde{O}(n^{-3/2})^4$ that depends on the disparity of the stratas (a measure of the problem complexity), and which corresponds to a stationary regime where the budget n is large compared to this complexity. (ii) A distribution-free bound $\tilde{O}(n^{-4/3})$ that does not depend on the the disparity of the stratas, and corresponds to a transitory regime where n is small compared to the complexity. The characterization of those two regimes and the fact that the corresponding excess error rates differ enlightens the fact that a finite-time analysis is very relevant for this problem.

The rest of the paper is organized as follows. In Section 2 we formalize the problem and introduce the notations used throughout the paper. Section 3 introduces the MC-UCB algorithm and reports performance bounds. We then discuss in Section 4 about the parameters of the algorithm and its performances. In Section 5 we report numerical experiments that

¹The target is defined in [Subsection 5.5] of (Rubinstein and Kroese, 2008) and later in this paper, see Equation 4.

²This is explained in (Rubinstein and Kroese, 2008) and will be formulated precisely later.

³Note that we consider a sampling strategy based on UCBs on the standard deviations of the arms whereas the so-called *UCB algorithm* of Auer et al. (2002), in the usual multi-armed bandit setting, computes UCBs on the mean rewards of the arms.

⁴The notation $\tilde{O}(\cdot)$ corresponds to $O(\cdot)$ up to logarithmic factors.

illustrate our method on the problem of pricing Asian options as introduced in (Glasserman et al., 1999). Finally, Section 6 concludes the paper and suggests future works.

2 Preliminaries

The allocation problem mentioned in the previous section is formalized as a K -armed bandit problem where each arm (stratum) $k = 1, \dots, K$ is characterized by a distribution ν_k with mean value μ_k and variance σ_k^2 . At each round $t \geq 1$, an allocation strategy (or algorithm) \mathcal{A} selects an arm k_t and receives a sample drawn from ν_{k_t} independently of the past samples. Note that a strategy may be adaptive, i.e., the arm selected at round t may depend on past observed samples. Let $\{w_k\}_{k=1, \dots, K}$ denote a known set of positive weights which sum to 1. For example in the setting of stratified sampling for Monte-Carlo, this would be the probability mass in each stratum. The goal is to define a strategy that estimates as precisely as possible $\mu = \sum_{k=1}^K w_k \mu_k$ using a total budget of n samples.

Let us write $T_{k,t} = \sum_{s=1}^t \mathbb{I}\{k_s = k\}$ the number of times arm k has been pulled up to time t , and $\hat{\mu}_{k,t} = \frac{1}{T_{k,t}} \sum_{s=1}^{T_{k,t}} X_{k,s}$ the empirical estimate of the mean μ_k at time t , where $X_{k,s}$ denotes the sample received when pulling arm k for the s -th time.

After n rounds, the algorithm \mathcal{A} returns the empirical estimate $\hat{\mu}_{k,n}$ of all the arms. Note that in the case of a deterministic strategy, the expected quadratic estimation error of the weighted mean μ as estimated by the weighted average $\hat{\mu}_n = \sum_{k=1}^K w_k \hat{\mu}_{k,n}$ satisfies:

$$\mathbb{E}\left[(\hat{\mu}_n - \mu)^2\right] = \mathbb{E}\left[\left(\sum_{k=1}^K w_k (\hat{\mu}_{k,n} - \mu_k)\right)^2\right] = \sum_{k=1}^K w_k^2 \mathbb{E}_{\nu_k}\left[(\hat{\mu}_{k,n} - \mu_k)^2\right].$$

We thus use the following measure for the performance of any algorithm \mathcal{A} :

$$L_n(\mathcal{A}) = \sum_{k=1}^K w_k^2 \mathbb{E}\left[(\mu_k - \hat{\mu}_{k,n})^2\right]. \quad (1)$$

The goal is to define an allocation strategy that minimizes the global loss defined in Equation 1. If the variance of the arms were known in advance, one could design an optimal static⁵ allocation strategy \mathcal{A}^* by pulling each arm k proportionally to the quantity $w_k \sigma_k$. Indeed, if arm k is pulled a deterministic number of times $T_{k,n}^*$, then

$$L_n(\mathcal{A}^*) = \sum_{k=1}^K w_k^2 \frac{\sigma_k^2}{T_{k,n}^*}. \quad (2)$$

By choosing $T_{k,n}^*$ such as to minimize L_n under the constraint that $\sum_{k=1}^K T_{k,n}^* = n$, the optimal static allocation (up to rounding effects) of algorithm \mathcal{A}^* is to pull each arm k ,

$$T_{k,n}^* = \frac{w_k \sigma_k}{\sum_{i=1}^K w_i \sigma_i} n, \quad (3)$$

times, and achieves a global performance

$$L_n(\mathcal{A}^*) = \frac{\Sigma_w^2}{n}, \quad (4)$$

where $\Sigma_w = \sum_{i=1}^K w_i \sigma_i$. In the following, we write $\lambda_k = \frac{T_{k,n}^*}{n} = \frac{w_k \sigma_k}{\Sigma_w}$ the optimal allocation proportion for arm k and $\lambda_{\min} = \min_{1 \leq k \leq K} \lambda_k$. Note that a small λ_{\min} means a large disparity of the $w_k \sigma_k$ and, as explained later, provides for the algorithm we build in Section 3 a characterization of the hardness of a problem.

However, in the setting considered here, the σ_k are unknown, and thus the optimal allocation is out of reach. A possible allocation is the uniform strategy \mathcal{A}^u , i.e., such that $T_k^u = \frac{w_k}{\sum_{i=1}^K w_i} n$. Its performance is

$$L_n(\mathcal{A}^u) = \sum_{k=1}^K w_k \sum_{k=1}^K \frac{w_k \sigma_k^2}{n} = \frac{\Sigma_{w,2}}{n},$$

⁵Static means that the number of pulls allocated to each arm does not depend on the received samples.

where $\Sigma_{w,2} = \sum_{k=1}^K w_k \sigma_k^2$. Note that by Cauchy-Schwartz's inequality, we have $\Sigma_w^2 \leq \Sigma_{w,2}$ with equality if and only if the (σ_k) are all equal. Thus \mathcal{A}^* is always at least as good as \mathcal{A}^u . In addition, since $\sum_i w_i = 1$, we have $\Sigma_w^2 - \Sigma_{w,2} = -\sum_k w_k (\sigma_k - \Sigma_w)^2$. The difference between those two quantities is the weighted quadratic variation of the σ_k around their weighted mean Σ_w . In other words, it is the variance of the $(\sigma_k)_{1 \leq k \leq K}$. As a result the gain of \mathcal{A}^* compared to \mathcal{A}^u grow with the disparity of the σ_k .

We would like to do better than the uniform strategy by considering an adaptive strategy \mathcal{A} that would estimate the σ_k at the same time as it tries to implement an allocation strategy as close as possible to the optimal allocation algorithm \mathcal{A}^* . This introduces a natural trade-off between the exploration needed to improve the estimates of the variances and the exploitation of the current estimates to allocate the pulls nearly-optimally.

In order to assess how well \mathcal{A} solves this trade-off and manages to sample according to the true standard deviations *without knowing them in advance*, we compare its performance to that of the optimal allocation strategy \mathcal{A}^* . For this purpose we define the notion of *regret* of an adaptive algorithm \mathcal{A} as the difference between the performance loss incurred by the algorithm and the optimal algorithm:

$$R_n(\mathcal{A}) = L_n(\mathcal{A}) - L_n(\mathcal{A}^*). \quad (5)$$

The *regret* indicates how much we loose in terms of expected quadratic estimation error by not knowing in advance the standard deviations (σ_k) . Note that since $L_n(\mathcal{A}^*) = \frac{\Sigma_w^2}{n}$, a consistent strategy i.e., asymptotically equivalent to the optimal strategy, is obtained whenever its regret is neglectable compared to $1/n$.

3 Allocation based on Monte Carlo Upper Confidence Bound

3.1 The algorithm

In this section, we introduce our adaptive algorithm for the allocation problem, called *Monte Carlo Upper Confidence Bound* (MC-UCB). The algorithm computes a high-probability bound on the standard deviation of each arm and samples the arms proportionally to their bounds times the corresponding weights. The MC-UCB algorithm, \mathcal{A}_{MC-UCB} , is described in Figure 1. It requires three parameters as inputs: c_1 and c_2 which are related to the shape of the distributions (see Assumption 1), and δ which defines the *confidence level* of the bound. In Subsection 4.2, we discuss a way to reduce the number of parameters from three to one. The amount of exploration of the algorithm can be adapted by properly tuning these parameters.

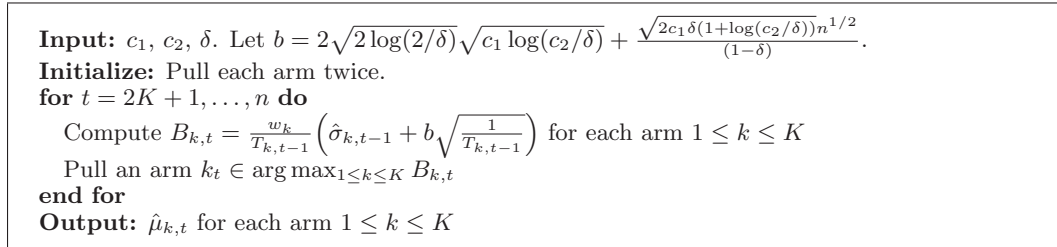


Figure 1: The pseudo-code of the MC-UCB algorithm. The empirical standard deviations $\hat{\sigma}_{k,t-1}$ are computed using Equation 6.

The algorithm starts by pulling each arm twice in rounds $t = 1$ to $2K$. From round $t = 2K+1$ on, it computes an upper confidence bound $B_{k,t}$ on the standard deviation σ_k , for each arm k , and then pulls the one with largest $B_{k,t}$. The upper bounds on the standard deviations are built by using Theorem 10 in (Maurer and Pontil, 2009)⁶ and based on the empirical standard deviation $\hat{\sigma}_{k,t-1}$:

$$\hat{\sigma}_{k,t-1}^2 = \frac{1}{T_{k,t-1} - 1} \sum_{i=1}^{T_{k,t-1}} (X_{k,i} - \hat{\mu}_{k,t-1})^2, \quad (6)$$

⁶We could also have used the variant reported in (Audibert et al., 2009).

where $X_{k,i}$ is the i -th sample received when pulling arm k , and $T_{k,t-1}$ is the number of pulls allocated to arm k up to time $t-1$. After n rounds, MC-UCB returns the empirical mean $\hat{\mu}_{k,n}$ for each arm $1 \leq k \leq K$.

3.2 Regret analysis of MC-UCB

Before stating the main results of this section, we state the assumption that the distributions are sub-Gaussian, which includes e.g., Gaussian or bounded distributions. See (Buldygin and Kozachenko, 1980) for more precisions.

Assumption 1 *There exist $c_1, c_2 > 0$ such that for all $1 \leq k \leq K$ and any $\epsilon > 0$,*

$$\mathbb{P}_{X \sim \nu_k}(|X - \mu_k| \geq \epsilon) \leq c_2 \exp(-\epsilon^2/c_1). \quad (7)$$

We provide two analyses, a *distribution-dependent* and a *distribution-free*, of MC-UCB, which are respectively interesting in two *regimes*, i.e., stationary and transitory *regimes*, of the algorithm. We will comment on this later in Section 4.

