

---

# Bayesian bandits: balancing the exploration-exploitation tradeoff via double sampling

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       Reinforcement learning studies how to balance exploration and exploitation in  
2       real-world systems, optimizing interactions with the world while simultaneously  
3       learning how the world works. One general class of algorithms for such learning is  
4       the multi-armed bandit setting (in which sequential interactions are independent  
5       and identically distributed) and the related contextual bandit case, in which the  
6       distribution depends on different information or ‘context’ presented with each  
7       interaction. Thompson sampling, though introduced in the 1930s, has recently  
8       been shown to perform well and to enjoy provable optimality properties, while  
9       at the same time permitting generative, interpretable modeling. In a Bayesian  
10      setting, prior knowledge is incorporated and the computed posteriors naturally  
11      capture the full state of knowledge. In several application domains, for example in  
12      health and medicine, each interaction with the world can be expensive and invasive,  
13      whereas drawing samples from the model is relatively inexpensive. Exploiting this  
14      viewpoint, we develop a double-sampling technique driven by the uncertainty in  
15      the learning process. It extends and out-performs (in the sense of regret) Thompson  
16      sampling. We illustrate using multi-armed bandit examples both with and without  
17      context.

## 18   1 Introduction

19   In a plethora of problems in science and engineering, one needs to decide which action to take next  
20   based on partial information about the options available: a doctor must prescribe a medicine to a  
21   patient, a manager must allocate resources to competing projects, an ad serving algorithm must decide  
22   where to place ads, etc. In practice, the underlying properties of each choice are only partially known  
23   at the time of the decision, but one hopes that the understanding of the caveats involved will improve  
24   as time passes.

25   This set of problems has an illustrative gambling analogy, where a person facing a row of slot  
26   machines needs to devise its playing strategy (policy): which arms to play and in which order.  
27   Statisticians have studied this abstraction under the name of the multi-armed bandit problem for  
28   decades, e.g., in the seminal works by Robbins [1956, 1952]. Since then, it has played an important  
29   role in many fields across science and engineering.

30   Several algorithms have been proposed to overcome the exploration-exploitation tradeoff in such  
31   problems, mostly based on heuristics, on upper confidence bounds, or on the Gittins index. From  
32   the former, the  $\epsilon$ -greedy approach (randomly pick an arm with probability  $\epsilon$ , otherwise be greedy)  
33   has become very popular, due to its simplicity while nonetheless retaining often good performance  
34   (Auer et al. [2002]). In the latter case, Gittins [1979] formulated a method based on computing  
35   the optimal strategy for some types of bandits, where geometrically discounted future rewards are  
36   considered. There are several difficulties inherent to the exact computation of the Gittins index and

thus, approximations have been developed as well (Brezzi and Lai [2002]). These and other intrinsic challenges of the method have limited its applicability (Sutton and Barto [1998]). Lai and Robbins [1985] and Lai [1987] introduced another class of algorithms, based on upper confidence intervals of the expected reward of each arm, for which strong theoretical guarantees were proved. Nevertheless, these algorithms might be far from optimal in the presence of dependent and more general reward distributions (Scott [2010]).

More recently, the problem has re-emerged both from a practical (importance in e-commerce and web applications, e.g. Li et al. [2010]) and a theoretical (research on probability matching algorithms and their regret bounds, e.g. Agrawal and Goyal [2011] and Maillard et al. [2011]) point of view. Contributing to this revival was the observation that one of the oldest heuristics to address the exploration-exploitation trade-off, i.e., Thompson [1935, 1933] sampling, has been empirically proven to perform satisfactorily (see Chapelle and Li [2011] and Scott [2015] for details). Contemporaneously, theoretical study established several performance bounds, both for problems with and without context (Agrawal and Goyal [2012a,b], Korda et al. [2013], Russo and Roy [2014, 2016]).

In this work, we are interested in the randomized probability matching approach, as it connects to the Bayesian learning paradigm which readily facilitates modeling and algorithm development in both sequential and batch processing scenarios. Specifically, we establish how one can extract more information about the environment by casting it as a Bayesian sequential learning problem, so that better informed decisions can be made, leading to a lower regret than in the Thompson sampling approach. Our motivation is cases where sampling from the model is inexpensive relative to interactions with the world, which may be expensive or invasive or, as in the medical application domain, both. We propose a technique that is based on Monte Carlo sampling (to approximate otherwise unsolvable integrals) and a sampling-based arm-selection policy. The policy is driven by the uncertainty in the learning process, as it favors exploitation when certain about the properties of each arm, exploring otherwise.

We formally introduce the problem in Section 2, before providing all the details of our proposed method, double sampling, in Section 3. The performance of double sampling is compared to the Thompson sampling approach in Section 4, and we conclude with final remarks in Section 5.

## 2 Problem formulation

We mathematically formulate the multi-armed bandit problem as follows. Let  $a \in \{1, \dots, A\}$  indicate the arms of the bandit (possible actions to take) and  $f_a(y|\theta)$  the stochastic reward distribution of each arm. For every time instant, the observed reward  $y_t$  is independently drawn from the reward distribution corresponding to the played arm. We denote  $a_t$  as the arm played at time instant  $t$ ;  $a_{1:t} \equiv (a_1, \dots, a_t)$  refers to the sequence of arms played, and similarly  $y_{1:t} \equiv (y_1, \dots, y_t)$  refers to the sequence of observed rewards up to time  $t$ .

