
Stochastic Rising Bandits: A Best Arm Identification Approach

Alessandro Montenegro
Politecnico di Milano, Milan
alessandro.montenegro@polimi.it

Marco Mussi
Politecnico di Milano, Milan
marco.mussi@polimi.it

Francesco Trovó
Politecnico di Milano, Milan
francesco1.trovo@polimi.it

Marcello Restelli
Politecnico di Milano, Milan
marcello.restelli@polimi.it

Alberto Maria Metelli
Politecnico di Milano, Milan
albertomaria.metelli@polimi.it

Abstract

Stochastic Rising Bandits (SRBs) model sequential decision-making problems in which the expected reward of the available options increases every time they are selected. This setting captures a wide range of scenarios in which the available options are *learning entities* whose performance improves (in expectation) over time. While previous works addressed the regret minimization problem, this paper focuses on the *fixed-budget Best Arm Identification* (BAI) problem for SRBs. In this scenario, given a fixed budget of rounds, we are asked to provide a recommendation about the best option at the end of the identification process. We propose two algorithms to tackle the above-mentioned setting, namely R-UCBE, which resorts to a UCB-like approach, and R-SR, which employs a successive reject procedure. Then, we prove that, with a sufficiently large budget, they provide guarantees on the probability of properly identifying the optimal option at the end of the learning process. Furthermore, we derive a lower bound on the error probability, matched by our R-SR (up to logarithmic factors), and illustrate how the need for a sufficiently large budget is unavoidable in the SRB setting. Finally, we numerically validate the proposed algorithms in synthetic and real-world environments and compare them with the currently available BAI strategies.

1 Introduction

Multi-Armed Bandits (MAB, Lattimore and Szepesvári, 2020) are a well-known framework that effectively solves learning problems requiring sequential decisions. Given a time horizon, the learner chooses, at each round, a single option (a.k.a. arm) and observes the corresponding noisy reward, which is a realization of an unknown distribution. The MAB problem is commonly studied in two flavours: *regret minimization* (Auer et al., 2002) and *best arm identification* (Bubeck et al., 2009). In regret minimization, the goal is to control the cumulative loss w.r.t. the optimal arm over a time horizon. Conversely, in best arm identification, the goal is to provide a recommendation about the best arm at the end of the time horizon. Specifically, we are interested in the fixed-budget scenario, where we seek to minimize the error probability of recommending the wrong arm at the end of the time budget, no matter the loss incurred during learning.

This work focuses on the *Stochastic Rising Bandits* (SRB), a specific instance of the *rested* bandit (Tekin and Liu, 2012) setting in which the expected reward of an arm increases according to the number of times it has been pulled. Online learning in such a scenario has been recently addressed from a regret minimization perspective by Metelli et al. (2022), in which the authors provide no-regret algorithms for the SRB setting in both the rested and restless cases. The SRB setting models several real-world scenarios where arms improve their performance over time. A classic example is the so-called *Combined Algorithm Selection and Hyperparameter optimization* (CASH, Thornton et al., 2013; Kotthoff et al., 2017; Erickson et al., 2020; Li et al., 2020; Zöller and Huber, 2021), a problem of paramount importance in *Automated Machine Learning* (AutoML, Feurer et al., 2015; Yao et al., 2018; Hutter et al., 2019; Mussi et al., 2023). In CASH, the goal is to identify the *best learning algorithm* together with the *best hyperparameter* configuration for a given ML task (e.g., classification or regression). In this problem, every arm represents a hyperparameter tuner acting on a specific learning algorithm. A pull corresponds to a unit of time/computation in which we improve (on average) the hyperparameter configuration (via the tuner) for the corresponding learning algorithm. CASH was handled in a bandit *Best Arm Identification* (BAI) fashion in Li et al. (2020) and Cella et al. (2021). The former handles the problem by considering rising rested bandits with *deterministic* rewards, failing to represent the intrinsic uncertain nature of such processes. Instead, the latter, while allowing stochastic rewards, assumes that the expected rewards evolve according to a *known* parametric functional class, whose parameters have to be learned.

Original Contributions In this paper, we address the design of algorithms to solve the BAI task in the rested SRB setting when a *fixed budget* is provided.¹ More specifically, we are interested in algorithms guaranteeing a sufficiently large probability of recommending the arm with the largest expected reward *at the end* of the time budget (as if only this arm were pulled from the beginning). The main contributions of the paper are summarized as follows:²

- We propose two *algorithms* to solve the BAI problem in the SRB setting: R-UCBE (an optimistic approach, Section 4) and R-SR (a phases-based rejection algorithm, Section 5). First, we introduce specifically designed estimators required by the algorithms (Section 3). Then, we provide guarantees on the error probability of the misidentification of the best arm.
- We derive the first error probability *lower bound* for the SRB setting, matched by our R-SR algorithm up to logarithmic factors, which highlights the complexity of the problem and the need for a sufficiently large time budget (Section 6).
- Finally, we conduct *numerical simulations* on synthetically generated data and a real-world online best model selection problem. We compare the proposed algorithms with the ones available in the bandit literature to tackle the SRB problem (Section 8).

2 Problem Formulation

In this section, we revise the Stochastic Rising Bandits (SRB) setting (Heidari et al., 2016; Metelli et al., 2022). Then, we formulate our best arm identification problem, introduce the definition of error probability, and provide a preliminary characterization of the problem.

Setting We consider a rested Multi-Armed Bandit problem $\{p_i\}_{i \in [K]}$ with a finite number of arms K .³ Let $T \in \mathbb{N}$ be the time budget of the learning process. At every round $t \in [T]$, the agent selects an arm $I_t \in [K]$, plays it, and observes a reward $X_t \sim p_{I_t, t}$, where $p_{I_t, t}$ is the reward distribution of the chosen arm I_t at round t and depends on the number of pulls performed so far $N_{i, t} = \sum_{s=1}^t \mathbb{1}_{\{I_s = i\}}$ (i.e., rested). The rewards are stochastic, formally $X_t \sim p_{I_t, t}$, where μ_{I_t} is the expected reward of arm I_t and ϵ_t is a zero-mean σ -subgaussian noise, conditioned to the past.⁴ As customary in the bandit literature, we assume that the rewards are bounded in expectation, formally $\mu_i \in [0, 1]$ for all $i \in [K]$. As in (Metelli et al., 2022), we focus on a particular family of rested bandits in which the expected rewards are *non-decreasing* and *concave* in expectation.

¹We focus on the rested setting only and, thus, from now on, we will omit “rested” in the setting name.

²Additional motivating examples are discussed in Appendix A. The proofs of all the statements in this work are provided in Appendix C.

³Let $y, z \in \mathbb{R}$, we denote with Jy, z : y, z , and with Jy, z : y, z .

⁴A zero-mean random variable X is σ -subgaussian if it holds $E_X e^{sX} \leq e^{\frac{\sigma^2 s^2}{2}}$ for every $s \in \mathbb{R}$.

Assumption 2.1 (Non-decreasing and concave expected rewards). Let \mathcal{M} be a rested MAB, defining $\mu_i(n) = \mathbb{E}[r_i(n)]$, for every $n \in \mathbb{N}$ and every arm $i \in [K]$ the rewards are non-decreasing and concave, formally:

$$\text{Non-decreasing: } \mu_i(n) \leq \mu_i(n+1); \quad \text{Concave: } \mu_i(n+1) - \mu_i(n) \leq \mu_i(n) - \mu_i(n-1).$$

Intuitively, the $\mu_i(n)$ represents the *increment* of the real process $r_i(n)$ evaluated at the n^{th} pull. Notice that concavity emerges in several settings, such as the best model selection and economics, representing the decreasing marginal returns (Lehmann et al., 2001; Heidari et al., 2016).⁵

Learning Problem The goal of BAI in the SRB setting is to select the arm providing the largest expected reward with a large enough probability given a fixed budget $T \in \mathbb{N}$. Unlike the stationary BAI problem (Audibert et al., 2010), in which the optimal arm is not changing, in this setting, we need to decide *when* to evaluate the optimality of an arm. We define optimality by considering the largest expected reward at time T . Formally, given a time budget T , the optimal arm $i^*(T) \in [K]$, which we assume unique, satisfies:

$$i^*(T) = \arg \max_{i \in [K]} \mu_i(T);$$

where we highlighted the dependence on T as, with different values of the budget, $i^*(T)$ may change. Let $i \in [K] \setminus \{i^*(T)\}$ be a suboptimal arm, we define the suboptimality gap as $\Delta_i(T) = \mu_{i^*(T)}(T) - \mu_i(T)$. We employ the notation $i^{\#}(T) \in [K]$ to denote the i^{th} best arm at time T (arbitrarily breaking ties), i.e., we have $\mu_{i^{\#}(T)}(T) \geq \mu_{i^{\#}(T)+1}(T)$. Given an algorithm \mathcal{A} that recommends $\hat{i}^*(T) \in [K]$ at the end of the learning process, we measure its performance with the *error probability*, i.e., the probability of recommending a suboptimal arm at the end of the time budget T :

$$e_{T,\mathcal{A}} = \mathbb{P}_{\mathcal{M}}[\hat{i}^*(T) \neq i^*(T)];$$

Problem Characterization We now provide a characterization of a specific class of polynomial functions to upper bound the increments $\mu_i(n)$.

Assumption 2.2 (Bounded $\mu_i(n)$). Let \mathcal{M} be a rested MAB, there exist $c_i \geq 0$ and $\beta_i \leq 1$ such that for every arm $i \in [K]$ and number of pulls $n \in \mathbb{N}$: $\mu_i(n) \leq c_i n^{\beta_i}$.

We anticipate that, even if our algorithms will not require such an assumption, it will be used for deriving the lower bound and for providing more human-readable error probability guarantees. Furthermore, we observe that our Assumption 2.2 is fulfilled by a strict superset of the functions employed in Cella et al. (2021).

3 Estimators

In this section, we introduce the estimators of the arm expected reward employed by the proposed algorithms.⁶ A visual representation of such estimators is provided in Figure 1.

Let $\rho_i(t) \in [0, 1]$ be the fraction of samples collected up to the current time t we use to build estimators of the expected reward. We employ an *adaptive arm-dependent window size* $h_i(N_{i,t}) = \lfloor \rho_i(t) N_{i,t} \rfloor$ to include the most recent samples collected only, avoiding the use of samples that are no longer representative. We define the set of the last $h_i(N_{i,t})$ rounds in which the i^{th} arm was pulled as:

$$T_{i,t} = \{t - P \in [K] : I_{i,t} \in [N_{i,t} - h_i(N_{i,t}), N_{i,t}]\};$$

Furthermore, the set of the pairs of rounds $(t, t-1)$ belonging to the sets of the last and second-last $h_i(N_{i,t})$ -wide windows of the i^{th} arm is defined as:

$$S_{i,t} = \left\{ (t, t-1) \in [K] \times [K] : I_{i,t} \in [N_{i,t} - h_i(N_{i,t}), N_{i,t}] \text{ and } I_{i,t-1} \in [N_{i,t-1} - h_i(N_{i,t-1}), N_{i,t-1}] \right\};$$

In the following, we design a *pessimistic* estimator and an *optimistic* estimator of the expected reward of each arm at the end of the budget time T , i.e., $\mu_{i^*(T)}$.⁷

⁵A discussion on the non-learnability of the problem when concavity assumption does not hold is provided in Appendix E.

⁶The estimators are adaptations of those presented by Metelli et al. (2022) to handle a fixed time budget T .

⁷Naïvely computing the estimators from their definition requires $\mathcal{O}(\rho_i h_i(N_{i,t}))$ number of operations. An efficient way to incrementally update them, using $\mathcal{O}(1)$ operations, is provided in Appendix B.

Pessimistic Estimator The *pessimistic* estimator $\hat{\mu}_i^p N_{i;t-1} q$ is a negatively biased estimate of $\mu_i p T q$ obtained assuming that the function $\mu_i p q$ remains constant up to time T . This corresponds to the minimum admissible value under Assumption 2.1 (due to the *Non-decreasing* constraint). This estimator is an average of the last $h p N_{i;t-1} q$ observed rewards collected from the i^{th} arm, formally:

$$\hat{\mu}_i^p N_{i;t-1} q : \frac{1}{h p N_{i;t-1} q} \sum_{\tau \in \mathcal{P}_{T_i;t}} x_{\tau} : \quad (1)$$

The estimator enjoys the following concentration property.

Lemma 3.1 (Concentration of $\hat{\mu}_i^p$). *Under Assumption 2.1, for every $a_i \geq 0$, simultaneously for every arm $i \in \mathcal{K}$ and number of pulls $n \in \mathbb{N}^+$, with probability at least $1 - 2TK e^{-a_i/2}$ it holds that:*

$$\hat{\mu}_i^p n q \leq \hat{\mu}_i^p n q \leq \mu_i^p n q \leq \mu_i^p n q \leq \hat{\mu}_i^p n q ;$$

where $\hat{\mu}_i^p n q : \frac{a_i}{h p n q}$ and $\hat{\mu}_i^p n q : \frac{1}{2} p 2 T n h p n q 1 q ; \mu_i^p n q h p n q 1 q ;$

As supported by intuition, we observe that the estimator is affected by a negative bias that is represented by $\hat{\mu}_i^p n q$ that vanishes as $n \rightarrow \infty$ under Assumption 2.1 with a rate that depends on the increment functions $\mu_i p q$. Considering also the term $\hat{\mu}_i^p n q$ and recalling that $h p n q \leq O p n q$, under Assumption 2.2, the overall concentration rate is $O p n^{-1/2} c T n q$.

Optimistic Estimator The *optimistic* estimator $\hat{\mu}_i^T p N_{i;t-1} q$ is a positively biased estimation of $\mu_i p T q$ obtained assuming that function $\mu_i p q$ linearly increases up to time T . This corresponds to the maximum value admissible under Assumption 2.1 (due to the *Concavity* constraint). The estimator is constructed by adding to the pessimistic estimator $\hat{\mu}_i^p N_{i;t-1} q$ an estimate of the increment occurring in the next step up to T . The latter uses the last $2h p N_{i;t-1} q$ samples to obtain an upper bound of such growth thanks to the concavity assumption, formally:

$$\hat{\mu}_i^T p N_{i;t-1} q : \hat{\mu}_i^p N_{i;t-1} q + \frac{p T}{j q ; k q \mathcal{P}_{S_{i;t}}} \sum_{j,k} \frac{x_j x_k}{h p N_{i;t-1} q^2} : \quad (2)$$

The estimator displays the following concentration guarantee.

Lemma 3.2 (Concentration of $\hat{\mu}_i^T$). *Under Assumption 2.1, for every $a_i \geq 0$, simultaneously for every arm $i \in \mathcal{K}$ and number of pulls $n \in \mathbb{N}^+$, with probability at least $1 - 2TK e^{-a_i/10}$ it holds that:*

$$\hat{\mu}_i^T p n q \leq \hat{\mu}_i^T p n q \leq \mu_i^T p n q \leq \hat{\mu}_i^T p n q ;$$

where $\hat{\mu}_i^T p n q : \frac{p T}{n h p n q 1 q} \frac{a_i}{h p n q^3}$ and $\hat{\mu}_i^T p n q : \frac{1}{2} p 2 T n h p n q 1 q ; \mu_i^T p n q 2 h p n q 1 q ;$

Differently from the pessimistic estimation, the optimistic one displays a positive vanishing bias $\hat{\mu}_i^T p n q$. Under Assumption 2.2, we observe that the overall concentration rate is $O p T n^{-3/2} c T n q$.

4 Optimistic Algorithm: Rising Upper Confidence Bound Exploration

In this section, we introduce and analyze Rising Upper Confidence Bound Exploration (R-UCBE) an *optimistic* error probability minimization algorithm for the SRB setting with a fixed budget. The algorithm explores by means of a UCB-like approach and, for this reason, makes use of the optimistic estimator $\hat{\mu}_i^T$ plus a bound to account for the uncertainty of the estimation. In R-UCBE, the choice of considering the optimistic estimator is natural and obliged since the pessimistic estimator is affected by negative bias and cannot be used to deliver optimistic estimates.

Algorithm The algorithm, whose pseudo-code is reported in Algorithm 1, requires as input an exploration parameter $a \in \mathbb{R}^+$, the window size $\lfloor \frac{a}{2} \rfloor$, the time budget T , and the number of

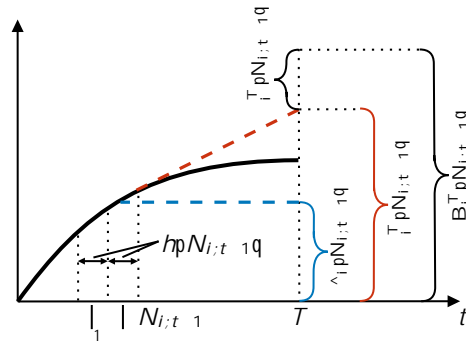


Figure 1: Graphical representation of the pessimistic $\hat{\mu}_i^p N_{i;t-1} q$ and the optimistic $\hat{\mu}_i^T p N_{i;t-1} q$ estimators.