A *distribution-dependent* result: We now report the first bound on the regret of MC-UCB algorithm. The proof is reported in (Carpentier and Munos, 2011). and relies on upper- and lower-bounds on $T_{k,t} - T_{k,t}^*$, i.e., the difference in the number of pulls of each arm compared to the optimal allocation (see Lemma 3).

Theorem 1 *Under Assumption 1 and if we choose c_2 such that $c_2 \geq 2Kn^{-5/2}$, the regret of MC-UCB run with parameter $\delta = n^{-7/2}$ with $n \geq 4K$ is bounded as*

$$R_n(\mathcal{A}_{MC-UCB}) \leq \frac{\log(n)c_1(c_2+2)}{n^{3/2}\lambda_{\min}^{3/2}} \left(112\Sigma_w + 6K\right) + \frac{19}{\lambda_{\min}^3 n^2} \left(K\Sigma_w^2 + 720c_1(c_2+1)\log(n)^2\right).$$

Note that this result crucially depends on the smallest proportion λ_{\min} which is a measure of the disparity of the standard deviations times their weight. For this reason we refer to it as “distribution-dependent” result.

A *distribution-free* result: Now we report our second regret bound that does not depend on λ_{\min} but whose rate is poorer. The proof is reported in (Carpentier and Munos, 2011) and relies on other upper- and lower-bounds on $T_{k,t} - T_{k,t}^*$ detailed in Lemma 4.

Theorem 2 *Under Assumption 1 and if we choose c_2 such that $c_2 \geq 2Kn^{-5/2}$, the regret of MC-UCB run with parameter $\delta = n^{-7/2}$ with $n \geq 4K$ is bounded as*

$$R_n(\mathcal{A}_{MC-UCB}) \leq \frac{200\sqrt{c_1}(c_2+2)\Sigma_w K}{n^{4/3}} \log(n) + \frac{365}{n^{3/2}} \left(129c_1(c_2+2)^2 K^2 \log(n)^2 + K\Sigma_w^2\right).$$

This bound does not depend on $1/\lambda_{\min}$. Note that the bound is not entirely distribution free since Σ_w appears. But it can be proved using Assumption 1 that $\Sigma_w^2 \leq c_1 c_2$. This is obtained at the price of the slightly worse rate $\tilde{O}(n^{-4/3})$.

4 Discussion on the results

4.1 Distribution-free versus distribution-dependent

Theorem 1 provides a regret bound of order $\tilde{O}(\lambda_{\min}^{-5/2} n^{-3/2})$, whereas Theorem 2 provides a bound in $\tilde{O}(n^{-4/3})$ independently of λ_{\min} . Hence, for a given problem i.e., a given λ_{\min} , the distribution-free result of Theorem 2 is more informative than the distribution-dependent result of Theorem 1 in the *transitory regime*, that is to say when n is small compared to λ_{\min}^{-1} . The distribution-dependent result of Theorem 1 is better in the *stationary regime* i.e., for large n . This distinction reminds us of the difference between distribution-dependent and distribution-free bounds for the UCB algorithm in usual multi-armed bandits⁷.

⁷The distribution dependent bound is in $O(K \log n / \Delta)$, where Δ is the difference between the mean value of the two best arms, and the distribution-free bound is in $O(\sqrt{nK \log n})$ as explained in (Auer et al., 2002, Audibert and Bubeck, 2009).

Although we do not have a lower bound on the regret yet, we believe that the rate $n^{-3/2}$ cannot be improved for general distributions. As explained in the proof in Appendix B of (Carpentier and Munos, 2011), this rate is a direct consequence of the high probability bounds on the estimates of the standard deviations of the arms which are in $O(1/\sqrt{n})$, *and those bounds are tight*. A natural question is whether there exists an algorithm with a regret of order $\tilde{O}(n^{-3/2})$ without any dependence in λ_{\min}^{-1} . Although we do not have an answer to this question, we can say that our algorithm MC-UCB does not satisfy this property. In Appendix D.1 of (Carpentier and Munos, 2011), we give a simple example where $\lambda_{\min} = 0$ and for which the rate of MC-UCB cannot be better than $\tilde{O}(n^{-4/3})$. This shows that our analysis of MC-UCB is tight.

The problem dependent upper bound is similar to the one provided for GAFS-WL in (Grover, 2009). We however expect that GAFS-WL has for some problems a sub-optimal behavior: it is possible to find cases where $R_n(\mathcal{A}_{GAFS-WL}) = \Omega(1/n)$, see Appendix D.1 of (Carpentier and Munos, 2011). Note however that when there is an arm with 0 standard deviation, GAFS-WL is likely to perform better than MC-UCB, as it will only sample this arm $O(\sqrt{n})$ times while MC-UCB samples it $\tilde{O}(n^{2/3})$ times.

4.2 The parameters of the algorithm

Our algorithm takes three parameters as input, namely c_1 , c_2 and δ , but we only use a combination of them in the algorithm, with the introduction of $b = 2\sqrt{2\log(2/\delta)}\sqrt{c_1\log(c_2/\delta)} + \frac{\sqrt{2c_1\delta(1+\log(c_2/\delta))}n^{1/2}}{(1-\delta)}$. For practical use of the method, it is enough to tune the algorithm with a single parameter b . By the choice of the value assigned to δ in the two theorems, b should be chosen of order $c\log(n)$, where c can be interpreted as a high probability bound on the range of the samples. We thus simply require a rough estimate of the magnitude of the samples. Note that in the case of bounded distributions, b can be chosen as $b = 4\sqrt{\frac{5}{2}}c\sqrt{\log(n)}$ where c is a true bound on the variables. This result is easy to deduce by simplifying Lemma 1 in Appendix A of (Carpentier and Munos, 2011) for the case of bounded variables.

5 Numerical experiment: Pricing of an Asian option

We consider the pricing problem of an Asian option introduced in (Glasserman et al., 1999) and later considered in (Kawai, 2010, Etoré and Jourdain, 2010). This uses a Black-Schole model with strike C and maturity T . Let $(W(t))_{0 \leq t \leq 1}$ be a Brownian motion that is discretized at d equidistant times $\{i/d\}_{1 \leq i \leq d}$, which defines the vector $W \in \mathbb{R}^d$ with components $W_i = W(i/d)$. The discounted payoff of the Asian option is defined as a function of W , by:

$$F(W) = \exp(-rT) \max \left[\frac{1}{d} \sum_{i=1}^d S_0 \exp \left[\left(r - \frac{1}{2} s_0^2 \right) \frac{iT}{d} + s_0 \sqrt{T} W_i \right] - C, 0 \right], \quad (8)$$

where S_0 , r , and s_0 are constants, and the price is defined by the expectation $p = \mathbb{E}_W F(W)$.

We want to estimate the price p by Monte-Carlo simulations (by sampling on $W = (W_i)_{1 \leq i \leq d}$). In order to reduce the variance of the estimated price, we can stratify the space of W . Glasserman et al. (1999) suggest to stratify according to a one dimensional projection of W , i.e., by choosing a projection vector $u \in \mathbb{R}^d$ and define the strata as the set of W such that $u \cdot W$ lies in intervals of \mathbb{R} . They further argue that the best direction for stratification is to choose $u = (0, \dots, 0, 1)$, i.e., to stratify according to the last component W_d of W . Thus we sample W_d and then conditionally sample W_1, \dots, W_{d-1} according to a Brownian Bridge as explained in (Kawai, 2010). Note that this choice of stratification is also intuitive since W_d has the largest exponent in the payoff (8), and thus the highest volatility. Kawai (2010) and Etoré and Jourdain (2010) also use the same direction of stratification.

Like in (Kawai, 2010) we consider 5 strata of equal weight. Since W_d follows a $\mathcal{N}(0, 1)$, the strata correspond to the 20-percentile of a normal distribution. The left plot of Figure 2 represents the cumulative distribution function of W_d and shows the strata in terms of

percentiles of W_d . The right plot represents, in dot line, the curve $\mathbb{E}[F(W)|W_d = x]$ versus $\mathbb{P}(W_d < x)$ parameterized by x , and the box plot represents the expectation and standard deviations of $F(W)$ conditioned on each stratum. We observe that this stratification produces an important heterogeneity of the standard deviations per stratum, which indicates that a stratified sampling would be profitable compared to a crude Monte-Carlo sampling.

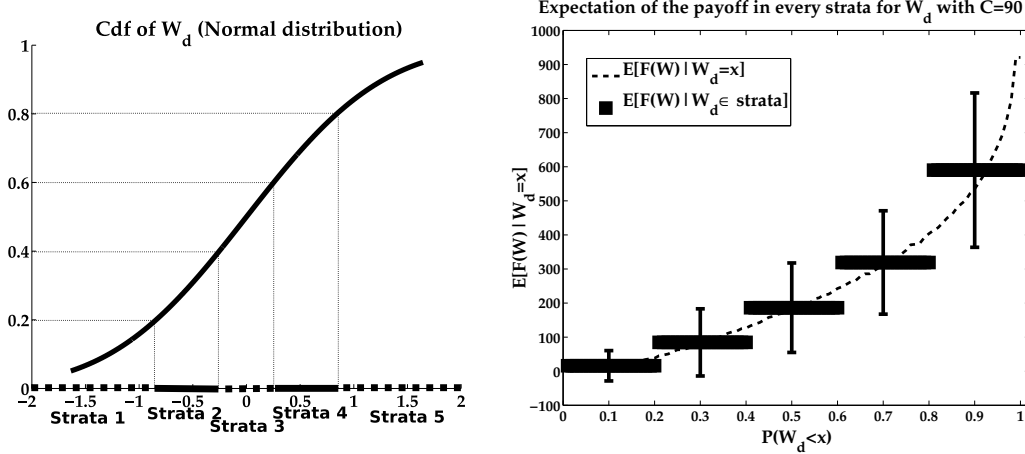


Figure 2: Left: Cdf of W_d and the definition of the strata. Right: expectation and standard deviation of $F(W)$ conditioned on each stratum for a strike $C = 90$.

We choose the same numerical values as Kawai (2010): $S_0 = 100$, $r = 0.05$, $s_0 = 0.30$, $T = 1$ and $d = 16$. Note that the strike C of the option has a direct impact on the variability of the strata. Indeed, the larger C , the more probable $F(W) = 0$ for strata with small W_d , and thus, the smaller λ_{\min} .

Our two main competitors are the SSAA algorithm of Etoré and Jourdain (2010) and GAFS-WL of Grover (2009). We did not compare to (Kawai, 2010) which aims at minimizing the computational time and not the loss considered here⁸. SSAA works in K_r rounds of length N_k where, at each round, it allocates proportionally to the empirical standard deviations computed in the previous rounds. Etoré and Jourdain (2010) report the asymptotic consistency of the algorithm whenever $\frac{k}{N_k}$ goes to 0 when k goes to infinity. Since their goal is not to obtain a finite-time performance, they do not mention how to calibrate the length and number of rounds in practice. We choose the same parameters as in their numerical experiments (Section 3.2.2 of (Etoré and Jourdain, 2010)) using 3 rounds. In this setting where we know the budget n at the beginning of the algorithm, GAFS-WL pulls each arm $a\sqrt{n}$ times and then pulls at time $t + 1$ the arm k_{t+1} that maximizes $\frac{w_k \hat{\sigma}_{k,t}}{T_{k,t}}$. We set $a = 1$.