In the multi-armed bandit setting one must decide, based on observed rewards  $y_{1:t}$  and actions  $a_{1:t}$ , which arm to play next in order to maximize rewards. Due to the stochastic nature of the rewards, their expectation under the arm's distribution is the most common metric used. We denote each arm's expected reward as  $\mu_a(\theta) = \mathbb{E}_a\{y|\theta\}$ , which is parameterized by the arm-dependent parameters  $\theta$ . When the properties of the arms (i.e., their parameters) are known, one can readily determine the optimal selection policy, i.e.,

$$a^*(\theta) = \underset{a}{\operatorname{argmax}} \mu_a(\theta) . \quad (1)$$

However, the optimal solution for the multi-armed bandit is only computable in closed form in very few special cases (Bellman [1956], Gittins [1979]), and it fails to generalize to more realistic reward distributions and scenarios (Scott [2010]). The biggest challenge occurs when the parameters are unknown, as one might end up playing the wrong arm forever if incomplete learning occurs (Brezzi and Lai [2000]). Practitioners have often turned to heuristics to overcome these issues, and amongst them, the randomized probability matching, i.e., playing each arm in proportion to its probability of being optimal, is a particularly appealing one.

Given the parameters  $\theta$ , the expected reward of each arm is deterministic and, thus, one must pick the arm with the maximum expected reward

$$\Pr [a = a_{t+1}^* | a_{1:t}, y_{1:t}, \theta] = \Pr [a = a_{t+1}^* | \theta] = I_a(\theta) , \quad (2)$$

87 where we use the indicator function

$$I_a(\theta) = \begin{cases} 1, & \mu_a(\theta) = \max\{\mu_1(\theta), \dots, \mu_A(\theta)\} \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

88 Under random probability matching, the aim is to compute the probability of a given arm  $a$  being  
89 optimal for the next time instant,  $p_{a,t+1} \in [0, 1]$ , even with unknown parameters. Mathematically,

$$p_{a,t+1} \equiv \Pr[a = a_{t+1}^* | a_{1:t}, y_{1:t}] = \Pr[\mu_a(\theta) = \max\{\mu_1(\theta), \dots, \mu_A(\theta)\} | a_{1:t}, y_{1:t}] \quad (4)$$

90 In this case, there is uncertainty about the unknown properties of the arms, as Eqn. 4 is parameterized  
91 by  $\theta$ , while Eqn. 3 is not. In order to compute a solution to this problem, recasting it as a Bayesian  
92 learning problem is of great help, as it allows for direct connection with the randomized probability  
93 matching technique (Scott [2010]). The randomized probability matching approach has shown to be  
94 easy to implement, efficient and broadly applicable.

### 95 3 Proposed method: double sampling

96 The multi-armed bandit problem consists of two separate but intertwined tasks: (1) learning about the  
97 properties of the arms, and (2) deciding what arm to play next. The problem is sequential in nature,  
98 as one makes a decision on which arm to play and learns from the observed reward, one observation  
99 at a time. We cast the multi-armed bandit problem as a sequential Bayesian learning task (Bernardo  
100 and Smith [2009]). By doing so, we capture the full state of knowledge about the world at every  
101 time instant. We incorporate any available prior information to the learning process, and update our  
102 knowledge about the unknown parameters  $\theta$ , as we sequentially play arms and observe rewards. This  
103 learning can be done both sequentially or in batches, as Bayesian posterior updates are computable  
104 for both cases. However, the solution to the probability matching equation is analytically intractable,  
105 so we approximate it via Monte Carlo sampling. For balancing the exploration-exploitation tradeoff,  
106 we propose a sampling-based probability matching technique. The proposed arm-selection policy  
107 is a function of the uncertainty in the learning process. The intuition is that we exploit only when  
108 certain about the properties of the arms, while we keep exploring otherwise. We elaborate on the  
109 foundations of the proposed method in the following sections, before presenting the algorithm in 1.

#### 110 3.1 The Bayesian multi-armed bandit

111 We are interested in computing, after playing arms  $a_{1:t}$  and observing rewards  $y_{1:t}$ , the probability of  
112 each arm  $a$  being optimal for the next time instant, i.e.,

$$p_{a,t+1} \equiv \Pr[a = a_{t+1}^* | a_{1:t}, y_{1:t}] = \Pr[\mu_a(\theta) = \max\{\mu_1(\theta), \dots, \mu_A(\theta)\} | a_{1:t}, y_{1:t}] \quad (5)$$

113 In practice, one needs to account for the lack of knowledge of each arm's properties, i.e., the unknown  
114 parameters  $\theta$ . We can do so by following the Bayesian methodology, where the parameters are  
115 considered to be another set of random variables. The uncertainty over the parameters can be  
116 accounted for by marginalizing over their probability distribution. Specifically, we marginalize over  
117 the posterior of the parameters after observing rewards and actions up to  $t$

$$p_{a,t+1} \equiv \Pr[a = a_{t+1}^* | a_{1:t}, y_{1:t}] = f(a = a_{t+1}^* | a_{1:t}, y_{1:t}) = \int f(a | a_{1:t}, y_{1:t}, \theta) f(\theta | a_{1:t}, y_{1:t}) d\theta \quad (6)$$

118 Given a prior for the parameters  $f(\theta)$  and the per-arm reward distribution  $f_a(y|\theta)$ , one can compute  
119 the posterior of each arm's parameters by

$$f(\theta | a_{1:t}, y_{1:t}) \propto f_{a_t}(y_t | \theta) f(\theta | a_{1:t-1}, y_{1:t-1}) = \left[ \prod_{\tau=1}^t f_{a_\tau}(y_\tau | \theta) \right] f(\theta) \quad (7)$$

120 This posterior provides information (with uncertainty) about the state of the arm. In this work, we  
121 focus on the parameter posterior updates for two of the most studied (the Bernoulli and the contextual  
122 linear Gaussian) bandits. Note that the updates (provided later in subsection 4.1 and subsection 4.2)

can be written in both sequential and batch forms. This flexibility is of great help in many practical scenarios, as one can learn from historic data, as well as process data as it comes.

