

---

# Best Arm Identification for Stochastic Rising Bandits

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

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

## 18 1 Introduction

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

29 This work focuses on the *Stochastic Rising Bandits* (SRB), a specific instance of the *rested* bandit  
30 (Tekin and Liu, 2012) setting in which the expected reward of an arm increases according to the  
31 number of times it has been pulled. Online learning in such a scenario has been recently addressed  
32 from a regret minimization perspective by Metelli et al. (2022), in which the authors provide no-  
33 regret algorithms for the SRB setting in both the rested and restless cases. The SRB setting models  
34 several real-world scenarios where arms improve their performance over time. A classic example is  
35 the so-called *Combined Algorithm Selection and Hyperparameter optimization* (CASH, Thornton  
36 et al., 2013; Kotthoff et al., 2017; Erickson et al., 2020; Li et al., 2020; Zöllner and Huber, 2021), a  
37 problem of paramount importance in *Automated Machine Learning* (AutoML, Feurer et al., 2015;  
38 Yao et al., 2018; Hutter et al., 2019; Mussi et al., 2023). In CASH, the goal is to identify the *best*  
39 *learning algorithm* together with the *best hyperparameter* configuration for a given ML task (e.g.,

40 classification or regression). In this problem, every arm represents a hyperparameter tuner acting  
 41 on a specific learning algorithm. A pull corresponds to a unit of time/computation in which we  
 42 improve (on average) the hyperparameter configuration (via the tuner) for the corresponding learning  
 43 algorithm. CASH was handled in a bandit *Best Arm Identification* (BAI) fashion in Li et al. (2020)  
 44 and Cella et al. (2021). The former handles the problem by considering rising rested bandits with  
 45 *deterministic* rewards, failing to represent the intrinsic uncertain nature of such processes. Instead,  
 46 the latter, while allowing stochastic rewards, assumes that the expected rewards evolve according to a  
 47 *known* parametric functional class, whose parameters have to be learned.<sup>1</sup>

48 **Original Contributions** In this paper, we address the design of algorithms to solve the BAI task  
 49 in the rested SRB setting when a *fixed budget* is provided.<sup>2</sup> More specifically, we are interested in  
 50 algorithms guaranteeing a sufficiently large probability of recommending the arm with the largest  
 51 expected reward *at the end* of the time budget (as if only this arm were pulled from the beginning).  
 52 The main contributions of the paper are summarized as follows:<sup>3</sup>

- 53 • We propose two *algorithms* to solve the BAI problem in the SRB setting: R-UCBE (an optimistic  
 54 approach, Section 4) and R-SR (a phases-based rejection algorithm, Section 5). First, we intro-  
 55 duce specifically designed estimators required by the algorithms (Section 3). Then, we provide  
 56 guarantees on the error probability of the misidentification of the best arm.
- 57 • We derive the first error probability *lower bound* for the SRB setting, matched by our R-SR  
 58 algorithm up to logarithmic factors, which highlights the complexity of the problem and the need  
 59 for a sufficiently large time budget (Section 6).
- 60 • Finally, we conduct *numerical simulations* on synthetically generated data and a real-world online  
 61 best model selection problem. We compare the proposed algorithms with the ones available in the  
 62 bandit literature to tackle the SRB problem (Section 7).

## 63 2 Problem Formulation

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

67 **Setting** We consider a rested Multi-Armed Bandit problem  $\nu = (\nu_i)_{i \in \llbracket K \rrbracket}$  with a finite number  
 68 of arms  $K$ .<sup>4</sup> Let  $T \in \mathbb{N}$  be the time budget of the learning process. At every round  $t \in \llbracket T \rrbracket$ , the  
 69 agent selects an arm  $I_t \in \llbracket K \rrbracket$ , plays it, and observes a reward  $x_t \sim \nu_{I_t}(N_{I_t, t})$ , where  $\nu_{I_t}(N_{I_t, t})$   
 70 is the reward distribution of the chosen arm  $I_t$  at round  $t$  and depends on the number of pulls  
 71 performed so far  $N_{i, t} := \sum_{\tau=1}^t \mathbb{1}\{I_\tau = i\}$  (i.e., rested). The rewards are stochastic, formally  
 72  $x_t := \mu_{I_t}(N_{I_t, t}) + \eta_t$ , where  $\mu_{I_t}(\cdot)$  is the expected reward of arm  $I_t$  and  $\eta_t$  is a zero-mean  $\sigma^2$ -  
 73 subgaussian noise, conditioned to the past.<sup>5</sup> As customary in the bandit literature, we assume that  
 74 the rewards are bounded in expectation, formally  $\mu_i(n) \in [0, 1], \forall i \in \llbracket K \rrbracket, n \in \llbracket T \rrbracket$ . As in (Metelli  
 75 et al., 2022), we focus on a particular family of rested bandits in which the expected rewards are  
 76 monotonically *non-decreasing* and *concave* in expectation.

77 **Assumption 2.1** (Non-decreasing and concave expected rewards). *Let  $\nu$  be a rested MAB, defining*  
 78  $\gamma_i(n) := \mu_i(n+1) - \mu_i(n)$ , *for every  $n \in \mathbb{N}$  and every arm  $i \in \llbracket K \rrbracket$  the rewards are non-decreasing*  
 79 *and concave, formally:*

$$\text{Non-decreasing: } \gamma_i(n) \geq 0, \quad \text{Concave: } \gamma_i(n+1) \leq \gamma_i(n).$$

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

<sup>1</sup>A complete discussion of the related works is available in Appendix A. Additional motivating examples are discussed in Appendix B.

<sup>2</sup>We focus on the rested setting only and, thus, from now on, we will omit “rested” in the setting name.

<sup>3</sup>The proofs of all the statements in this work are provided in Appendix D.

