

---

# Exploiting Structure of Uncertainty for Efficient Matroid Semi-Bandits

---

Pierre Perrault<sup>1,2</sup> Vianney Perchet<sup>2,3</sup> Michal Valko<sup>1,4</sup>

## Abstract

We improve the efficiency of algorithms for stochastic *combinatorial semi-bandits*. In most interesting problems, state-of-the-art algorithms take advantage of structural properties of rewards, such as *independence*. However, while being optimal in terms of asymptotic regret, these algorithms are inefficient. In our paper, we first reduce their implementation to a specific *submodular maximization*. Then, in case of *matroid* constraints, we design adapted approximation routines, thereby providing the first efficient algorithms that rely on reward structure to improve regret bound. In particular, we improve the state-of-the-art efficient gap-free regret bound by a factor  $\sqrt{m}/\log m$ , where  $m$  is the maximum action size. Finally, we show how our improvement translates to more general *budgeted combinatorial semi-bandits*.

## 1. Introduction

Stochastic bandits model sequential decision-making in which an *agent* selects an arm (a decision) at each round and observes a realization of the corresponding unknown reward distribution. The goal is to maximize the expected cumulative reward, or equivalently, to minimize the *expected regret*, defined as the difference between the expected cumulative reward achieved by an oracle algorithm always selecting the optimal arm and that achieved by the agent. To accomplish this objective, the agent must trade-off between *exploration* (gaining information about reward distributions) and *exploitation* (using greedily the information collected so far) as it was already discovered by Robbins (1952). Bandits have been applied to many fields such as mechanism design (Mohri & Munoz, 2014), search advertising (Tran-Thanh et al., 2014), and personalized recommendation (Li et al., 2010). We improve the computational efficiency (i.e., the

time and space complexity) of their *combinatorial* generalization, in which the agent selects at each round a *subset* of arms, that we refer to as an *action* in the rest of the paper (Cesa-Bianchi & Lugosi, 2012; Gai et al., 2012; Audibert et al., 2014; Chen et al., 2014).

Different kinds of feedback provided by the environment are possible. First, with *bandit feedback* (also called full bandit or opaque feedback), the agent only observes the *total* reward associated to the selected action. Second, with *semi-bandit feedback*, the agent observes the *partial* reward of each base arm in the selected action. Finally, with *full information feedback*, the agent observes the partial reward of all arms. We give results for semi-bandit feedback only.

There are two main questions that come up with combinatorial (semi)-bandits: 1° *How can the stochastic structure of the reward vector be exploited to reduce the regret?* and 2° *Can algorithms be efficient?* Combes et al. (2015) answer the first question assuming that reward distributions are *mutually independent*. Later, Degenne & Perchet (2016) generalize the algorithm of Combes et al. (2015) to a larger class of *sub-Gaussian* rewards by exploiting the covariance structure of the arms. They also show the optimality of proposed algorithms, in particular, that an upper bound on their regret matches the asymptotic gap dependent lower bound of this class. However, algorithms of Combes et al. (2015) and Degenne & Perchet (2016) are computationally inefficient. The second question is studied by Kveton et al. (2015), who give an efficient algorithm based on the UCB algorithm of Auer et al. (2002). While being efficient, the algorithm of Kveton et al. (2015) assumes the worst case class of *arbitrary correlated*<sup>1</sup> rewards, i.e., it does not exploit any properties of rewards and therefore does not match the lower bound of Degenne & Perchet (2016).

On the other hand, efficient algorithms for matroid (Whitney, 1935) semi-bandits exist (Kveton et al., 2014; Talebi & Proutiere, 2016) and their regret bounds match the asymptotic gap dependent lower bound, which is the same for both sub-Gaussian and arbitrary correlated rewards. Among these algorithms, the state-of-the-art gap-free regret bound is of order  $\mathcal{O}(\sqrt{nmT \log T})$ , where  $T$  is the number of rounds,  $n$  is the number of base arms, and  $m$  is the maximum action size (Kveton et al., 2014).

---

<sup>1</sup>SequeL team, INRIA Lille - Nord Europe <sup>2</sup>CMLA, ENS Paris-Saclay <sup>3</sup>Criteo AI Lab <sup>4</sup>now with DeepMind. Correspondence to: P. Perrault <pierre.perrault@inria.fr>.

<sup>1</sup>Any dependence can exist between rewards.

**Our contributions** In this paper, we show how algorithms of Combes et al. (2015) and Degenne & Perchet (2016) can be *efficiently* approximated for matroid. This improves the bound of Kveton et al. (2014) by a factor  $\sqrt{m}/\log(m)$  for the class of *sub-Gaussian* rewards. We first locate the source of inefficiency of these algorithms: At each round, they have to solve a *submodular maximization* problem. We then provide efficient, adapted LOCALSEARCH and GREEDY-based algorithms that exploit the submodularity to give approximation guarantees on the regret upper bound. These algorithms can be of independent interest. We also extend our approximation techniques to more challenging *budgeted combinatorial semi-bandits* via binary search methods and exhibit the same improvement for this setting as well.