As mentioned in Subsection 4.2, an advantage of our algorithm is that it requires a single parameter to tune. We chose $b = 1000 \log(n)$ where 1000 is a high-probability range of the variables (see right plot of Figure 2). Table 5 reports the performance of MC-UCB, GAFS-WL, SSAA, and the uniform strategy, for different values of strike C i.e., for different values of λ_{\min}^{-1} and $\Sigma_{w,2}/\Sigma_w^2 = \frac{\sum_k w_k \sigma_k^2}{(\sum_k w_k \sigma_k)^2}$. The total budget is $n = 10^5$. The results are averaged on 50000 trials. We notice that MC-UCB outperforms SSAA, the uniform strategy, and GAFS-WL strategy. Note however that, in the case of GAFS-WL strategy, the small gain could come from the fact that there are more parameters in MC-UCB, and that we were thus able to adjust them (even if we kept the same parameters for the three values of C).

In the left plot of Figure 3, we plot the rescaled regret $R_n n^{3/2}$, averaged over 50000 trials, as a function of n , where n ranges from 50 to 5000. The value of the strike is $C = 120$. Again, we notice that MC-UCB performs better than Uniform and SSAA because it adapts

⁸In that paper, the computational costs for each stratum vary, i.e. it is faster to sample in some strata than in others, and the aim of their paper is to minimize the global computational cost while achieving a given performance.

C	$\frac{1}{\lambda_{\min}}$	$\Sigma_{w,2}/\Sigma_w^2$	Uniform	SSAA	GAFS-WL	MC-UCB
60	6.18	1.06	$2.52 \cdot 10^{-2}$	$5.87 \cdot 10^{-3}$	$8.25 \cdot 10^{-4}$	$7.29 \cdot 10^{-4}$
90	15.29	1.24	$3.32 \cdot 10^{-2}$	$6.14 \cdot 10^{-3}$	$8.58 \cdot 10^{-4}$	$8.07 \cdot 10^{-4}$
120	744.25	3.07	$3.56 \cdot 10^{-2}$	$6.22 \cdot 10^{-3}$	$9.89 \cdot 10^{-4}$	$9.28 \cdot 10^{-4}$

Table 1: Characteristics of the distributions (λ_{\min}^{-1} and $\Sigma_{w,2}/\Sigma_w^2$) and regret of the Uniform, SSAA, and MC-UCB strategies, for different values of the strike C .

faster to the distributions of the strata. But it performs very similarly to GAFS-WL. In addition, it seems that the regret of Uniform and SSAA grows faster than the rate $n^{3/2}$, whereas MC-UCB, as well as GAFS-WL, grow with this rate. The right plot focuses on the MC-UCB algorithm and rescales the y -axis to observe the variations of its rescaled regret more accurately. The curve grows first and then stabilizes. This could correspond to the two regimes discussed previously.

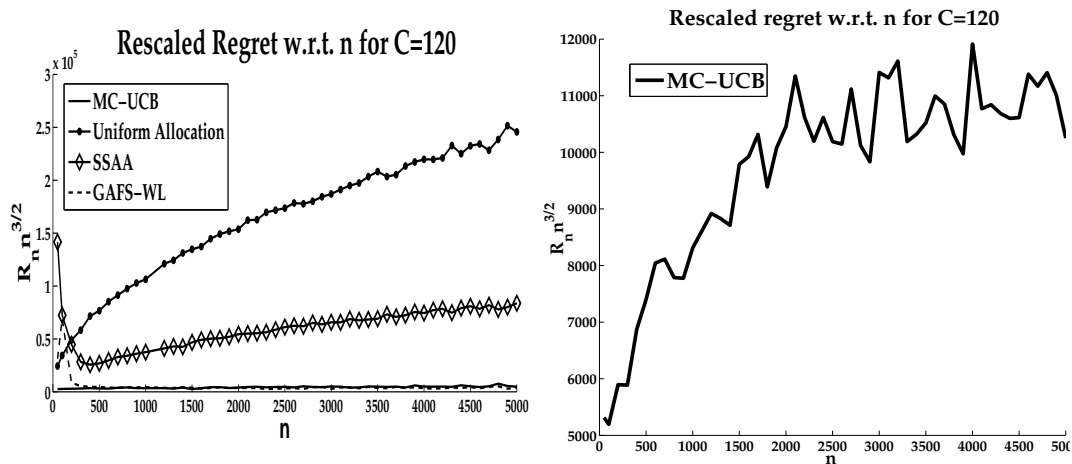


Figure 3: Left: Rescaled regret ($R_n n^{3/2}$) of the Uniform, SSAA, and MC-UCB strategies. Right: zoom on the rescaled regret for MC-UCB that illustrates the two regimes.

6 Conclusions

We provided a finite-time analysis for stratified sampling for Monte-Carlo in the case of fixed strata. We reported two bounds: (i) a distribution dependent bound $\tilde{O}(n^{-3/2}\lambda_{\min}^{-5/2})$ which is of interest when n is large compared to a measure of disparity λ_{\min}^{-1} of the standard deviations (*stationary regime*), and (ii) a distribution free bound in $\tilde{O}(n^{-4/3})$ which is of interest when n is small compared to λ_{\min}^{-1} (*transitory regime*).

Possible directions for future work include: (i) making the MC-UCB algorithm anytime (i.e. not requiring the knowledge of n), (ii) investigating whether there exists an algorithm with $\tilde{O}(n^{-3/2})$ regret without dependency on λ_{\min}^{-1} , and (iii) deriving distribution-dependent and distribution-free lower-bounds for this problem.

Acknowledgements

We thank András Antos for several comments that helped us to improve the quality of the paper. This research was partially supported by Region Nord-Pas-de-Calais Regional Council, French ANR EXPLO-RA (ANR-08-COSI-004), the European Communitys Seventh Framework Programme (FP7/2007-2013) under grant agreement 231495 (project ComplACS), and by Pascal-2.

References

- András Antos, Varun Grover, and Csaba Szepesvári. Active learning in heteroscedastic noise. *Theoretical Computer Science*, 411:2712–2728, June 2010.
- B. Arouna. Adaptive monte carlo method, a variance reduction technique. *Monte Carlo Methods and Applications*, 10(1):1–24, 2004.
- J.Y. Audibert and S. Bubeck. Minimax policies for adversarial and stochastic bandits. In *22nd annual conference on learning theory*, 2009.
- J.Y. Audibert, R. Munos, and Cs. Szepesvári. Exploration-exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19):1876–1902, 2009.
- P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.
- V.V. Buldygin and Y.V. Kozachenko. Sub-gaussian random variables. *Ukrainian Mathematical Journal*, 32(6):483–489, 1980.
- A. Carpentier and R. Munos. Finite-time analysis of stratified sampling for monte carlo. Technical Report inria-00636924, INRIA, 2011.
- A. Carpentier, A. Lazaric, M. Ghavamzadeh, R. Munos, and P. Auer. Upper-confidence-bound algorithms for active learning in multi-armed bandits. In *Algorithmic Learning Theory*, pages 189–203. Springer, 2011.
- Pierre Etoré and Benjamin Jourdain. Adaptive optimal allocation in stratified sampling methods. *Methodol. Comput. Appl. Probab.*, 12(3):335–360, September 2010.
- Pierre Etoré, Gersende Fort, Benjamin Jourdain, and Éric Moulines. On adaptive stratification. *Ann. Oper. Res.*, 2011. to appear.
- P. Glasserman. *Monte Carlo methods in financial engineering*. Springer Verlag, 2004. ISBN 0387004513.
- P. Glasserman, P. Heidelberger, and P. Shahabuddin. Asymptotically optimal importance sampling and stratification for pricing path-dependent options. *Mathematical Finance*, 9(2):117–152, 1999.
- V. Grover. Active learning and its application to heteroscedastic problems. *Department of Computing Science, Univ. of Alberta, MSc thesis*, 2009.
- R. Kawai. Asymptotically optimal allocation of stratified sampling with adaptive variance reduction by strata. *ACM Transactions on Modeling and Computer Simulation (TOMACS)*, 20(2):1–17, 2010. ISSN 1049-3301.
- A. Maurer and M. Pontil. Empirical bernstein bounds and sample-variance penalization. In *Proceedings of the Twenty-Second Annual Conference on Learning Theory*, pages 115–124, 2009.
- S.I. Resnick. *A probability path*. Birkhäuser, 1999.
- R.Y. Rubinstein and D.P. Kroese. *Simulation and the Monte Carlo method*. Wiley-interscience, 2008. ISBN 0470177942.

Supplementary material for the paper : Finite-Time Analysis of Stratified Sampling for Monte Carlo

A Main tools

A.1 The main tool: a high probability bound on the standard deviations

Upper bound on the standard deviation: The upper confidence bounds $B_{k,t}$ used in the MC-UCB algorithm is motivated by Theorem 10 in (Maurer and Pontil, 2009) (a variant of this result is also reported in (Audibert et al., 2009)). We extend this result to sub-Gaussian random variables.

Lemma 1 *Let Assumption 1 hold and $n \geq 2$. Define the following event*

$$\xi = \xi_{K,n}(\delta) = \bigcap_{1 \leq k \leq K, 2 \leq t \leq n} \left\{ \left| \sqrt{\frac{1}{t-1} \sum_{i=1}^t \left(X_{k,i} - \frac{1}{t} \sum_{j=1}^t X_{k,j} \right)^2} - \sigma_k \right| \leq 2a \sqrt{\frac{\log(2/\delta)}{t}} \right\}, \quad (9)$$

where $a = \sqrt{2c_1 \log(c_2/\delta)} + \frac{\sqrt{c_1 \delta(1+c_2+\log(c_2/\delta))}}{(1-\delta)\sqrt{2\log(2/\delta)}} n^{1/2}$. Then $\Pr(\xi) \geq 1 - 2nK\delta$.

Note that the first term in the absolute value in Equation 9 is the empirical standard deviation of arm k computed as in Equation 6 for t samples. The event ξ plays an important role in the proofs of this section and a number of statements will be proved on this event.