Even if analytical expressions for the parameter posteriors are available for many models of interest, computing the probability of any given arm being optimal is analytically intractable, due to the nonlinearities induced by the indicator function

$$p_{a,t+1} \equiv f(a = a_{t+1}^* | a_{1:t}, y_{1:t}) = \int f(a | a_{1:t}, y_{1:t}, \theta) f(\theta | a_{1:t}, y_{1:t}) d\theta = \int I_a(\theta) f(\theta | a_{1:t}, y_{1:t}) d\theta . \quad (8)$$

### 3.2 Monte-Carlo integration

We harness the power of Monte Carlo sampling to compute the otherwise analytically intractable integral for  $p_{a,t+1} \in [0, 1]$  in Eqn. 8. We obtain its Monte Carlo approximation as follows:

1. Draw  $M$  parameter samples from the updated posterior distribution

$$\theta^{(m)} \sim f(\theta | a_{1:t}, y_{1:t}), \quad m = \{1, \dots, M\} . \quad (9)$$

2. For each parameter sample  $\theta^{(m)}$ , compute the expected reward and determine the best arm

$$a_{t+1}^*(\theta^{(m)}) = \underset{a}{\operatorname{argmax}} \mu_a(\theta^{(m)}) . \quad (10)$$

3. Define the random measure approximation to the probability distribution in Eqn. ?? as

$$f_M(a = a_{t+1}^* | a_{1:t}, y_{1:t}) = \frac{1}{M} \sum_{m=1}^M \delta \left( a - a_{t+1}^*(\theta^{(m)}) \right) , \quad (11)$$

where  $\delta(\cdot)$  denotes the Dirac delta function.

4. Estimate the first- and second-order sufficient statistics of the distribution

$$\begin{cases} \hat{p}_{a,t+1} = \frac{1}{M} \sum_{m=1}^M I_a(\theta^{(m)}) , \\ \hat{\sigma}_{a,t+1}^2 = \frac{1}{M} \sum_{m=1}^M (I_a(\theta^{(m)}) - \hat{p}_{a,t+1})^2 . \end{cases} \quad (12)$$

5. Estimate which is the optimal arm and with what probability

$$\begin{cases} \hat{a}_{t+1}^* = \operatorname{argmax}_a \hat{p}_{a,t+1} , \\ \hat{p}_{a,t+1}^* = \max_a \hat{p}_{a,t+1} . \end{cases} \quad (13)$$

### 3.3 Sampling-based policy

Given the available information at any given time  $t$ , one needs to decide which arm to play next. A randomized probability matching technique would pick the next arm  $a$  with probability  $p_{a,t+1}$ . On the contrary, a greedy approach would choose the arm with the highest probability of being optimal (i.e.,  $p_{a,t+1}^*$ ). We present an alternative sampling-based probability matching arm-selection policy that automatically finds a balance between these two cases. We rely on the Monte Carlo approximation to Eqn. 8, and leverage the estimated sufficient statistics in Eqn. 12 to balance the exploration-exploitation tradeoff.

We draw candidate arm samples from the random measure in Eqn. 11, and automatically adjust the probability matching technique according to the accuracy of this approximation. We account for the uncertainty of the learning process in the arm-selection policy, dynamically balancing exploration and exploitation. The number of candidate arm samples drawn is instrumental for our sampling based policy. We automatically adjust its value according to the uncertainty on the optimality of each arm (i.e.,  $\hat{\sigma}_{a,t+1}^2$ ).

The number of candidate arm samples to draw is inversely proportional to the probability of not picking the optimal arm. We denote this probability as  $p_{FA}$ , which is computed for each arm as

$$p_{FA}^{(a)} = \Pr(p_{a,t+1} > p_{a,t+1}^*) = 1 - F_{p_{a,t+1}}(p_{a,t+1}^*) , \quad (14)$$

where  $p_{a,t+1}^* = \max_a p_{a,t+1}$ . Since we can not exactly evaluate  $p_{a,t+1}$ , we resort to our Monte Carlo estimates in Eqn. 12. The true cumulative density function  $F_{p_{a,t+1}}(\cdot)$  is analytically intractable as well, but we approximate it with a truncated Gaussian

$$F_{p_{a,t+1}}(p_{a,t+1}^*) \approx \Phi_{[0,1]} \left( \frac{\hat{p}_{a,t+1} - \hat{p}_{a,t+1}^*}{\hat{\sigma}_{a,t+1}} \right). \quad (15)$$

The proposed sampling policy proceeds as follows:

1. Determine  $N_{t+1}$ , the number of candidate arm samples to draw

$$N_{t+1} \propto \log \left( \frac{1}{p_{FA}} \right), \quad p_{FA} = \frac{1}{K-1} \sum_{a \neq \hat{a}_{t+1}^*} p_{FA}^{(a)}, \quad p_{FA}^{(a)} \approx 1 - \Phi_{[0,1]} \left( \frac{\hat{p}_{a,t+1} - \hat{p}_{a,t+1}^*}{\hat{\sigma}_{a,t+1}} \right). \quad (16)$$

2. Draw  $N_{t+1}$  candidate arm samples from the random measure in Eqn. 11

$$\hat{a}_{t+1}^{(n)} \sim \text{Cat}(\hat{p}_{a,t+1}), \quad n = 1, \dots, N_{t+1}. \quad (17)$$

3. Pick the most probable optimal arm, given drawn candidate arm samples  $\hat{a}_{t+1}^{(n)}$

$$a_{t+1} = \text{Mode}(\hat{a}_{t+1}^{(n)}), \quad n = 1, \dots, N_{t+1}. \quad (18)$$

By allowing for  $N_{t+1}$  to be adjusted based upon the uncertainty of the learning process, we balance the exploration-exploitation tradeoff. The proposed sampling policy reduces to a probabilistic matching regime when uncertain about the arms, (i.e.,  $N_t \approx 1$ ), but favors exploitation ( $N_t \gg 1$ ) when the probability of picking a suboptimal arm is low. In other words, double sampling exploits only when confident about the learned probabilities ( $\hat{\sigma}_{a,t+1} \rightarrow 0, N \gg 1$ ), and picks the arm with the highest probability  $\hat{p}_{a,t+1}$ . However, for  $N_t \approx 1$ , a randomized probability matching is in play. Note that Thompson sampling is a special case of double sampling, when  $N_t = 1, \forall t$ .