**Related work** Efficiency in combinatorial bandits, or more generally, in linear bandits with a large set of arms is an open problem. Some methods, such as *convex hull representation*, (Koolen et al., 2010), *hashing* (Jun et al., 2017), or *DAG-encoding* (Sakaue et al., 2018) reduce the algorithmic complexity. For semi-bandit feedback, in the adversarial case, Neu & Bartók (2013) proposed an efficient implementation via *geometric resampling*. In the stochastic case, many efficient Bayesian algorithms exist (see e.g., Russo & Van Roy, 2013; Russo & Roy, 2016), although they are only shown to be optimal for *Bayesian regret*.<sup>2</sup>

## 2. Background

We denote the set of arms by  $[n] \triangleq \{1, 2, \dots, n\}$ , we typeset vectors in bold and indicate components with indices, i.e.,  $\mathbf{a} = (a_i)_{i \in [n]} \in \mathbb{R}^n$ . We let  $\mathcal{P}([n]) \triangleq \{A, A \subset [n]\}$  be the power set of  $[n]$ . Let  $\mathbf{e}_i \in \mathbb{R}^n$  denote the  $i^{\text{th}}$  canonical unit vector. The *incidence vector* of any set  $A \in \mathcal{P}([n])$  is

$$\mathbf{e}_A \triangleq \sum_{i \in A} \mathbf{e}_i.$$

The above definition allows us to represent a subset of  $[n]$  as an element of  $\{0, 1\}^n$ . We denote the Minkowski sum of two sets  $Z, Z' \subset \mathbb{R}^n$  as  $Z + Z' \triangleq \{z + z', z \in Z, z' \in Z'\}$ , and  $Z + z' \triangleq Z + \{z'\}$ . Let  $\mathcal{A} \subset \mathcal{P}([n])$  be a set of *actions*. We define the maximum possible cardinality of an element of  $\mathcal{A}$  as  $m \triangleq \max\{|A|, A \in \mathcal{A}\}$ .

### 2.1. Stochastic Combinatorial Semi-Bandits

In combinatorial semi-bandits, an agent selects an action  $A_t \in \mathcal{A}$  at each round  $t \in \mathbb{N}^*$ , and receives a reward  $\mathbf{e}_{A_t}^\top \mathbf{X}_t$ , where  $\mathbf{X}_t \in \mathbb{R}^n$  is an unknown random vector of rewards. The successive reward vectors  $(\mathbf{X}_t)_{t \geq 1}$  are i.i.d., with an unknown mean  $\boldsymbol{\mu}^* \triangleq \mathbb{E}[\mathbf{X}] \in \mathbb{R}^n$ , where  $\mathbf{X} = \mathbf{X}_1$ . After

<sup>2</sup>A setting where arm distributions depend on a random parameter drawn from a known prior, and expectation of the regret is also taken with respect to this parameter.

selecting an action  $A_t$  in round  $t$ , the agent observes the partial reward of each individual arm in  $A_t$ . The goal of the agent is to minimize the expected regret, defined with  $A^* \in \arg \max_{A \in \mathcal{A}} \mathbf{e}_A^\top \boldsymbol{\mu}^*$  as

$$R_T \triangleq \mathbb{E} \left[ \sum_{t=1}^T (\mathbf{e}_{A^*} - \mathbf{e}_{A_t})^\top \mathbf{X}_t \right].$$

For any action  $A \in \mathcal{A}$ , we define its gap as the difference  $\Delta(A) \triangleq (\mathbf{e}_{A^*} - \mathbf{e}_A)^\top \boldsymbol{\mu}^*$ . We then rewrite the expected cumulative regret as  $R_T = \mathbb{E} \left[ \sum_{t=1}^T \Delta(A_t) \right]$ . Finally, we define  $\Delta \triangleq \min_{A \in \mathcal{A}, \Delta(A) > 0} \Delta(A)$ .

Combinatorial semi-bandits have been introduced by Cesa-Bianchi & Lugosi (2012). More recently, different algorithms have been proposed (Talebi et al., 2013; Combes et al., 2015; Kveton et al., 2015; Degenne & Perchet, 2016), depending whether the random vector  $\mathbf{X}$  satisfies specific properties. Some of these properties commonly assumed are a subset of the following ones:

- (i)  $X_1, \dots, X_n \in \mathbb{R}$  are *mutually independent*,
- (ii)  $X_1, \dots, X_n \in \mathbb{R}$  are *arbitrary correlated*,
- (iii)  $\mathbf{X} \in [-1, 1]^n$ ,
- (iv)  $\mathbf{X} \in \mathbb{R}^n$  is *multivariate sub-Gaussian*,  
i.e.,  $\mathbb{E} \left[ e^{\boldsymbol{\lambda}^\top (\mathbf{X} - \boldsymbol{\mu}^*)} \right] \leq e^{\|\boldsymbol{\lambda}\|_2^2 / 2}$ ,  $\forall \boldsymbol{\lambda} \in \mathbb{R}^n$ ,
- (v)  $\mathbf{X} \in \mathbb{R}^n$  is *component-wise sub-Gaussian*,  
i.e.,  $\mathbb{E} \left[ e^{\lambda_i (X_i - \mu_i^*)} \right] \leq e^{\lambda_i^2 / 2}$ ,  $\forall i \in [n], \forall \boldsymbol{\lambda} \in \mathbb{R}^n$ .