Proof:

Step 1. Truncating sub-Gaussian variables. We want to characterize the mean and variance of the variables $X_{k,t}$ given that $|X_{k,t} - \mu_k| \leq \sqrt{c_1 \log(c_2/\delta)}$. For any positive random variable Y and any $b \geq 0$, $\mathbb{E}(Y \mathbb{I}\{Y > b\}) = \int_b^\infty \mathbb{P}(Y > \epsilon) d\epsilon + b\mathbb{P}(Y > b)$. If we take $b = c_1 \log(c_2/\delta)$ and use Assumption 1, we obtain:

$$\begin{aligned} \mathbb{E}\left[|X_{k,t} - \mu_k|^2 \mathbb{I}\{|X_{k,t} - \mu_k|^2 > b\}\right] &= \int_b^{+\infty} \mathbb{P}(|X_{k,t} - \mu_k|^2 > \epsilon) d\epsilon + b\mathbb{P}(|X_{k,t} - \mu_k|^2 > b) \\ &\leq \int_b^{+\infty} c_2 \exp(-\epsilon/c_1) d\epsilon + bc_2 \exp(-b/c_1) \\ &\leq c_1 \delta + c_1 \log(c_2/\delta) \delta \\ &\leq c_1 \delta (1 + \log(c_2/\delta)). \end{aligned}$$

We have $\mathbb{E}\left[|X_{k,t} - \mu_k|^2 \mathbb{I}\{|X_{k,t} - \mu_k|^2 > b\}\right] + \mathbb{E}\left[|X_{k,t} - \mu_k|^2 \mathbb{I}\{|X_{k,t} - \mu_k|^2 \leq b\}\right] = \sigma_k^2$, which, combined with the previous equation, implies that

$$\begin{aligned} \left| \mathbb{E}\left[|X_{k,t} - \mu_k|^2 \mid |X_{k,t} - \mu_k|^2 \leq b\right] - \sigma_k^2 \right| &= \frac{\left| \mathbb{E}\left[\left((X_{k,t} - \mu_k)^2 - \sigma_k^2\right) \mathbb{I}\{|X_{k,t} - \mu_k|^2 > b\}\right] \right|}{\mathbb{P}\left(|X_{k,t} - \mu_k|^2 \leq b\right)} \\ &\leq \frac{c_1 \delta (1 + \log(c_2/\delta)) + \delta \sigma_k^2}{1 - \delta}. \end{aligned} \quad (10)$$

Note also that Cauchy-Schwartz inequality implies

$$\begin{aligned} \left| \mathbb{E} \left[\left(X_{k,t} - \mu_k \right) \mathbb{I} \{ |X_{k,t} - \mu_k|^2 > b \} \right] \right| &\leq \sqrt{\mathbb{E} \left[(X_{k,t} - \mu_k)^2 \mathbb{I} \{ |X_{k,t} - \mu_k|^2 > b \} \right]} \\ &\leq \sqrt{c_1 \delta (1 + \log(c_2/\delta))}. \end{aligned}$$

Now, notice that $\mathbb{E} \left[X_{k,t} \mathbb{I} \{ |X_{k,t} - \mu_k|^2 > b \} \right] + \mathbb{E} \left[X_{k,t} \mathbb{I} \{ |X_{k,t} - \mu_k|^2 \leq b \} \right] = \mu_k$, which, combined with the previous result and using $n \geq K \geq 2$, implies that

$$|\tilde{\mu}_k - \mu_k| = \frac{\left| \mathbb{E} \left[\left(X_{k,t} - \mu_k \right) \mathbb{I} \{ |X_{k,t} - \mu_k|^2 > b \} \right] \right|}{\mathbb{P} \left(|X_{k,t} - \mu_k|^2 \leq b \right)} \leq \frac{\sqrt{c_1 \delta (1 + \log(c_2/\delta))}}{1 - \delta}, \quad (11)$$

$$\text{where } \tilde{\mu}_k \stackrel{\text{def}}{=} \mathbb{E} \left[X_{k,t} \mid |X_{k,t} - \mu_k|^2 \leq b \right] = \frac{\mathbb{E} \left[X_{k,t} \mathbb{I} \{ |X_{k,t} - \mu_k|^2 \leq b \} \right]}{\mathbb{P} \left(|X_{k,t} - \mu_k|^2 \leq b \right)}.$$

We note $\tilde{\sigma}_k^2 \stackrel{\text{def}}{=} \mathbb{V} \left[X_{k,t} \mid |X_{k,t} - \mu_k|^2 \leq b \right] = \mathbb{E} \left[|X_{k,t} - \mu_k|^2 \mid |X_{k,t} - \mu_k|^2 \leq b \right] - (\mu_k - \tilde{\mu}_k)^2$. From Equations 10 and 11, we derive

$$\begin{aligned} |\tilde{\sigma}_k^2 - \sigma_k^2| &\leq \left| \mathbb{E} \left[|X_{k,t} - \mu_k|^2 \mid |X_{k,t} - \mu_k|^2 \leq b \right] - \sigma_k^2 \right| + |\tilde{\mu}_k - \mu_k|^2 \\ &\leq \frac{c_1 \delta (1 + \log(c_2/\delta)) + \delta \sigma_k^2}{1 - \delta} + \frac{c_1 \delta (1 + \log(c_2/\delta))}{(1 - \delta)^2} \\ &\leq \frac{2c_1 \delta (1 + \log(c_2/\delta)) + \delta \sigma_k^2}{(1 - \delta)^2}, \end{aligned}$$

from which we deduce, because $\sigma_k^2 \leq c_1 c_2$

$$|\tilde{\sigma}_k - \sigma_k| \leq \frac{\sqrt{2c_1 \delta (1 + c_2 + \log(c_2/\delta))}}{1 - \delta}. \quad (12)$$

Step 2. Application of large deviation inequalities.

Let $\xi_1 = \xi_{1,K,n}(\delta)$ be the event:

$$\xi_1 = \bigcap_{1 \leq k \leq K, 1 \leq t \leq n} \left\{ |X_{k,t} - \mu_k| \leq \sqrt{c_1 \log(c_2/\delta)} \right\}.$$

Under Assumption 1, using a union bound, we have that the probability of this event is at least $1 - nK\delta$.

We now recall Theorem 10 of (Maurer and Pontil, 2009):

Theorem 1 (Maurer and Pontil (2009)) *Let (X_1, \dots, X_t) be $t \geq 2$ i.i.d. random variables of variance σ^2 and mean μ and such that $\forall i \leq t, X_i \in [a, a + c]$. Then with probability at least $1 - \delta$:*

$$\left| \sqrt{\frac{1}{t-1} \sum_{i=1}^t \left(X_i - \frac{1}{t} \sum_{j=1}^t X_j \right)^2} - \sigma \right| \leq 2c \sqrt{\frac{\log(2/\delta)}{t-1}}.$$

On ξ_1 , the $\{X_{k,i}\}_{i, 1 \leq k \leq K, 1 \leq i \leq t}$ are t i.i.d. bounded random variables with standard deviation $\tilde{\sigma}_k$.

Let $\xi_2 = \xi_{2,K,n}(\delta)$ be the event:

$$\xi_2 = \bigcap_{1 \leq k \leq K, 1 \leq t \leq n} \left\{ \left| \sqrt{\frac{1}{t-1} \sum_{i=1}^t \left(X_{k,i} - \frac{1}{t} \sum_{j=1}^t X_{k,j} \right)^2} - \tilde{\sigma}_k \right| \leq 2\sqrt{c_1 \log(c_2/\delta)} \sqrt{\frac{\log(2/\delta)}{t-1}} \right\}.$$

Using Theorem 10 of (Maurer and Pontil, 2009) and a union bound, we deduce that $\Pr(\xi_1 \cap \xi_2) \geq 1 - 2nK\delta$.

Now, from Equation 12, we have on $\xi_1 \cap \xi_2$, for all $1 \leq k \leq K$, $2 \leq t \leq n$:

$$\begin{aligned} \left| \sqrt{\frac{1}{t-1} \sum_{i=1}^t \left(X_{k,i} - \frac{1}{t} \sum_{j=1}^t X_{k,j} \right)^2} - \sigma_k \right| &\leq 2\sqrt{c_1 \log(c_2/\delta)} \sqrt{\frac{\log(2/\delta)}{t-1}} \\ &\quad + \frac{\sqrt{2c_1\delta(1+c_2+\log(c_2/\delta))}}{1-\delta} \\ &\leq 2\sqrt{2c_1 \log(c_2/\delta)} \sqrt{\frac{\log(2/\delta)}{t}} \\ &\quad + \frac{\sqrt{2c_1\delta(1+c_2+\log(c_2/\delta))}}{1-\delta}, \end{aligned}$$

from which we deduce Lemma 1 (since $\xi_1 \cap \xi_2 \subseteq \xi$ and $2 \leq t \leq n$). \square

We deduce the following corollary when the number of samples $T_{k,t}$ are random.

Corollary 1 *For any $k = 1, \dots, K$ and $t = 2K, \dots, n$, let $\{X_{k,i}\}_i$ be n i.i.d. random variables drawn from ν_k , satisfying Assumption 1. Let $T_{k,t}$ be any random variable taking values in $\{2, \dots, n\}$. Let $\hat{\sigma}_{k,t}^2$ be the empirical variance computed from Equation 6. Then, on the event ξ , we have:*

$$|\hat{\sigma}_{k,t} - \sigma_k| \leq 2a \sqrt{\frac{\log(2/\delta)}{T_{k,t}}}. \quad (13)$$

A.2 Other important properties

A stopping time problem: We now draw a connection between the adaptive sampling and stopping time problems. We report the following proposition which is a type of Wald's Theorem for variance (see e.g. Resnick (1999)).

Proposition 1 *Let $\{\mathcal{F}_t\}$ be a filtration and X_t a \mathcal{F}_t -adapted sequence of i.i.d. random variables with variance σ^2 . Assume that \mathcal{F}_t and the σ -algebra generated by $\{X_i : i \geq t+1\}$ are independent and T is a stopping time w.r.t. \mathcal{F}_t with a finite expected value. If $\mathbb{E}[X_1^2] < \infty$ then*

$$\mathbb{E} \left[\left(\sum_{i=1}^T X_i - T \mu \right)^2 \right] = \mathbb{E}[T] \sigma^2. \quad (14)$$

Bound on the regret outside of ξ . The next lemma provides a bound for the loss whenever the event ξ does not hold.

Lemma 2 *Let Assumption 1 holds. Then for every arm k :*

$$\mathbb{E} \left[|\hat{\mu}_{k,n} - \mu_k|^2 \mathbb{I} \{ \xi^C \} \right] \leq 2c_1 n^2 K \delta (1 + \log(c_2/2nK\delta)).$$

Proof: Since the arms have sub-Gaussian distribution, for any $1 \leq k \leq K$ and $1 \leq t \leq n$, we have

$$\mathbb{P}(|X_{k,t} - \mu_k|^2 \geq \epsilon) \leq c_2 \exp(-\epsilon/c_1),$$

and thus by setting $\epsilon = c_1 \log(c_2/2nK\delta)^9$, we obtain

$$\mathbb{P}(|X_{k,t} - \mu_k|^2 \geq c_1 \log(c_2/2nK\delta)) \leq 2nK\delta.$$

⁹Note that we need to choose c_2 such that $c_2 \geq 2nK\delta = 2Kn^{-5/2}$ if $\delta = n^{-7/2}$.