## 4 Evaluation: Thompson sampling vs. double sampling

We now provide empirical evidence of our double sampling's lower regret relative to Thompson sampling. We consider both discrete and continuous multi-armed bandits, and the contextual setting as well. We compare the performance of our algorithm, as presented in 1, to that of Thompson sampling, for both the Bernoulli and the contextual linear Gaussian bandits. As is standard in the literature, we measure performance in the cumulative regret sense, i.e.,

$$R_t = \sum_{\tau=0}^t \mathbb{E} \{ (y_\tau^* - y_\tau) \} = \sum_{\tau=0}^t \mu_\tau^* - \bar{y}_\tau, \quad (27)$$

where for each time instant  $t$ ,  $\mu_t^*$  denotes the expected reward of the optimal arm and  $\bar{y}_t$  the empirical mean of the observed rewards.

### 4.1 Bernoulli bandits

Bernoulli bandits are well suited for applications with binary rewards (i.e., success or failure of an action). The rewards of each arm are modeled as independent draws from a Bernoulli distribution with success probabilities  $\theta_a$ , i.e.,

$$f_a(y|\theta) = \theta_a^y (1 - \theta_a)^{(1-y)}. \quad (28)$$

For this reward distribution, the posterior parameter update can be computed using the conjugate prior distribution  $f(\theta_a|\alpha_{a,0}, \beta_{a,0}) = \text{Beta}(\theta_a|\alpha_{a,0}, \beta_{a,0})$ . After observing actions  $a_{1:t}$  and rewards  $y_{1:t}$ , the posterior parameter distribution follows an updated Beta distribution

$$f(\theta_a|a_{1:t}, y_{1:t}, \alpha_{a,0}, \beta_{a,0}) = f(\theta_a|\alpha_{a,t}, \beta_{a,t}) = \text{Beta}(\theta_a|\alpha_{a,t}, \beta_{a,t}),$$

with  $\begin{cases} \text{sequential updates} & \begin{cases} \alpha_{a,t} = \alpha_{a,t-1} + y_t \cdot \mathbb{1}[a_t = a] \\ \beta_{a,t} = \beta_{a,t-1} + (1 - y_t) \cdot \mathbb{1}[a_t = a] \end{cases} \\ \text{batch updates} & \begin{cases} \alpha_{a,t} = \alpha_{a,0} + \sum_{t|a_t=a} y_t \\ \beta_{a,t} = \beta_{a,0} + \sum_{t|a_t=a} (1 - y_t). \end{cases} \end{cases} \quad (29)$

---

**Algorithm 1** Bayesian Double Sampling algorithm

---

**Require:** Number of arms  $A$

**Require:** Prior over model parameters  $f(\theta)$

**Require:** Per-arm reward distributions  $f_a(y|\theta)$

**Require:** Horizon  $T$

$D = \emptyset$

**for**  $t = 1, \dots, T$  **do**

Draw  $M$  parameter samples from the updated posterior distribution

$$\theta_{t+1}^{(m)} \sim f(\theta|a_{1:t}, y_{1:t}), \quad m = \{1, \dots, M\}. \quad (19)$$

If applicable, receive context  $x_{t+1}$

**for**  $a = 1, \dots, A$  **do**

Compute the expected reward for each parameter sample  $\theta_{t+1}^{(m)}$

$$\mu_{a,t+1}(\theta_{t+1}^{(m)}) = \mu_a(x_{t+1}, \theta_{t+1}^{(m)}) = \mathbb{E}_a\{y|x_{t+1}, \theta_{t+1}^{(m)}\} \quad (20)$$

Compute the first- and second-order sufficient statistics of each arm being optimal

$$\begin{cases} \hat{p}_{a,t+1} = \frac{1}{M} \sum_{m=1}^M I_a(\theta_{t+1}^{(m)}), \\ \hat{\sigma}_{a,t+1}^2 = \frac{1}{M} \sum_{m=1}^M \left( I_a(\theta_{t+1}^{(m)}) - \hat{p}_{a,t+1} \right)^2, \end{cases} \quad (21)$$

where  $I_a(\theta_{t+1}^{(m)}) = \begin{cases} 1, & \mu_a(\theta_{t+1}^{(m)}) = \max\{\mu_1(\theta_{t+1}^{(m)}), \dots, \mu_A(\theta_{t+1}^{(m)})\}, \\ 0, & \text{otherwise} \end{cases}$ .

**end for**

Estimate the optimal arm and its probability of being optimal

$$\begin{cases} \hat{a}_{t+1}^* = \operatorname{argmax}_a \hat{p}_{a,t+1}, \\ \hat{p}_{a,t+1}^* = \max_a \hat{p}_{a,t+1}. \end{cases} \quad (22)$$

**for**  $a = 1, \dots, A$  **do**

Compute the estimated probability of each arm not being the optimal arm

$$\begin{aligned} \hat{p}_{FA}^{(a)} &= Pr(\hat{p}_{a,t+1} > \hat{p}_{a,t+1}^*) = 1 - F_{\hat{p}_{a,t+1}}(\hat{p}_{a,t+1}^*), \\ F_{\hat{p}_{a,t+1}}(\hat{p}_{a,t+1}^*) &\approx \Phi_{[0,1]} \left( \frac{\hat{p}_{a,t+1} - \hat{p}_{a,t+1}^*}{\hat{\sigma}_{a,t+1}} \right). \end{aligned} \quad (23)$$