### 2.2. Lower Bounds

Combining some of the above properties, we consider different classes of possible distributions for  $\mathbf{X}$ . In Table 1, we show two existing gap-dependent lower bounds on  $R_T$  that depend on the respective class. They are valid for at least one distribution of  $\mathbf{X}$  belonging to the corresponding class, one combinatorial structure  $\mathcal{A} \subset \mathcal{P}([n])$ , and for any consistent algorithm (Lai & Robbins, 1985), for which the regret on any problem verifies  $R_T = o(T^a)$  as  $T \rightarrow \infty$  for all  $a > 0$ . Table 1 suggests that a tighter regret rate can be reached with some prior knowledge on the random vector  $\mathbf{X}$ .

## 3. (In)efficiency of Existing Algorithms

In this section, we discuss the efficiency of existing algorithms matching the lower bounds in Table 1. We consider that an algorithm is *efficient* as soon as the time and space complexity for each round  $t$  is polynomial in  $n$  and polylog-

Table 1. Gap-dependent lower bounds proved on different classes of possible distributions for  $\mathbf{X}$ .

Class of possible reward distributions	Gap-dependent lower bound
$(i) + (iii)$ $\Rightarrow (i) + (v)$ $\Rightarrow (iv)$	$\Omega\left(\frac{n \log T}{\Delta}\right)$ Combes et al., 2015
$(ii) + (iii)$ $\Rightarrow (ii) + (v)$	$\Omega\left(\frac{nm \log T}{\Delta}\right)$ Kveton et al., 2015

arithmetic<sup>3</sup> in  $t$ . Notice that the per-round complexity depends substantially on  $\mathcal{A}$ . We assume  $\mathcal{A}$  is such that linear optimization problems on  $\mathcal{A}$  — of the form  $\max_{A \in \mathcal{A}} \mathbf{e}_A^\top \boldsymbol{\delta}$  for some  $\boldsymbol{\delta} \in \mathbb{R}^n$  — can be solved efficiently. As a consequence, an agent knowing  $\boldsymbol{\mu}^*$  can efficiently compute  $A^*$ . Assuming efficient linear maximization is crucial (cf. Neu & Bartók, 2013; Combes et al., 2015; Kveton et al., 2015; Degenne & Perchet, 2016). Without this assumption, e.g., for  $\mathcal{A}$  being dominating sets in a graph, even the offline problem cannot be solved efficiently, and we would have to consider the notion of approximation regret instead, as was done by Chen et al. (2013).

### 3.1. A Generic Algorithm

As mentioned above, the action  $A_t$  is selected based on the feedback received up to round  $t - 1$ . A common statistic computed from this feedback is the *empirical average* of each arm  $i \in [n]$ , defined as

$$\bar{\mu}_{i,0} = 0, \quad \forall t \geq 2, \quad \bar{\mu}_{i,t-1} = \frac{\sum_{u \in [t-1]} \mathbb{I}\{i \in A_u\} X_{i,u}}{N_{i,t-1}},$$

where  $\forall t \geq 1$ ,  $N_{i,t-1} \triangleq \sum_{u \in [t-1]} \mathbb{I}\{i \in A_u\}$ . Many combinatorial semi-bandit algorithms, in particular, those listed in Table 2, can be seen as a special case of Algorithm 1 for different confidence regions  $\mathcal{C}_t$  around  $\bar{\boldsymbol{\mu}}_{t-1}$ .

---

#### Algorithm 1 Generic confidence-region-based algorithm.

---

At each round  $t$  :

- Find a confidence region  $\mathcal{C}_t \subset \mathbb{R}^n$ .
- Solve the bilinear program

$$(\boldsymbol{\mu}_t, A_t) \in \arg \max_{\boldsymbol{\mu} \in \mathcal{C}_t, A \in \mathcal{A}} \mathbf{e}_A^\top \boldsymbol{\mu}.$$

Play  $A_t$ .

---

<sup>3</sup>In streaming settings with near real-time requirements, it is imperative to have algorithms that can run with a complexity that stay almost constant across rounds.

We further assume that  $\mathcal{C}_t$  is defined through some parameters  $p, r \in \{1, \infty\}$ , and some functions  $g_{i,t}, i \in [n]$  by

$$\mathcal{C}_t \triangleq [-r, r]^n \cap \left( \bar{\boldsymbol{\mu}}_{t-1} + \left\{ \boldsymbol{\delta} \in \mathbb{R}^n, \|(g_{i,t}(\boldsymbol{\delta}))_i\|_p \leq 1 \right\} \right),$$