We thus know that

$$\begin{aligned}
& \max_{\Omega/\mathbb{P}(\Omega)=2nK\delta} \mathbb{E}[|X_{k,t} - \mu_k|^2 \mathbb{I}\{\Omega\}] \\
& \leq \int_{c_1 \log(c_2/2nK\delta)}^{\infty} c_2 \exp(-\epsilon/c_1) d\epsilon + c_1 \log(c_2/2nK\delta) \mathbb{P}(\Omega) \\
& = 2c_1 nK\delta (1 + \log(c_2/2nK\delta)).
\end{aligned}$$

Since the event ξ^C has a probability at most $2nK\delta$, for any $1 \leq k \leq K$ and $1 \leq t \leq n$, we have

$$\mathbb{E}[|X_{k,t} - \mu_k|^2 \mathbb{I}\{\xi^C\}] \leq \max_{\Omega/\mathbb{P}(\Omega)=2nK\delta} \mathbb{E}[|X_{k,t} - \mu_k|^2 \mathbb{I}\{\Omega\}] \leq 2c_1 nK\delta (1 + \log(c_2/2nK\delta)).$$

The claim follows from the fact that $\mathbb{E}[|\hat{\mu}_{k,n} - \mu_k|^2 \mathbb{I}\{\xi^C\}] \leq \sum_{t=1}^n \mathbb{E}[|X_{k,t} - \mu_k|^2 \mathbb{I}\{\xi^C\}] \leq 2c_1 n^2 K\delta (1 + \log(c_2/2nK\delta))$. \square

A.3 Technical inequalities

Upper and lower bound on a : If $\delta = n^{-7/2}$, with $n \geq 4K \geq 8$

$$\begin{aligned}
a &= \sqrt{2c_1 \log(c_2/\delta)} + \frac{\sqrt{c_1 \delta (1 + c_2 + \log(c_2/\delta))}}{(1 - \delta) \sqrt{2 \log(2/\delta)}} n^{1/2} \\
&\leq \sqrt{7c_1 (c_2 + 1) \log(n)} + \frac{1}{n^{3/2}} \sqrt{c_1 (2 + c_2)} \\
&\leq 2\sqrt{2c_1 (c_2 + 2) \log(n)}.
\end{aligned}$$

We also have by just keeping the first term and choosing c_2 such that $c_2 \geq e\delta = en^{-7/2}$

$$\begin{aligned}
a &= \sqrt{2c_1 \log(c_2/\delta)} + \frac{\sqrt{c_1 \delta (1 + c_2 + \log(c_2/\delta))}}{(1 - \delta) \sqrt{2 \log(2/\delta)}} n^{1/2} \\
&\geq \sqrt{2c_1} \geq \sqrt{c_1}.
\end{aligned}$$

Lower bound on $c(\delta)$ when $\delta = n^{-7/2}$: Since the arms have sub-Gaussian distribution, for any $1 \leq k \leq K$ and $1 \leq t \leq n$, we have

$$\mathbb{P}(|X_{k,t} - \mu_k|^2 \geq \epsilon) \leq c_2 \exp(-\epsilon/c_1),$$

We then have

$$\mathbb{E}[|X_{k,t} - \mu_k|^2] \leq \int_0^{\infty} c_2 \exp(-\epsilon/c_1) d\epsilon = c_2 c_1$$

We then have $\Sigma_w \leq \sqrt{c_2 c_1}$.

If $\delta = n^{-7/2}$, we obtain by using the lower bound on a that

$$\begin{aligned}
c(\delta = n^{-7/2}) &= \left(\frac{2a \sqrt{\log(2/\delta)}}{\Sigma_w + 4a \sqrt{\log(2/\delta)}} \frac{1}{K} \right)^{2/3} \\
&= \left(\frac{1}{2K} - \frac{1}{2K} \frac{\Sigma_w}{\Sigma_w + 4a \sqrt{\log(2/\delta)}} \right)^{2/3} \\
&\geq \left(\frac{1}{2K} - \frac{1}{2K} \frac{\Sigma_w}{\Sigma_w + 4\sqrt{c_1} \log(n)} \right)^{2/3} \\
&\geq \left(\frac{1}{2K} \right)^{2/3} \left(\frac{\sqrt{c_1}}{\Sigma_w + \sqrt{c_1}} \right)^{2/3} \geq \left(\frac{1}{2K} \right)^{2/3} \left(\frac{1}{\sqrt{c_2} + 1} \right)^{2/3},
\end{aligned}$$

by using $\Sigma_w \leq \sqrt{c_2 c_1}$ for the last step.

Upper bound on $\mathbb{E}[|\hat{\mu}_{k,n} - \mu_k|^2 \mathbb{I}\{\xi^C\}]$ when $\delta = n^{-7/2}$: We get from Lemma 2 when $\delta = n^{-7/2}$ and when choosing c_2 such that $c_2 \geq 2enK\delta = 2Ken^{-5/2}$

$$\begin{aligned} \mathbb{E}[|\hat{\mu}_{k,n} - \mu_k|^2 \mathbb{I}\{\xi^C\}] &\leq 2c_1 n^2 K \delta (1 + \log(c_2/2nK\delta)) \\ &\leq 2c_1 K (1 + \frac{5}{2}(c_2 + 1) \log(n)) n^{-3/2} \\ &\leq 6c_1 K (c_2 + 1) \log(n) n^{-3/2}. \end{aligned}$$

B Proof of Theorem 1

In this section, we first provide the proof for the following Lemma and then use the result to prove Theorem 1.

Lemma 3 *Let Assumption 1 hold. Let $0 < \delta \leq 1$ be arbitrary and $n \geq 4K$. The difference between the allocation $T_{p,n}$ implemented by the MC-UCB algorithm described in Figure 1 and the optimal allocation rule $T_{p,n}^*$ has the following upper and lower bounds, on ξ (and thus with probability at least $1 - 2nK\delta$), for any arm $1 \leq p \leq K$:*

$$-12a\lambda_p \frac{\sqrt{\log(2/\delta)}}{\Sigma_w \lambda_{\min}^{3/2}} \sqrt{n} - 4K\lambda_p \leq T_{p,n} - T_{p,n}^* \leq 12a \frac{\sqrt{\log(2/\delta)}}{\Sigma_w \lambda_{\min}^{3/2}} \sqrt{n} + 4K. \quad (15)$$

where $a = \sqrt{2c_1 \log(c_2/\delta)} + \frac{\sqrt{c_1 \delta (1+c_2+\log(c_2/\delta))}}{(1-\delta)\sqrt{2\log(2/\delta)}} n^{1/2}$.

In Equation 15, the difference $T_{p,n} - T_{p,n}^*$ is bounded with $\tilde{O}(\sqrt{n})$. This is directly linked to the parametric rate of convergence of the estimation of σ_k , which is of order $1/\sqrt{n}$. Note that Equation 15 also shows the inverse dependency on the smallest proportion λ_{\min} .

Proof: [Lemma 3] The proof consists of the following three main steps.

Step 1. Properties of the algorithm. Recall the definition of the upper bound used in MC-UCB when $t > 2K$:

$$B_{q,t+1} = \frac{w_q}{T_{q,t}} \left(\hat{\sigma}_{q,t} + 2a \sqrt{\frac{\log(2/\delta)}{T_{q,t}}} \right), \quad 1 \leq q \leq K.$$

From Corollary 1, we obtain the following upper and lower bounds for $B_{q,t+1}$ on ξ :

$$\frac{w_q \sigma_q}{T_{q,t}} \leq B_{q,t+1} \leq \frac{w_q}{T_{q,t}} \left(\sigma_q + 4a \sqrt{\frac{\log(2/\delta)}{T_{q,t}}} \right). \quad (16)$$

Let $t+1 > 2K$ be the time at which a given arm k is pulled for the last time, i.e., $T_{k,t} = T_{k,n} - 1$ and $T_{k,(t+1)} = T_{k,n}$. Note that as $n \geq 4K$, there is at least one arm k such that this happens, i.e. such that it is pulled after the initialization phase. Since \mathcal{A}_{MC-UCB} chooses to pull arm k at time $t+1$, we have for any arm p

$$B_{p,t+1} \leq B_{k,t+1}. \quad (17)$$

From Equation 16 and the fact that $T_{k,t} = T_{k,n} - 1$, we obtain

$$B_{k,t+1} \leq \frac{w_k}{T_{k,t}} \left(\sigma_k + 4a \sqrt{\frac{\log(2/\delta)}{T_{k,t}}} \right) = \frac{w_k}{T_{k,n} - 1} \left(\sigma_k + 4a \sqrt{\frac{\log(2/\delta)}{T_{k,n} - 1}} \right). \quad (18)$$

Using the lower bound in Equation 16 and the fact that $T_{p,t} \leq T_{p,n}$, we may lower bound $B_{p,t+1}$ as

$$B_{p,t+1} \geq \frac{w_p \sigma_p}{T_{p,t}} \geq \frac{w_p \sigma_p}{T_{p,n}}. \quad (19)$$

Combining Equations 17, 18, and 19, we obtain

$$\frac{w_p \sigma_p}{T_{p,n}} \leq \frac{w_k}{T_{k,n} - 1} \left(\sigma_k + 4a \sqrt{\frac{\log(2/\delta)}{T_{k,n} - 1}} \right). \quad (20)$$

Note that at this point there is no dependency on t , and thus, the probability that Equation 20 holds for any p and for any k such that arm k is pulled after the initialization phase, i.e., such that $T_{k,n} > 2$, is at least $1 - 2nK\delta$ (probability of event ξ).

Step 2. Lower bound on $T_{p,n}$. If an arm p is under-pulled compared to its optimal allocation *without taking into account the initialization phase*, i.e., $T_{p,n} - 2 < \lambda_p(n - 2K)$, then from the constraint $\sum_k (T_{k,n} - 2) = n - 2K$ and the definition of the optimal allocation, we deduce that there exists at least another arm k that is over-pulled compared to its optimal allocation *without taking into account the initialization phase*, i.e., $T_{k,n} - 2 > \lambda_k(n - 2K)$. Note that for this arm, $T_{k,n} - 2 > \lambda_k(n - 2K) \geq 0$, so we know that this specific arm is pulled at least once *after* the initialization phase and that it satisfies Equation 20. Using the definition of the optimal allocation $T_{k,n}^* = nw_k \sigma_k / \Sigma_w$, and the fact that $T_{k,n} \geq \lambda_k(n - 2K) + 2$, Equation 20 may be written as for any arm p

$$\begin{aligned} \frac{w_p \sigma_p}{T_{p,n}} &\leq \frac{w_k}{T_{k,n}^*} \frac{n}{(n - 2K)} \left(\sigma_k + 4a \sqrt{\frac{\log(2/\delta)}{\lambda_k(n - 2K) + 1}} \right) \\ &\leq \frac{\Sigma_w}{n} + \frac{4K \Sigma_w}{n^2} + 8\sqrt{2}a \frac{\sqrt{\log(2/\delta)}}{n^{3/2} \lambda_k^{3/2}}, \end{aligned}$$

because $n \geq 4K$. The previous Equation, combined with the fact that $\lambda_k \geq \lambda_{\min}$, may be written as

$$\frac{w_p \sigma_p}{T_{p,n}} \leq \frac{\Sigma_w}{n} + 12a \frac{\sqrt{\log(2/\delta)}}{n^{3/2} \lambda_{\min}^{3/2}} + \frac{4K \Sigma_w}{n^2}. \quad (21)$$

By rearranging Equation 21, we obtain the lower bound on $T_{p,n}$:

$$T_{p,n} \geq \frac{w_p \sigma_p}{\frac{\Sigma_w}{n} + 12a \frac{\sqrt{\log(2/\delta)}}{n^{3/2} \lambda_{\min}^{3/2}} + \frac{4K \Sigma_w}{n^2}} \geq T_{p,n}^* - 12a \lambda_p \frac{\sqrt{\log(2/\delta)}}{\Sigma_w \lambda_{\min}^{3/2}} \sqrt{n} - 4K \lambda_p, \quad (22)$$

where in the second inequality we use $1/(1+x) \geq 1-x$ (for $x > -1$). Note that the lower bound holds on ξ for any arm p .