Algorithm 1: R-UCBE.

Input : Time budget T , Number of arms K ,
Window size w , Exploration parameter a

- 1 Initialize $N_{i;0} = 0$,
- 2 $B_i^T \leftarrow \frac{a}{w}$ $\forall i \in \{1, \dots, K\}$
- 3 **for** $t \in \{1, \dots, T\}$ **do**
- 4 Compute $I_t \leftarrow \arg \max_{i \in \{1, \dots, K\}} B_i^T \frac{a}{N_{i;t}}$
- 5 Pull arm I_t and observe X_t
- 6 $N_{I_t;t} \leftarrow N_{I_t;t-1} + 1$
- 7 $B_{I_t;t} \leftarrow \frac{a}{N_{I_t;t}}$ $\forall i = I_t$
- 8 Update $\bar{I}_t \leftarrow \frac{a}{N_{I_t;t}}$
- 9 Update $\bar{I}_t \leftarrow \frac{a}{N_{I_t;t}}$
- 10 Compute $B_{I_t;t}^T \leftarrow \frac{a}{N_{I_t;t}}$ $\bar{I}_t \leftarrow \frac{a}{N_{I_t;t}}$ $\bar{I}_t \leftarrow \frac{a}{N_{I_t;t}}$
- 11 **end**
- 12 Recommend $\hat{P} \leftarrow \arg \max_{i \in \{1, \dots, K\}} B_i^T \frac{a}{N_{i;T}}$

Algorithm 2: R-SR.

Input : Time budget T , Number of arms K ,
Window size w

- 1 Initialize $t \leftarrow 1$; $N_0 = 0$; $X_0 = \{1, \dots, K\}$
- 2 **for** $j \in \{1, \dots, K\}$ **do**
- 3 **for** $i \in X_{j-1}$ **do**
- 4 **for** $l \in \{1, \dots, N_j\}$ **do**
- 5 Pull arm i and observe X_t
- 6 $t \leftarrow t + 1$
- 7 **end**
- 8 Update $\hat{\mu}_i \leftarrow \frac{a}{N_j}$
- 9 **end**
- 10 Define $\bar{I}_j \leftarrow \arg \min_{i \in X_{j-1}} \hat{\mu}_i \frac{a}{N_j}$
- 11 Update $X_j = X_{j-1} \setminus \bar{I}_j$
- 12 **end**
- 13 Recommend $\hat{P} \leftarrow X_{K-1}$ (unique)

First of all, we notice that the error probability $e_{T,R-UCBE}$ presented in Theorem 4.2 holds under the condition that the time budget T fulfills Equation (8). We defer a more detailed discussion on this condition to Remark 5.1, where we show that the existence of a finite value of T fulfilling Equation (8) is ensured under mild conditions.

Let us remark that term $H_1; \frac{a}{w}$ characterizes the complexity of the SRB setting, corresponding to term H_1 of Audibert et al. (2010) for the classical BAI problem when $w = 2$. As expected, in the small- w regime (i.e., $w \leq 3\{2\}$), looking at the dependence of $H_1; \frac{a}{w}$ on w , we realize that the complexity of a problem decreases as the parameter w increases. Indeed, the larger w , the faster the expected reward reaches a stationary behavior. Nevertheless, even in the large- w regime (i.e., $w \geq 3\{2\}$), the complexity of the problem is governed by $H_1; 2\{3\} \frac{a}{w}$, leading to an error probability larger than the corresponding one for BAI in standard bandits (Audibert et al., 2010). This can be explained by the fact that R-UCBE uses the optimistic estimator that, as shown in Section 3, enjoys a slower concentration rate compared to the standard sample mean, even for stationary bandits.

This two-regime behavior has an interesting interpretation when comparing Corollary 4.2 with Theorem 4.1. Indeed, $w = 3\{2\}$ is the break-even threshold in which the two terms of the l.h.s. of Equation (6) have the same convergence rate. Specifically, the term $\frac{a}{w}$ takes into account the expected rewards growth (i.e., the bias in the estimators), while $\frac{a}{w}$ considers the uncertainty in the estimations of the R-UCBE algorithm (i.e., the variance). Intuitively, when the expected reward function displays a slow growth (i.e., $\frac{a}{w} \propto cn$ with $w \geq 3\{2\}$), the bias term $\frac{a}{w}$ dominates the variance term $\frac{a}{w}$ and the value of a changes accordingly. Conversely, when the variance term $\frac{a}{w}$ is the dominant one (i.e., $\frac{a}{w} \propto cn$ with $w \leq 3\{2\}$), the threshold a is governed by the estimation uncertainty, being the bias negligible.

As common in optimistic algorithms for BAI (Audibert et al., 2010), setting a theoretically sound value of exploration parameter a (i.e., computing a), requires additional knowledge of the setting, namely the complexity index $H_1; \frac{a}{w}$.⁸ In the next section, we propose an algorithm that relaxes this requirement.

5 Phase-Based Algorithm: Rising Successive Rejects

In this section, we introduce the **Rising Successive Rejects** (R-SR), a phase-based solution inspired by the one proposed by Audibert et al. (2010), which overcomes the drawback of R-UCBE of requiring knowledge of $H_1; \frac{a}{w}$.

Algorithm R-SR, whose pseudo-code is reported in Algorithm 2, takes as input the time budget T and the number of arms K . At first, it initializes the set of the active arms X_0 with all the available arms (Line 1). This set will contain the arms that are still eligible candidates to be recommended. The entire process proceeds through $K - 1$ phases. More specifically, during the j^{th} phase, the arms

⁸We defer the empirical study of the sensitivity of a to Section 8.

still remaining in the active arms set X_{j-1}^R are played (Line 5) for $N_j - N_{j-1}$ times each, where:

$$N_j := \frac{1}{\overline{\log_p K q}} \frac{T - K}{K - 1} \frac{1}{j}; \quad (9)$$

and $\overline{\log_p K q} := \frac{1}{2} \sum_{i=2}^K \frac{1}{i}$. At the end of each phase, the arm with the smallest value of the pessimistic estimator $\hat{\mu}_{i|N_j} q$ is discarded from the set of active arms (Line 11). At the end of the $pK - 1$ qth phase, the algorithm recommends the (unique) arm left in X_{K-1} (Line 13).

It is worth noting that R-SR makes use of the pessimistic estimator $\hat{\mu}_{i|N} q$. Even if both estimators defined in Section 3 are viable for R-SR, the choice of using the pessimistic estimator is justified by its better concentration rate $O(pn^{-1/2}q)$ compared to that of the optimistic estimator $O(pTn^{-3/2}q)$, being $n \propto T$ (see Section 3).

Note that the phase lengths are the ones adopted by Audibert et al. (2010). This choice allows us to provide theoretical results without requiring domain knowledge (still under a large enough budget). An optimized version of N_j may be derived assuming full knowledge of the gaps $\mu_i p T q$, but, unfortunately, such a hypothetical approach would have similar drawbacks as R-UCBE.

Bound on the Error Probability of R-SR The following theorem provides the guarantee on the error probability for the R-SR algorithm.

Theorem 5.1. *Under Assumptions 2.1 and 2.2, if the time budget T satisfies:*

$$T \geq 2^{-1} c^{-1} \overline{\log_p K q}^{-1} \max_{i \in J_{2:K}} i^{-1} \mu_i p T q^{-1}; \quad (10)$$

then, the error probability of R-SR is bounded by:

$$e_T p R\text{-SR} q \leq \frac{K p K - 1 q}{2} \exp \left(- \frac{T - K}{\overline{\log_p K q} H_2 p T q} \right);$$

where $H_2 p T q := \max_{i \in J_{2:K}} i^{-1} \mu_i p T q^{-2}$ and $\overline{\log_p K q} := \frac{1}{2} \sum_{i=2}^K \frac{1}{i}$.

Similar to the R-UCBE, the complexity of the problem is characterized by term $H_2 p T q$ that, for the standard MAB setting, reduces to the H_2 term of Audibert et al. (2010). Furthermore, when the condition of Equation (10) on the time budget T is satisfied, the error probability coincides with that of the SR algorithm for standard MABs (apart for constant terms). The following remark elaborates on the conditions of Equations (8) and (10) about the minimum requested time budget.

Remark 5.1 (About the minimum time budget T). *To satisfy the e_T bounds presented in Corollary 4.2 and Theorem 5.1, R-UCBE and R-SR require the conditions provided by Equations (8) and (10) about the time budget T , respectively. First, let us notice that if the suboptimal arms converge to an expected reward different from that of the optimal arm as $T \rightarrow \infty$, it is always possible to find a finite value of $T \rightarrow \infty$ such that these conditions are fulfilled. Formally, assume that there exists $T_0 \rightarrow \infty$ and that for every $T \geq T_0$ we have that for all suboptimal arms $i \in I p T q$ it holds that $\mu_i p T q \neq \mu^* p T q$. In such a case, the l.h.s. of Equations (8) and (10) are upper bounded by a function of $\mu^* p T q$ and are independent on T . Instead, if a suboptimal arm converges to the same expected reward as the optimal arm when $T \rightarrow \infty$, the identification problem is more challenging and, depending on the speed at which the two arms converge as a function of T , might slow down the learning process arbitrarily. This should not surprise as the BAI problem becomes non-learnable even in standard (stationary) MABs when multiple optimal arms are present (Heide et al., 2021).*

6 Lower Bound

In this section, we investigate the complexity of the BAI problem for SRBs with a fixed budget.

Minimum time budget T We show that, under Assumptions 2.1 and 2.2, any algorithm requires a minimum time budget T to be guaranteed to identify the optimal arm, even in a deterministic setting.

Theorem 6.1. *For every algorithm A , there exists a deterministic SRB satisfying Assumptions 2.1 and 2.2 with $K \geq 8^{1/p} - 1$ such that the optimal arm $i^* p T q$ cannot be identified for some time budgets T unless:*

$$T \geq H_{1,1/p} \mu_{i^*} p T q^{-1} \frac{1}{\mu_{i^*} p T q^{-1}}; \quad (11)$$

	Error Probability $e_{T,p,q}$	Time Budget T
SRB	$\frac{1}{4} \exp\left(-\frac{8T}{i \cdot pTq \cdot \frac{1}{i \cdot pTq}}\right)$	$\frac{1}{i \cdot pTq \cdot \frac{1}{i \cdot pTq}}$
R-UCBE	$2TK \exp\left(-\frac{a}{10}\right)$	$\begin{cases} c^{\frac{1}{3}} p^{\frac{1}{3}} 2^{\frac{1}{3}} q^{\frac{1}{3}} \frac{1}{i \cdot pTq \cdot \frac{1}{i \cdot pTq}} \cdot pK \cdot 1q^{\frac{1}{3}} & \text{if } P \in \{1, 3\} \{2q\} \\ c^{\frac{2}{3}} p^{\frac{2}{3}} 2^{\frac{2}{3}} q^{\frac{2}{3}} \frac{1}{i \cdot pTq \cdot \frac{1}{i \cdot pTq}} \cdot pK \cdot 1q^{\frac{2}{3}} & \text{if } P \in \{2, 3\} \{8q\} \end{cases}$
R-SR	$\frac{KpK}{2} 1q \exp\left(-\frac{T}{8 \cdot \frac{1}{\log pKq \max_{i \in [p,K]} i \cdot pTq}}\right)$	$2^{\frac{1}{3}} c^{\frac{1}{3}} \log pKq^{\frac{1}{3}} \max_{i \in [p,K]} i \cdot pTq^{\frac{1}{3}}$

Table 1: Bounds on the time budget and error probability: lower for the setting and upper for the algorithms.

Theorem 6.1 formalizes the intuition that any of the suboptimal arms must be pulled a sufficient number of times to ensure that, if pulled further, it cannot become the optimal arm. Indeed, in the proof of Theorem 6.1, we show that each suboptimal arm $i \in [p, Tq]$ has to be pulled at least an expected number of times of order $\text{Er}N_{i,pTq} \asymp \frac{1}{i \cdot pTq} \log \frac{1}{i \cdot pTq}$ (see Equation 36). It is worth comparing this bound on the time budget with the corresponding conditions on the minimum time budget requested by Equations (8) and (10) for R-UCBE and R-SR, respectively. Regarding R-UCBE, we notice that the minimum admissible time budget in the small- p regime is of order $H_{1,1} \frac{1}{pTq} \log \frac{1}{pTq}$ which is larger than term $H_{1,1} \frac{1}{pTq}$ of Equation (11).⁹ Similarly, in the large- p regime (i.e., $p \geq 3\{2\}$), the R-UCBE requirement is of order $H_{1,2} \frac{1}{pTq} \log \frac{1}{pTq} \asymp H_{1,2} \frac{1}{pTq}$ which is larger than the term of Theorem 6.1 since $\frac{1}{pTq} \log \frac{1}{pTq} \geq \frac{1}{pTq}$. Concerning R-SR, it is easy to show that $H_{1,1} \frac{1}{pTq} \log \frac{1}{pTq} \asymp \max_{i \in [p,K]} i \cdot pTq \log \frac{1}{i \cdot pTq}$, apart from logarithmic terms, by means of the argument provided by (Audibert et al., 2010, Section 6.1). Thus, up to logarithmic terms, Equation (10) provides a tight condition on the minimum budget.

Error Probability Lower Bound We now present a lower bound on the error probability that every algorithm performing BAI in the SRB setting suffers.

Theorem 6.2. *For every algorithm A run with a time budget T fulfilling Equation (11), there exists a SRB satisfying Assumptions 2.1 and 2.2 with $K \asymp 8 \log \frac{1}{i \cdot pTq}$ such that the error probability is lower bounded by:*

$$e_{T,p,q} \asymp \frac{1}{4} \exp\left(-\frac{8T}{2H_{1,2} pTq}\right); \text{ where } H_{1,2} pTq \asymp \frac{1}{i \cdot pTq}.$$

Some comments are in order. First, we stated the lower bound for the case in which the minimum time budget satisfies the inequality of Theorem 6.1, which is a necessary condition for identifying the optimal arm. Second, the lower bound on the error probability matches, up to logarithmic factors, that of our R-SR, suggesting the superiority of this algorithm compared to R-UCBE. Finally, provided that the identifiability condition of Equation (11), such a result corresponds to that of the standard (stationary) MABs (Audibert et al., 2010; Kaufmann et al., 2016). A summary of all the bounds provided in the paper is presented in Table 1.

7 Related Works

In this section, we summarize the relevant literature related both to the works focusing on the best arm identification problem and rested bandits. The SRB setting was proposed by Heidari et al. (2016) and analyzed by Metelli et al. (2022) from the regret minimization perspective.

Best Arm Identification in Stochastic Rising Bandits As highlighted in Section 1, the works mostly related to ours are the ones by Li et al. (2020) and Cella et al. (2021). They both focus on the BAI problem in the rested setting, given a fixed-budget. More specifically, Li et al. (2020) consider rising rested bandits in which the reward function of each arm increases as it is pulled. However, they limit to deterministic arms and, thus, fail to deal with the intrinsic stochasticity of the real-world

⁹See Lemma C.11.

processes they want to model. Instead, Cella et al. (2021) deal with the problem of identifying the arm with the smallest loss in a setting where the losses incurred by selecting an arm decrease over time. It is easy to show that such a setting can be transformed straightforwardly in the SRB one. However, the authors develop two algorithms whose theoretical guarantees hold under the assumption that the expected loss follows a specific known parametric functional form, whose parameters are to be estimated. This constitutes a major limitation to the presented work since checking such an assumption is not feasible in real-world settings.

Best Arm Identification The pure exploration and BAI problems have been first introduced by Bubeck et al. (2009), while algorithms able to learn in such a setting have been provided by Audibert et al. (2010). The work by Gabillon et al. (2012) proposes a unified approach to deal with stochastic best arm identification problems by having either a fixed budget or fixed confidence. However, the stochastic algorithms developed in this line of research only provide theoretical guarantees in settings where the expected reward is stationary over the pulls. Abbasi-Yadkori et al. (2018) propose a method able to handle both the stochastic and adversarial cases, but they do not make explicit use of the properties (e.g., increasing nature) of the expected reward. Finally, (Garivier and Kaufmann, 2016; Kaufmann et al., 2016; Carpentier and Locatelli, 2016) analyze the problem of BAI from the lower bound perspective.

Rested Bandits Bandit settings in which the evolution of an arm reward depends on the number of times the arm has been pulled, such as the one analyzed in our paper, are generally referred to as *rested*. A first general formulation of the rested bandit setting appeared in the work by Tekin and Liu (2012) and was further discussed by Mintz et al. (2020) and Seznec et al. (2020). In these works, the evolution of the expected reward of each arm is regulated by a Markovian process that is assumed to visit the same state multiple times. This is not the case for the rising bandits, where the arm expected rewards continuously increase over the time budget. Finally, a specific instance of the rested bandits is constituted by the *rotting* bandits (Levine et al., 2017; Seznec et al., 2019, 2020), in which the expected payoff for a given arm decreases with the number of pulls. However, as pointed out by Metelli et al. (2022), techniques developed for this setting cannot be directly translated into ours, due to the inherently different nature of the problem.