where  $g_{i,t} = 0$  if  $N_{i,t-1} = 0$  and, otherwise, is convex, strictly decreasing on  $[-r - \bar{\mu}_{i,t-1}, 0]$  and strictly increasing on  $[0, r - \bar{\mu}_{i,t-1}]$  such that  $g_{i,t}(0) = 0$ . Typically,  $r = 1$  under assumption (iii) and  $r = \infty$  otherwise. Table 2 lists variants of Algorithm 1, with the corresponding reward class under which they can be used. Each of these algorithms is matching the lower bound corresponding to the respective reward class considered in Table 1, i.e.,  $R_T$  is a ‘big  $\mathcal{O}$ ’ of the lower bound, up to a polylogarithmic factor in  $m$  (Degenne & Perchet, 2016). Notice that THOMPSONSAMPLING is not an instance of Algorithm 1. However, we are not aware of any tight analysis: The one by Wang & Chen (2018) matches the lower bound  $nm \log(T)/\Delta$ , but only for mutually independent rewards, where the lower bound is  $n \log(T)/\Delta$ . The regret upper bound of algorithms in Table 1 with  $p = 1$  have an additive constant term w.r.t.  $T$  but exponential in  $n$ , which can be replaced with a different analysis to get either:

- an exponential term in  $m$  plus a term of order  $1/\Delta^2$ ,
- a term of order  $1/\Delta^2$  — by changing  $\log(t) + m$  to  $\log(t) + m \log \log(t)$  in the algorithm,
- or can be removed — by changing  $\log(t) + m$  to  $\log(t) + n \log \log(t)$  in the algorithm.

On one hand, the arbitrary correlated case can be considered as solved, since the matching lower bound algorithm CUCB (Kveton et al., 2015) is efficient.<sup>4</sup> On the other hand, considering the reward class given by the first line of Table 1, the known algorithms that match the lower bound are inefficient.<sup>5</sup> We further discuss the efficiency of Algorithm 1 in the following subsection.

### 3.2. Submodular Maximization

In Algorithm 1, only  $A_t$  needs to be computed. It is a maximizer over  $\mathcal{A}$  of the set function

$$\begin{aligned} \mathcal{P}([n]) &\rightarrow \mathbb{R} \\ A &\mapsto \max_{\boldsymbol{\mu} \in \mathcal{C}_t} \mathbf{e}_A^\top \boldsymbol{\mu}. \end{aligned} \quad (1)$$

We can easily evaluate the function (1) above for some set  $A \in \mathcal{P}([n])$ , since it only requires solving a linear optimization problem on the convex<sup>6</sup> set  $\mathcal{C}_t$ . In Proposition 1, we

<sup>4</sup>In Theorem 1, we recover that  $A_t$  is computed by Algorithm 1 by solving a linear optimization problem.

<sup>5</sup> $\mathcal{A}$  may have up to  $2^n$  elements.

<sup>6</sup> $\mathcal{C}_t$  is convex since functions  $g_{i,t}$  are convex.

Table 2. Some combinatorial semi-bandit algorithms and their properties.

Class	Algorithm	$p$	$g_{i,t}$ (for $N_{i,t-1} \geq 1$ , up to universal factor)	Efficient?
(ii) + (v)	CUCB (Kveton et al., 2015)	$\infty$	$\delta \mapsto \frac{\delta^2 N_{i,t-1}}{\log(t)}$	Yes
(i) + (iii)	ESCB-KL (Combes et al., 2015)	1	$\delta \mapsto \frac{\text{kl}\left(\frac{1+\bar{\mu}_{i,t-1}}{2}, \frac{1+\bar{\mu}_{i,t-1}+\delta}{2}\right) N_{i,t-1}}{\log(t) + m}$	No
(iv)	ESCB (Combes et al., 2015), OLS-UCB (Degenne & Perchet, 2016)	1	$\delta \mapsto \frac{\delta^2 N_{i,t-1}}{\log(t) + m}$	No

show that in some cases, the evaluation can be even simpler. However, maximizing the function (1) over a combinatorial set  $\mathcal{A}$  is not straightforward. Before studying this function more closely, Definition 1 recalls some well-known properties that can be satisfied by a set function  $F : \mathcal{P}([n]) \rightarrow \mathbb{R}$ .

**Definition 1.** A set function  $F$  is:

- *normalized*, if  $F(\emptyset) = 0$ ,
- *linear (or modular)* if  $F(A) = \mathbf{e}_A^\top \boldsymbol{\delta} + b$ , for some  $\boldsymbol{\delta} \in \mathbb{R}^n$ ,  $b \in \mathbb{R}$ ,
- *non-decreasing* if  $F(A) \leq F(B) \forall A \subset B \subset [n]$ ,
- *submodular* if for all  $A, B \subset [n]$ ,

$$F(A \cup B) + F(A \cap B) \leq F(A) + F(B).$$

Equivalently,  $F$  is submodular if for all  $A \subset B \subset [n]$ , and  $i \notin B$ ,  $F(A \cup \{i\}) - F(A) \geq F(B \cup \{i\}) - F(B)$ .