Step 3. Upper bound on $T_{p,n}$. Using Equation 22 and the fact that $\sum_k T_{k,n} = n$, we obtain

$$T_{p,n} = n - \sum_{k \neq p} T_{k,n} \leq \left(n - \sum_{k \neq p} T_{k,n}^* \right) + \sum_{k \neq p} \left(12a \lambda_p \frac{\sqrt{\log(2/\delta)}}{\Sigma_w \lambda_{\min}^{3/2}} \sqrt{n} + 4K \lambda_p \right).$$

And we deduce because $\sum_{k \neq p} \lambda_k \leq 1$

$$T_{p,n} \leq T_{p,n}^* + 12a \frac{\sqrt{\log(2/\delta)}}{\Sigma_w \lambda_{\min}^{3/2}} \sqrt{n} + 4K. \quad (23)$$

The lemma follows by combining the lower and upper bounds in Equations 22 and 23. \square

We are now ready to prove Theorem 1.

Theorem 1 Under Assumption 1 and if c_2 is chosen such that $c_2 \geq 2Kn^{-5/2}$, the regret of MC-UCB run with parameter $\delta = n^{-7/2}$ with $n \geq 4K$ is bounded as

$$R_n(\mathcal{A}_{MC-UCB}) \leq \frac{\log(n)}{n^{3/2}\lambda_{\min}^{3/2}} \left(112\Sigma_w \sqrt{c_1(c_2 + 2)} + 6c_1(c_2 + 2)K \right) + \frac{19}{\lambda_{\min}^3 n^2} \left(K\Sigma_w^2 + 720c_1(c_2 + 1) \log(n)^2 \right).$$

Proof: [Theorem 1] The proof consists of the following two steps.

Step 1. $T_{k,n}$ is a stopping time. Consider an arm k . At each time step $t + 1$, the MC-UCB algorithm decides which arm to pull according to the current values of the upper-bounds $\{B_{k,t+1}\}_k$. Thus for any arm k , $T_{k,(t+1)}$ depends only on the values $\{T_{k,t}\}_k$ and $\{\hat{\sigma}_{k,t}\}_k$. So by induction, $T_{k,(t+1)}$ depends on the sequence $\{X_{k,1}, \dots, X_{k,T_{k,t}}\}$, and on the samples of the other arms (which are independent of the samples of arm k). We deduce that $T_{k,n}$ is a stopping time adapted to the process $(X_{k,t})_{t \leq n}$.

Step 2. Regret bound. By definition, the loss of the algorithm writes

$$L_n = \sum_{k=1}^K w_k^2 \mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \right] = \sum_{k=1}^K w_k^2 \mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \mathbb{I}\{\xi\} \right] + \sum_{k=1}^K w_k^2 \mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \mathbb{I}\{\xi^C\} \right].$$

Using the definition of $\hat{\mu}_{k,n}$ and Proposition 1 we bound the first term as

$$\sum_{k=1}^K w_k^2 \mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \mathbb{I}\{\xi\} \right] \leq \sum_{k=1}^K w_k^2 \frac{\sigma_k^2 \mathbb{E}[T_{k,n}]}{\underline{T}_{k,n}^2}, \quad (24)$$

where $\underline{T}_{k,n}$ is the lower bound on $T_{k,n}$ on the event ξ .

Note that as $\sum_k T_{k,n} = n$, we also have $\sum_k \mathbb{E}[T_{k,n}] = n$.

Using Equation 24 and Equation 21 for $w_k \sigma_k / \underline{T}_{k,n}$ (which is equivalent to using a lower bound on $T_{k,n}$ on the event ξ), we obtain

$$\sum_{k=1}^K w_k^2 \frac{\sigma_k^2 \mathbb{E}[T_{k,n}]}{\underline{T}_{k,n}^2} \leq \sum_{k=1}^K \left(\frac{\Sigma_w}{n} + 12a \frac{\sqrt{\log(2/\delta)}}{n^{3/2}\lambda_{\min}^{3/2}} + \frac{4K\Sigma_w}{n^2} \right)^2 \mathbb{E}[T_{k,n}]. \quad (25)$$

Equation 25 may be bounded using the fact that $\sum_k \mathbb{E}[T_{k,n}] = n$ as

$$\begin{aligned} \sum_{k=1}^K w_k^2 \frac{\sigma_k^2 \mathbb{E}[T_{k,n}]}{\underline{T}_{k,n}^2} &\leq \left(\frac{\Sigma_w}{n} + 12a \frac{\sqrt{\log(2/\delta)}}{n^{3/2}\lambda_{\min}^{3/2}} + \frac{4K\Sigma_w}{n^2} \right)^2 n \\ &\leq \left(\left(\frac{\Sigma_w}{n} \right)^2 + 24a\Sigma_w \frac{\sqrt{\log(2/\delta)}}{n^{5/2}\lambda_{\min}^{3/2}} + \frac{8K\Sigma_w^2}{n^3} + 288a^2 \frac{\log(2/\delta)}{n^3\lambda_{\min}^3} + \frac{8K^2\Sigma_w^2}{n^4} \right) n \\ &= \frac{\Sigma_w^2}{n} + 24a\Sigma_w \frac{\sqrt{\log(2/\delta)}}{n^{3/2}\lambda_{\min}^{3/2}} + \frac{8K\Sigma_w^2}{n^2} + 288a^2 \frac{\log(2/\delta)}{n^2\lambda_{\min}^3} + \frac{8K^2\Sigma_w^2}{n^3} \\ &\leq \frac{\Sigma_w^2}{n} + 24a\Sigma_w \frac{\sqrt{\log(2/\delta)}}{n^{3/2}\lambda_{\min}^{3/2}} + \frac{16}{\lambda_{\min}^3 n^2} \left(K\Sigma_w^2 + 18a^2 \log(2/\delta) \right). \end{aligned}$$

From Lemma 2, we have $\mathbb{E}\left[(\hat{\mu}_{k,n} - \mu_k)^2 \mathbb{I}\{\xi^C\}\right] \leq 2c_1 n^2 K \delta (1 + \log(c_2/2nK\delta))$. Thus using the previous equation, we deduce

$$\begin{aligned} L_n &\leq \frac{\Sigma_w^2}{n} + 24a\Sigma_w \frac{\sqrt{\log(2/\delta)}}{n^{3/2}\lambda_{\min}^{3/2}} + \frac{16}{\lambda_{\min}^3 n^2} \left(K\Sigma_w^2 + 18a^2 \log(2/\delta)\right) + 2c_1 n^2 K \delta (1 + \log(c_2/2nK\delta)) \\ &\leq \frac{\Sigma_w^2}{n} + 54a\Sigma_w \frac{\sqrt{\log(n)}}{n^{3/2}\lambda_{\min}^{3/2}} + \frac{16}{\lambda_{\min}^3 n^2} \left(K\Sigma_w^2 + 90a^2 \log(n)\right) + 6c_1 K (c_2 + 1) \log(n) n^{-3/2} \\ &\leq \frac{\Sigma_w^2}{n} + \frac{\log(n)}{n^{3/2}\lambda_{\min}^{3/2}} \left(112\Sigma_w \sqrt{c_1(c_2 + 2)} + 6c_1(c_2 + 2)K\right) + \frac{19}{\lambda_{\min}^3 n^2} \left(K\Sigma_w^2 + 720c_1(c_2 + 1) \log(n)^2\right). \end{aligned}$$

where we use $a \leq 2\sqrt{2c_1(c_2 + 2)\log(n)}$ and $\mathbb{E}[|\hat{\mu}_{k,n} - \mu_k|^2 \mathbb{I}\{\xi^C\}] \leq 6c_1 K (c_2 + 1) \log(n) n^{-3/2}$. Those bounds are made explicit in A.3.

The Theorem follows by expressing the regret. \square

C Proof of Theorem 2

Again, we first state and prove the following Lemma and then use this result to prove Theorem 2.

Lemma 4 *Let Assumption 1 hold. For any $0 < \delta \leq 1$ and for $n \geq 4K$, the algorithm MC-UCB satisfies on ξ , and thus with probability at least $1 - 2nK\delta$, for any arm p ,*

$$T_{p,n} \geq T_{p,n}^* - \left(24aK^{2/3} \frac{1}{\Sigma_w} \lambda_q \sqrt{\frac{\log(2/\delta)}{c(\delta)}} n^{2/3} + 12K\lambda_q\right), \quad (26)$$

and

$$T_{p,n} \leq T_{p,n}^* + \left(24aK^{2/3} \frac{1}{\Sigma_w} \sqrt{\frac{\log(2/\delta)}{c(\delta)}} n^{2/3} + 12K\Sigma_w\right), \quad (27)$$

where $c(\delta) = \left(\frac{2a\sqrt{\log(2/\delta)}}{\Sigma_w + 4a\sqrt{\log(2/\delta)}} \frac{1}{K}\right)^{2/3}$ and $a = \sqrt{2c_1 \log(c_2/\delta)} + \frac{\sqrt{c_1\delta(1+c_2+\log(c_2/\delta))}}{(1-\delta)\sqrt{2\log(2/\delta)}} n^{1/2}$.

Unlike the bounds proved in Lemma 3, the difference between $T_{p,n}$ and $T_{p,n}^*$ is bounded by $\tilde{O}(n^{2/3})$ without any inverse dependency on λ_{\min} .

Proof:

Step 1. Lower bound of order $\tilde{O}(n^{2/3})$. Let k be the index of an arm such that $T_{k,n} \geq \frac{n}{K}$ (this implies $T_{k,n} \geq 3$ as $n \geq 4K$, and arm k is thus pulled after the initialization) and let $t + 1 \leq n$ be the last time at which it was pulled¹⁰, i.e., $T_{k,t} = T_{k,n} - 1$ and $T_{k,t+1} = T_{k,n}$. From Equation 13 and the fact that $T_{k,n} \geq \frac{n}{K}$, we obtain on ξ

$$B_{k,t} \leq \frac{w_k}{T_{k,t}} \left(\sigma_k + 4a\sqrt{\frac{\log(2/\delta)}{T_{k,t}}}\right) \leq \frac{K\left(\Sigma_w + 4a\sqrt{\log(2/\delta)}\right)}{n}, \quad (28)$$

where the second inequality follows from the facts that $T_{k,t} \geq 1$, $w_k\sigma_k \leq \Sigma_w$, and $w_k \leq \sum_k w_k = 1$. Since at time $t + 1$ the arm k has been pulled, then for any arm q , we have

$$B_{q,t} \leq B_{k,t}. \quad (29)$$

¹⁰Note that such an arm always exists for any possible allocation strategy given the constraint $n = \sum_q T_{q,n}$.