8 Numerical Validation

In this section, we provide a numerical validation of R-UCBE and R-SR. We compare them with state-of-the-art bandit baselines designed for stationary and non-stationary BAI in a synthetic setting, and we evaluate the sensitivity of R-UCBE to its exploration parameter a . Additional details about the experiments presented in this section are available in Appendix F. Additional experimental results on both synthetic settings and in a real-world experiment are available in Appendix G. The code to reproduce the experiments can be found at <https://github.com/MontenegroAlessandro/BestArmIdSRB>.

Baselines We compare our algorithms against a wide range of solutions for BAI:

- RR: uniformly pulls all the arms until the budget ends in a *round-robin* fashion and, in the end, makes a recommendation based on the empirical mean of their reward over the collected samples;
- RR-SW: makes use of the same exploration strategy as RR to pull arms but makes a recommendation based on the empirical mean over the last $\frac{T}{K}$ collected samples from an arm.¹⁰
- UCB-E and SR (Audibert et al., 2010): algorithms for the stationary BAI problem;
- Prob-1 (Abbasi-Yadkori et al., 2018): an algorithm dealing with the adversarial BAI setting;
- ETC and Rest-Sure (Cella et al., 2021): algorithms developed for the decreasing loss BAI setting, that we converted through a linear transformation of the reward.

The hyperparameters required by the above methods have been set as prescribed in the original papers. For both our algorithms and RR-SW, we set $\epsilon = 0.25$.

Setting To assess the quality of the recommendation $\hat{\mu}_{pTq}$ provided by our algorithms, we consider a synthetic SRB setting with $K = 5$ and $\epsilon = 0.01$. Figure 2 shows the evolution of the expected values of the arms w.r.t. the number of pulls. In this setting, the optimal arm changes depending on whether $T \leq 185$ or $T > 185$.¹¹ Thus, when the time budget is close to that value, the problem is more challenging since the optimal and second-best arms expected rewards are close to each other. For this reason, the BAI algorithms are less likely to provide a correct recommendation

¹⁰The formal description of this baseline, as well as its theoretical analysis, is provided in Appendix D.

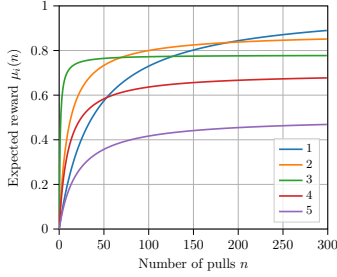


Figure 2: Expected values for the arms of the synthetic setting.

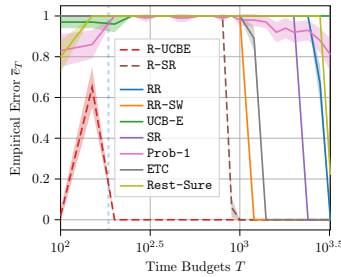


Figure 3: Empirical error rate for the synthetically generated setting (100 runs, mean 95% c.i.).

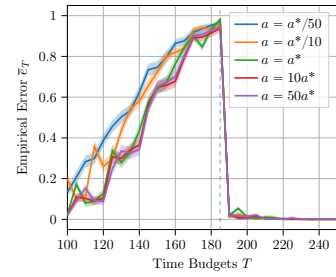


Figure 4: Empirical error rate for the R-UCBE at different a (1000 runs, mean 95% c.i.).

than for time budgets for which the two expected rewards are well separated. We compare the analyzed algorithms A in terms of empirical error $e_{T,p}A$ (the smaller, the better), i.e., the empirical counterpart of $e_{T,p}A$ averaged over 100 runs, considering time budgets $T \in \{100, 3200\}$.

Results The empirical error probability provided by the analyzed algorithms in the synthetically generated setting is presented in Figure 3. We report with a dashed vertical blue line at $T = 185$, i.e., the budgets after which the optimal arm no longer changes. Before such a budget, all the algorithms provide large errors (i.e., $e_{T,p}A \geq 0.2$). However, R-UCBE outperforms the others by a large margin, suggesting that an optimistic estimator might be advantageous when the time budget is small. Shortly after $T = 185$, R-UCBE starts providing the correct suggestion consistently. R-SR begins to identify the optimal arm (i.e., with $e_{T,p}R-SR \leq 0.05$) for time budgets $T \geq 1000$. Nonetheless, both algorithms perform significantly better than the baseline algorithms used for comparison.

Sensitivity Analysis for the Exploration Parameter of R-UCBE We perform a sensitivity analysis on the exploration parameter a of R-UCBE. Such a parameter should be set to a value less or equal to a^* , and the computation of the latter is challenging. We tested the sensitivity of R-UCBE to this hyperparameter by looking at the error probability for $a \in \{a^*/50, a^*/10, a^*, 10a^*, 50a^*\}$. Figure 4 shows the empirical errors of R-UCBE with different parameters a , where the blue dashed vertical line denotes the last time the optimal arm changes over the time budget. It is worth noting how, even in this case, we have two significantly different behaviors before and after such a time. Indeed, if $T \approx 185$, we have that misspecification with larger values than a^* does not significantly impact the performance of R-UCBE, while smaller values slightly decrease the performance. Conversely, for $T \leq 185$ learning with different values of a seems not to impact the algorithm performance significantly. This corroborates the previous results about the competitive performance of R-UCBE.

9 Discussion and Conclusions

This paper introduces the BAI problem with a fixed budget for the Stochastic Rising Bandits setting. Notably, such setting models many real-world scenarios in which the reward of the available options increases over time, and the interest is on the recommendation of the one having the largest expected rewards after the time budget has elapsed. In this setting, we presented two algorithms, namely R-UCBE and R-SR providing theoretical guarantees on the error probability. R-UCBE is an optimistic algorithm requiring an exploration parameter whose optimal value requires prior information on the setting. Conversely, R-SR is a phase-based solution that only requires the time budget to run. We established lower bounds for the error probability an algorithm suffers in such a setting, which is matched by our R-SR, up to logarithmic factors. Furthermore, we showed how a requirement on the minimum time budget is unavoidable to ensure the identifiability of the optimal arm. Finally, we validate the performance of the two algorithms in both synthetically generated and real-world settings. A possible future line of research is to derive an algorithm balancing the tradeoff between theoretical guarantees on the e_T and the chance of providing such guarantees with lower time budgets.

Acknowledgements

This paper is supported by PNRR-PE-AI FAIR project funded by the NextGeneration EU program.

References

- Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2):235–256, 2002.
- Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in multi-armed bandits problems. In *Proceedings of the Algorithmic Learning Theory Conference (ALT)*, volume 5809, pages 23–37, 2009.
- Cem Tekin and Mingyan Liu. Online learning of rested and restless bandits. *IEEE Transaction on Information Theory*, 58(8):5588–5611, 2012.
- Alberto Maria Metelli, Francesco Trovò, Matteo Pirola, and Marcello Restelli. Stochastic rising bandits. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 15421–15457, 2022.
- Chris Thornton, Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Auto-weka: Combined selection and hyperparameter optimization of classification algorithms. In *Proceedings of the ACM Conference of Knowledge Discovery and Data Mining (SIGKDD)*, pages 847–855, 2013.
- Lars Kotthoff, Chris Thornton, Holger H Hoos, Frank Hutter, and Kevin Leyton-Brown. Auto-weka 2.0: Automatic model selection and hyperparameter optimization in weka. *Journal of Machine Learning Research*, 18:1–5, 2017.
- Nick Erickson, Jonas Mueller, Alexander Shirkov, Hang Zhang, Pedro Larroy, Mu Li, and Alexander Smola. Autogluon-tabular: Robust and accurate automl for structured data. *arXiv preprint arXiv:2003.06505*, 2020.
- Yang Li, Jiawei Jiang, Jinyang Gao, Yingxia Shao, Ce Zhang, and Bin Cui. Efficient automatic CASH via rising bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4763–4771, 2020.
- Marc-André Zöller and Marco F Huber. Benchmark and survey of automated machine learning frameworks. *Journal of Artificial Intelligence Research*, 70:409–472, 2021.
- Matthias Feurer, Aaron Klein, Katharina Eggensperger, Jost Tobias Springenberg, Manuel Blum, and Frank Hutter. Efficient and robust automated machine learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2962–2970, 2015.
- Quanming Yao, Mengshuo Wang, Yuqiang Chen, Wenyuan Dai, Yu-Feng Li, Wei-Wei Tu, Qiang Yang, and Yang Yu. Taking human out of learning applications: A survey on automated machine learning. *arXiv preprint arXiv:1810.13306*, 2018.
- Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren. *Automated machine learning: methods, systems, challenges*. Springer Nature, 2019.
- Marco Mussi, Davide Lombarda, Alberto Maria Metelli, Francesco Trovò, and Marcello Restelli. Arlo: A framework for automated reinforcement learning. *Expert Systems with Applications*, 224: 119883, 2023.
- Leonardo Cella, Massimiliano Pontil, and Claudio Gentile. Best model identification: A rested bandit formulation. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 139, pages 1362–1372, 2021.
- Hoda Heidari, Michael J. Kearns, and Aaron Roth. Tight policy regret bounds for improving and decaying bandits. In *Proceeding of the International Joint Conference on Artificial Intelligence (IJCAI)*, pages 1562–1570, 2016.
- Benny Lehmann, Daniel Lehmann, and Noam Nisan. Combinatorial auctions with decreasing marginal utilities. In *ACM Proceedings of the Conference on Electronic Commerce (EC)*, pages 18–28, 2001.

- Jean-Yves Audibert, Sébastien Bubeck, and Rémi Munos. Best arm identification in multi-armed bandits. In *Proceedings of the Conference on Learning Theory (COLT)*, pages 41–53, 2010.
- Rianne De Heide, James Cheshire, Pierre Ménard, and Alexandra Carpentier. Bandits with many optimal arms. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 22457–22469, 2021.
- Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm identification in multi-armed bandit models. *Journal of Machine Learning Research*, 17:1:1–1:42, 2016.
- Victor Gabillon, Mohammad Ghavamzadeh, and Alessandro Lazaric. Best arm identification: A unified approach to fixed budget and fixed confidence. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 3221–3229, 2012.
- Yasin Abbasi-Yadkori, Peter L. Bartlett, Victor Gabillon, Alan Malek, and Michal Valko. Best of both worlds: Stochastic & adversarial best-arm identification. In *Proceedings of the Conference on Learning Theory (COLT)*, volume 75, pages 918–949, 2018.
- Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence. In *Proceedings of the Conference on Learning Theory (COLT)*, volume 49, pages 998–1027, 2016.
- Alexandra Carpentier and Andrea Locatelli. Tight (lower) bounds for the fixed budget best arm identification bandit problem. In *Proceedings of the Conference on Learning Theory (COLT)*, volume 49, pages 590–604, 2016.
- Yonatan Mintz, Anil Aswani, Philip Kaminsky, Elena Flowers, and Yoshimi Fukuoka. Nonstationary bandits with habituation and recovery dynamics. *Operations Research*, 68(5):1493–1516, 2020.
- Julien Seznec, Pierre Ménard, Alessandro Lazaric, and Michal Valko. A single algorithm for both restless and rested rotting bandits. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 108, pages 3784–3794, 2020.
- Nir Levine, Koby Crammer, and Shie Mannor. Rotting bandits. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 3074–3083, 2017.
- Julien Seznec, Andrea Locatelli, Alexandra Carpentier, Alessandro Lazaric, and Michal Valko. Rotting bandits are no harder than stochastic ones. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 89, pages 2564–2572, 2019.

A Additional Motivating Examples

In this appendix, we provide two additional motivating examples to better understand and appreciate the SRB setting.

Selection of Athletes for Competitions Consider the role of a professional trainer for a team, having several athletes (i.e., our arms) to train in order to increase their performances. The final goal is to select a single athlete to represent the team in a competition. The performances of athletes increase when the trainer properly follows them. However, a trainer can follow just one athlete at a time. The trainer can be modeled as an agent performing best-arm (athlete) identification, and the athletes represent the arms that increase their payoffs (i.e., performance) when pulled (i.e., when the trainer follows them).

Online Best Model Selection Suppose we have to choose among a set of algorithms to maximize a given index (e.g., accuracy) over a training set. In this setting, we expect that all the algorithms progressively increase (on average) the index value and converge to their optimum value with different convergence rates. Therefore, we want to identify which candidate algorithm is the most likely to reach optimal performances, given the budget, and assign the available resources (e.g., computational power or samples). In summary, this problem reduces to the identification, with the largest probability, of the algorithm that converges faster to the optimum. A real-world example of such a scenario is provided in Figure 8.

B Estimators Efficient Update

In this appendix, we describe how to implement an efficient version (i.e., fully online) of the estimators we presented in the main paper. We resort to the update developed by Metelli et al. (2022). This update provides a way to achieve an $\mathcal{O}(1)$ computational complexity at each step for the update of the estimates for the pessimistic estimator $\hat{\mu}_i$ and optimistic estimator $\hat{\nu}_i$.

With a slight abuse of notation, only in this appendix, for the sake of simplicity, we denote with $X_{i;n}$ the reward collected at the n^{th} pull from the arm i and with $h_{i;t} = \mathbb{1}_{\{N_{i;t} = t\}}$ the window size. Differently, from what we use in the paper, here the reward subscript indicates the arm i and the number of pulls of that arm n instead of the total number of pulls t we used in the definition of X_t .

More specifically, the *pessimistic* estimator $\hat{\mu}_i$ can be written as:

$$\hat{\mu}_i = \frac{\bar{a}_i}{h_{i;t}}$$

where the quantity \bar{a}_i is updated as follows:

$$\bar{a}_i \leftarrow \begin{cases} \bar{a}_i - X_{i;N_{i;t}} + X_{i;N_{i;t} - h_{i;t}} & \text{if } h_{i;t} = h_{i;t-1} \\ \bar{a}_i & \text{otherwise} \end{cases}$$

and $\bar{a}_i = 0$ as the algorithm starts.

Instead, the *optimistic* estimator $\hat{\nu}_i$ is updated using:

$$\hat{\nu}_i = \frac{1}{h_{i;t}} \left(\bar{a}_i + \frac{\tau \bar{a}_i - \bar{b}_i}{h_{i;t}} + \frac{\tau \bar{c}_i - \bar{d}_i}{h_{i;t}} \right)$$

Where the quantity \bar{a}_i is defined and updated above and \bar{b}_i , \bar{c}_i , and \bar{d}_i are updated as follows:

$$\begin{aligned} \bar{b}_i &\leftarrow \begin{cases} \bar{b}_i - X_{i;N_{i;t}} + h_{i;t} X_{i;N_{i;t} - 2h_{i;t}} & \text{if } h_{i;t} = h_{i;t-1} \\ \bar{b}_i - X_{i;N_{i;t}} + 2h_{i;t} & \text{otherwise} \end{cases} \\ \bar{c}_i &\leftarrow \begin{cases} \bar{c}_i - N_{i;t} X_{i;N_{i;t}} + \rho N_{i;t} h_{i;t} X_{i;N_{i;t} - h_{i;t}} & \text{if } h_{i;t} = h_{i;t-1} \\ \bar{c}_i - N_{i;t} X_{i;N_{i;t}} & \text{otherwise} \end{cases} \\ \bar{d}_i &\leftarrow \begin{cases} \bar{d}_i - N_{i;t} X_{i;N_{i;t}} + h_{i;t} \rho N_{i;t} h_{i;t} X_{i;N_{i;t} - 2h_{i;t}} & \text{if } h_{i;t} = h_{i;t-1} \\ \bar{d}_i - \rho N_{i;t} h_{i;t} X_{i;N_{i;t} - 2h_{i;t}} & \text{otherwise} \end{cases} \end{aligned}$$

Similarly to what is presented above, the quantities are initialized to 0 as the algorithms start.

C Proofs and Derivations

In this appendix, we provide all the proofs omitted in the main paper. For the sake of generality, we will provide the derivations for a generic choice of the window size of the estimators $h_{i;t}$ which depends on the arm $i \in \{1, \dots, K\}$ and on the round $t \in \{1, \dots, T\}$. When needed, we will particularize the choice for the case in which the window size depends on the number of pulls only $h_{i;t} = h_{p} N_{i;t-1}^{-1}$.