The function (1) is clearly normalized, and it can be decomposed into two set functions in the following way,

$$\forall A \subset [n], \max_{\boldsymbol{\mu} \in \mathcal{C}_t} \mathbf{e}_A^\top \boldsymbol{\mu} = \mathbf{e}_A^\top \bar{\boldsymbol{\mu}}_{t-1} + \max_{\boldsymbol{\delta} \in \mathcal{C}_t - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta}.$$

The linear part  $A \mapsto \mathbf{e}_A^\top \bar{\boldsymbol{\mu}}_{t-1}$  is efficiently maximized alone, we thus focus on the other part,  $A \mapsto \max_{\boldsymbol{\delta} \in \mathcal{C}_t - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta}$ , usually called an *exploration bonus*. It aims to compensate for the negative selection bias of the first term. We define

$$\begin{aligned} \mathcal{C}_t^+ &\triangleq [-r, r]^n \cap \left( \bar{\boldsymbol{\mu}}_{t-1} + \left\{ \boldsymbol{\delta} \in \mathbb{R}_+^n, \|(g_{i,t}(\boldsymbol{\delta}))_i\|_p \leq 1 \right\} \right) \\ &= \mathcal{C}_t \cap \left\{ \boldsymbol{\mu} \in \mathbb{R}^n, \boldsymbol{\mu} \geq \bar{\boldsymbol{\mu}}_{t-1} \right\} \end{aligned}$$

and rewrite  $A \mapsto \max_{\boldsymbol{\delta} \in \mathcal{C}_t - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta}$  through Lemma 1.

**Lemma 1.** For all  $A \in \mathcal{P}([n])$ ,  $\max_{\boldsymbol{\delta} \in \mathcal{C}_t - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta} = \max_{\boldsymbol{\delta} \in \mathcal{C}_t^+ - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta}$ .

The lemma holds as  $\left\{ (\boldsymbol{\delta}_i^+)_i, \boldsymbol{\delta} \in \mathcal{C}_t - \bar{\boldsymbol{\mu}}_{t-1} \right\} \subset \mathcal{C}_t - \bar{\boldsymbol{\mu}}_{t-1}$ . As a corollary, this set function is non-negative, and non-decreasing. It can be written in closed form under additional assumptions, see Proposition 1 and Example 1.

**Proposition 1.** Let  $A \in \mathcal{P}([n])$ ,  $t \in \mathbb{N}^*$ ,  $p = 1$ . Assume that for all  $i \in A$ ,  $g_{i,t}$  has a strictly increasing, continuous derivative  $g'_{i,t}$  defined on  $[0, r - \bar{\mu}_{i,t-1}]$ . For  $i \in A$ , let

$$f_i(\lambda) \triangleq \begin{cases} g_{i,t}^{\prime-1}(1/\lambda) & \text{if } 1/\lambda < g'_{i,t}(r - \bar{\mu}_{i,t-1}), \\ r - \bar{\mu}_{i,t-1} & \text{otherwise,} \end{cases}$$

defined for  $\lambda \geq 0$ . Then, the smallest  $\lambda^*$  satisfying

$$\mathbf{e}_A^\top (g_{i,t}(f_i(\lambda^*)))_i \leq 1$$

is such that

$$(\boldsymbol{\delta}_i^*)_i \triangleq (\mathbb{I}\{i \in A\} f_i(\lambda^*))_i \in \arg \max_{\boldsymbol{\delta} \in \mathcal{C}_t^+ - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta}.$$

The proof of Proposition 1 can be found in Appendix A. An important use-case example of Proposition 1 is the following

**Example 1.** Let  $A \in \mathcal{P}([n])$ ,  $t \in \mathbb{N}^*$ . If for all  $i \in [n]$ ,  $g_{i,t} = (\cdot)^2 \alpha_{i,t}$  for some  $\alpha_{i,t} > 0$ , and  $r = \infty$ ,  $p = 1$ , then

$$\max_{\boldsymbol{\delta} \in \mathcal{C}_t^+ - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta} = \sqrt{\mathbf{e}_A^\top \left( \frac{1}{\alpha_{i,t}} \right)_i}.$$

Indeed, since the maximizer  $\boldsymbol{\delta}^*$  lies at the boundary,  $\max_{\boldsymbol{\delta} \in \mathcal{C}_t^+ - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta} = \max_{\boldsymbol{\delta} \in \mathbb{R}_+^n, \sum_i \alpha_{i,t} \delta_i^2 = 1} \mathbf{e}_A^\top \boldsymbol{\delta}$ , and from the first-order optimality condition we deduce that  $\mathbf{e}_A = 2\lambda^* (\alpha_{i,t} \boldsymbol{\delta}_i^*)_i$ , i.e.,  $\boldsymbol{\delta}_i^* = \mathbb{I}\{i \in A\} / 2\lambda^* \alpha_{i,t}$ , where  $\lambda^*$  is necessarily  $\frac{1}{2} \sqrt{\mathbf{e}_A^\top (1/\alpha_{i,t})_i}$ . We thus recover the ESCB's exploration bonus for  $\alpha_{i,t} = N_{i,t-1} / \log t$ .

**Remark 1.** The proof of Proposition 1 follows the same technique as the proof of Theorem 4 by Combes et al. (2015) for developing the computation of the ESCB-KL exploration bonus.