From the definition of $B_{q,t}$, and also using the fact that $T_{q,t} \leq T_{q,n}$, we deduce on ξ that

$$B_{q,t} \geq 2aw_q \frac{\sqrt{\log(2/\delta)}}{T_{q,t}^{3/2}} \geq 2aw_q \frac{\sqrt{\log(2/\delta)}}{T_{q,n}^{3/2}}. \quad (30)$$

Combining Equations 28–30, we obtain on ξ

$$2aw_q \frac{\sqrt{\log(2/\delta)}}{T_{q,n}^{3/2}} \leq \frac{K(\Sigma_w + 4a\sqrt{\log(2/\delta)})}{n}.$$

Finally, this implies on ξ that for any q ,

$$T_{q,n} \geq \left(\frac{2aw_q \sqrt{\log(2/\delta)}}{\Sigma_w + 4a\sqrt{\log(2/\delta)}} \frac{n}{K} \right)^{2/3}. \quad (31)$$

In order to simplify the notation, in the following we define

$$c(\delta) = \left(\frac{2a\sqrt{\log(2/\delta)}}{\Sigma_w + 4a\sqrt{\log(2/\delta)}} \frac{1}{K} \right)^{2/3},$$

thus the lower bound on $T_{q,n}$ on ξ writes $T_{q,n} \geq w_q^{2/3} c(\delta) n^{2/3}$.

Step 2. Properties of the algorithm. We follow a similar analysis to Step 1 of the proof of Lemma 3. We first recall the definition of $B_{q,t+1}$ used in the MC-UCB algorithm

$$B_{q,t+1} = \frac{w_q}{T_{q,t}} \left(\hat{\sigma}_{q,t} + 2a\sqrt{\frac{\log(2/\delta)}{T_{q,t}}} \right).$$

Using Corollary 1 it follows that, on ξ

$$\frac{w_q \sigma_q}{T_{q,t}} \leq B_{q,t+1} \leq \frac{w_q}{T_{q,t}} \left(\sigma_q + 4a\sqrt{\frac{\log(2/\delta)}{T_{q,t}}} \right). \quad (32)$$

Let $t+1 \geq 2K+1$ be the time at which an arm q is pulled for the last time, that is $T_{q,t} = T_{q,n} - 1$. Note that there is at least one arm such that this happens as $n \geq 4K$. Since at $t+1$ arm q is chosen, then for any other arm p , we have

$$B_{p,t+1} \leq B_{q,t+1}. \quad (33)$$

From Equation 32 and $T_{q,t} = T_{q,n} - 1$, we obtain on ξ

$$B_{q,t+1} \leq \frac{w_q}{T_{q,t}} \left(\sigma_q + 4a\sqrt{\frac{\log(2/\delta)}{T_{q,t}}} \right) = \frac{w_q}{T_{q,n} - 1} \left(\sigma_q + 4a\sqrt{\frac{\log(2/\delta)}{T_{q,n} - 1}} \right). \quad (34)$$

Furthermore, since $T_{p,t} \leq T_{p,n}$, then on ξ

$$B_{p,t+1} \geq \frac{w_p \sigma_p}{T_{p,t}} \geq \frac{w_p \sigma_p}{T_{p,n}}. \quad (35)$$

Combining Equations 33–35, we obtain on ξ

$$\frac{w_p \sigma_p}{T_{p,n}} (T_{q,n} - 1) \leq w_q \left(\sigma_q + 4a\sqrt{\frac{\log(2/\delta)}{T_{q,n} - 1}} \right).$$

Summing over all q such that the previous Equation is verified, i.e. such that $T_{q,n} \geq 3$, on both sides, we obtain on ξ

$$\frac{w_p \sigma_p}{T_{p,n}} \sum_{q|T_{q,n} \geq 3} (T_{q,n} - 1) \leq \sum_{q|T_{q,n} \geq 3} w_q \left(\sigma_q + 4a\sqrt{\frac{\log(2/\delta)}{T_{q,n} - 1}} \right).$$

This implies

$$\frac{w_p \sigma_p}{T_{p,n}}(n - 3K) \leq \sum_{q=1}^K w_q \left(\sigma_q + 4a \sqrt{\frac{\log(2/\delta)}{T_{q,n} - 1}} \right). \quad (36)$$

Step 3. Lower bound. Plugging Equation 31 in Equation 36,

$$\begin{aligned} \frac{w_p \sigma_p}{T_{p,n}}(n - 3K) &\leq \sum_q w_q \left(\sigma_q + 4a \sqrt{\frac{\log(2/\delta)}{T_{q,n} - 1}} \right) \\ &\leq \sum_q w_q \left(\sigma_q + 4a \sqrt{\frac{2 \log(2/\delta)}{w_q^{2/3} c(\delta) n^{2/3}}} \right) \\ &\leq \Sigma_w + \sum_q 4a w_q^{2/3} \sqrt{2 \frac{\log(2/\delta)}{c(\delta) n^{2/3}}} \leq \Sigma_w + 6a K^{2/3} \sqrt{\frac{\log(2/\delta)}{c(\delta) n^{2/3}}}, \end{aligned}$$

on ξ , since $\sum_q w_q^{2/3} \leq K^{2/3}$ by Jensen inequality and because $T_{q,n} - 1 \geq \frac{T_{q,n}}{2}$ (as $T_{q,n} \geq 2$). Finally as $n \geq 4K$, we obtain on ξ the following bound

$$\frac{w_p \sigma_p}{T_{p,n}} \leq \frac{\Sigma_w}{n} + 24a K^{2/3} \sqrt{\frac{\log(2/\delta)}{c(\delta)}} n^{-4/3} + \frac{12K \Sigma_w}{n^2}. \quad (37)$$

We now invert the bound and obtain on ξ the final lower-bound on $T_{p,n}$ as follows:

$$T_{p,n} \geq \frac{w_p \sigma_p}{\frac{\Sigma_w}{n} + 24a K^{2/3} \sqrt{\frac{\log(2/\delta)}{c(\delta)}} n^{-4/3} + \frac{12K \Sigma_w}{n^2}} \geq T_{p,n}^* - 24a K^{2/3} \frac{1}{\Sigma_w} \lambda_p \sqrt{\frac{\log(2/\delta)}{c(\delta)}} n^{2/3} - 12K \lambda_p.$$

Note that the above lower bound holds with high probability for any arm p .

Step 4. Upper bound. An upper bound on $T_{p,n}$ on ξ follows by using $T_{p,n} = n - \sum_{q \neq p} T_{q,n}$ and the previous lower bound, that is

$$\begin{aligned} T_{p,n} &\leq n - \sum_{q \neq p} T_{q,n}^* + \sum_{q \neq p} \left(12K \lambda_q + 24a K^{2/3} \frac{1}{\Sigma_w} \lambda_q \sqrt{\frac{\log(2/\delta)}{c(\delta)}} n^{2/3} \right) \\ &\leq T_{p,n}^* + \left(24a K^{2/3} \frac{1}{\Sigma_w} \sqrt{\frac{\log(2/\delta)}{c(\delta)}} n^{2/3} + 12K \right), \end{aligned}$$

□

because $\sum_{q \neq p} \lambda_q \leq 1$.

We are now ready to prove Theorem 2.

Theorem 2 *Under Assumption 1 and by ensuring that $c_2 \geq 2Kn^{-5/2}$, the regret of MC-UCB run with parameter $\delta = n^{-7/2}$ with $n \geq 4K$ is bounded as*

$$R_n(\mathcal{A}_{MC-UCB}) \leq \frac{200\sqrt{c_1}(c_2 + 2)\Sigma_w K}{n^{4/3}} \log(n) + \frac{365}{n^{3/2}} \left(129c_1(c_2 + 2)^2 K^2 \log(n)^2 + K \Sigma_w^2 \right).$$

Proof: [Theorem 2]

We decompose the loss on ξ and its complement:

$$L_n = \sum_{k=1}^K w_k^2 \mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \right] = \sum_{k=1}^K w_k^2 \mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \mathbb{I}\{\xi\} \right] + \sum_{k=1}^K w_k^2 \mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \mathbb{I}\{\xi^C\} \right].$$

Using the definition of $\hat{\mu}_{k,n}$ and Proposition 1 we bound the first term as

$$\sum_{k=1}^K w_k^2 \mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \mathbb{I} \{ \xi \} \right] \leq \sum_{k=1}^K w_k^2 \frac{\sigma_k^2 \mathbb{E}[T_{k,n}]}{\underline{T}_{k,n}}, \quad (38)$$

where $\underline{T}_{k,n}$ is the lower bound on $T_{k,n}$ on ξ .

Note also that as $\sum_k T_{k,n} = n$, we also have $\sum_k \mathbb{E}[T_{k,n}] = n$. Using Equation 38 and Equation 37 which provides an upper bound on ξ on $\frac{w_k \sigma_k}{T_{k,n}}$ (and thus a lower bound on ξ on $T_{k,n}$), we deduce

$$\sum_{k=1}^K w_k^2 \mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \mathbb{I} \{ \xi \} \right] \leq \sum_{k=1}^K \left(\frac{\Sigma_w}{n} + 24aK^{2/3} \sqrt{\frac{\log(2/\delta)}{c(\delta)}} n^{-4/3} + \frac{12K\Sigma_w}{n^2} \right)^2 \mathbb{E}[T_{k,n}]. \quad (39)$$

Using the fact that $\sum_k \mathbb{E}[T_{k,n}] = n$, Equation 39 may be rewritten as

$$\begin{aligned} \sum_{k=1}^K w_k^2 \mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \mathbb{I} \{ \xi \} \right] &\leq \left(\frac{\Sigma_w}{n} + 24aK^{2/3} \sqrt{\frac{\log(2/\delta)}{c(\delta)}} n^{-4/3} + \frac{12K\Sigma_w}{n^2} \right)^2 n \\ &\leq \left(\left(\frac{\Sigma_w}{n} \right)^2 + \frac{48\Sigma_w a K^{2/3}}{n^{7/3}} \sqrt{\frac{\log(2/\delta)}{c(\delta)}} \right. \\ &\quad \left. + \frac{12K\Sigma_w^2}{n^3} + \frac{1152a^2 K^{4/3}}{n^{8/3}} \frac{\log(2/\delta)}{c(\delta)} + \frac{288K^2 \Sigma_w^2}{n^4} \right) n \\ &= \frac{\Sigma_w^2}{n} + \frac{48\Sigma_w a K^{2/3}}{n^{4/3}} \sqrt{\frac{\log(2/\delta)}{c(\delta)}} \\ &\quad + \frac{12K\Sigma_w^2}{n^2} + \frac{1152a^2 K^{4/3}}{n^{5/3}} \frac{\log(2/\delta)}{c(\delta)} + \frac{288K^2 \Sigma_w^2}{n^3} \\ &\leq \frac{\Sigma_w^2}{n} + \frac{48\Sigma_w a K^{2/3}}{n^{4/3}} \sqrt{\frac{\log(2/\delta)}{c(\delta)}} + \frac{300}{n^2} \left(4a^2 K^{4/3} \frac{\log(2/\delta)}{c(\delta)} + K\Sigma_w^2 \right). \end{aligned}$$