**end for**

Compute  $N_{t+1}$ , the number of candidate arm samples to draw

$$N_{t+1} \propto \log \left( \frac{1}{\hat{p}_{FA}} \right), \quad \hat{p}_{FA} = \frac{1}{A-1} \sum_{a \neq \hat{a}_{t+1}^*} \hat{p}_{FA}^{(a)}. \quad (24)$$

Draw  $N_{t+1}$  candidate arm samples

$$\hat{a}_{t+1}^{(n)} \sim \operatorname{Cat}(\hat{p}_{a,t+1}), \quad n = 1, \dots, N_{t+1}. \quad (25)$$

Play arm

$$a_{t+1} = \operatorname{Mode}(\hat{a}_{t+1}^{(n)}), \quad n = 1, \dots, N_{t+1}. \quad (26)$$

Observe reward  $y_{t+1}$

$D = D \cup \{x_{t+1}, a_{t+1}, y_{t+1}\}$

**end for**

---

182 We dissect the key aspects of double sampling in a particular realization of a three-armed Bernoulli  
 183 bandit with parameters  $\theta = (0.4 \ 0.7 \ 0.8)$ . The sequential Bayesian learning process is illustrated in  
 184 Fig. 1a, where we show the evolution of the probability of each arm being optimal (i.e., the Monte  
 185 Carlo approximation in Eqn. 11). For all results to follow, Monte Carlo integration is attained with  
 186  $M = 1000$  samples, as larger  $M$ s do not significantly improve regret performance. In Fig. 1b, we  
 187 illustrate how double sampling is *automatically* adjusted according to the learning process uncertainty,  
 188 via the number of arm samples to draw (i.e.,  $N_{t+1}$  in Eqn. 16).

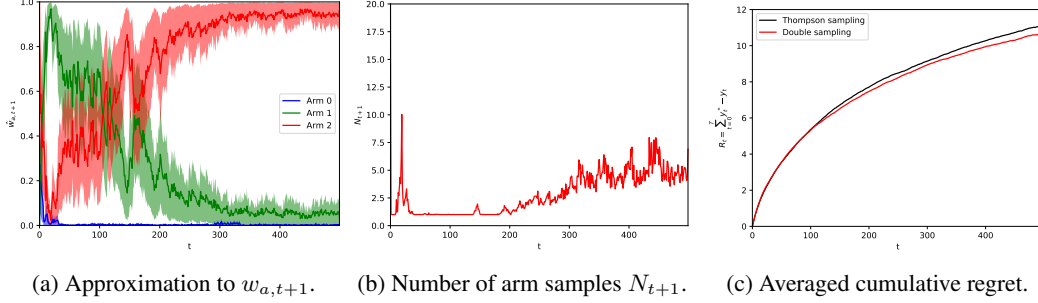


Figure 1: Empirical results with double sampling.

189 For the case shown in Fig. 1, arm 0 (with  $\theta_0 = 0.4$ ) is quickly discarded as being optimal, while the  
 190 decision over which of the other two arms is better requires further learning. For some time ( $t < 200$ ),  
 191 there is high uncertainty about the properties of both arms (high variance in Fig. 1a). Thus, double  
 192 sampling favors exploration ( $N_{t+1} \approx 1$  in Fig. 1b), until the uncertainty about which arm is best  
 193 is reduced. Once the algorithm becomes more certain about the better reward properties of arm 2  
 194 ( $t > 200$ ), the sampling approach gradually favors a greedier policy ( $N_{t+1} > 1$ ). As a result, the  
 195 cumulative regret of our proposed technique is reduced when compared to the Thompson sampling  
 196 approach (see averaged cumulative regrets in Fig. 1c). All the averaged results presented in this paper  
 197 are computed over 5000 realizations of the same set of parameters.

198 In Fig. 2, we show two antagonistic examples of the performance of double sampling. Our algorithm  
 199 provides a clear benefit over Thompson sampling in Fig. 2a, though it performs only slightly better  
 200 than Thompson sampling in Fig. 2b. Such a cumulative regret difference is explained by the unknown  
 201 nature of the bandit’s arms. When the properties of the arms are very similar to each other, our  
 202 algorithm resorts to a Thompson sampling-like policy ( $N_{t+1} \approx 1$ ), yielding near-equivalent sampling  
 203 (and thus regret). However, when the learning process is certain about the properties of the arms,  
 204 double sampling can considerably reduce the method’s cumulative regret. Mathematically, the  
 205 difference in arm properties can be captured by the divergence between their respective reward  
 206 distributions. By computing the minimum Kullback-Leibler (KL) divergence between arms, we get  
 207 an insight onto how “difficult” a multi-armed bandit problem is.

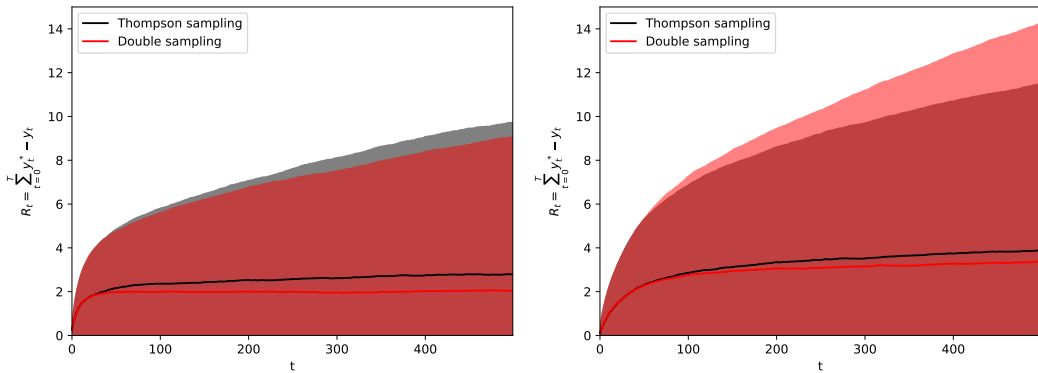


Figure 2: Averaged cumulative regret comparison (standard deviation shown as shaded region).