<sup>4</sup>Let  $y, z \in \mathbb{N}$ , we denote with  $\llbracket z \rrbracket := \{1, \dots, z\}$ , and with  $\llbracket y, z \rrbracket := \{y, \dots, z\}$ .

<sup>5</sup>A zero-mean random variable  $x$  is  $\sigma^2$ -subgaussian if it holds  $\mathbb{E}_x[e^{\xi x}] \leq e^{\frac{\sigma^2 \xi^2}{2}}$  for every  $\xi \in \mathbb{R}$ .

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

$$i^*(T) := \arg \max_{i \in \llbracket K \rrbracket} \mu_i(T),$$

where we highlighted the dependence on  $T$  as, with different values of the budget,  $i^*(T)$  may  
 change. Let  $i \in \llbracket K \rrbracket \setminus \{i^*(T)\}$  be a suboptimal arm, we define the suboptimality gap as  $\Delta_i(T) :=$   
 $\mu_{i^*(T)}(T) - \mu_i(T)$ . We employ the notation  $(i) \in \llbracket K \rrbracket$  to denote the  $i^{\text{th}}$  best arm at time  $T$  (arbitrarily  
 breaking ties), i.e., we have  $\Delta_{(2)}(T) \leq \dots \leq \Delta_{(K)}(T)$ . Given an algorithm  $\mathfrak{A}$  that recommends  
 $\hat{I}^*(T) \in \llbracket K \rrbracket$  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  $T$ :

$$e_T(\mathfrak{A}) := \mathbb{P}_{\mathfrak{A}}(\hat{I}^*(T) \neq i^*(T)).$$

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

91 **Assumption 2.2** (Bounded  $\gamma_i(n)$ ). *Let  $\nu$  be a rested MAB, there exist  $c > 0$  and  $\beta > 1$  such that for*  
 92 *every arm  $i \in \llbracket K \rrbracket$  and number of pulls  $n \in \llbracket 0, T \rrbracket$  it holds that  $\gamma_i(n) \leq cn^{-\beta}$ .*

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

### 97 3 Estimators

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

Let  $\varepsilon \in (0, 1/2)$  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*  $h(N_{i,t-1}) := \lfloor \varepsilon N_{i,t-1} \rfloor$   
 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  $h(N_{i,t-1})$  rounds in which the  $i^{\text{th}}$  arm was pulled as:

$$\mathcal{T}_{i,t} := \{\tau \in \llbracket T \rrbracket : I_\tau = i \wedge N_{i,\tau} = N_{i,t-1} - l, l \in \llbracket 0, h(N_{i,t-1}) - 1 \rrbracket\}.$$

100 Furthermore, the set of the pairs of rounds  $\tau$  and  $\tau'$  belonging to the sets of the last and second-last  
 101  $h(N_{i,t-1})$ -wide windows of the  $i^{\text{th}}$  arm is defined as:

$$\mathcal{S}_{i,t} := \{(\tau, \tau') \in \llbracket T \rrbracket \times \llbracket T \rrbracket : I_\tau = I_{\tau'} = i \wedge N_{i,\tau} = N_{i,t-1} - l, \\ N_{i,\tau'} = N_{i,\tau} - h(N_{i,t-1}), l \in \llbracket 0, h(N_{i,t-1}) - 1 \rrbracket\}.$$

102 In the following, we design a *pessimistic* estimator and an *optimistic* estimator of the expected reward  
 103 of each arm at the end of the budget time  $T$ , i.e.,  $\mu_i(T)$ .<sup>7</sup>

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

$$\hat{\mu}_i(N_{i,t-1}) := \frac{1}{h(N_{i,t-1})} \sum_{\tau \in \mathcal{T}_{i,t}} x_\tau. \quad (1)$$

108 The estimator enjoys the following concentration property.

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

<sup>7</sup>Naïvely computing the estimators from their definition requires  $\mathcal{O}(h(N_{i,t-1}))$  number of operations. An efficient way to incrementally update them, using  $\mathcal{O}(1)$  operations, is provided in Appendix C.

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

$$\hat{\beta}_i(n) - \hat{\zeta}_i(n) \leq \hat{\mu}_i(n) - \mu_i(n) \leq \hat{\beta}_i(n), \quad (2)$$

111 where  $\hat{\beta}_i(n) := \sigma \sqrt{\frac{a}{h(n)}}$  and  $\hat{\zeta}_i(n) := \frac{1}{2}(2T - n + h(n) - 1) \gamma_i(n - h(n) + 1)$ .

112 As supported by intuition, we observe that the estimator  
 113 is affected by a negative bias that is represented by  
 114  $\hat{\zeta}_i(n)$  that vanishes as  $n \rightarrow \infty$  under Assumption 2.1  
 115 with a rate that depends on the increment functions  
 116  $\gamma_i(\cdot)$ . Considering also the term  $\hat{\beta}_i(n)$  and recalling  
 117 that  $h(n) = \mathcal{O}(n)$ , under Assumption 2.2, the overall  
 118 concentration rate is  $\mathcal{O}(n^{-1/2} + cTn^{-\beta})$ .