C.1 Proofs of Section 3

Lemma C.1. *Under Assumption 2.1, for every $i \in \{1, \dots, K\}$, $j; k \in \{1, \dots, K\}$ with $k \neq j$, it holds that:*

$$r_{i;t}^{p,j,q} \leq \frac{r_{i;t}^{p,j,q} - r_{i;t}^{p,k,q}}{j - k}.$$

Proof. Using Assumption 2.1, we get:

$$r_{i;t}^{p,j,q} - \frac{1}{j-k} \sum_{l=k}^{j-1} r_{i;t}^{p,l,q} \leq r_{i;t}^{p,j,q} - \frac{1}{j-k} \sum_{l=k}^{j-1} r_{i;t}^{p,l,q} = \frac{1}{j-k} \sum_{l=k}^{j-1} (r_{i;t}^{p,j,q} - r_{i;t}^{p,l,q}) = \frac{r_{i;t}^{p,j,q} - r_{i;t}^{p,k,q}}{j-k},$$

where the first inequality comes from the concavity of the expected reward function (Assumption 2.1), and the second equality comes from the definition of increment. \square

Lemma C.2. *For every arm $i \in \{1, \dots, K\}$, every round $t \in \{1, \dots, T\}$, and window width $1 \leq h_{i;t} \leq t N_{i;t-1}^{-1} \{2u\}$, let us define:*

$$r_{i;t}^{T,p} N_{i;t}^{-1} : \frac{1}{h_{i;t-1} N_{i;t-1}^{-1} h_{i;t-1}} \sum_{l=1}^{N_{i;t-1}} (r_{i;t}^{p,l,q} - r_{i;t}^{p,l,q} \frac{h_{i;t-1}}{h_{i;t}});$$

otherwise if $h_{i;t} = 0$, we set $r_{i;t}^{T,p} N_{i;t}^{-1} : \emptyset$. Then, $r_{i;t}^{T,p} N_{i;t}^{-1} \leq r_{i;t}^{p,T,q}$ and, if $N_{i;t-1} \geq 2$, it holds that:

$$r_{i;t}^{T,p} N_{i;t}^{-1} - r_{i;t}^{p,T,q} \leq \frac{1}{2} p \sum_{l=1}^{N_{i;t-1}} (h_{i;t-1}^{-1} - h_{i;t}^{-1}) r_{i;t}^{p,l,q} \leq \frac{1}{2} p N_{i;t-1}^{-1} (h_{i;t-1}^{-1} - h_{i;t}^{-1}):$$

Proof. Following the derivation provided above, we have for every $l \in \{2, \dots, N_{i;t-1}\}$:

$$r_{i;t}^{p,l,q} - r_{i;t}^{p,l,q} \frac{h_{i;t-1}}{h_{i;t}} \leq r_{i;t}^{p,l,q} - r_{i;t}^{p,l,q} \frac{h_{i;t-1}}{h_{i;t}} \quad (12)$$

$$\leq r_{i;t}^{p,l,q} - r_{i;t}^{p,l,q} \frac{h_{i;t-1}}{h_{i;t}} \quad (13)$$

where Equation (12) follows from Assumption 2.1, and Equation (13) is obtained from Lemma C.1. By averaging over the most recent $1 \leq h_{i;t} \leq t N_{i;t-1}^{-1} \{2u\}$ pulls, we get:

$$r_{i;t}^{p,T,q} - \frac{1}{h_{i;t-1} N_{i;t-1}^{-1} h_{i;t-1}} \sum_{l=1}^{N_{i;t-1}} (r_{i;t}^{p,l,q} - r_{i;t}^{p,l,q} \frac{h_{i;t-1}}{h_{i;t}}) = r_{i;t}^{T,p} N_{i;t}^{-1}:$$

For the bias bound, when $N_{i;t-1} \geq 2$, we retrieve:

$$r_{i;t}^{T,p} N_{i;t}^{-1} - r_{i;t}^{p,T,q} \leq \frac{1}{h_{i;t-1} N_{i;t-1}^{-1} h_{i;t-1}} \sum_{l=1}^{N_{i;t-1}} (h_{i;t-1}^{-1} - h_{i;t}^{-1}) r_{i;t}^{p,l,q} \leq r_{i;t}^{p,T,q} \quad (14)$$

$$\begin{aligned}
&\leq \frac{1}{h_{i;t}} \sum_{l=1}^{N_{i;t}-1} \sum_{j=1}^{l-1} \frac{pT}{h_{i;t}} \left| \frac{X_{i;l}}{h_{i;t}} - \frac{X_{j;l}}{h_{i;t}} \right| \\
&\leq \frac{1}{h_{i;t}} \sum_{l=1}^{N_{i;t}-1} \sum_{j=1}^{l-1} \frac{pT}{h_{i;t}} \left| \frac{X_{i;l}}{h_{i;t}} - \frac{X_{j;l}}{h_{i;t}} \right| \\
&\leq \frac{1}{h_{i;t}} \sum_{l=1}^{N_{i;t}-1} \sum_{j=1}^{l-1} \frac{pT}{h_{i;t}} \left| \frac{X_{i;l}}{h_{i;t}} - \frac{X_{j;l}}{h_{i;t}} \right| \\
&\leq \frac{1}{h_{i;t}} \sum_{l=1}^{N_{i;t}-1} \sum_{j=1}^{l-1} \frac{pT}{h_{i;t}} \left| \frac{X_{i;l}}{h_{i;t}} - \frac{X_{j;l}}{h_{i;t}} \right|
\end{aligned} \tag{15}$$

$$\leq \frac{1}{2} p^2 T \sum_{l=1}^{N_{i;t}-1} \sum_{j=1}^{l-1} \frac{1}{h_{i;t}} \left| \frac{X_{i;l}}{h_{i;t}} - \frac{X_{j;l}}{h_{i;t}} \right| \tag{16}$$

where Equation (14) follows from Assumption 2.1 applied as $\frac{X_{i;l}}{h_{i;t}} \leq \frac{X_{i;l}}{h_{i;t}}$, Equation (15) follows from Assumption 2.1 and bounding $\frac{1}{h_{i;t}} \sum_{j=1}^{l-1} \frac{1}{h_{i;t}} \left| \frac{X_{i;l}}{h_{i;t}} - \frac{X_{j;l}}{h_{i;t}} \right| \leq \frac{1}{h_{i;t}}$, and Equation (16) is derived still from Assumption 2.1, $\frac{X_{i;l}}{h_{i;t}} \leq \frac{X_{i;l}}{h_{i;t}}$ and computing the summation. \square

Lemma C.3. For every arm $i \in \mathcal{K}$, every round $t \in \mathcal{T}$, and window width $1 \leq h \leq t$, let us define:

$$\begin{aligned}
\bar{X}_{i;t} &:= \frac{1}{h} \sum_{l=1}^{N_{i;t}-1} \sum_{j=1}^{l-1} \frac{X_{i;l}}{h} - \frac{X_{j;l}}{h} \\
\bar{X}_{i;t} &:= \frac{1}{h} \sum_{l=1}^{N_{i;t}-1} \sum_{j=1}^{l-1} \frac{X_{i;l}}{h} - \frac{X_{j;l}}{h}
\end{aligned}$$

where $X_{i;l}$ denotes the reward collected from arm i when pulled for the l -th time. Otherwise, if $h_{i;t} = 0$, we set $\bar{X}_{i;t} := \emptyset$ and $\bar{X}_{i;t} := \emptyset$. Then, if the window size depends on the number of pulls only $h_{i;t} = h_{i;t}$, it holds that:

$$\mathbb{P} \left(\bar{X}_{i;t} - \bar{X}_{i;t} \geq \frac{a}{10} \right) \leq 2T \exp \left(-\frac{a}{10} \right)$$

Proof. Before starting the proof, it is worth noting that under the event $h_{i;t} = 0$, it holds that $\bar{X}_{i;t} = \bar{X}_{i;t} = \emptyset$. Thus, under the convention that $\emptyset \leq \emptyset = 0$, then $0 \leq \bar{X}_{i;t}$ is not satisfied. For this reason, we need to perform our analysis under the event $h_{i;t} \geq 1$.

The first thing to do is to remove the dependence on the number of pulls that, in a generic time instant, represents a random variable. So, we can write:

$$\begin{aligned}
\mathbb{P} \left(\bar{X}_{i;t} - \bar{X}_{i;t} \geq \frac{a}{10} \right) &\leq \mathbb{P} \left(\bar{X}_{i;t} - \bar{X}_{i;t} \geq \frac{a}{10} \right) \\
&\leq \mathbb{P} \left(\bar{X}_{i;t} - \bar{X}_{i;t} \geq \frac{a}{10} \right) \\
&\leq \mathbb{P} \left(\bar{X}_{i;t} - \bar{X}_{i;t} \geq \frac{a}{10} \right)
\end{aligned} \tag{17}$$

where Equation (17) follows from a union bound over the possible values of $N_{i;t}$.

Now that we have a fixed value of n , consider a generic time t in which arm i has been pulled. We will observe a reward X_t composed by the mean of the process $\frac{X_{i;t}}{h_{i;t}}$ plus some noise. The noise will be equal to $\frac{X_t}{h_{i;t}} - \frac{X_{i;t}}{h_{i;t}}$, i.e., as the difference (not known) between the observed value for the arm i at time t and its real value at the same time. Let us rewrite the quantity to be bounded as follows for every n :

$$\begin{aligned}
h_{i;n} \left(\frac{X_t}{h_{i;n}} - \frac{X_{i;n}}{h_{i;n}} \right) &= \sum_{l=1}^n \left(\frac{X_{i;l}}{h_{i;n}} - \frac{X_{i;l}}{h_{i;n}} \right) \\
&= \sum_{l=1}^n \left(\frac{X_{i;l}}{h_{i;n}} - \frac{X_{i;l}}{h_{i;n}} \right)
\end{aligned}$$

$$\sum_{l=1}^n \sum_{h_{i;n}=1}^{l-1} \frac{T-l}{h_{i;n}} \rho_{l,q} \quad \sum_{l=1}^n \sum_{h_{i;n}=1}^{l-1} \frac{T-l}{h_{i;n}} \rho_{l,q};$$

Here, notice that all the quantities $\rho_{l,q}$ and $\rho_{l,q} / h_{i;n}$ are independent since the number of pulls l is fully determined by n and $h_{i;n}$, that now are non-random quantities.

Now, we apply the Azuma-Hoeffding's inequality of Lemma C.5 from Metelli et al. (2022) for weighted sums of subgaussian martingale difference sequences. To this purpose, we compute the sum of the square weights:

$$\begin{aligned} & \sum_{l=1}^n \sum_{h_{i;n}=1}^{l-1} \frac{T-l}{h_{i;n}} \quad \sum_{l=1}^n \sum_{h_{i;n}=1}^{l-1} \frac{T-l}{h_{i;n}} \\ & \propto \sum_{h_{i;n}=1}^{n-1} \frac{T-n+h_{i;n}-1}{h_{i;n}} \quad \sum_{h_{i;n}=1}^{n-1} \frac{T-n+h_{i;n}-1}{h_{i;n}} \\ & \propto \frac{5pT-n+h_{i;n}-1q}{h_{i;n}}. \end{aligned}$$

Given the previous argument, we have, for a fixed n :

$$\begin{aligned} & \mathbb{P} \left| \sum_{l=1}^n \sum_{h_{i;n}=1}^{l-1} \frac{T-l}{h_{i;n}} \rho_{l,q} \right| \geq \sum_{l=1}^n \sum_{h_{i;n}=1}^{l-1} \frac{T-l}{h_{i;n}} \rho_{l,q} \\ & \propto \mathbb{P} \left| \sum_{l=1}^n \sum_{h_{i;n}=1}^{l-1} \frac{T-l}{h_{i;n}} \rho_{l,q} \right| \geq \sum_{l=1}^n \sum_{h_{i;n}=1}^{l-1} \frac{T-l}{h_{i;n}} \rho_{l,q} \geq \sum_{h_{i;n}=1}^{n-1} \frac{T-n+h_{i;n}-1}{h_{i;n}} \rho_{h_{i;n},q} \\ & \propto 2 \exp \left[- \frac{h_{i;n}^2 \sum_{l=1}^n \sum_{h_{i;n}=1}^{l-1} \frac{T-l}{h_{i;n}} \rho_{l,q}^2}{2 \sum_{h_{i;n}=1}^{n-1} \frac{5pT-n+h_{i;n}-1q}{h_{i;n}}} \right] \\ & = 2 \exp \left[- \frac{a}{10} \right]. \end{aligned}$$

By replacing the obtained result into Equation (17) we get:

$$\sum_{n \in \mathcal{J}_0; T \setminus \mathcal{K}; h_{i;n} \neq 1} 2 \exp \left[- \frac{a}{10} \right] \leq \sum_{n=1}^t 2 \exp \left[- \frac{a}{10} \right] \leq 2T \exp \left[- \frac{a}{10} \right].$$

□

Lemma C.4. For every arm $i \in \mathcal{J} \setminus \mathcal{K}$, every round $t \in \mathcal{J} \setminus \mathcal{K}$, and window width $1 \leq h_{i;t} \leq t \vee N_{i;t-1} \vee 2u$, let us define:

$$\rho_{N_{i;t},q} := \frac{1}{h_{i;t}} \sum_{l=N_{i;t-1}+1}^{N_{i;t}} \rho_{l,q};$$

otherwise, if $h_{i;t} = 0$, we set $\rho_{N_{i;t},q} := \emptyset$. Then, $\sum_{l=1}^t \rho_{N_{i;t},q} \leq \sum_{l=1}^t \rho_{l,q}$ and, if $N_{i;t-1} \neq 2$, it holds that:

$$\sum_{l=1}^t \rho_{l,q} - \rho_{N_{i;t},q} \leq \frac{1}{2} p 2T \sum_{l=N_{i;t-1}+1}^{N_{i;t}} \frac{1}{h_{i;t}} \rho_{l,q} - \sum_{l=N_{i;t-1}+1}^{N_{i;t}} \frac{1}{h_{i;t}} \rho_{l,q};$$

Proof. Following the derivation provided above, we have for every $l \in \{2, \dots, N_{i;T-1}\}$:

$$\sum_{l=1}^t \rho_{l,q} - \rho_{l,q} \leq \sum_{j=1}^{t-1} \rho_{j,q} \leq \sum_{j=1}^t \rho_{j,q}. \quad (18)$$

Thus, by averaging over the most recent $1 \leq h_{i,t} \leq tN_{i,t-1} \{2u\}$ pulls, we get:

$$\begin{aligned}
 & \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j} \\
 & \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j} \\
 & \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j} \\
 & \frac{1}{2} \sum_{j=1}^{N_{i,t-1}} X_{i,j}
 \end{aligned}$$

where we used Assumption 2.1. □

Lemma C.5. For every arm $i \in \mathcal{K}$, every round $t \in \mathcal{T}$, and window width $1 \leq h_{i,t} \leq tN_{i,t-1} \{2u\}$, let us define:

$$\begin{aligned}
 \hat{X}_{i,t} &: \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j} \\
 \hat{X}_{i,t} &: \frac{1}{h_{i,t}} \sum_{j=1}^{N_{i,t-1}} X_{i,j}
 \end{aligned}$$

where $X_{i,j}$ denotes the reward collected from arm i when pulled for the j -th time. Otherwise, if $h_{i,t} = 0$, we set $\hat{X}_{i,t} : \emptyset$ and $\hat{X}_{i,t} : \emptyset$. Then, if the window size depends on the number of pulls only $h_{i,t} = h_{i,t-1} \{2u\}$, it holds that:

$$\mathbb{P} \left(\left| \hat{X}_{i,t} - \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j} \right| \geq \frac{a}{2} \right) \leq 2T \exp \left(-\frac{a^2}{2} \right)$$

Proof. Before starting the proof, it is worth noting that under the event $h_{i,t} = 0$, it holds that $\hat{X}_{i,t} = \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j}$. Thus, under the convention that $\emptyset = 0$, then $0 \leq \hat{X}_{i,t}$ is not satisfied. For this reason, we need to perform our analysis under the event $h_{i,t} \neq 0$.

The first thing to do is to remove the dependence on the number of pulls that, in a generic time instant, represents a random variable. So, we can write:

$$\begin{aligned}
 \mathbb{P} \left(\left| \hat{X}_{i,t} - \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j} \right| \geq \frac{a}{2} \right) & \\
 & \leq \mathbb{P} \left(\left| \hat{X}_{i,t} - \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j} \right| \geq \frac{a}{2} \right) \\
 & \leq \mathbb{P} \left(\left| \hat{X}_{i,t} - \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j} \right| \geq \frac{a}{2} \right)
 \end{aligned} \tag{19}$$

where Equation (19) follows from a union bound over the possible values of $N_{i,t}$.

Now that we have a fixed value of n , consider a generic time t in which arm i has been pulled. We will observe a reward X_t composed by the mean of the process $\frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j}$ plus some noise. The noise will be equal to $X_t - \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j}$, i.e., as the difference (not known) between the observed value for the arm i at time t and its real value at the same time. Let us rewrite the quantity to be bounded as follows, for every n :

$$X_t - \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j}$$

Here we can note that all the quantities $\frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j}$ and $X_t - \frac{1}{h_{i,t-1} N_{i,t-1} h_{i,t-1}} \sum_{j=1}^{N_{i,t-1}} X_{i,j}$ are independent since the number of pulls l is fully determined by n and $h_{i,n}$, that now are non-random quantities.