208 We now proceed to show the relative difference between the averaged expected rewards of our  
 209 algorithm and Thompson sampling, over many parameterizations of different Bernoulli multi-armed  
 210 bandits. We evaluate

$$\Delta_t = \frac{R_t^{(DS)}}{R_t^{(TS)}} - 1, \quad (30)$$

211 where  $R_t^{(DS)}$  denotes the regret of the proposed double sampling approach at time  $t$ , and  $R_t^{(TS)}$ , that  
 212 of Thompson sampling.

213 We show in Fig. 3 average results over 5000 realizations of Bernoulli bandits with  $K=2$  and  $K=3$   
 214 arms, for all per-arm parameter  $\theta_a$  combinations in the range  $[0, 1]$  with grid size 0.05. Note that  
 215 the KL metric may map many parameter combinations to the same point in Fig. 3. When the best  
 216 arms are very similar (small KL), our performance is comparable to that of the Thompson sampling  
 217 technique. On the contrary, double sampling performs considerably better when it is certain about the  
 218 learned arm parameters (with regret reductions around 40% for many cases).

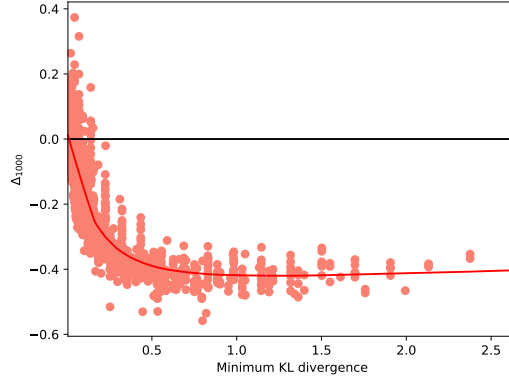


Figure 3: Relative average expected reward differences at  $t = 1000$ .

## 219 4.2 Contextual linear Gaussian bandits

220 Another set of well studied bandits are those with continuous reward functions. It is interesting to  
 221 consider the contextual case as well, where the reward distribution of each arm is dependent on a  
 222 time-varying  $d$ -dimensional context vector. The contextual linear Gaussian bandit model is suited for  
 223 these scenarios, where the expected reward of each arm is linearly dependent on the context  $x \in \mathbb{R}^d$ ,  
 224 and the idiosyncratic parameters of the bandit  $\theta \equiv \{w, \sigma\}$ ,

$$f_a(y|x, \theta) = \mathcal{N}(y|x^\top w_a, \sigma_a^2). \quad (31)$$

225 For this reward distribution, the posterior can be computed with the Normal Inverse Gamma conjugate  
 226 prior distribution

$$\begin{aligned} f(w_a, \sigma_a^2 | u_{a,0}, V_{a,0}, \alpha_{a,0}, \beta_{a,0}) &= \text{NIG}(w_a, \sigma_a^2 | u_{a,0}, V_{a,0}, \alpha_{a,0}, \beta_{a,0}) \\ &= \mathcal{N}(w_a | u_{a,0}, \sigma_a^2 V_{a,0}) \cdot \Gamma^{-1}(\sigma_a^2 | \alpha_{a,0}, \beta_{a,0}), \end{aligned} \quad (32)$$



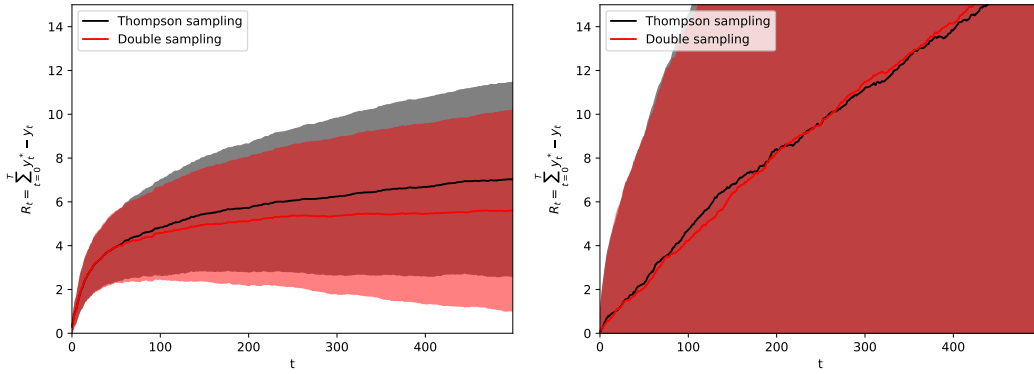
that yields, given previous actions  $a_{1:t}$ , contexts  $x_{1:t}$  and rewards  $y_{1:t}$ , the following posterior

$$f(w_a, \sigma_a^2 | a_{1:t}, y_{1:t}, u_{a,0}, V_{a,0}, \alpha_{a,0}, \beta_{a,0}) = \text{NIG}(w_a, \sigma_a^2 | u_{a,t}, V_{a,t}, \alpha_{a,t}, \beta_{a,t}),$$