119 **Optimistic Estimator** The *optimistic* estimator  
 120  $\check{\mu}_i^T(N_{i,t-1})$  is a positively biased estimation of  $\mu_i(T)$   
 121 obtained assuming that function  $\mu_i(\cdot)$  linearly increas-  
 122 es up to time  $T$ . This corresponds to the  
 123 maximum value admissible under Assumption 2.1  
 124 (due to the *Concavity* constraint). The estimator is  
 125 constructed by adding to the pessimistic estimator  
 126  $\hat{\mu}_i(N_{i,t-1})$  an estimate of the increment occurring  
 127 in the next step up to  $T$ . The latter uses the last  
 128  $2h(N_{i,t-1})$  samples to obtain an upper bound of such  
 129 growth thanks to the concavity assumption, formally:

$$\check{\mu}_i^T(N_{i,t-1}) := \hat{\mu}_i(N_{i,t-1}) + \sum_{(j,k) \in \mathcal{S}_{i,t}} (T-j) \frac{x_j - x_k}{h(N_{i,t-1})^2}. \quad (3)$$

130 The estimator displays the following concentration guarantee.

131 **Lemma 3.2** (Concentration of  $\check{\mu}_i^T$ ). *Under Assumption 2.1, for every  $a > 0$ , simultaneously for every*  
 132 *arm  $i \in \llbracket K \rrbracket$  and number of pulls  $n \in \llbracket 0, T \rrbracket$ , with probability at least  $1 - 2TK e^{-a/10}$  it holds that:*

$$\check{\beta}_i^T(n) \leq \check{\mu}_i^T(n) - \mu_i(n) \leq \check{\beta}_i^T(n) + \check{\zeta}_i^T(n), \quad (4)$$

133 where  $\check{\beta}_i^T(n) := \sigma \cdot (T - n + h(n) - 1) \sqrt{\frac{a}{h(n)^3}}$  and  $\check{\zeta}_i^T(n) := \frac{1}{2}(2T - n + h(n) - 1) \gamma_i(n - 2h(n) + 1)$ .

134 Differently from the pessimistic estimation, the optimistic one displays a positive vanishing bias  
 135  $\check{\zeta}_i^T(n)$ . Under Assumption 2.2, we observe that the overall concentration rate is  $\mathcal{O}(Tn^{-3/2} + cTn^{-\beta})$ .

## 136 4 Optimistic Algorithm: Rising Upper Confidence Bound Exploration

137 In this section, we introduce and analyze Rising Upper Confidence Bound Exploration  
 138 (R-UCBE) an *optimistic* error probability minimization algorithm for the SRB setting with a fixed  
 139 budget. The algorithm explores by means of a UCB-like approach and, for this reason, makes use of  
 140 the optimistic estimator  $\check{\mu}_i^T$  plus a bound to account for the uncertainty of the estimation.<sup>8</sup>

141 **Algorithm** The algorithm, whose pseudo-code is reported in Algorithm 1, requires as input an  
 142 exploration parameter  $a \geq 0$ , the window size  $\varepsilon \in (0, 1/2)$ , the time budget  $T$ , and the number of  
 143 arms  $K$ . At first, it initializes to zero the counters  $N_{i,0}$ , and sets to  $+\infty$  the upper bounds  $B_i^T(N_{i,0})$   
 144 of all the arms (Line 2). Subsequently, at each time  $t \in \llbracket T \rrbracket$ , the algorithm selects the arm  $I_t$  with the  
 145 largest upper confidence bound (Line 4):

$$I_t \in \arg \max_{i \in \llbracket K \rrbracket} B_i^T(N_{i,t-1}) := \check{\mu}_i^T(N_{i,t-1}) + \check{\beta}_i^T(N_{i,t-1}), \quad (5)$$

$$\text{with: } \check{\beta}_i^T(N_{i,t-1}) := \sigma \cdot (T - N_{i,t-1} + h(N_{i,t-1}) - 1) \sqrt{\frac{a}{h(N_{i,t-1})^3}}, \quad (6)$$

<sup>8</sup>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.

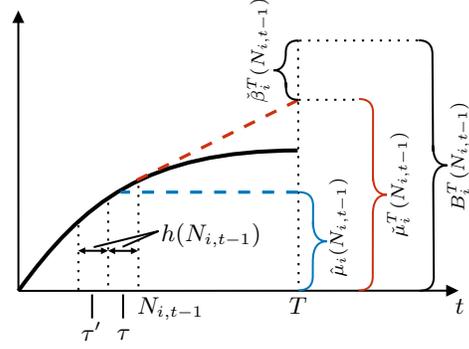


Figure 1: Graphical representation of the pessimistic  $\hat{\mu}_i(N_{i,t-1})$  and the optimistic  $\check{\mu}_i^T(N_{i,t-1})$  estimators.

146 where  $\tilde{\beta}_i^T(N_{i,t-1})$  represents the exploration bonus (a graphical representation is reported in Figure 1).  
 147 Once the arm is chosen, the algorithm plays it and observes the feedback  $x_t$  (Line 5). Then, the  
 148 optimistic estimate  $\tilde{\mu}_{I_t}^T(N_{I_t,t})$  and the exploration bonus  $\tilde{\beta}_{I_t}^T(N_{I_t,t})$  of the selected arm  $I_t$  are updated  
 149 (Lines 8-9). The procedure is repeated until the algorithm reaches the time budget  $T$ . The final  
 150 recommendation of the best arm is performed using the last computed values of the bounds  $B_i^T(N_{i,T})$ ,  
 151 returning the arm  $\hat{I}^*(T)$  corresponding to the largest upper confidence bound (Line 12).

152 **Bound on the Error Probability of R-UCBE** We now provide bounds on the error probability for  
 153 R-UCBE. We start with a general analysis that makes no assumption on the increments  $\gamma_i(\cdot)$  and, then,  
 154 we provide a more explicit result under Assumption 2.2. The general result is formalized as follows.