From Lemma 2, we have $\mathbb{E} \left[(\hat{\mu}_{k,n} - \mu_k)^2 \mathbb{I} \{ \xi^C \} \right] \leq 2c_1 n^2 K \delta (1 + \log(c_2/2nK\delta))$. Thus using the last equation and the fact that $\delta = n^{-7/2}$, the loss is bounded as

$$\begin{aligned} L_n &\leq \frac{\Sigma_w^2}{n} + \frac{48\Sigma_w a K^{2/3}}{n^{4/3}} \sqrt{\frac{\log(2/\delta)}{c(\delta)}} + \frac{300}{n^2} \left(4a^2 K^{4/3} \frac{\log(2/\delta)}{c(\delta)} + K\Sigma_w^2 \right) + 2c_1 n^2 K \delta (1 + \log(c_2/2nK\delta)) \\ &\leq \frac{\Sigma_w^2}{n} + \frac{96\Sigma_w a K}{n^{4/3}} \sqrt{\log(n)} (\sqrt{c_2} + 1)^{1/3} + \frac{300}{n^2} \left(16a^2 K^2 \log(n) (\sqrt{c_2} + 1)^{2/3} + K\Sigma_w^2 \right) \\ &\quad + 6c_1 K (c_2 + 1) \log(n) n^{-3/2} \\ &\leq \frac{\Sigma_w^2}{n} + \frac{200\sqrt{c_1(c_2+2)}\Sigma_w K}{n^{4/3}} \log(n) (\sqrt{c_2} + 1)^{1/3} \\ &\quad + \frac{365}{n^{3/2}} \left(16a^2 K^2 \log(n) (\sqrt{c_2} + 1)^{2/3} + K\Sigma_w^2 + c_1(c_2+2)K \log(n) \right) \\ &\leq \frac{\Sigma_w^2}{n} + \frac{200\sqrt{c_1(c_2+2)}\Sigma_w K}{n^{4/3}} \log(n) + \frac{365}{n^{3/2}} \left(129c_1(c_2+2)^2 K^2 \log(n)^2 + K\Sigma_w^2 \right). \end{aligned}$$

where we use $a \leq 2\sqrt{2c_1(c_2+2)\log(n)}$, $c(\delta) = \left(\frac{1}{2K}\right)^{2/3} \left(\frac{1}{\sqrt{c_2+1}}\right)^{2/3}$ and $\mathbb{E} \left[|\hat{\mu}_{k,n} - \mu_k|^2 \mathbb{I} \{ \xi^C \} \right] \leq 6c_1 K (c_2 + 1) \log(n) n^{-3/2}$. Those bound are made explicit in A.3.

□

D Comments on problem independent bounds for MC-UCB and GAFS-WL

D.1 Note on the problem independent bound for MC-UCB

An interesting question is whether it is possible to obtain a regret bound of order $n^{-3/2}$ without the dependency on λ_{\min}^{-1} . We provide a simple example that demonstrates that MC-UCB does not possess this property.

Consider a problem with $K = 2$ arms with $\sigma_1 = 1$ and $\sigma_2 = 0$ and $w_1 = w_2 = 0.5$. The optimal allocation strategy for this problem is $T_{1,n}^* = n - 1$, $T_{2,n}^* = 1$ (we need only one sample of the second arm to estimate its mean). Since $\lambda_{\min} = 0$ the bound in Theorem 1 is meaningless (although MC-UCB is still able to minimize the regret as demonstrated by Theorem 2). Indeed, the definition of the upper-confidence bound on the standard deviation forces the algorithm to pull each arm at least $\tilde{O}(n^{2/3})$ times, including those arms with zero variance. Hence in this example, arm 2 will be pulled $\tilde{O}(n^{2/3})$ times, which results in under-pulling arm 1 by the same amount, and thus, in worsening its estimation. It can be shown that the resulting regret is $\tilde{O}(n^{-4/3})$, which still decreases to zero faster than $1/n$ but with a poorer rate. A sketch of the proof of this argument is as follows. Using the definition of $B_{k,t}$ and Equation 18 (see Appendix B), since $\hat{\sigma}_2 = 0$, we have that at any time $t + 1 > 2$

$$B_{1,t+1} \leq \frac{1}{2T_{1,t}}(1+b) \quad \text{and} \quad B_{2,t+1} = \frac{1}{2T_{2,t}}\left(b\sqrt{\frac{1}{T_{2,t}}}\right). \quad (40)$$

Let $t+1 \leq n$ be the last time that arm 1 was pulled, i.e., $T_{1,t} = T_{1,n} - 1$ and $B_{1,t+1} \geq B_{2,t+1}$. From Equation 40, we have

$$B_{2,t+1} = \frac{b}{2T_{2,t}^{3/2}} \leq B_{1,t+1} \leq \frac{1}{2T_{1,n} - 1}(1+b). \quad (41)$$

Now consider the two possible cases: **1)** $T_{1,n} \leq n/2$, in which case obviously $T_{2,n} \geq n/2$ and **2)** $T_{1,n} > n/2$, in this case Equation 41 implies that $T_{2,n} \geq T_{2,t} = \tilde{O}(n^{2/3})$. Thus in both cases, we may write $T_{2,n} = \tilde{O}(n^{2/3})$, which indicates that arm 2 (resp. arm 1) is over-sampled (resp. under-sampled) by a number of pulls of order $\tilde{O}(n^{2/3})$. By following the same arguments as in the proof of Theorem 2, we deduce that the regret in this case is of order $\tilde{O}(n^{-4/3})$. Note that this poorer rate is the result of over-sampling the arm with the smaller variance (and as a consequence under-sampling at least one arm with a larger variance).

Thus in the case of an arm with 0 standard deviation, the regret of MC-UCB is at least $\tilde{O}(n^{-4/3})$.

D.2 Note for a problem independent bound for GAFS-WL

Let $n \geq 4$ be the budget. We face a two-arms bandit problem with $w_1 = w_2 = \frac{1}{2}$ and such that (i) the distribution of the first arm is a Bernoulli of parameter $p = \frac{1}{n^{1/2+\epsilon}}$ with ϵ such that $1/6 > \epsilon > 0$ and that (ii) the distribution of the second arm is such that $\sigma_2 = 1$ and bounded by c .

Note that

$$\frac{1}{2n^{1/4+\epsilon/2}} \leq \sigma_1 \leq \frac{1}{n^{1/4+\epsilon/2}} \quad \text{and} \quad \sigma_2 = 1,$$

because $\sigma_1 = \sqrt{p(1-p)}$ and that thus

$$L_n^* \leq \frac{(1 + n^{-1/4-\epsilon/2})^2}{4n} \leq \frac{1 + 3n^{-1/4-\epsilon/2}}{4n} \leq \frac{1}{4n} + \frac{1}{n^{5/4+\epsilon/2}}.$$

We run algorithm GAFS-WL on that problem. Note that algorithm GAFS-WL pull each arm $\lfloor a\sqrt{n} \rfloor$ times and then pull the arms according to $\frac{w_k \hat{\sigma}_{k,t}}{T_{k,t}}$.

We call $\{X_{p,u}\}_{p=1,2;u=1,\dots,n}$ the samples of the arms.

Note that:

$$\begin{aligned} \mathbb{P}\left(X_{1,1} = 0, \dots, X_{1,\lfloor a\sqrt{n} \rfloor} = 0\right) &\geq \left(1 - \frac{1}{n^{1/2+\epsilon}}\right)^{a\sqrt{n}} \\ &\geq \left(1 - \frac{an^{-\epsilon}}{a\sqrt{n}}\right)^{a\sqrt{n}} \\ &\geq (1 - an^{-\epsilon}) \exp(-an^{-\epsilon}) \geq (1 - an^{-\epsilon})^2. \end{aligned}$$

Note on the other hand, that $\mathbb{P}(|\hat{\sigma}_{2,a\sqrt{n}} - 1| \geq \frac{2\sqrt{\log(2/\delta)}}{\sqrt{an^{1/4}}}) \leq \delta$. This means that with probability at least $1 - 2\exp(-a\sqrt{n}/4)$, we have $\hat{\sigma}_{2,a\sqrt{n}} > 0$.

The probability that $\hat{\sigma}_{1,a\sqrt{n}} = 0$ goes to 1 when n goes to $+\infty$. The probability that $\hat{\sigma}_{2,a\sqrt{n}} > 0$ goes to 1 when n goes to $+\infty$. This means that the probability that GAFS-WL stops pulling arm 1 after $a\sqrt{n}$ pulls goes to 1 when n goes to $+\infty$, and arm 1 is under-pulled if $\epsilon < 1/2$ (it should be pulled $n^{3/4-\epsilon/2}$).

Note that on the event such that $(X_{1,1} = 0, \dots, X_{1,\lfloor a\sqrt{n} \rfloor} = 0)$, we know that $\hat{\mu}_{1,a\sqrt{n}} = 0$. Note also that we know that as arm 2 is gaussian, we have $\mathbb{E}(\hat{\mu}_{2,n} - \mu_2)^2 \leq \frac{1}{4n}$. The performance of GAFS-WL then verifies

$$\begin{aligned} L_n(\mathcal{A}_{GAFS-WL}) &\geq \frac{1}{4n} + \mathbb{P}(\hat{\sigma}_{1,a\sqrt{n}} = 0)\mathbb{P}(\hat{\sigma}_{2,a\sqrt{n}} > 0) \left(n^{-1/2-\epsilon}\right)^2 \\ &\geq \frac{1}{4n} + (1 - 2\exp(-a\sqrt{n}/4))(1 - an^{-\epsilon})^2 \left(n^{-1-2\epsilon}\right) \\ &\geq \frac{1}{4n} + \left(1 - \frac{8}{a\sqrt{n}}\right) \left(1 - 2\frac{a}{n^\epsilon}\right) \frac{1}{n^{1+2\epsilon}} \\ &\geq \frac{1}{4n} + \frac{1}{n^{1+2\epsilon}} - \frac{8}{an^{3/2+2\epsilon}} - \frac{2a}{n^{1+3\epsilon}} \\ &\geq \frac{1}{4n} + \frac{1}{n^{1+2\epsilon}} - \frac{10 \max(a, 1/a)}{n^{1+3\epsilon}}, \end{aligned}$$

where the last line is obtained using the fact that $\epsilon < 1/6$.

The loss thus verifies

$$\begin{aligned} R_n(\mathcal{A}_{GAFS-WL}) &\geq \frac{1}{n^{1+2\epsilon}} - \frac{10 \max(a, 1/a)}{n^{1+3\epsilon}} - \frac{1}{n^{5/4+\epsilon/2}} \\ &\geq \frac{1}{n^{1+2\epsilon}} - \frac{11 \max(a, 1/a)}{n^{1+3\epsilon}}, \end{aligned}$$

again because $\epsilon < 1/6$. This implies that for n such that $n \geq \left(\frac{11 \max(a, 1/a)}{2}\right)^{1/\epsilon}$, we have

$$R_n(\mathcal{A}_{GAFS-WL}) \geq \frac{1}{2n^{1+2\epsilon}},$$

with ϵ arbitrarily close to 0.