Now, we apply the Azuma-Hoeffding's inequality of Lemma C.5 from Metelli et al. (2022) for sums of subgaussian martingale difference sequences. For a fixed n , we have:

$$\begin{aligned} \mathbb{P} \left[\left| \hat{\mu}_i(n) - \mu_i \right| \geq \frac{a}{2} \right] &\leq \mathbb{P} \left[\sum_{t=1}^n \left| \tilde{\mu}_i(t) - \mu_i \right| \geq \frac{a}{2} \right] \\ &\leq 2 \exp \left(-\frac{h_{i;n} \sum_{t=1}^n \mu_i^2}{2 \frac{a^2}{4}} \right) \\ &\leq 2 \exp \left(-\frac{a}{2} \right). \end{aligned}$$

By replacing the obtained result into Equation (19) we get:

$$\sum_{i \in \mathcal{J}; T_k: h_{i;n} \geq 1} 2 \exp \left(-\frac{a}{2} \right) \leq \sum_{i \in \mathcal{J}} 2 \exp \left(-\frac{a}{2} \right) \leq 2T \exp \left(-\frac{a}{2} \right).$$

□

Lemma 3.1 (Concentration of $\hat{\mu}_i$). *Under Assumption 2.1, for every $a \geq 0$, simultaneously for every arm $i \in \mathcal{J}$ and number of pulls $n \geq \frac{2T}{a}$, with probability at least $1 - 2TK e^{-a/2}$ it holds that:*

$$\hat{\mu}_i(n) - \hat{\mu}_i(n) \leq \mu_i(n) - \mu_i(n) \leq \hat{\mu}_i(n);$$

where $\hat{\mu}_i(n) := \frac{a}{2n}$ and $\hat{\mu}_i(n) := \frac{1}{2} \left(\frac{2T}{n} - h_{i;n} \right) \mathbb{1}_{\left\{ \frac{2T}{n} - h_{i;n} \geq 1 \right\}}$.

Proof. The proof simply combines Lemmas C.4 and C.5 and a union bound over the arms. □

Lemma 3.2 (Concentration of $\tilde{\mu}_i$). *Under Assumption 2.1, for every $a \geq 0$, simultaneously for every arm $i \in \mathcal{J}$ and number of pulls $n \geq \frac{2T}{a}$, with probability at least $1 - 2TK e^{-a/2}$ it holds that:*

$$\tilde{\mu}_i(n) - \tilde{\mu}_i(n) \leq \mu_i(n) - \tilde{\mu}_i(n) \leq \tilde{\mu}_i(n);$$

where $\tilde{\mu}_i(n) := \frac{a}{2n} - h_{i;n} \mathbb{1}_{\left\{ \frac{a}{2n} - h_{i;n} \geq 1 \right\}}$ and $\tilde{\mu}_i(n) := \frac{1}{2} \left(\frac{2T}{n} - h_{i;n} \right) \mathbb{1}_{\left\{ \frac{2T}{n} - h_{i;n} \geq 1 \right\}}$.

Proof. The proof simply combines Lemmas C.2 and C.3 and a union bound over the arms. □

C.2 Proofs of Section 4

In this appendix, we provide the proofs we have omitted in the main paper for what concerns the theoretical results about R-UCBE. All the lemma below are assuming that the strategy we use for selecting the arm is R-UCBE.

Let us define the *good event* corresponding to the scenario in which all (over the rounds and over the arms) the bounds B_i^T hold for the projection up to time T of the real reward expected value μ_i , formally:

$$\mathcal{G} := \left\{ \forall i \in \mathcal{J}; \forall t \in \mathcal{T}; \left| \tilde{\mu}_i(t) - \mu_i \right| \leq \frac{a}{2} \right\};$$

where $\tilde{\mu}_i(t)$ is the deterministic counterpart of $\tilde{\mu}_i(t)$ considering the expected payoffs μ_i instead of the realizations, formally:

$$\tilde{\mu}_i(t) := \frac{1}{h_{i;t}} \sum_{s=1}^{N_{i;t}-1} \mu_i \mathbb{1}_{\left\{ \mu_i \geq \frac{a}{2} \right\}} - \frac{a}{2} \mathbb{1}_{\left\{ \mu_i < \frac{a}{2} \right\}};$$

Lemma C.6. *Under Assumption 2.1 and assuming that the good event holds, the maximum number of pulls $N_{i;T}$ of a sub-optimal arm ($i \neq \arg \max_{i \in \mathcal{J}} \mu_i$) performed by the R-UCBE is upper bounded by the*

maximum integer $y_i p a q$ which satisfies the following condition:

$$T - i p t p 1 - 2 \lceil y_i p a q u \rceil - 2 T \leq \frac{C}{t \lceil y_i p a q u \rceil^3} \leq i p T q.$$

Proof. In the following, we will use $\tilde{y}_i^T p N_{i;t} - 1 q$ to bound the bias introduced by $\tilde{y}_i^T p N_{i;t} - 1 q$ and, subsequently, to find a number of pulls such that the algorithm cannot suggest pulling a suboptimal arm. Using Lemma C.4, we have that $\tilde{y}_i^T p N_{i;t} - 1 q \leq \frac{1}{2} p 2 T - 2 N_{i;t} - 1 h_{i;t} - 1 q - i p N_{i;t} - 1 2 h_{i;t} - 1 q$ with $N_{i;t} - 1 \leq 2$, it holds that:

$$\tilde{y}_i^T p N_{i;t} - 1 q - i p T q \leq \frac{1}{2} p 2 T - 2 N_{i;t} - 1 h_{i;t} - 1 q - i p N_{i;t} - 1 2 h_{i;t} - 1 q. \quad (20)$$

Let us assume that, at round t , the R-UCBE algorithm pulls the arm $i^* p J K K$ such that $i^* - i^* p T q$. From now on, to avoid weighing down the notation, we will omit the dependence of the optimal arm $i^* p T q$ on the budget T , simply denoting it as i^* , and the window size will be denoted as $h_{i;t} = h p N_{i;t} - 1 q$. By construction, the algorithm chooses the arm with the largest upper confidence bound $B_i^T p N_{i;t} - 1 q$. Thus, we have that $B_i^T p N_{i;t} - 1 q \leq B_{i^*}^T p N_{i;t} - 1 q$. Now, we want to identify the minimum number of pulls such that this event no longer occurs, assuming that the good event holds. We have that, if we pull such an arm $i^* - i^*$, it holds:

$$\begin{aligned} B_i^T p N_{i;t} - 1 q &\leq B_{i^*}^T p N_{i;t} - 1 q \\ B_i^T p N_{i;t} - 1 q - B_{i^*}^T p N_{i;t} - 1 q &\leq 0 \\ i p T q - B_i^T p N_{i;t} - 1 q - B_{i^*}^T p N_{i;t} - 1 q &\leq - i p T q. \end{aligned}$$

Using the definition of $i p T q$ and the definition of the upper confidence bound $B_i^T p N_{i;t} - 1 q$ in Equation (3) for i and i^* , we have:

$$i p T q - i p T q - \tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q - \tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q \leq - i p T q.$$

Given Assumption 2.1 we have that $\tilde{y}_i^T p N_{i;t} - 1 q \leq \tilde{y}_{i^*}^T p N_{i;t} - 1 q$, and, therefore, we have:

$$\tilde{y}_i^T p N_{i;t} - 1 q - i p T q - \tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q - \tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q \leq - i p T q;$$

and, since under the good event, it holds that $\tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q - \tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q \leq 0$, we have:

$$i p T q - \tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q - \tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q \leq - i p T q$$

$$i p T q - \tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q - \tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q \leq - i p T q$$

where we added and subtracted $\tilde{y}_i^T p N_{i;t} - 1 q$ in the last equation. Under the good event, we can upper bound $|p D q| = |\tilde{y}_i^T p N_{i;t} - 1 q - \tilde{y}_{i^*}^T p N_{i;t} - 1 q| - \tilde{y}_i^T p N_{i;t} - 1 q$:

$$\tilde{y}_i^T p N_{i;t} - 1 q - i p T q - 2 \tilde{y}_i^T p N_{i;t} - 1 q \leq - i p T q.$$

Using Equation (20), and substituting the definition of $\tilde{y}_i^T p N_{i;t} - 1 q$ provided in Equation (4), we have:

$$\frac{1}{2} p 2 T - 2 N_{i;t} - 1 h_{i;t} - 1 q - i p N_{i;t} - 1 2 h_{i;t} - 1 q$$

$$\leq 2 \tilde{y}_i^T p N_{i;t} - 1 q - \frac{C}{h_{i;t}^3} \leq - i p T q$$

$$\sum_{i \in \mathcal{P}} \frac{2TK}{N_{i,t}} \leq \frac{C}{a} \quad (21)$$

where we used the definition of $h_{i,t} = \frac{2TK}{N_{i,t}}$ and the fact that $N_{i,t-1} \leq N_{i,t} - 1$ since at time t the algorithm pulls the i -th arm.

This concludes the proof. \square

Theorem 4.1. *Under Assumption 2.1, let a be the largest positive value of a satisfying:*

$$\sum_{i \in \mathcal{P}} \frac{2TK}{y_i(a)} \leq 1; \quad (5)$$

where for every $i \in \mathcal{P}$, $y_i(a)$ is the largest integer for which it holds:

$$\sum_{i \in \mathcal{P}} \frac{2TK}{y_i(a)} \leq \frac{C}{a} \quad (6)$$

If a exists, then for every $a \in (0; a]$ the error probability of R-UCBE is bounded by:

$$e_{\mathcal{T}}(R\text{-UCBE}) \leq 2TK \exp\left(-\frac{a}{10}\right); \quad (7)$$

Proof. From the definition of the error probability, we have:

$$e_{\mathcal{T}}(R\text{-UCBE}) = \mathbb{P}\left(\hat{p} \neq p^* \mid \mathcal{E}\right) = \mathbb{P}\left(\bigcup_{i \in \mathcal{P}} \mathcal{E}_i \mid \mathcal{E}\right);$$

Therefore, we need to evaluate the probability that the R-UCBE algorithm would pull a suboptimal arm in the $\mathcal{T} - 1$ round. Given that Assumption 2.1 and that each suboptimal arms have been pulled a number of times $N_{i;\mathcal{T}}$ (and cannot be pulled more) at the end of the time budget \mathcal{T} , under the good event \mathcal{E} , we are guaranteed to recommend the optimal arm if:

$$\sum_{i \in \mathcal{P}} \frac{2TK}{N_{i;\mathcal{T}}} \leq 1; \quad (22)$$

If Equation (22) holds, a suboptimal arm can be selected by R-UCBE for the next round $\mathcal{T} - 1$ only if the good event \mathcal{E} does not hold $e_{\mathcal{T}}(R\text{-UCBE}) = \mathbb{P}(\mathcal{E}^c)$, where we denote with \mathcal{E}^c the complementary of event \mathcal{E} . This probability is upper bounded by Lemma C.5 as:

$$e_{\mathcal{T}}(R\text{-UCBE}) = \mathbb{P}(\mathcal{E}^c) \leq 2TK \exp\left(-\frac{a}{10}\right);$$

We now derive a condition for a in order to make Equation (22) hold. Thanks to Lemma C.6 we know that $N_{i;\mathcal{T}} \geq y_i(a)$ where $y_i(a)$ is the maximum integer such that:

$$\sum_{i \in \mathcal{P}} \frac{2TK}{y_i(a)} \leq \frac{C}{a};$$

From this condition, we observe that $y_i(a)$ is an increasing function of a . Therefore, we can select a in the interval $(0; a]$, where a is the maximum value of a such that:

$$\sum_{i \in \mathcal{P}} \frac{2TK}{y_i(a)} \leq 1; \quad (23)$$

Note that we are not guaranteed that such a value of $a > 0$ exists. In such a case, we cannot provide meaningful guarantees on the error probability of R-UCBE. \square

Corollary 4.2. Under Assumptions 2.1 and 2.2, if the time budget T satisfies:

$$T \geq \begin{cases} \frac{c^2 p_1}{4} \left(\frac{T-1}{H_{1,1}(pTq)} \right)^2 \left(\frac{pK-1}{pTq} \right) & \text{if } P \geq \frac{3}{2}; \mathcal{S} \\ \frac{c^2 p_1}{4} \left(\frac{T-3}{H_{1,2}(3pTq)} \right)^2 \left(\frac{pK-1}{pTq} \right) & \text{if } P < \frac{3}{2}; \mathcal{S} \end{cases} \quad (8)$$

there exists a $\alpha > 0$ defined as:

$$\alpha = \begin{cases} \frac{T-1}{4} \left(\frac{pK-1}{pTq} \right) & \text{if } P \geq \frac{3}{2}; \mathcal{S} \\ \frac{T-3}{4} \left(\frac{pK-1}{pTq} \right) & \text{if } P < \frac{3}{2}; \mathcal{S} \end{cases}$$

where $H_{1,1}(pTq) = \frac{1}{pTq}$ for $\alpha > 0$. Then, for every $a > 0$, the error probability of R-UCBE is bounded by:

$$e_{T, R-UCBE} \leq 2TK \exp \left(-\frac{a}{10} \right)$$

Proof. We recall that Assumption 2.2 states that all the increment functions $\mu_{i,pTq}$ are such that $\mu_{i,pTq} \leq cn$. We use such a fact to provide an explicit solution for the optimal value of a . From Theorem 4.1 and using the fact that $\mu_{i,pTq} \leq cn$, we have that Equation (6) becomes:

$$\frac{Tc}{p_1 - 2\alpha} \leq \frac{2T}{t} \frac{a^{\frac{1}{2}}}{y^{\frac{3}{2}}} \leq \mu_{i,pTq} \quad (24)$$

Or, more restrictively:

$$\frac{Tc p_1 - 2\alpha}{p_1 - 1} \leq \frac{2T}{p_1 - 1} \frac{a^{\frac{1}{2}}}{1q^{\frac{3}{2}}} \leq \mu_{i,pTq}$$

Let us solve Equation (24) by analyzing separately the two cases in which one of the two terms in the l.h.s. of such equation become prevalent.

Case 1: $P \geq \frac{3}{2}; \mathcal{S}$ In this branch, we can upper bound the left-side part of the inequality in Equation (24) by:

$$\frac{Tc p_1 - 2\alpha}{p_1 - 1q^{\frac{3}{2}}} \leq \frac{2T}{p_1 - 1q^{\frac{3}{2}}} \frac{a^{\frac{1}{2}}}{1q^{\frac{3}{2}}} \leq \mu_{i,pTq}$$

Thus, we can derive:

$$y_{i,pTq} \leq 1 + \frac{Tc p_1 - 2\alpha}{p_1 - 1q^{\frac{3}{2}}} \frac{2T}{p_1 - 1q^{\frac{3}{2}}} \frac{a^{\frac{1}{2}}}{1q^{\frac{3}{2}}} \quad (25)$$

Using the above value in Equation (23), provides:

$$T \leq \frac{pK-1}{pTq} \left(\frac{Tc p_1 - 2\alpha}{p_1 - 1q^{\frac{3}{2}}} \frac{2T}{p_1 - 1q^{\frac{3}{2}}} \frac{a^{\frac{1}{2}}}{1q^{\frac{3}{2}}} \right)^{\frac{2}{3}} \leq \frac{1}{p_1 - 1q^{\frac{3}{2}}} \quad \alpha > 0$$

$$T \leq \frac{pK-1}{pTq} \left(\frac{Tc p_1 - 2\alpha}{p_1 - 1q^{\frac{3}{2}}} \frac{2T}{p_1 - 1q^{\frac{3}{2}}} \frac{a^{\frac{1}{2}}}{1q^{\frac{3}{2}}} \right)^{\frac{2}{3}} H_{1,2}(3pTq) \quad \alpha > 0$$

$$a \leq \frac{\frac{pT^{1(3)} - T^{2(3)} pK - 1q^{\frac{3}{2}}}{pH_{1,2}(3pTq)^{\frac{3}{2}}} \left(p_1 - 2\alpha \right)^2}{4 - 2\alpha^3}$$

$$a = \frac{\frac{\rho T^{1-\zeta} \rho K 1 q q^{\frac{3}{2}}}{\rho H_{1,2}(\zeta) \rho T q q^{\frac{3}{2}}} \rho p 1 2'' q^2}{4 2'' 3};$$

where the last expression is obtained by observing that $T \neq 1$ and for obtaining a more manageable expression, under the assumption that $\frac{\rho T^{1-\zeta} \rho K 1 q q^{\frac{3}{2}}}{\rho H_{1,2}(\zeta) \rho T q q^{\frac{3}{2}}} \rho p 1 2'' q \neq 0$.