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

$$T - \sum_{i \neq i^*(T)} y_i(a) \geq 1, \quad (7)$$

156 where for every  $i \in \llbracket K \rrbracket$ ,  $y_i(a)$  is the largest integer for which it holds:

$$\underbrace{T\gamma_i(\lfloor (1-2\varepsilon)y \rfloor)}_{(A)} + \underbrace{2T\sigma\sqrt{\frac{a}{[\varepsilon y]^3}}}_{(B)} \geq \Delta_i(T). \quad (8)$$

157 If  $a^*$  exists, then for every  $a \in [0, a^*]$  the error probability of R-UCBE is bounded by:

$$e_T(\text{R-UCBE}) \leq 2TK \exp\left(-\frac{a}{10}\right). \quad (9)$$

158 Some comments are in order. First,  $a^*$  is defined implicitly, depending on the constants  $\sigma$ ,  $T$ , the  
 159 increments  $\gamma_i(\cdot)$ , and the suboptimality gaps  $\Delta_i(T)$ . In principle, there might exist no  $a^* > 0$   
 160 fulfilling condition in Equation (7) (this can happen, for instance, when the budget  $T$  is not large  
 161 enough), and, in such a case, we are unable to provide theoretical guarantees on the error probability  
 162 of R-UCBE. Second, the result presented in Theorem 4.1 holds for generic increasing and concave  
 163 expected reward functions. This result shows that, as expected, the error probability decreases when  
 164 the exploration parameter  $a$  increases. However, this behavior stops when we reach the threshold  $a^*$ .  
 165 Intuitively, the value of  $a^*$  sets the maximum amount of exploration we should use for learning.

166 Under Assumption 2.2, i.e., using the knowledge on the increment  $\gamma_i(\cdot)$  upper bound, we derive a  
 167 result providing conditions on the time budget  $T$  under which  $a^*$  exists and an explicit value for  $a^*$ .

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

$$T \geq \begin{cases} \left( c^{\frac{1}{\beta}} (1-2\varepsilon)^{-1} (H_{1,1/\beta}(T)) + (K-1) \right)^{\frac{\beta}{\beta-1}} & \text{if } \beta \in (1, 3/2) \\ \left( c^{\frac{2}{3}} (1-2\varepsilon)^{-\frac{2}{3}\beta} (H_{1,2/3}(T)) + (K-1) \right)^3 & \text{if } \beta \in [3/2, +\infty) \end{cases}, \quad (10)$$

169 there exists  $a^* > 0$  defined as:

$$a^* = \begin{cases} \frac{c^3}{4\sigma^2} \left( \left( \frac{T^{1-1/\beta} - (K-1)}{H_{1,1/\beta}(T)} \right)^\beta - c(1-2\varepsilon)^{-\beta} \right)^2 & \text{if } \beta \in (1, 3/2) \\ \frac{c^3}{4\sigma^2} \left( \left( \frac{T^{1/3} - (K-1)}{H_{1,2/3}(T)} \right)^{3/2} - c(1-2\varepsilon)^{-\beta} \right)^2 & \text{if } \beta \in [3/2, +\infty) \end{cases},$$

170 where  $H_{1,\eta}(T) := \sum_{i \neq i^*(T)} \frac{1}{\Delta_i^\eta(T)}$  for  $\eta > 0$ . Then, for every  $a \in [0, a^*]$ , the error probability of  
 171 R-UCBE is bounded by:

$$e_T(\text{R-UCBE}) \leq 2TK \exp\left(-\frac{a}{10}\right).$$

172 First of all, we notice that the error probability  $e_T(\text{R-UCBE})$  presented in Theorem 4.2 holds under  
 173 the condition that the time budget  $T$  fulfills Equation (10). We defer a more detailed discussion  
 174 on this condition to Remark 5.1, where we show that the existence of a finite value of  $T$  fulfilling  
 175 Equation (10) is ensured under mild conditions.

176 Let us remark that term  $H_{1,\eta}(T)$  characterizes the complexity of the SRB setting, corresponding to  
 177 term  $H_1$  of Audibert et al. (2010) for the classical BAI problem when  $\eta = 2$ . As expected, in the  
 178 small- $\beta$  regime (i.e.,  $\beta \in (1, 3/2]$ ), looking at the dependence of  $H_{1,1/\beta}(T)$  on  $\beta$ , we realize that

---

**Algorithm 1: R-UCBE.**

---

**Input :** Time budget  $T$ , Number of arms  $K$ ,  
Window size  $\varepsilon$ , Exploration parameter  $a$

- 1 Initialize  $N_{i,0} = 0$ ,
- 2  $B_i^T(0) = +\infty, \forall i \in \llbracket K \rrbracket$
- 3 **for**  $t \in \llbracket T \rrbracket$  **do**
- 4     Compute  $I_t \in \arg \max_{i \in \llbracket K \rrbracket} B_i^T(N_{i,t-1})$
- 5     Pull arm  $I_t$  and observe  $x_t$
- 6      $N_{I_t,t} \leftarrow N_{I_t,t-1} + 1$
- 7      $N_{i,t} \leftarrow N_{i,t-1}, \forall i \neq I_t$
- 8     Update  $\check{\mu}_{I_t}^T(N_{I_t,t})$
- 9     Update  $\check{\beta}_{I_t}^T(N_{I_t,t})$
- 10    Compute  $B_{I_t}^T(N_{I_t,t}) = \check{\mu}_{I_t}^T(N_{I_t,t}) + \check{\beta}_{I_t}^T(N_{I_t,t})$
- 11 **end**
- 12 Recommend  $\hat{I}^*(T) \in \arg \max_{i \in \llbracket K \rrbracket} B_i^T(N_{i,T})$

---



---

**Algorithm 2: R-SR.**

---

**Input :** Time budget  $T$ , Number of arms  $K$ ,  
Window size  $\varepsilon$

- 1 Initialize  $t \leftarrow 1, N_0 = 0, \mathcal{X}_0 = \llbracket K \rrbracket$
- 2 **for**  $j \in \llbracket K - 1 \rrbracket$  **do**
- 3     **for**  $i \in \mathcal{X}_{j-1}$  **do**
- 4         **for**  $l \in \llbracket N_{j-1} + 1, N_j \rrbracket$  **do**
- 5             Pull arm  $i$  and observe  $x_t$
- 6              $t \leftarrow t + 1$
- 7         **end**
- 8         Update  $\hat{\mu}_i(N_j)$
- 9     **end**
- 10    Define  $\bar{I}_j \in \arg \min_{i \in \mathcal{X}_{j-1}} \hat{\mu}_i(N_j)$
- 11    Update  $\mathcal{X}_j = \mathcal{X}_{j-1} \setminus \{\bar{I}_j\}$
- 12 **end**
- 13 Recommend  $\hat{I}^*(T) \in \mathcal{X}_{K-1}$  (unique)