Example 1 is a specific case where the exploration bonus  $A \mapsto \max_{\boldsymbol{\delta} \in \mathcal{C}_t^+ - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta}$  has a particularly simple form: It is the square root of a non-decreasing linear set function. Such a set function is known to be submodular (Stobbe & Krause, 2010). This interesting property helps for maximizing the function (1). In Theorem 1, we prove that  $A \mapsto \max_{\boldsymbol{\delta} \in \mathcal{C}_t^+ - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta}$  is in fact always submodular.

**Theorem 1.** *The following two properties hold.*

- For  $p = \infty$ ,  $A \mapsto \max_{\delta \in \mathcal{C}_t^+ - \bar{\mu}_{t-1}} \mathbf{e}_A^\top \delta$  is linear.
- For  $p = 1$ ,  $A \mapsto \max_{\delta \in \mathcal{C}_t^+ - \bar{\mu}_{t-1}} \mathbf{e}_A^\top \delta$  is submodular.

The proof is deferred to Appendix B and uses a result on *polymatroids* by He et al. (2012). Theorem 1 first implies the efficiency of any variant of Algorithm 1 with  $p = \infty$ , since it reduces to optimizing a linear set function over  $\mathcal{A}$ . Theorem 1 also yields that when the reward class is strengthened to target the tighter lower bound  $n \log(T)/\Delta$ , Algorithm 1 reduces to maximizing a submodular set function over  $\mathcal{A}$  (the sum of a linear and a submodular function is submodular). Submodular maximization problems have been applied in machine learning before (see e.g., Krause & Golovin, 2011; Bach, 2011), however, maximizing a submodular function  $F$ , even for  $\mathcal{A} = \{A, |A| \leq m\}$  and  $F$  non-decreasing, is NP-Hard in general (Schrijver, 2008), with an approximation factor of  $1 + 1/(e - 1)$  by the GREEDY algorithm (Nemhauser et al., 1978). This is problematic as the typical analysis is based on controlling with high probability the error  $\Delta(A_t)$  at round  $t$  by the quantity  $2 \max_{\delta \in \mathcal{C}_t - \bar{\mu}_{t-1}} \mathbf{e}_{A_t}^\top (|\delta_i|)_i$ . More precisely, since  $\mu^*$  belongs with high probability to the confidence region  $\mathcal{C}_t$ ,  $\mu^* - \bar{\mu}_{t-1}$  belongs with high probability to  $\mathcal{C}_t - \bar{\mu}_{t-1}$ . Under this event, and for  $\kappa \geq 1$ , a  $\kappa$ -approximation algorithm for maximizing the function (1) would only guarantee the following:

$$\begin{aligned}
 \Delta(A_t) &= (\mathbf{e}_{A^*} - \mathbf{e}_{A_t})^\top \mu^* \\
 &\leq \max_{\mu \in \mathcal{C}_t} \mathbf{e}_{A^*}^\top \mu - \mathbf{e}_{A_t}^\top \mu^* \\
 &= \max_{\mu \in \mathcal{C}_t^+} \mathbf{e}_{A^*}^\top \mu - \mathbf{e}_{A_t}^\top \mu^* \\
 &\leq \kappa \max_{\mu \in \mathcal{C}_t^+} \mathbf{e}_{A_t}^\top \mu - \mathbf{e}_{A_t}^\top \mu^* \quad (2) \\
 &= \kappa \max_{\delta \in \mathcal{C}_t^+ - \bar{\mu}_{t-1}} \mathbf{e}_{A_t}^\top \delta \\
 &\quad + \mathbf{e}_{A_t}^\top (\bar{\mu}_{t-1} - \mu^*) + (\kappa - 1) \mathbf{e}_{A_t}^\top \bar{\mu}_{t-1} \quad (3) \\
 &\leq (\kappa + 1) \max_{\delta \in \mathcal{C}_t - \bar{\mu}_{t-1}} \mathbf{e}_{A_t}^\top (|\delta_i|)_i + (\kappa - 1) \mathbf{e}_{A_t}^\top \bar{\mu}_{t-1}.
 \end{aligned}$$

If  $\kappa \neq 1$ , the term  $(\kappa - 1) \mathbf{e}_{A_t}^\top \bar{\mu}_{t-1}$  gives linear regret bounds. In the next section, with a stronger assumption on  $\mathcal{A}$  (but for which submodular maximization is still NP-Hard), we show that both parts of the objective can have different approximation factors. More precisely, we show how to approximate the linear part with factor 1, and the submodular part with a constant factor  $\kappa > 1$ . Then, (2) can be replaced by

$$\kappa \cdot \max_{\mu \in \mathcal{C}_t^+ - \bar{\mu}_{t-1}} \mathbf{e}_{A_t}^\top \mu + 1 \cdot \mathbf{e}_{A_t}^\top \bar{\mu}_{t-1} - \mathbf{e}_{A_t}^\top \mu^*.$$

Therefore, in (3), the extra  $(\kappa - 1) \mathbf{e}_{A_t}^\top \bar{\mu}_{t-1}$  term is removed.

## 4. Efficient Algorithms for Matroid Constraints

In this section, we will consider additional structure on  $\mathcal{A}$ , using the notion of matroid, recalled below.