This implies a constraint on the minimum time budget T , which explicit form for the case $P \in \frac{3}{2}; \mathcal{B}$ is provided in Lemma C.7

Case 2: $P \in 1; \frac{3}{2}$ In this case, we enforce the more restrictive condition:

$$\frac{T c p 1 2'' q}{\rho y 1 q} = \frac{2 T'' \frac{3}{2} a^{\frac{1}{2}}}{\rho y 1 q} \neq \rho p T q;$$

the value for the number of pulls is:

$$y_i \rho a q \propto 1 = \frac{T c p 1 2'' q}{\rho p T q} = \frac{2 T'' \frac{3}{2} a^{\frac{1}{2}}}{\rho p T q} \quad (26)$$

and the value for a becomes:

$$T \rho K 1 q = T c p 1 2'' q = 2 T'' \frac{3}{2} a^{\frac{1}{2}} \frac{1}{\rho p T q} \sum_{i=1}^N \frac{1}{\rho p T q} \rho p T q \rho H_{1,1}(\zeta) \rho p T q \rho p T q \rho p T q$$

$$a = \frac{\frac{\rho T^{1-\zeta} \rho K 1 q q^{\frac{3}{2}}}{\rho H_{1,1}(\zeta) \rho T q q^{\frac{3}{2}}} \rho p 1 2'' q^2}{4 2'' 3}$$

$$a = \frac{\frac{\rho T^{1-\zeta} \rho K 1 q q^{\frac{3}{2}}}{\rho H_{1,1}(\zeta) \rho T q q^{\frac{3}{2}}} \rho p 1 2'' q^2}{4 2'' 3};$$

where the last expression is obtained by observing that $T \neq 1$ and for obtaining a more convenient expression, under the assumption that $\frac{\rho T^{1-\zeta} \rho K 1 q q^{\frac{3}{2}}}{\rho H_{1,1}(\zeta) \rho T q q^{\frac{3}{2}}} \rho p 1 2'' q \neq 0$.

Also here, this implies a constraint on the minimum time budget T for the case $P \in 1; \frac{3}{2}$, which explicit form is provided in Lemma C.7 \square

Lemma C.7. Under Assumptions 2.1 and 2.2, the minimum time budget T for which the theoretical guarantees of R-UCBE hold is:

$$T \neq \begin{cases} c^1 p 1 2'' q^{-1} \rho H_{1,1}(\zeta) \rho p T q q \rho K 1 q^{-1} & \text{if } P \in 1; \frac{3}{2} \mathcal{A} \\ c^{\frac{2}{3}} p 1 2'' q^{\frac{2}{3}} \rho H_{1,2}(\zeta) \rho p T q q \rho K 1 q^{-1} & \text{if } P \in \frac{3}{2}; \mathcal{B} \end{cases}$$

and $H_{1,1}(\zeta) \rho p T q = \sum_{i=1}^N \frac{1}{\rho p T q}$ for $\alpha = 1$.

Proof. Given Corollary 4.2, we want to find the values of T such that a value of $a \in \mathbb{R}_0^+$ should exist. This implies having $a \neq 0$. Given the value of α , we can derive a lower bound for the time budget T .

Case 1: $P \geq \frac{3}{2}; \delta \leq \frac{1}{2}$:

$$\frac{pT^{1/3} \sqrt{pK} \sqrt{1/q}}{pH_{1,2}(pTq)^{3/2}} \leq c p^{1/2} q \leq 0:$$

From this, it follows:

$$T \leq c^{2/3} p^{1/2} q^{2/3} pH_{1,2}(pTq) \sqrt{pK} \sqrt{1/q}^3 :$$

Case 2: $P \leq 1; \frac{3}{2} \leq \delta \leq 1$:

$$\frac{pT^{1-\delta} \sqrt{pK} \sqrt{1/q}}{pH_{1,1}(pTq)} \leq c p^{1/2} q \leq 0:$$

From this, it follows:

$$T \leq c^{1/\delta} p^{1/2} q^{1-\delta} pH_{1,1}(pTq) \sqrt{pK} \sqrt{1/q}^{1-\delta} :$$

□

C.3 Proofs of Section 5

In this appendix, we provide the proofs we have omitted in the main paper for what concerns the theoretical results about R-SR. We recall that with a slight abuse of notation, as done in Section 5, we denote with $p_{i,q}pTq$ the i^{th} gap rearranged in increasing order, i.e., we have $p_{i,q}pTq \leq p_{j,q}pTq$ for $i \leq j$.

Lemma C.8. For every arm $i \in [K]$ and every round $t \in [T]$, let us define:

$$\bar{p}_{i,t} = \frac{1}{h_{i,t}} \sum_{s=1}^{N_{i,t}} p_{i,q}^s$$

if $N_{i,t} \leq 2$, then $\bar{p}_{i,t} \leq p_{i,t}$, and if $h_{i,t} \leq N_{i,t}/2$, it holds that:

$$p_{i,t} \leq \bar{p}_{i,t} \leq T^{-1} N_{i,t}^{-2} \sum_{s=1}^{N_{i,t}} p_{i,q}^s \quad (27)$$

Proof. The proof follows trivially from Lemma C.2. □

Lemma C.9 (Lower Bound for the Time Budget for R-SR). Under Assumptions 2.1 and 2.2, the R-SR algorithm is s.t. the minimum value for the horizon T ensuring that $\sum_{j \in [K]} p_{j,t} \leq 1$ and $\sum_{i \in [K]} \bar{p}_{i,t} \leq 1$:

$$T \geq \sum_{j \in [K]} N_j \sqrt{1/q} \leq \frac{\sqrt{pK} \sqrt{1/q} pTq}{2},$$

is:

$$T \geq c^{1-\delta} 2^{1-\delta} \log pK q^{-\delta} \max_{i \in [2;K]} i^{-\delta} p_{i,q}^{-1} pTq^{\delta} :$$

Proof. First, using Assumption 2.2, we derive an upper bound on the bias between $p_{i,t}$ and $\bar{p}_{i,t}$ (r.h.s. of Equation 27), where N_j is a generic time corresponding to the end of a phase of the R-SR algorithm:

$$T^{-1} N_j^{-2} \sum_{s=1}^{N_j} p_{i,q}^s \leq c T^{-1} N_j^{-2} \sum_{s=1}^{N_j} p_{i,q}^s :$$

Substituting the definition of N_j into the above equation, we get:

$$\begin{aligned} T^{-1} N_j^{-2} \sum_{s=1}^{N_j} p_{i,q}^s &\leq T^{-1} \frac{1}{\log pK q} \frac{T}{K-1} \frac{1}{j} \leq 1 \\ &\leq T^{-1} \frac{1}{\log pK q} \frac{T}{K-1} \frac{1}{j} \leq 1 \end{aligned} \quad (28)$$

$$\propto T \frac{T}{\overline{\log p K q} p K - 1 j q} \quad ; \quad (29)$$

Requiring that, for a generic N_j , the maximum possible bias is lower than a fraction of the suboptimality gap of arm $K - 1 - j$:

$$\begin{aligned} c T t N_j \{2u\} &\propto \frac{p K - 1 j q p T q}{2} \\ c T \frac{T}{2 \overline{\log p K q} p K - 1 j q} &\propto \frac{p K - 1 j q p T q}{2} \\ T^{-1} &\propto \frac{p K - 1 j q p T q}{c 2^1 \overline{\log p K q} p K - 1 j q} \\ T &\preceq \frac{p K - 1 j q p T q}{c^{1-1} 2^{1-1} \overline{\log p K q} p K - 1 j q^{1-1}}. \end{aligned}$$

Requiring that the above condition holds for all the phases $j \in \{K - 1, \dots, 1\}$ we have:

$$\begin{aligned} T &\preceq \max_{j \in \{K-1, \dots, 1\}} \frac{p K - 1 j q p T q}{c^{1-1} 2^{1-1} \overline{\log p K q} p K - 1 j q^{1-1}} \\ &\preceq c^{1-1} 2^{1-1} \overline{\log p K q}^{-1} \max_{j \in \{K-1, \dots, 1\}} \frac{p K - 1 j q p T q}{p K - 1 j q} \\ &\preceq c^{1-1} 2^{1-1} \overline{\log p K q}^{-1} \max_{j \in \{K-1, \dots, 1\}} p K - 1 j q \frac{p K - 1 j q p T q}{p K - 1 j q} \\ &\preceq c^{1-1} 2^{1-1} \overline{\log p K q}^{-1} \max_{i \in \{2, \dots, K\}} i^{-1} p i q^{-1} p T q \quad ; \end{aligned}$$

□

Theorem 5.1. Under Assumptions 2.1 and 2.2, if the time budget T satisfies:

$$T \preceq 2^{-1} c^{-1} \overline{\log p K q}^{-1} \max_{i \in \{2, \dots, K\}} i^{-1} p i q^{-1} p T q^{-1} \quad ; \quad (10)$$

then, the error probability of R-SR is bounded by:

$$e_T p_{R-SR} \propto \frac{K p K - 1 q}{2} \exp \left(- \frac{T}{8} \frac{K}{\overline{\log p K q} H_2 p T q} \right) ;$$

where $H_2 p T q = \max_{i \in \{2, \dots, K\}} i^{-1} p i q^{-1} p T q^{-2}$ and $\overline{\log p K q} = \frac{1}{2} \sum_{i=2}^K \frac{1}{i}$.

Proof. The R-SR algorithm makes an error when at the end of a phase j the optimal arm has a pessimistic estimator $\hat{p}_1 p N_j q$ is smallest among the arms, formally:

$$\begin{aligned} e_T p_{R-SR} &\propto \mathbb{P} \left(\bigcup_{j=1}^{K-1} D_j \mid \bigcap_{j=1}^{K-1} \{K-1 - j; K\} : \hat{p}_1 p N_j q < \hat{p}_{i q} p N_j q \right) \\ &\propto \sum_{j=1}^{K-1} \mathbb{P} \left(D_j \mid \bigcap_{j=1}^{K-1} \{K-1 - j; K\} : \hat{p}_1 p N_j q < \hat{p}_{i q} p N_j q \right) \\ &\propto \sum_{j=1}^{K-1} \sum_{i=K-1-j}^K \mathbb{P} \left(\hat{p}_1 p N_j q < \hat{p}_{i q} p N_j q \right) ; \end{aligned}$$

where we use a union bound over the phases and over the arms still in the available arm set X_{j-1} in each phase. Let us focus on $\hat{\mu}_{p1q} \hat{\mu}_{Nj} \propto \hat{\mu}_{p1q} \hat{\mu}_{Nj}$. We have that the optimal arm has a smaller pessimistic estimator than the l^{th} one when:

$$\begin{aligned} & \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \hat{\mu}_{p1q} \hat{\mu}_{Nj} \\ & \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq 0 \\ & \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \\ & \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \\ & \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \end{aligned} \quad (30)$$

$$\begin{aligned} & \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \\ & \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \\ & \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \end{aligned} \quad (31)$$

$$\begin{aligned} & \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \\ & \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \\ & \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \end{aligned} \quad (32)$$

where we added $\hat{\mu}_{p1q} \hat{\mu}_{Nj}$ to derive Equation (30), and added $-\hat{\mu}_{p1q} \hat{\mu}_{Nj}$ to derive Equation (31), we used the results in Lemma C.8 and from the fact that the reward function is increasing. Since we are with a time budget T satisfying Theorem C.9, we have that:

$$T \hat{\mu}_{p1q} \hat{\mu}_{Nj} - 1q \propto \frac{\hat{\mu}_{p1q} \hat{\mu}_{Nj}}{2} \quad (33)$$

Substituting into Equation (32) the above, we have:

$$\hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \frac{\hat{\mu}_{p1q} \hat{\mu}_{Nj}}{2} ;$$

and the error probability becomes:

$$e_{T \text{ pR-SRq}} \propto \sum_{j=1}^{K-1} \sum_{i=K-1}^K \mathbb{P} \left(\hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} - \hat{\mu}_{p1q} \hat{\mu}_{Nj} \not\leq \frac{\hat{\mu}_{p1q} \hat{\mu}_{Nj}}{2} \right) ;$$

For the previous argumentation, we apply the Azuma-Hoeffding's inequality to the latter probability:

$$\begin{aligned} e_{T \text{ pR-SRq}} & \propto \sum_{j=1}^{K-1} \sum_{i=K-1}^K \exp \left(-\frac{\hat{\mu}_{p1q} \hat{\mu}_{Nj}^2}{2} \right) \\ & \propto \sum_{j=1}^{K-1} j \exp \left(-\frac{\hat{\mu}_{p1q} \hat{\mu}_{Nj}^2}{8} \right) \end{aligned} ;$$

Now, given that:

$$\begin{aligned} \frac{\hat{\mu}_{p1q} \hat{\mu}_{Nj}^2}{8} & \not\leq \frac{T}{8} \frac{K}{\log p K q p K - 1} \frac{K}{j q p K - 1} \\ & \not\leq \frac{T}{8} \frac{K}{\log p K q H_2 p T q} ; \end{aligned}$$

we finally derive the following:

$$e_{T \text{ pR-SRq}} \propto \frac{K p K - 1 q}{2} \exp \left(-\frac{T}{8} \frac{K}{\log p K q H_2 p T q} \right) ;$$

which concludes the proof. □

C.4 Proofs of Section 6

In this appendix, we provide the proofs of the lower bound on the error probability presented in Section 6.

Theorem 6.1. For every algorithm A , there exists a deterministic SRB satisfying Assumptions 2.1 and 2.2 with $K \asymp 8^{1/p-1} q$ such that the optimal arm $i^* \in \mathcal{P}^*$ cannot be identified for some time budgets T unless:

$$T \not\asymp H_{1,1/p-1}(\mathcal{P}^*) \asymp \frac{1}{i^* \Delta_i^{1/p-1}}. \quad (11)$$

Proof. We define for every suboptimal arm $i \in \mathcal{P} \setminus \mathcal{P}^*$ the suboptimality gap reached at $T \in \mathbb{N}$ as $\Delta_i \in (0, 1/2]$. We consider the base instance (see Figure 5) in which the (deterministic) reward functions are defined for $i \in \mathcal{P} \setminus \mathcal{P}^*$ and $n \in \mathbb{N}$ as:

$$r_i(p, n) = \frac{1}{2} - \frac{1}{n} \quad ;$$

$$r_i(p, n) = \min \left\{ \frac{1}{2} - \frac{1}{n}, \frac{1}{2} - \frac{1}{n} \right\} \quad ; \quad i \in \mathcal{P} \setminus \mathcal{P}^* ;$$

Clearly, r_i fulfills Assumption 2.1 and it is simple to show that also Assumption 2.2 is satisfied. Indeed, by first-order Taylor expansion:

$$r_i(p, n) - r_i(p^*, n) = \frac{1}{n} - \frac{1}{n} \asymp \sup_{x \in \mathcal{P} \setminus \mathcal{P}^*} \frac{B}{Bx} - \frac{1}{2} \asymp \frac{1}{2} - \frac{1}{n} \quad ; \quad (34)$$

$$r_i(p, n) - r_i(p^*, n) = \frac{1}{n} - \frac{1}{n} \asymp \sup_{x \in \mathcal{P} \setminus \mathcal{P}^*} \frac{B}{Bx} - \frac{1}{2} \asymp \frac{1}{2} - \frac{1}{n}$$

Thus, we simply take $c = 1$ in Assumption 2.2. Let us define n_i the number of pulls in which arm $i \in \mathcal{P} \setminus \mathcal{P}^*$ reaches the stationary behavior:

$$\frac{1}{2} - \frac{1}{n_i} = \frac{1}{2} - \frac{1}{n_i} \asymp \frac{1}{2} - \frac{1}{n_i} \quad ;$$

A sufficient condition on the time budget so that the optimal arm is 1 (i.e., $i^* \in \mathcal{P}^*$) is given by $T \not\asymp T^*$, where T^* is the point in which the curve of the optimal arm intersects that of any of the suboptimal arms $i \in \mathcal{P} \setminus \mathcal{P}^*$:

$$\frac{1}{2} - \frac{1}{T^*} = \frac{1}{2} - \frac{1}{n_i} \quad ; \quad T^* = \max_{i \in \mathcal{P} \setminus \mathcal{P}^*} \frac{1}{2 - \frac{1}{n_i}}$$

Consider now the regime in which $T \asymp T^*$. We proceed by contradiction. Suppose that there exists an algorithm A that identifies the optimal arm such that on the bandit \mathcal{B} and that the suboptimal arm $i \in \mathcal{P} \setminus \mathcal{P}^*$ has an expected number of pulls satisfying:

$$\mathbb{E} N_{i, p, T} \asymp n_i \quad ; \quad (35)$$

Consider now the alternative bandit \mathcal{B}^* constructed from \mathcal{B} by keeping all the arms unaltered, except for arm i that is made optimal:

$$r_i(p, n) = \frac{1}{2} - \frac{1}{n} \quad ;$$

$$j \in \mathcal{P} \setminus \mathcal{K} \text{ or } j \in \mathcal{K} \setminus \mathcal{J}.$$