---

179 the complexity of a problem decreases as the parameter  $\beta$  increases. Indeed, the larger  $\beta$ , the faster  
180 the expected reward reaches a stationary behavior. Nevertheless, even in the large- $\beta$  regime (i.e.,  
181  $\beta > 3/2$ ), the complexity of the problem is governed by  $H_{1,2/3}(T)$ , leading to an error probability  
182 larger than the corresponding one for BAI in standard bandits (Audibert et al., 2010). This can be  
183 explained by the fact that R-UCBE uses the optimistic estimator that, as shown in Section 3, enjoys a  
184 slower concentration rate compared to the standard sample mean, even for stationary bandits.

185 This two-regime behavior has an interesting interpretation when comparing Corollary 4.2 with  
186 Theorem 4.1. Indeed,  $\beta = 3/2$  is the break-even threshold in which the two terms of the l.h.s. of  
187 Equation (8) have the same convergence rate. Specifically, the term (A) takes into account the  
188 expected rewards growth (i.e., the bias in the estimators), while (B) considers the uncertainty in  
189 the estimations of the R-UCBE algorithm (i.e., the variance). Intuitively, when the expected reward  
190 function displays a slow growth (i.e.,  $\gamma_i(n) \leq cn^{-\beta}$  with  $\beta < 3/2$ ), the bias term (A) dominates  
191 the variance term (B) and the value of  $a^*$  changes accordingly. Conversely, when the variance term  
192 (B) is the dominant one (i.e.,  $\gamma_i(n) \leq cn^{-\beta}$  with  $\beta > 3/2$ ), the threshold  $a^*$  is governed by the  
193 estimation uncertainty, being the bias negligible.

194 As common in optimistic algorithms for BAI (Audibert et al., 2010), setting a theoretically sound  
195 value of exploration parameter  $a$  (i.e., computing  $a^*$ ), requires additional knowledge of the setting,  
196 namely the complexity index  $H_{1,\eta}(T)$ .<sup>9</sup> In the next section, we propose an algorithm that relaxes this  
197 requirement.

## 198 5 Phase-Based Algorithm: Rising Successive Rejects

199 In this section, we introduce the Rising Successive Rejects (R-SR), a phase-based solution  
200 inspired by the one proposed by Audibert et al. (2010), which overcomes the drawback of R-UCBE of  
201 requiring knowledge of  $H_{1,\eta}(T)$ .

202 **Algorithm** R-SR, whose pseudo-code is reported in Algorithm 2, takes as input the time budget  $T$   
203 and the number of arms  $K$ . At first, it initializes the set of the active arms  $\mathcal{X}_0$  with all the available  
204 arms (Line 1). This set will contain the arms that are still eligible candidates to be recommended.  
205 The entire process proceeds through  $K - 1$  phases. More specifically, during the  $j^{\text{th}}$  phase, the arms  
206 still remaining in the active arms set  $\mathcal{X}_{j-1}$  are played (Line 5) for  $N_j - N_{j-1}$  times each, where:

$$N_j := \left\lceil \frac{1}{\log(K)} \frac{T - K}{K + 1 - j} \right\rceil, \quad (11)$$

207 and  $\overline{\log}(K) := \frac{1}{2} + \sum_{i=2}^K \frac{1}{i}$ . At the end of each phase, the arm with the smallest value of the  
208 pessimistic estimator  $\hat{\mu}_i(N_j)$  is discarded from the set of active arms (Line 11). At the end of the  
209  $(K - 1)^{\text{th}}$  phase, the algorithm recommends the (unique) arm left in  $\mathcal{X}_{K-1}$  (Line 13).

---

<sup>9</sup>We defer the empirical study of the sensitivity of  $a$  to Section 7.

210 It is worth noting that R-SR makes use of the pessimistic estimator  $\hat{\mu}_i(n)$ . Even if both estimators  
 211 defined in Section 3 are viable for R-SR, the choice of using the pessimistic estimator is justified  
 212 by its better concentration rate  $\mathcal{O}(n^{-1/2})$  compared to that of the optimistic estimator  $\mathcal{O}(Tn^{-3/2})$ ,  
 213 being  $n \leq T$  (see Section 3).

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

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

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

$$T \geq 2^{\frac{\beta+1}{\beta-1}} c^{\frac{1}{\beta-1}} \overline{\log}(K)^{\frac{\beta}{\beta-1}} \max_{i \in \llbracket 2, K \rrbracket} \left\{ i^{\frac{\beta}{\beta-1}} \Delta_{(i)}(T)^{-\frac{1}{\beta-1}} \right\}, \quad (12)$$

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

$$e_T(\text{R-SR}) \leq \frac{K(K-1)}{2} \exp\left(-\frac{\varepsilon}{8\sigma^2} \cdot \frac{T-K}{\overline{\log}(K)H_2(T)}\right),$$

221 where  $H_2(T) := \max_{i \in \llbracket K \rrbracket} \{i \Delta_{(i)}(T)^{-2}\}$  and  $\overline{\log}(K) = \frac{1}{2} + \sum_{i=2}^K \frac{1}{i}$ .