**Definition 2.** *A matroid is a pair  $([n], \mathcal{I})$ , where  $\mathcal{I}$  is a family of subsets of  $[n]$ , called the independent sets, with the following properties:*

- *The empty set is independent, i.e.,  $\emptyset \in \mathcal{I}$ .*
- *Every subset of an independent set is independent, i.e., for all  $A \in \mathcal{I}$ , if  $A' \subset A$ , then  $A' \in \mathcal{I}$ .*
- *If  $A$  and  $B$  are two independent sets, and  $|A| > |B|$ , then there exists  $x \in A \setminus B$  such that  $B \cup \{x\} \in \mathcal{I}$ .*

Matroids generalize the notion of linear independence. A maximal (for the inclusion) independent set is called *basis* and all bases have the same cardinality  $m$ , which is called the *rank* of the matroid (Whitney, 1935). Many combinatorial problems such as building a spanning tree for network routing (Oliveira & Pardalos, 2005) can be expressed as a linear optimization on a matroid (see Edmonds & Fulkerson, 1965 or Perfect, 1968, for other examples).

Let  $\mathcal{I} \in \mathcal{P}([n])$  be such that  $([n], \mathcal{I})$  forms a matroid. Let  $\mathcal{B} \subset \mathcal{I}$  be the set of bases of the matroid  $([n], \mathcal{I})$ . In the following, we may assume that  $\mathcal{A}$  is either  $\mathcal{I}$  or  $\mathcal{B}$ . We also assume that an independence oracle is available, i.e., given an input  $A \subset [n]$ , it returns TRUE if  $A \in \mathcal{I}$  and FALSE otherwise. Maximizing a linear set function  $L$  on  $\mathcal{A}$  is efficient, and it can be done as follows (Edmonds, 1971): Let  $\sigma$  be a permutation of  $[n]$  and  $j$  an integer such that  $j = m$  in case  $\mathcal{A} = \mathcal{B}$  and otherwise,  $j$  satisfies

$$\ell_1 \geq \dots \geq \ell_j \geq 0 \geq \ell_{j+1} \geq \dots \geq \ell_n,$$

where  $\ell_i = L(\{\sigma(i)\}) \forall i \in [n]$ . The optimal  $S$  is built greedily starting from  $S = \emptyset$ , and for  $i \in [j]$ ,  $\sigma(i)$  is added to  $S$  only if  $S \cup \{\sigma(i)\} \in \mathcal{I}$ .

Matroid bandits with  $\mathcal{A} = \mathcal{B}$  has been studied by Kveton et al. (2014); Talebi & Proutiere (2016). In this case, the two lower bounds in Table 1 coincide to  $\Omega(n \log(T)/\Delta)$ , and CUCB reaches it, with the following gap-free upper bound:  $R_T(\text{CUCB}) = \mathcal{O}(\sqrt{nmT \log T})$ . Assuming sub-Gaussian rewards to use any Algorithm of Table 2 with  $p = 1$  would tighten (Degenne & Perchet, 2016) this gap-free upper bound to  $\mathcal{O}(\sqrt{n \log^2 mT \log T})$ . Notice, due to the  $\sqrt{\log T}$  factor, this does not contradict the  $\Omega(\sqrt{nmT})$  gap-free lower bound for multi-play bandits.

In the rest of this section, we provide efficient approximation routines to maximize the function (1) on  $\mathcal{A} = \mathcal{I}$  and  $\mathcal{B}$  without having the extra undesirable term  $(\kappa - 1) \mathbf{e}_{A_t}^\top \bar{\mu}_{t-1}$ , that a standard  $\kappa$ -approximation algorithm would suffer.

Therefore, using these routines in Algorithm 1 do not alter its regret upper bound.

Let  $L$  be a normalized, linear set function, that will correspond to the linear part  $A \mapsto \mathbf{e}_A^\top \bar{\boldsymbol{\mu}}_{t-1}$ ; and let  $F$  denote a normalized, non-decreasing, submodular function, that will correspond to the submodular part  $A \mapsto \max_{\boldsymbol{\delta} \in \mathcal{C}_t^+ - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta}$ . Unless stated otherwise, we further assume that  $F$  is positive (for  $A \neq \emptyset$ ). This is a mild assumption as it holds for  $A \mapsto \max_{\boldsymbol{\delta} \in \mathcal{C}_t^+ - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta}$  in the unbounded case, i.e., if (iii) is not assumed and  $r = \infty$ . If (iii) is true, then adding an extra  $\mathbf{e}_A^\top \left( \frac{1}{N_{i,t-1}^2} \right)_i$  term will recover positivity and increase the regret upper bound by only an additive constant. In the following subsections, we will provide algorithms that efficiently outputs  $S$  such that

$$L(S) + \kappa F(S) \geq L(O) + F(O), \quad \forall O \in \mathcal{A}, \quad (4)$$

with some appropriate approximation factor  $\kappa \geq 1$ . It is possible to efficiently output  $S_1$  and  $S_2$  such that we get  $L(S_1) \geq L(O_1)$  and  $\kappa F(S_2) \geq F(O_2)$  for any  $O_1, O_2 \in \mathcal{A}$ . Although we can take  $O_1 = O_2$ ,  $S_1$  and  $S_2$  are not necessarily equal, and (4) is not straightforward. The last subsection is we apply this approach to *budgeted matroid semi-bandits*.