Clearly the bandit $\mathcal{P} \setminus \mathcal{K}$ fulfills Assumption 2.1 and, with calculations similar to those in Equation (34), we conclude that it satisfies Assumption 2.2 with $c = 1$. A sufficient condition on T for which arm i is optimal in bandit $\mathcal{P} \setminus \mathcal{K}$ is that $T \geq T_2$ in which the curve of arm i intersects that of the arms j such that $j \in \mathcal{P} \setminus \mathcal{K}$:

$$\frac{1}{2} \mu_i - \frac{1}{2} \mu_{\mathcal{P} \setminus \mathcal{K}} \leq \frac{1}{2} \mu_j \iff T \geq \max_{j \in \mathcal{P} \setminus \mathcal{K}} \frac{1}{2} \frac{\mu_i - \mu_j}{\mu_i - \mu_{\mathcal{P} \setminus \mathcal{K}}}.$$

Thus, we take:

$$T_2 := \frac{1}{2} \frac{\mu_i - \mu_{\mathcal{P} \setminus \mathcal{K}}}{\mu_i - \mu_{\mathcal{P} \setminus \mathcal{K}}}.$$

Clearly, for $T \geq T_2$ since all the suboptimality gaps are at most $\frac{1}{2}$. Thus, we continue in the regime $T \geq T_2$. Since $\mu_i > \mu_j$ if $n_i > n_j$, it follows that under condition (35), algorithm A cannot distinguish between the two bandits and, consequently, cannot identify the optimal arm on bandit $\mathcal{P} \setminus \mathcal{K}$. Thus, it must follow, from the contradiction, that:

$$\mathbb{E} n_{i, \mathcal{P} \setminus \mathcal{K}} \geq n_i.$$

By summing over $i \in \mathcal{P} \setminus \mathcal{K}$, we obtain:

$$\sum_{i \in \mathcal{P} \setminus \mathcal{K}} \mathbb{E} n_{i, \mathcal{P} \setminus \mathcal{K}} \geq \sum_{i \in \mathcal{P} \setminus \mathcal{K}} n_i \geq \sum_{i \in \mathcal{P} \setminus \mathcal{K}} \frac{1}{2} \frac{\mu_i - \mu_{\mathcal{P} \setminus \mathcal{K}}}{\mu_i - \mu_{\mathcal{P} \setminus \mathcal{K}}} \geq \frac{1}{4} \sum_{i \in \mathcal{P} \setminus \mathcal{K}} \frac{1}{\mu_i - \mu_{\mathcal{P} \setminus \mathcal{K}}}. \quad (36)$$

Thus, we have found an interval $T \in [T_2, T_3]$ in which identification cannot be performed. Notice that it is simple to enforce that $T_3 \leq T$ with a sufficiently large number of arms $K \geq 2^{1/T_3}$.

To conclude, we need to relate μ_i with $\mu_{i, \mathcal{P} \setminus \mathcal{K}}$ and $\mu_{i, \mathcal{K}}$. We perform the computation for both the instances $\mathcal{P} \setminus \mathcal{K}$ and \mathcal{K} , in the regime $T \geq 2^{1/T_3}$. Let us start with $\mathcal{P} \setminus \mathcal{K}$:

$$\mu_{i, \mathcal{P} \setminus \mathcal{K}} \leq \frac{1}{2} \mu_{\mathcal{P} \setminus \mathcal{K}} + \frac{1}{2} \mu_i \iff \mu_{i, \mathcal{P} \setminus \mathcal{K}} \leq \frac{\mu_i}{2}; \quad i \in \mathcal{P} \setminus \mathcal{K}$$

We move to \mathcal{K} :

$$\begin{aligned} \mu_{i, \mathcal{K}} &\leq \frac{1}{2} \mu_{\mathcal{K}} + \frac{1}{2} \mu_i \iff \mu_{i, \mathcal{K}} \leq \frac{\mu_i}{2}; \\ \mu_{i, \mathcal{K}} &\leq \frac{1}{2} \mu_{\mathcal{K}} + \frac{1}{2} \mu_i \iff \mu_{i, \mathcal{K}} \leq \frac{\mu_i}{2}; \quad i \in \mathcal{K} \end{aligned}$$

Thus, a necessary condition for the correct identification of the optimal arm is:

$$T \geq \frac{1}{4} \frac{\mu_i - \mu_{\mathcal{P} \setminus \mathcal{K}}}{\mu_i - \mu_{\mathcal{P} \setminus \mathcal{K}}} \geq \frac{1}{8} \frac{1}{\mu_i - \mu_{\mathcal{P} \setminus \mathcal{K}}} \geq 2^{1/T_3} T.$$

Similarly, with $K \geq 8^{1/T_3}$, we can enforce $2^{1/T_3} T \leq 2^{1/T_3} T$. \square

Theorem 6.2. For every algorithm A run with a time budget T fulfilling Equation (11), there exists a SRB satisfying Assumptions 2.1 and 2.2 with $K \geq 8^{1/T_3}$ such that the error probability is lower bounded by:

$$e_{\mathcal{P} \setminus \mathcal{K}} \geq \frac{1}{4} \exp \left(-\frac{8T}{2H_{1,2, \mathcal{P} \setminus \mathcal{K}}} \right); \quad \text{where } H_{1,2, \mathcal{P} \setminus \mathcal{K}} := \frac{1}{\mu_i - \mu_{\mathcal{P} \setminus \mathcal{K}}}.$$

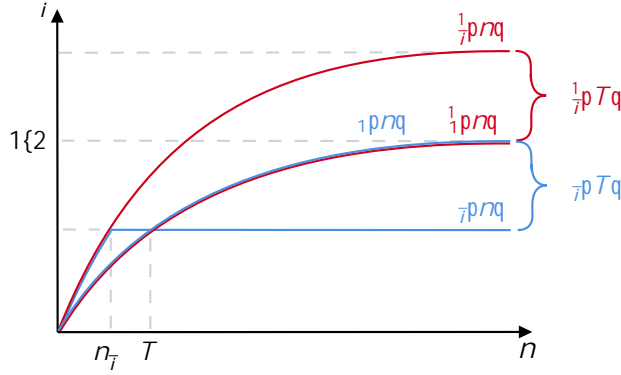


Figure 5: Instances and of SRB used in Theorem 6.1 and Theorem 6.2.

Proof. The proof imports the technique from (Kaufmann et al., 2016, Theorem 16 and 17). We consider the Gaussian bandit with variance σ^2 equal for all the arms and the expected reward μ_i as in the base instance of proof of Theorem 6.1 (see Figure 5). Let us define by convention $\mu_1 = \frac{1}{2}$. Let A be an algorithm, it is simple to show that there exists an arm $i \in \{1, \dots, K\}$, such that:

$$\mathbb{E} r_{i,pTq} \asymp \frac{T}{H_2(i)};$$

where $H_2(i) = \sum_{k=1}^K \frac{1}{k^2}$. We consider two cases. Suppose that $i = 1$ and we construct the alternative Gaussian bandit with the same variance σ^2 and the expected rewards defined as follows:

$$\mu_{1,pTq} = \min\left\{\frac{1}{2}, \frac{1}{n}\right\}; \quad \mu_{i,pTq} = \frac{1}{2} - \frac{1}{n};$$

For T sufficiently large as in Theorem 6.1, while in bandit the optimal arm is 1, in bandit the optimal arm is 2. Instead, suppose that $i \neq 1$ and we construct the alternative Gaussian bandit with the same variance σ^2 and the expected rewards defined as follows:

$$\mu_{i,pTq} = \frac{1}{2} - \frac{1}{n}; \quad \mu_{1,pTq} = \frac{1}{n};$$

For T sufficiently large as in Theorem 6.1, while in bandit the optimal arm is 1, in bandit the optimal arm is i . Let us denote with $r_{i,tq}$ the distribution of the reward at time t for arm i . By the Bretagnolle-Huber's inequality, we obtain:

$$\begin{aligned} \max_{t \in [1, T]} \mathbb{P}(r_{i,tq} \neq \mu_{i,pTq}) &\leq \frac{1}{4} \exp\left(-\frac{1}{4} \sum_{t=1}^T \frac{1}{t}\right) \mathbb{E} \sum_{t=1}^T \frac{1}{t} \mathbb{E} \sum_{t=1}^T \frac{1}{t} \mathbb{E} \sum_{t=1}^T \frac{1}{t} \mathbb{E} \sum_{t=1}^T \frac{1}{t} \\ &\leq \frac{1}{4} \exp\left(-\frac{1}{4} \sum_{t=1}^T \frac{1}{t}\right) \mathbb{E} \sum_{t=1}^T \frac{1}{t} \mathbb{E} \sum_{t=1}^T \frac{1}{t} \mathbb{E} \sum_{t=1}^T \frac{1}{t} \mathbb{E} \sum_{t=1}^T \frac{1}{t} \\ &\leq \frac{1}{4} \exp\left(-\frac{1}{4} \sum_{t=1}^T \frac{1}{t}\right) \mathbb{E} r_{i,pTq} \frac{p^2 \sigma^2}{2} : \frac{1}{4} \exp\left(-\frac{2T}{2H_2}\right) \\ &\leq \frac{1}{4} \exp\left(-\frac{2T}{2H_2}\right); \end{aligned}$$

where we observed that for every $n \in \mathbb{N}$, we have $|\mu_{i,pTq} - \mu_{1,pTq}| \leq \frac{1}{n}$. To conclude, we relate H_2 with $H_{1,2pTq}$. Using an argument analogous to that of the last part of the proof Theorem 6.1 it is simple to observe that, for sufficiently large T , we have $H_2(i) \asymp 2 \mu_{i,pTq}$ (condition already satisfied

by requesting the bound on T of Theorem 6.1), from which we have:

$$H_2 \sum_{i=1}^K \sum_{t=1}^T x_t^2 \leq \sum_{i=1}^K \sum_{t=1}^T x_t^2 \leq \sum_{i=1}^K \sum_{t=1}^T x_t^2 \leq \frac{1}{4} \sum_{i=1}^K \sum_{t=1}^T x_t^2 \leq \frac{1}{4} H_{1,2} p T q$$

□

C.5 Auxiliary Lemmas

Lemma C.10 (Hoeffding-Azuma’s inequality for weighted martingales). *Let $F_1 \subseteq \dots \subseteq F_n$ be a filtration and X_1, \dots, X_n be real random variables such that X_t is F_t -measurable, $\mathbb{E} x_t | F_{t-1} = 0$ (i.e., a martingale difference sequence), and $\mathbb{E} \exp(x_t | F_{t-1}) \leq \exp(\frac{\sigma_t^2}{2})$ for any $t \in [n]$ (i.e., σ_t^2 -subgaussian). Let w_1, \dots, w_n be non-negative real numbers. Then, for every $\epsilon \in (0, 1]$ it holds that:*

$$\mathbb{P} \left(\sum_{t=1}^n w_t X_t \geq \epsilon \right) \leq 2 \exp \left(-\frac{\epsilon^2}{2 \sum_{t=1}^n w_t \sigma_t^2} \right)$$

Proof. For a complete demonstration of this statement, we refer to Lemma C.5 of Metelli et al. (2022). □

Lemma C.11. *Let $\epsilon \in (0, 1]$, then it holds that:*

$$H_{1,1}(\epsilon, p, T, q) \leq H_{1,1}(\epsilon, p, T, q)$$

Proof. We prove the equivalent statement, being $\epsilon \in (0, 1]$:

$$H_{1,1}(\epsilon, p, T, q)^{\epsilon} \leq H_{1,1}(\epsilon, p, T, q)$$

Recalling that the function x^{ϵ} is subadditive, being $\epsilon \in (0, 1]$, we have:

$$H_{1,1}(\epsilon, p, T, q)^{\epsilon} \leq \sum_{i=1}^K \frac{1}{p T q} \leq \sum_{i=1}^K \frac{1}{p T q} = H_{1,1}(\epsilon, p, T, q)$$

□

D Theoretical Analysis of a baseline: RR-SW

In this appendix, we provide the theoretical analysis for the algorithm Round Robin Sliding Window (RR-SW), as it represents the most intuitive baseline for this setting. First, we need to formalize the algorithm, whose pseudo-code is provided in Algorithm 3.

Algorithm 3: RR-SW.

Input : Time budget T , Number of arms K , Window size w

- 1 Initialize $t \leftarrow 1$
 - 2 Estimate $N \leftarrow \frac{T}{w}$
 - 3 **for** $i \in \{1, \dots, K\}$ **do**
 - 4 **for** $j \in \{1, \dots, N\}$ **do**
 - 5 Pull arm i and observe X_t
 - 6 $t \leftarrow t + 1$
 - 7 **end**
 - 8 Update \hat{p}_i
 - 9 **end**
 - 10 Recommend $\hat{p} = \arg \max_{i \in \{1, \dots, K\}} \hat{p}_i$
-

Algorithm The algorithm takes as input the time budget T and the number of arms K . Then, it computes the number of pulls $N = \frac{T}{K}$ we need to perform for each arm. After having computed the number of pulls, RR-SW plays all the arms N times in a round-robin fashion. After the N pulls, it estimates $\hat{\mu}_i$ using the last $\frac{T}{K}$ samples (i.e., the ones from $(i-1)N$ to N). Finally, it recommends i^* , corresponding to the one which the highest estimated $\hat{\mu}_i$.

Error probability bound Before presenting the error probability bound for RR-SW, we need to introduce Δ , which represents the minimum suboptimality gap at a given time budget T . It is actually the gap between the optimal arm and the first sub-optimal one. Formally: $\Delta = \min_{i \neq i^*} \mu_{i^*} - \mu_i$. Given this quantity, the error probability for the RR-SW algorithm can be bounded as follows.

Theorem D.1. Under Assumptions 2.1 and 2.2, considering a time budget T satisfying:

$$T \geq 2 \frac{1}{\Delta} c^{-1} \mu_{i^*} \ln \frac{K}{\Delta} \frac{1}{\Delta} \mu_{i^*} T; \quad (37)$$

the error probability of RR-SW is bounded by:

$$e_{\text{RR-SW}} \leq K \exp \left(-\frac{T}{8K^2} \Delta^2 \right);$$

Some comments are in order. First, it is worth noting how, as expected, by increasing the number of samples considered in the estimator, we reduce the error probability $e_{\text{RR-SW}}$ at the cost of a more strict constraint on the time budget T . This is due to the request that the arms must be already separated at the beginning of the window we use to estimate the $\hat{\mu}_i$. Second, the error probability scales as an (inverse) function only of the smallest suboptimality gap Δ .

D.1 Proofs

Before demonstrating Theorem D.1, we need to introduce the following technical lemma.

Lemma D.2 (Lower Bound for the Time Budget). Under Assumptions 2.1 and 2.2, the RR-SW algorithm is s.t. the minimum value for the horizon T ensuring that $\mathbb{P}(i^* \neq \hat{i}) \leq \epsilon$:

$$T \geq \frac{1}{\Delta} \mu_{i^*} \ln \frac{K}{\Delta} \frac{1}{\Delta} \mu_{i^*} T; \quad (38)$$

where $N = \frac{T}{K}$ and $\mu_{i^*} \geq \Delta$ is:

$$T \geq 2 \frac{1}{\Delta} c^{-1} \mu_{i^*} \ln \frac{K}{\Delta} \frac{1}{\Delta} \mu_{i^*} T;$$

Proof. First of all, we recall that $N = \frac{T}{K}$ is the number of times each arm has been pulled, considering K arms by running a round-robin procedure until we reach a time budget T . We consider the pessimistic estimator described in Section 3. Considering such an estimator and the RR-SW algorithm, which runs a round-robin procedure, what we get at the end of the time budget is a sliding-window estimator for the value of μ_{i^*} , which will include the last $\frac{T}{K}$ samples. In this lemma, we want to find the minimum value of the time budget T for which, at the first samples we consider, the real process of the arms are separated by at least $\frac{\Delta}{2}$. In this estimator, we consider samples in the range of $(i-1)N$ to N , so we need to ensure, given Assumption 2.1, that:

$$T \geq \frac{1}{\Delta} \mu_{i^*} \ln \frac{K}{\Delta} \frac{1}{\Delta} \mu_{i^*} T; \quad (38)$$

Given that, for Assumptions 2.1 and 2.2, it holds:

$$T \geq \frac{1}{\Delta} \mu_{i^*} \ln \frac{K}{\Delta} \frac{1}{\Delta} \mu_{i^*} T \quad \text{and} \quad T \geq \frac{1}{\Delta} \mu_{i^*} \ln \frac{K}{\Delta} \frac{1}{\Delta} \mu_{i^*} T$$

$$\alpha T c p1 \quad "q \frac{T}{K} \quad : \quad (39)$$

By introducing the term derived in Equation (39) into Equation (38) we obtain:

$$T c p1 \quad "q \frac{T}{K} \quad \alpha \frac{p_{2q} p T q}{2}.$$

This implies that the minimum time budget T which guarantees the initial condition of Equation (39) is:

$$T \geq 2^{-1} c^{-1} p1 \quad "q \quad K^{-1} \quad p_{2q}^{-1} p T q;$$

where $p_{2q} p T q$ is the minimum suboptimality gap ($p_{2q} p T q = \min_i |i p T q - i p T q_u$). \square

Now, we can find the error probability $e_T pRR-SWq$, which will hold for all the time budgets which satisfy the condition of Lemma D.2.