222 Similar to the R-UCBE, the complexity of the problem is characterized by term  $H_2(T)$  that, for the  
 223 standard MAB setting, reduces to the  $H_2$  term of Audibert et al. (2010). Furthermore, when the  
 224 condition of Equation (12) on the time budget  $T$  is satisfied, the error probability coincides with that  
 225 of the SR algorithm for standard MABs (apart for constant terms). The following remark elaborates  
 226 on the conditions of Equations (10) and (12) about the minimum requested time budget.

227 **Remark 5.1** (About the minimum time budget  $T$ ). *To satisfy the  $e_T$  bounds presented in Corollary 4.2  
 228 and Theorem 5.1, R-UCBE and R-SR require the conditions provided by Equations (10) and (12)  
 229 about the time budget  $T$ , respectively. First, let us notice that if the suboptimal arms converge to  
 230 an expected reward different from that of the optimal arm as  $T \rightarrow +\infty$ , it is always possible to  
 231 find a finite value of  $T < +\infty$  such that these conditions are fulfilled. Formally, assume that there  
 232 exists  $T_0 < +\infty$  and that for every  $T \geq T_0$  we have that for all suboptimal arms  $i \neq i^*(T)$  it holds  
 233 that  $\Delta_i(T) \geq \Delta_\infty > 0$ . In such a case, the l.h.s. of Equations (10) and (12) are upper bounded by  
 234 a function of  $\Delta_\infty$  and are independent on  $T$ . Instead, if a suboptimal arm converges to the same  
 235 expected reward as the optimal arm when  $T \rightarrow +\infty$ , the identification problem is more challenging  
 236 and, depending on the speed at which the two arms converge as a function of  $T$ , might slow down the  
 237 learning process arbitrarily. This should not surprise as the BAI problem becomes non-learnable  
 238 even in standard (stationary) MABs when multiple optimal arms are present (Heide et al., 2021).*

## 239 6 Lower Bound

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

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

243 **Theorem 6.1.** *For every algorithm  $\mathfrak{A}$ , there exists a deterministic SRB satisfying Assumptions 2.1  
 244 and 2.2 such that the optimal arm  $i^*(T)$  cannot be identified for some time budgets  $T$  unless:*

$$T \geq H_{1,1/(\beta-1)}(T) = \sum_{i \neq i^*(T)} \frac{1}{\Delta_i(T)^{\frac{1}{\beta-1}}}. \quad (13)$$

245 Theorem 6.1 formalizes the intuition that any of the suboptimal arms must be pulled a sufficient  
 246 number of times to ensure that, if pulled further, it cannot become the optimal arm. It is worth  
 247 comparing this bound on the time budget with the corresponding conditions on the minimum  
 248 time budget requested by Equations (10) and (12) for R-UCBE and R-SR, respectively. Regarding  
 249 R-UCBE, we notice that the minimum admissible time budget in the small- $\beta$  regime is of order  
 250  $H_{1,1/\beta}(T)^{\beta/(\beta-1)}$  which is larger than term  $H_{1,1/(\beta-1)}(T)$  of Equation (13).<sup>10</sup> Similarly, in the

<sup>10</sup>See Lemma D.12.

	Error Probability $e_T(\cdot)$	Time Budget $T$
SRB	$\frac{1}{4} \exp\left(-\frac{8T}{\sigma^2 \sum_{i \neq i^*(T)} \frac{1}{\Delta_i^2(T)}}\right)$	$\sum_{i \neq i^*(T)} \frac{1}{\Delta_i(T)^{\beta-1}}$
R-UCBE	$2TK \exp\left(-\frac{a}{10}\right)$	$\begin{cases} \left(c^{\frac{1}{\beta}}(1-2\varepsilon)^{-1} \left(\sum_{i \neq i^*(T)} \frac{1}{\Delta_i^{1/\beta}(T)}\right) + (K-1)\right)^{\beta-1} & \text{if } \beta \in (1, 3/2) \\ \left(c^{\frac{2}{3}}(1-2\varepsilon)^{-\frac{2}{3}\beta} \left(\sum_{i \neq i^*(T)} \frac{1}{\Delta_i^{2/3}(T)}\right) + (K-1)\right)^3 & \text{if } \beta \in [3/2, +\infty) \end{cases}$
R-SR	$\frac{K(K-1)}{2} \exp\left(-\frac{\varepsilon}{8\sigma^2} \frac{T-K}{\log(K) \max_{i \in [K]} \{i\Delta_{(i)}^{-2}(T)\}}\right)$	$2^{\frac{1+\beta}{\beta-1}} c^{\frac{1}{\beta-1}} \log(K)^{\frac{\beta}{\beta-1}} \max_{i \in [2, K]} \{i^{\frac{\beta}{\beta-1}} \Delta_{(i)}(T)^{-\frac{1}{\beta-1}}\}$

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

251 large- $\beta$  regime (i.e.,  $\beta > 3/2$ ), the R-UCBE requirement is of order  $H_{1,2/3}(T)^3 \geq H_{1,2}(T)$  which  
252 is larger than the term of Theorem 6.1 since  $1/(\beta-1) < 2$ . Concerning R-SR, it is easy to show  
253 that  $H_{1,1/(\beta-1)}(T) \approx \max_{i \in [2, K]} i\Delta_{(i)}(T)^{-1/(\beta-1)}$ , apart from logarithmic terms, by means of  
254 the argument provided by (Audibert et al., 2010, Section 6.1). Thus, up to logarithmic terms,  
255 Equation (12) provides a tight condition on the minimum budget.

256 **Error Probability Lower Bound** We now present a lower bound on the error probability.