$$\text{with } \begin{cases} \text{sequential updates} & \begin{cases} V_{a,t}^{-1} = V_{a,t-1}^{-1} + x_t x_t^\top \cdot \mathbb{1}[a_t = a] \\ u_{a,t} = V_{a,t} (V_{a,t-1}^{-1} u_{a,t-1} + x_t y_t \cdot \mathbb{1}[a_t = a]) \\ \alpha_{a,t} = \alpha_{a,t-1} + \frac{\mathbb{1}[a_t = a]}{2} \\ \beta_{a,t} = \beta_{a,t-1} + \frac{\mathbb{1}[a_t = a](y_{t_a} - x_t^\top \theta_{a,t-1})^2}{2(1 + x_t^\top \Sigma_{a,t-1} x_t)} \end{cases} \\ \text{batch updates} & \begin{cases} V_{a,t}^{-1} = V_{a,0}^{-1} + x_{1:t|t_a} x_{1:t|t_a}^\top \\ u_{a,t} = V_{a,t} (V_{a,0}^{-1} u_{a,0} + x_{1:t|t_a} y_{1:t|t_a}) \\ \alpha_{a,t} = \alpha_{a,0} + \frac{|t_a|}{2} \\ \beta_{a,t} = \beta_{a,0} + \frac{(y_{1:t|t_a} y_{1:t|t_a}^\top + u_{a,0}^\top V_{a,0}^{-1} u_{a,0} - u_{a,t}^\top V_{a,t}^{-1} u_{a,t})}{2} \end{cases} \end{cases} \quad (33)$$

where  $t_a = \{t | a_t = a\}$  indicates the set of time instances when arm  $a$  is played.

We provide results in Fig. 4 for the two-armed contextual Gaussian bandit with uniformly distributed random uncorrelated context, i.e.,  $x_{i,t} \sim \mathcal{U}(0, 1)$ ,  $i \in \{1, \dots, d\}$ ,  $t \in \mathbb{N}$ . When the arms' reward distributions are comparable, our method performs similar to the Thompson sampling approach (see Fig. 4b). However, when the difference between arms is easier to learn (as in Fig. 4a), double sampling attains reduced cumulative regret.



(a) Averaged cumulative regret with  $A = 2$ ,  $w_0 = (0.4, 0.4)^\top$ ,  $w_1 = (0.8, 0.8)^\top$ ,  $\sigma_0 = \sigma_1 = 0.2$ . (b) Averaged cumulative regret with  $A = 2$ ,  $w_0 = (0.4, 0.4)^\top$ ,  $w_1 = (0.5, 0.5)^\top$ ,  $\sigma_0 = \sigma_1 = 1$ .

Figure 4: Averaged cumulative regret comparison (standard deviation shown as shaded region).

We elaborate on the power of the proposed double sampling technique by evaluating the relative difference between the averaged expected rewards of our algorithm and Thompson sampling, with respect to the Kullback-Leibler divergence. Fig. 5 contains average results over 5000 realizations of the two-dimensional contextual linear Gaussian bandit with two arms, for per-dimension parameter combinations in the range  $[-1, 1]$  with step size 0.1. We observe that double sampling can provide significant regret reductions for all the parameterizations of the two-armed contextual linear Gaussian bandit.

## 5 Conclusion

We have presented a new sampling-based probability matching technique for the multi-armed bandit setting. We formulated the problem as a Bayesian sequential learning one, and leveraged random sampling to overcome two of its main challenges: approximating the analytically unsolvable integrals, and automatically balancing the exploration-exploitation tradeoff. We empirically show that additional sampling from the model can provide improved regrets, which is in many application domains inexpensive in comparison with interacting with the world. Encouraged by these findings, we aim at extending this technique to other reward distributions and implementing it with datasets from domain application in which actions are expensive relative to samples drawn from the current model.

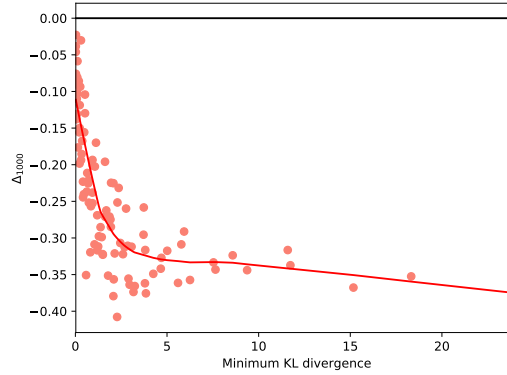


Figure 5: Relative average expected reward differences at  $t = 1000$ .

## References

- Shipra Agrawal and Navin Goyal. Analysis of Thompson Sampling for the multi-armed bandit problem. *CoRR*, abs/1111.1797, 2011. URL <http://arxiv.org/abs/1111.1797>.
- Shipra Agrawal and Navin Goyal. Thompson Sampling for Contextual Bandits with Linear Payoffs. *CoRR*, abs/1209.3352, 2012a. URL <http://arxiv.org/abs/1209.3352>.
- Shipra Agrawal and Navin Goyal. Further Optimal Regret Bounds for Thompson Sampling. *CoRR*, abs/1209.3353, 2012b. URL <http://arxiv.org/abs/1209.3353>.
- Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time Analysis of the Multiarmed Bandit Problem. *Machine Learning*, 47(2-3):235–256, May 2002. ISSN 0885-6125. doi: 10.1023/A:1013689704352. URL <http://dx.doi.org/10.1023/A:1013689704352>.
- Richard Bellman. A Problem in the Sequential Design of Experiments. *Sankhya: The Indian Journal of Statistics (1933 - 1960)*, 16(3/4):221–229, 1956. URL <http://www.jstor.org/stable/25048278>.
- José M. Bernardo and Adrian F.M. Smith. *Bayesian Theory*. Wiley Series in Probability and Statistics. Wiley, 2009. ISBN 9780470317716. doi: 10.1002/9780470316870. URL <http://onlinelibrary.wiley.com/book/10.1002/9780470316870>.
- Monica Brezzi and Tze Leung Lai. Incomplete Learning from Endogenous Data in Dynamic Allocation. *Econometrica*, 68(6):1511–1516, 2000. ISSN 1468 0262. doi: 10.1111/1468-0262.00170. URL <http://dx.doi.org/10.1111/1468-0262.00170>.
- Monica Brezzi and Tze Leung Lai. Optimal learning and experimentation in bandit problems. *Journal of Economic Dynamics and Control*, 27(1):87 – 108, 2002. ISSN 0165-1889. doi: [https://doi.org/10.1016/S0165-1889\(01\)00028-8](https://doi.org/10.1016/S0165-1889(01)00028-8). URL <http://www.sciencedirect.com/science/article/pii/S0165188901000288>.
- Olivier Chapelle and Lihong Li. An Empirical Evaluation of Thompson Sampling. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 2249–2257. Curran Associates, Inc., 2011. URL <http://papers.nips.cc/paper/4321-an-empirical-evaluation-of-thompson-sampling.pdf>.
- J. C. Gittins. Bandit Processes and Dynamic Allocation Indices. *Journal of the Royal Statistical Society. Series B (Methodological)*, 41(2):148–177, 1979. ISSN 00359246. URL <http://www.jstor.org/stable/2985029>.
- Nathaniel Korda, Emilie Kaufmann, and Rémi Munos. Thompson Sampling for 1-Dimensional Exponential Family Bandits. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 1448–1456. Curran Associates, Inc., 2013. URL <http://papers.nips.cc/paper/5110-thompson-sampling-for-1-dimensional-exponential-family-bandits.pdf>.