Theorem D.1. *Under Assumptions 2.1 and 2.2, considering a time budget T satisfying:*

$$T \geq 2^{-1} c^{-1} p1 \quad "q \quad K^{-1} \quad p_{2q}^{-1} p T q; \quad (37)$$

the error probability of RR-SW is bounded by:

$$e_T pRR-SWq \leq K \exp \left(- \frac{T}{8 K^2} p_{2q} p T q \right);$$

Proof. The RR-SW algorithm makes an error in predicting the best arm when, at the end of the process (at T total pulls), the optimal arm has a pessimistic estimator $\hat{\mu}_1 p N q$ that is not the highest among the arms (we consider w.l.o.g. that the best arm is the arm 1). Formally:

$$e_T pRR-SWq = P p D i P J K K : \hat{\mu}_1 p N q \neq \hat{\mu}_i p N q q \\ \leq \sum_{i P J K K} P p \hat{\mu}_1 p N q \neq \hat{\mu}_i p N q q;$$

Let us focus on a single arm i , where we want to upper bound the probability that $P p \hat{\mu}_1 p N q \neq \hat{\mu}_i p N q q$. Let us focus on the term inside the probability:

$$\hat{\mu}_i p N q \neq \hat{\mu}_1 p N q \\ \hat{\mu}_i p N q - \hat{\mu}_1 p N q \neq 0 \\ \frac{1}{T} \sum_{t=1}^T \mu_i p t q - \frac{1}{T} \sum_{t=1}^T \mu_1 p t q \neq 0 \quad (40)$$

$$\frac{1}{T} \sum_{t=1}^T \mu_i p t q - \frac{1}{T} \sum_{t=1}^T \mu_1 p t q \neq 0 \\ \frac{1}{T} \sum_{t=1}^T \mu_i p t q - \mu_i p N q \neq \frac{1}{T} \sum_{t=1}^T \mu_1 p t q - \mu_1 p N q \quad (41)$$

$$\frac{1}{T} \sum_{t=1}^T \mu_i p t q - \mu_i p N q \neq \frac{1}{T} \sum_{t=1}^T \mu_1 p t q - \mu_1 p N q \quad (42)$$

where we added $\mu_i p T q$ to derive Equation (40), and added $-\mu_1 p N q$ to derive Equation (41), we used the results in Lemma C.8 and from the fact that the reward function is increasing. Considering a time budget T satisfying Theorem D.2, and $\mu_i p T q \neq p_{2q} p T q; @ i P J K K$, we have that:

$$T \mu_1 p N q - \mu_1 q \leq \frac{p T q}{2}. \quad (43)$$

Equation (43) holds since we are considering a time budget T which satisfies a more restrictive condition (we are considering a time budget at which this separation already holds for $p1 \quad "q N$, so it also holds now).

	b	c	
Arm 1	37	1	1
Arm 2	10	0.88	1
Arm 3	1	0.78	1
Arm 4	10	0.7	1
Arm 5	20	0.5	1

Table 2: Numerical values of the parameters characterizing the functions for the synthetically generated setting.

F Experimental Details

In this section, we provide all the details about the presented experiments.

The payoff functions characterizing the arms shown in Figure 2 belong to the family:

$$F = \left\{ f_{p,q} : c \in [0, 1], \frac{b}{pb^c - q} \right\};$$

where $c \in [0, 1]$ and $b \in (0, \infty)$. Note that, by construction, all the functions laying in F satisfy the Assumptions 2.1 and 2.2. In particular, the largest value of b satisfying Assumption 2.2, for the setting presented in Section 8, is 1.3. In Table 2, we report the value of the parameters characterizing the function employed in the synthetically generated setting presented in the main paper.

F.1 Parameters Values for the Algorithms

This section provides a detailed view of the parameter values we employed in the presented experiments. More specifically, the parameters, which may still depend on the time budget T and on the number of arms K , are set as follows:

- UCB-E: for the exploration parameter a , we used the optimal value, i.e., the one that minimizes the upper bound of the error probability, as prescribed in Audibert et al. (2010), formally:

$$a = \frac{25pT - Kq}{36H_1};$$

where $H_1 = \sum_{i=1}^K \frac{1}{p_i q_i}$;

- R-UCBE: we used the value prescribed by Corollary 4.2 where we set the value 1.3;
- ETC and Rest-Sure: we set 0.8 and $U = 1$ as suggested by Cella et al. (2021).

F.2 Running Time

The code used for the results provided in this section has been run on an Intel(R) I7 9750H @ 2.6GHz CPU with 16 GB of *LPDDR4* system memory. The operating system was *MacOS 13.1*, and the experiments were run on *Python 3.10*. A run of R-UCBE over a time budget of $T = 3200$ takes 0.07 seconds (on average), while a run of R-SR takes 0.06 seconds (on average).

G Additional Experimental Results

In this section, we present additional results in terms of empirical error \bar{e}_T of R-UCBE, R-SR, and the other baselines presented in Section 8.

G.1 Challenging scenario

Here we test the algorithms on a challenging scenario in which we consider $K = 3$ arms whose increment changes *abruptly*. The setting is presented in Figure 6a. The results corresponding to

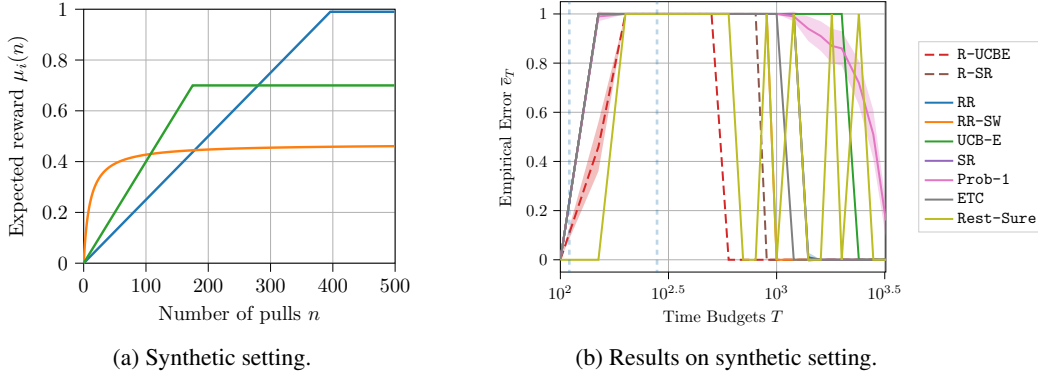


Figure 6: Challenging scenario in which the arm reward increment rate changes abruptly.

such a setting are presented in Figure 6b. In this case, the last time the optimal arm does not change anymore is $T = 280$. Similarly to the synthetic setting presented in the main paper, we have two different behaviors for time budgets $T \leq 280$ and $T > 280$. For short time budgets, the algorithm providing the best performance is Rest-Sure, and the second best is R-UCBE. Conversely, for time budgets $T > 550$ R-UCBE provides a correct suggestion in most of the cases providing an error $\mathbb{P}_{\mathcal{T}}\{R\text{-UCBE} \neq \text{opt}\} = 0.01$. Instead, Rest-Sure is not consistently providing reliable suggestions. This is allegedly due to the fact that such an algorithm has been designed to work in less general settings than the one we are tackling. Even in this case, the R-SR starts providing a small value for the error probability after R-UCBE does, at $T = 1000$. However, it is still better behaving than the other baseline algorithms. Note that Rest-Sure has a peculiar behavior. Indeed, it seems that even for large values of the time budget, it does not consistently suggest the optimal arm (i.e., the error probability does not go to zero). This is likely due to the nature of the parametric shape enforced by the algorithm, which may result in unpredictable behaviors when it does not reflect the nature of the real reward functions.

G.2 Sensitivity Analysis on the Noise Variance

In what follows, we report the analysis of the robustness of the analyzed algorithms as noise standard deviation σ changes in the collected samples. The setting we considered is the one described in Section 8. The results are provided in Figure 7. Let us focus on the performances of the R-UCBE algorithm. For small values of the standard deviation ($\sigma = 0.01$), we have the same behavior in terms of error probability, i.e., a progressive degradation of the performances for time budget $T = 150$. Indeed, at this time budget, the expected rewards of 3 arms are close to each other, and determining the optimal arm is a challenging problem. However, the performances are better or equal to all the other algorithms even at this point. Conversely, for values of the standard deviation $\sigma \geq 0.05$, the performance of R-UCBE starts to degrade, with behavior for $\sigma = 0.5$ which is constant w.r.t. the chosen time budget with a value of $\mathbb{P}_{\mathcal{T}}\{R\text{-UCBE} \neq \text{opt}\} = 0.8$. This suggests that such an algorithm suffers in the case the stochasticity of the problem is significant. Let us focus on R-SR. This algorithm does not change its performances w.r.t. changes in terms of σ . Indeed, only for $\sigma = 0.5$, we have that it does not provide an error probability close to zero for time budget $T > 1000$. However, excluding R-UCBE, we have that the R-SR algorithm is the best/close to the best performing algorithm. This is also true in the case of $\sigma = 0.5$, in which the R-UCBE fails in providing a reliable recommendation for the optimal arm with a large probability.

G.3 Real-world Experiment on IMDB dataset

Description We validate our algorithms and the baselines on an AutoML task, namely an *online best model selection* problem with a real-world dataset. We employ the IMDB dataset, made of 50,000 reviews of movies (scores from 0 to 10). We preprocessed the data as done by Metelli et al. (2022), and run the algorithms for time budgets $T \in \{500; 1000; \dots; 15000; 20000; 30000\}$. A graphical representation of the reward (in this case, represented by the accuracy) of the different models is presented in Figure 8. Since, in this case, we only had a single realization to estimate the error probability $\mathbb{P}_{\mathcal{T}}\{A \neq \text{opt}\}$, we report the success rate $\mathbb{P}_{\mathcal{T}}\{A = \text{opt}\}$ instead, i.e., the ratio between the number

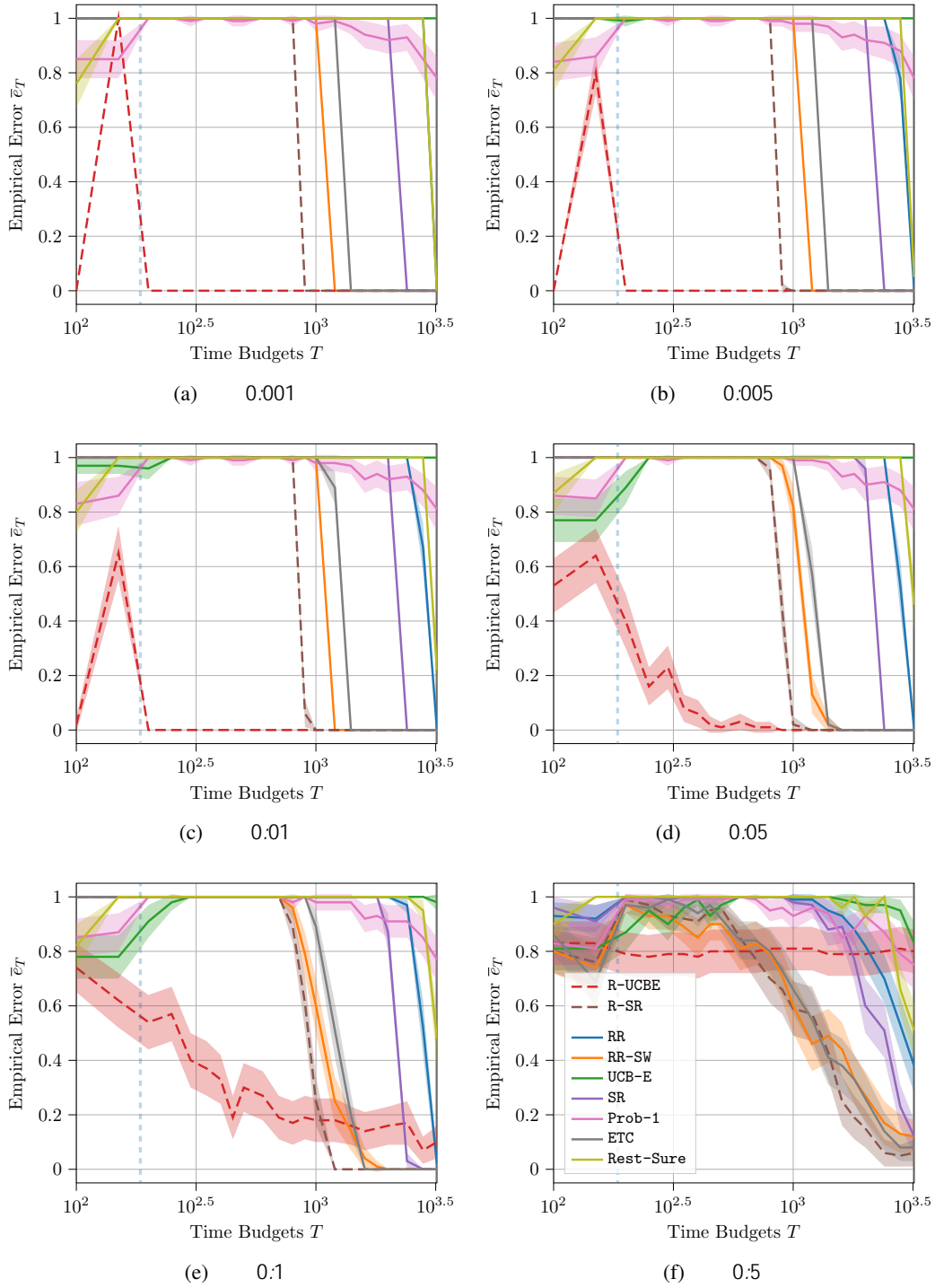


Figure 7: Empirical error probability for the synthetically generated setting, with different values of the noise standard deviation σ .

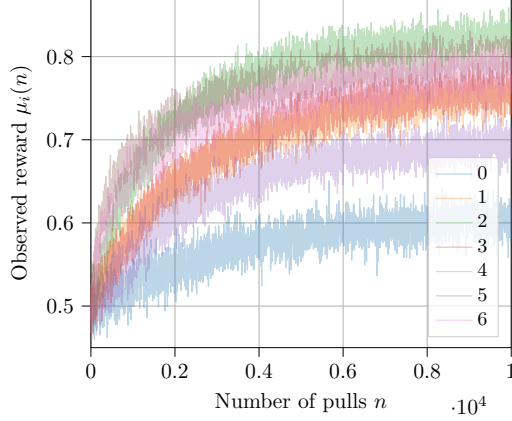


Figure 8: Rewards for the arms of the IMDB experiments.

	500	1000	2000	3000	4000	T		10000	15000	20000	30000	$\mathcal{R}\rho\mathcal{U}\mathcal{q}$	
Optimal Arm	5	5	2	2	2	2	2	2	2	2	2		
R-UCBE (ours)	6	5	2	5	2	2	2	2	2	2	2	9{11}	
R-SR (ours)	5	5	5	5	5	5	5	5	2	2	2	5{11}	
Algorithms	RR	5	5	5	5	5	5	5	5	5	5	2{11}	
	RR-SW	5	5	5	5	5	5	5	5	2	2	4{11}	
	SR	5	5	5	5	5	5	5	5	5	2	3{11}	
	UCB-E	5	5	5	5	5	5	5	5	5	5	2{11}	
	Prob-1	1	5	2	5	5	5	5	1	5	6	2	3{11}
	ETC	5	5	5	5	5	5	5	5	5	5	2	3{11}
	Rest-Sure	6	5	2	2	2	1	0	2	5	0	2	6{11}

Table 3: Optimal arm for different time budgets on the IMDB dataset (first row) and corresponding recommendations provided by the algorithms (second to last row). In the last column, we compute the corresponding success rate.

of times an algorithm provides a correct suggestion and the number of budget values we considered, formally defined as $\mathcal{R}\rho\mathcal{U}\mathcal{q} : \frac{1}{|T|} \sum_{T \in \mathcal{T}} \mathbb{1}_{\hat{u} = u}$ (the larger, the better).

Results The results are reported in Table 3. The algorithm with the largest success rate $\mathcal{R}\rho\mathcal{U}\mathcal{q}$ is the R-UCBE, while R-SR provides the third best success rate. Moreover, Rest-Sure, the only algorithm providing a success rate larger than R-SR, has issues with large time budgets since for $T \notin \{5000\}$ is able to provide only 2 correct guesses of the optimal arm over 6 attempts. Conversely, our algorithms progressively provide more and more correct guesses as the time budget T increases. The above results on a real-world dataset corroborate the evidence presented above that the proposed algorithms outperform state-of-the-art ones for the BAI problem in SRB.