257 **Theorem 6.2.** *For every algorithm  $\mathfrak{A}$  run with a time budget  $T$  fulfilling Equation (13), there exists a*  
258 *SRB satisfying Assumptions 2.1 and 2.2 such that the error probability is lower bounded by:*

$$e_T(\mathfrak{A}) \geq \frac{1}{4} \exp\left(-\frac{8T}{\sigma^2 H_{1,2}(T)}\right), \text{ where } H_{1,2}(T) := \sum_{i \neq i^*(T)} \frac{1}{\Delta_i^2(T)}.$$

259 Some comments are in order. First, we stated the lower bound for the case in which the minimum  
260 time budget satisfies the inequality of Theorem 6.1, which is a necessary condition for identifying the  
261 optimal arm. Second, the lower bound on the error probability matches, up to logarithmic factors,  
262 that of our R-SR, suggesting the superiority of this algorithm compared to R-UCBE. Finally, provided  
263 that the identifiability condition of Equation (13), such a result corresponds to that of the standard  
264 (stationary) MABs (Audibert et al., 2010; Kaufmann et al., 2016). A summary of all the bounds  
265 provided in the paper is presented in Table 1.

## 266 7 Numerical Validation

267 In this section, we provide a numerical validation of R-UCBE and R-SR. We compare them with  
268 state-of-the-art bandit baselines designed for stationary and non-stationary BAI in a synthetic setting,  
269 and we evaluate the sensitivity of R-UCBE to its exploration parameter  $a$ . Additional details about the  
270 experiments presented in this section are available in Appendix G. Additional experimental results on  
271 both synthetic settings and in a real-world experiment are available in Appendix H.<sup>11</sup>

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

- 273 • RR: uniformly pulls all the arms until the budget ends in a *round-robin* fashion and, in the end,  
274 makes a recommendation based on the empirical mean of their reward over the collected samples;
- 275 • RR-SW: makes use of the same exploration strategy as RR to pull arms but makes a recommendation  
276 based on the empirical mean over the last  $\frac{\varepsilon T}{K}$  collected samples from an arm.<sup>12</sup>
- 277 • UCB-E and SR (Audibert et al., 2010): algorithms for the stationary BAI problem;
- 278 • Prob-1 (Abbasi-Yadkori et al., 2018): an algorithm dealing with the adversarial BAI setting;
- 279 • ETC and Rest-Sure (Cella et al., 2021): algorithms developed for the decreasing loss BAI setting.<sup>13</sup>

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

<sup>11</sup>The code to run the experiments is available in the supplementary material. It will be published in a public repository conditionally to the acceptance of the paper.

<sup>12</sup>The formal description of this baseline, as well as its theoretical analysis, is provided in Appendix E.

<sup>13</sup>This problem is equivalent to ours, given a linear transformation of the reward.

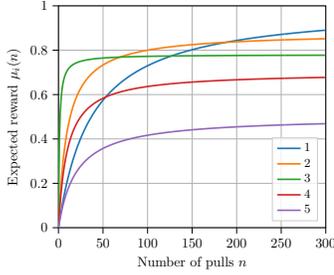


Figure 2: Expected values  $\mu_i(n)$  for the arms of the synthetic setting.

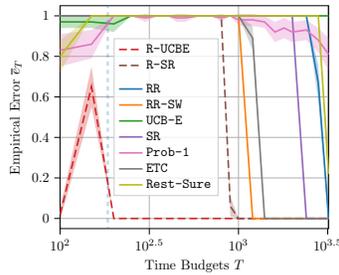


Figure 3: Empirical error rate for the synthetically generated setting (100 runs, mean  $\pm$  95% c.i.).

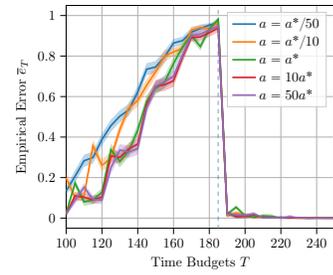


Figure 4: Empirical error rate for the R-UCBE at different  $a$  (1000 runs, mean  $\pm$  95% c.i.).

282 **Setting** To assess the quality of the recommendation  $\hat{I}^*(T)$  provided by our algorithms, we consider  
 283 a synthetic SRB setting with  $K = 5$  and  $\sigma = 0.01$ . Figure 2 shows the evolution of the expected  
 284 values of the arms w.r.t. the number of pulls. In this setting, the optimal arm changes depending  
 285 on whether  $T \in [1, 185]$  or  $T \in (185, +\infty)$ . Thus, when the time budget is close to that value, the  
 286 problem is more challenging since the optimal and second-best arms expected rewards are close to  
 287 each other. For this reason, the BAI algorithms are less likely to provide a correct recommendation  
 288 than for time budgets for which the two expected rewards are well separated. We compare the  
 289 analyzed algorithms  $\mathfrak{A}$  in terms of empirical error  $\bar{e}_T(\mathfrak{A})$  (the smaller, the better), i.e., the empirical  
 290 counterpart of  $e_T(\mathfrak{A})$  averaged over 100 runs, considering time budgets  $T \in [100, 3200]$ .