- 286 Tze Leung Lai. Adaptive Treatment Allocation and the Multi-Armed Bandit Problem. *The Annals*  
287 *of Statistics*, 15(3):1091–1114, 1987. ISSN 00905364. URL [http://www.jstor.org/stable/](http://www.jstor.org/stable/2241818)  
288 2241818.
- 289 Tze Leung Lai and Herbert Robbins. Asymptotically Efficient Adaptive Allocation Rules. *Advances*  
290 *in Applied Mathematics*, 6(1):4–22, mar 1985. ISSN 0196-8858. doi: 10.1016/0196-8858(85)  
291 90002-8. URL [http://dx.doi.org/10.1016/0196-8858\(85\)90002-8](http://dx.doi.org/10.1016/0196-8858(85)90002-8).
- 292 Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. A Contextual-Bandit Approach to  
293 Personalized News Article Recommendation. *CoRR*, abs/1003.0146, 2010. URL [http://arxiv.](http://arxiv.org/abs/1003.0146)  
294 [org/abs/1003.0146](http://arxiv.org/abs/1003.0146).
- 295 Odalric-Ambrym Maillard, Rémi Munos, and Gilles Stoltz. Finite-Time Analysis of Multi-armed  
296 Bandits Problems with Kullback-Leibler Divergences. In *Conference On Learning Theory*, 2011.  
297 URL <http://hal.inria.fr/inria-00574987/fr/>.
- 298 Herbert Robbins. Some aspects of the sequential design of experiments. *Bulletin of*  
299 *the American Mathematical Society*, (58):527–535, 1952. doi: [https://doi.org/10.1090/](https://doi.org/10.1090/S0002-9904-1952-09620-8)  
300 [S0002-9904-1952-09620-8](https://doi.org/10.1090/S0002-9904-1952-09620-8).
- 301 Herbert Robbins. A sequential decision procedure with a finite memory. *Proceedings of the National*  
302 *Academy of Science*, (42):920 – 923, 1956.
- 303 Daniel Russo and Benjamin Van Roy. Learning to optimize via posterior sampling. *Mathematics of*  
304 *Operations Research*, 39(4):1221–1243, 2014. doi: [http://pubsonline.informs.org/doi/abs/10.1287/](http://pubsonline.informs.org/doi/abs/10.1287/moor.2014.0650)  
305 [moor.2014.0650](http://pubsonline.informs.org/doi/abs/10.1287/moor.2014.0650).
- 306 Daniel Russo and Benjamin Van Roy. An information-theoretic analysis of thompson sampling. *The*  
307 *Journal of Machine Learning Research*, 17(1):2442–2471, 2016. URL [http://www.jmlr.org/](http://www.jmlr.org/papers/volume17/14-087/14-087.pdf)  
308 [papers/volume17/14-087/14-087.pdf](http://www.jmlr.org/papers/volume17/14-087/14-087.pdf).
- 309 Steven L. Scott. A modern Bayesian look at the multi-armed bandit. *Applied Stochastic Models in*  
310 *Business and Industry*, 26(6):639–658, 2010. ISSN 1526-4025. doi: 10.1002/asmb.874. URL  
311 <http://dx.doi.org/10.1002/asmb.874>.
- 312 Steven L. Scott. Multi-armed bandit experiments in the online service economy. *Applied Stochastic*  
313 *Models in Business and Industry*, 31:37–49, 2015. URL [http://onlinelibrary.wiley.com/](http://onlinelibrary.wiley.com/doi/10.1002/asmb.2104/abstract)  
314 [doi/10.1002/asmb.2104/abstract](http://onlinelibrary.wiley.com/doi/10.1002/asmb.2104/abstract). Special issue on actual impact and future perspectives on  
315 stochastic modelling in business and industry.
- 316 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press: Cam-  
317 bridge, MA, 1998. URL <https://mitpress.mit.edu/books/reinforcement-learning>.
- 318 William R. Thompson. On the Likelihood that One Unknown Probability Exceeds Another in View  
319 of the Evidence of Two Samples. *Biometrika*, 25(3/4):285–294, 1933. ISSN 00063444. URL  
320 <http://www.jstor.org/stable/2332286>.
- 321 William R. Thompson. On the Theory of Apportionment. *American Journal of Mathematics*, 57(2):  
322 450–456, 1935. ISSN 00029327, 10806377. URL <http://www.jstor.org/stable/2371219>.