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

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

## 310 8 Discussion and Conclusions

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

## 324 References

- 325 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 326 Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit  
327 problem. *Machine Learning*, 47(2):235–256, 2002.
- 328 Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. Pure exploration in multi-armed bandits problems.  
329 In *Proceedings of the Algorithmic Learning Theory (ALT)*, volume 5809, pages 23–37, 2009.
- 330 Cem Tekin and Mingyan Liu. Online learning of rested and restless bandits. *IEEE Transaction on*  
331 *Information Theory*, 58(8):5588–5611, 2012.
- 332 Alberto Maria Metelli, Francesco Trovò, Matteo Pirola, and Marcello Restelli. Stochastic rising  
333 bandits. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages  
334 15421–15457, 2022.
- 335 Chris Thornton, Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Auto-weka: Combined  
336 selection and hyperparameter optimization of classification algorithms. In *Proceedings of the ACM*  
337 *(SIGKDD)*, pages 847–855, 2013.
- 338 Lars Kotthoff, Chris Thornton, Holger H Hoos, Frank Hutter, and Kevin Leyton-Brown. Auto-weka  
339 2.0: Automatic model selection and hyperparameter optimization in weka. *Journal of Machine*  
340 *Learning Research*, 18:1–5, 2017.
- 341 Nick Erickson, Jonas Mueller, Alexander Shirkov, Hang Zhang, Pedro Larroy, Mu Li, and Alexander  
342 Smola. Autogluon-tabular: Robust and accurate automl for structured data. *arXiv preprint*  
343 *arXiv:2003.06505*, 2020.
- 344 Yang Li, Jiawei Jiang, Jinyang Gao, Yingxia Shao, Ce Zhang, and Bin Cui. Efficient automatic  
345 CASH via rising bandits. In *Proceedings of the Conference on Artificial Intelligence (AAAI)*, pages  
346 4763–4771, 2020.
- 347 Marc-André Zöller and Marco F Huber. Benchmark and survey of automated machine learning  
348 frameworks. *Journal of Artificial Intelligence Research*, 70:409–472, 2021.
- 349 Matthias Feurer, Aaron Klein, Katharina Eggenberger, Jost Tobias Springenberg, Manuel Blum, and  
350 Frank Hutter. Efficient and robust automated machine learning. In *Advances in Neural Information*  
351 *Processing Systems (NeurIPS)*, pages 2962–2970, 2015.
- 352 Quanming Yao, Mengshuo Wang, Yuqiang Chen, Wenyuan Dai, Yu-Feng Li, Wei-Wei Tu, Qiang  
353 Yang, and Yang Yu. Taking human out of learning applications: A survey on automated machine  
354 learning. *arXiv preprint arXiv:1810.13306*, 2018.
- 355 Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren. *Automated machine learning: methods, systems,*  
356 *challenges*. Springer Nature, 2019.
- 357 Marco Mussi, Davide Lombarda, Alberto Maria Metelli, Francesco Trovò, and Marcello Restelli.  
358 Arlo: A framework for automated reinforcement learning. *Expert Systems with Applications*, 224:  
359 119883, 2023.
- 360 Leonardo Cella, Massimiliano Pontil, and Claudio Gentile. Best model identification: A rested  
361 bandit formulation. In *Proceedings of the International Conference on Machine Learning (ICML)*,  
362 volume 139, pages 1362–1372, 2021.
- 363 Hoda Heidari, Michael J. Kearns, and Aaron Roth. Tight policy regret bounds for improving and  
364 decaying bandits. In *Proceeding of the International Joint Conference on Artificial Intelligence*  
365 *(AISTATS)*, pages 1562–1570, 2016.
- 366 Benny Lehmann, Daniel Lehmann, and Noam Nisan. Combinatorial auctions with decreasing  
367 marginal utilities. In *ACM Proceedings of the Conference on Electronic Commerce (EC)*, pages  
368 18–28, 2001.
- 369 Jean-Yves Audibert, Sébastien Bubeck, and Rémi Munos. Best arm identification in multi-armed  
370 bandits. In *Proceedings of the Conference on Learning Theory (COLT)*, pages 41–53, 2010.

- 371 Rianne De Heide, James Cheshire, Pierre Ménard, and Alexandra Carpentier. Bandits with many  
372 optimal arms. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 22457–  
373 22469, 2021.
- 374 Emilie Kaufmann, Olivier Cappé, and Aurélien Garivier. On the complexity of best-arm identification  
375 in multi-armed bandit models. *Journal of Machine Learning Research*, 17:1:1–1:42, 2016.
- 376 Yasin Abbasi-Yadkori, Peter L. Bartlett, Victor Gabillon, Alan Malek, and Michal Valko. Best of  
377 both worlds: Stochastic & adversarial best-arm identification. In *Proceedings of the Conference on  
378 Learning Theory (COLT)*, volume 75, pages 918–949, 2018.
- 379 Victor Gabillon, Mohammad Ghavamzadeh, and Alessandro Lazaric. Best arm identification: A  
380 unified approach to fixed budget and fixed confidence. In *Advances in Neural Information  
381 Processing Systems (NeurIPS)*, pages 3221–3229, 2012.
- 382 Aurélien Garivier and Emilie Kaufmann. Optimal best arm identification with fixed confidence. In  
383 *Proceedings of the Conference on Learning Theory (COLT)*, volume 49, pages 998–1027, 2016.
- 384 Alexandra Carpentier and Andrea Locatelli. Tight (lower) bounds for the fixed budget best arm iden-  
385 tification bandit problem. In *Proceedings of the 29th Conference on Learning Theory*, volume 49,  
386 pages 590–604, 2016.
- 387 Yonatan Mintz, Anil Aswani, Philip Kaminsky, Elena Flowers, and Yoshimi Fukuoka. Nonstationary  
388 bandits with habituation and recovery dynamics. *Operations Research*, 68(5):1493–1516, 2020.
- 389 Julien Seznec, Pierre Ménard, Alessandro Lazaric, and Michal Valko. A single algorithm for both  
390 restless and rested rotting bandits. In *Proceedings of the International Conference on Artificial  
391 Intelligence and Statistics (AISTATS)*, volume 108, pages 3784–3794, 2020.
- 392 Nir Levine, Koby Crammer, and Shie Mannor. Rotting bandits. In *Advances in Neural Information  
393 Processing Systems (NeurIPS)*, pages 3074–3083, 2017.
- 394 Julien Seznec, Andrea Locatelli, Alexandra Carpentier, Alessandro Lazaric, and Michal Valko.  
395 Rotting bandits are no harder than stochastic ones. In *Proceedings of the International Conference  
396 on Artificial Intelligence and Statistics*, volume 89, pages 2564–2572, 2019.