#### 4.1. Local Search Algorithm

In this subsection, we assume that  $\mathcal{A} = \mathcal{I}$ . In Algorithm 2, we provide a specific instance of LOCALSEARCH that we tailored to our needs to approximately maximize  $L + F$ .

It starts from the greedy solution  $S_{\text{init}} \in \arg \max_{A \in \mathcal{I}} L(A)$ . Then, Algorithm 2 repeatedly tries three basic operations in order to improve the current solution. Since every  $S \in \mathcal{A}$  can potentially be visited, only *significant* improvements are considered, i.e., improvements greater than  $\frac{\varepsilon}{m} F(S)$  for some input parameter  $\varepsilon > 0$ . The smaller  $\varepsilon$  is, the higher complexity will be. Notice the improvement threshold  $\frac{\varepsilon}{m} F(S)$  does not depend on  $L$ . In fact, this is crucial to ensure that the approximation factor of  $L$  is 1. However, this can increase the time complexity. To prevent this increase, the second essential ingredient is the initialization, where only  $L$  plays a role.

In Theorem 2, we state the approximation guarantees for Algorithm 2 and its time complexity. The proof of Theorem 2 is in Appendix C. For  $\mathcal{C}_t$  given by any algorithm of Table 2,  $F(A) = \max_{\boldsymbol{\delta} \in \mathcal{C}_t^+ - \bar{\boldsymbol{\mu}}_{t-1}} \mathbf{e}_A^\top \boldsymbol{\delta}$ , and  $\varepsilon = 1$ , the time complexity is bounded by  $\mathcal{O}(m^2 n \log(mt))$ , and is thus efficient. Another algorithm enjoying an improved time complexity is provided in the next subsection, in the case where  $\mathcal{A} = \mathcal{B}$ .

**Theorem 2.** *Algorithm 2 outputs  $S \in \mathcal{I}$  such that*

$$L(S) + 2(1 + \varepsilon)F(S) \geq L(O) + F(O), \quad \forall O \in \mathcal{I}.$$

**Algorithm 2** LOCALSEARCH for maximizing  $L + F$  on  $\mathcal{I}$ .

**Input:**  $L, F, \mathcal{I}, m, \varepsilon > 0$ .

**Initialization:**  $S_{\text{init}} \in \arg \max_{A \in \mathcal{I}} L(A)$ .

**if**  $S_{\text{init}} = \emptyset$  **then**

**if**  $\exists \{x\} \in \mathcal{I}$  such that  $(L + F)(\{x\}) > 0$  **then**

$S_0 \in \arg \max_{\{x\} \in \mathcal{I}, (L+F)(\{x\}) > 0} L(\{x\})$ .

**else**

**Output**  $\emptyset$

**end if**

**else**

$S_0 \leftarrow S_{\text{init}}$

**end if**

$S \leftarrow S_0$ .

Repeatedly perform one of the following local improvements **while** possible:

• **Delete an element:**

**if**  $\exists x \in S$  such that

$(L + F)(S \setminus \{x\}) > (L + F)(S) + \frac{\varepsilon}{m} F(S)$ ,

**then**  $S \leftarrow S \setminus \{x\}$ .

**end if**

• **Add an element:**

**if**  $\exists y \in [n] \setminus S, S \cup \{y\} \in \mathcal{I}$ , such that

$(L + F)(S \cup \{y\}) > (L + F)(S) + \frac{\varepsilon}{m} F(S)$ ,

**then**  $S \leftarrow S \cup \{y\}$ .

**end if**

• **Swap a pair of elements:**

**if**  $\exists (x, y) \in S \times [n] \setminus S, S \setminus \{x\} \cup \{y\} \in \mathcal{I}$ , such that

$(L + F)(S \setminus \{x\} \cup \{y\}) > (L + F)(S) + \frac{\varepsilon}{m} F(S)$

**then**  $S \leftarrow S \setminus \{x\} \cup \{y\}$  **end if**

**end while**

**Output:**  $S$ .

*Its complexity is  $\mathcal{O}\left(mn \log\left(\frac{\max_{A \in \mathcal{I}} F(A)}{F(S_0)}\right) / \log\left(1 + \frac{\varepsilon}{m}\right)\right)$ .*

Theorem 2 gives a parameter  $\kappa$  arbitrary close to 2 in (4).<sup>7</sup>

#### 4.2. Greedy Algorithm

In this section, we assume that  $\mathcal{A} = \mathcal{B}$ . This situation happens, for instance, under a non-negativity assumption on  $L$ , i.e., if we consider non-negative rewards  $X_i$ . We show that the standard GREEDY algorithm (Algorithm 3) improves over Algorithm 2 by exploiting this extra constraint, both in terms of the running time and the approximation factor. We state the result in Theorem 3 and prove it in Appendix D.

<sup>7</sup>We could design a different version of Algorithm 2, based on NON-OBLIVIOUSLOCALSEARCH (Filmus & Ward, 2012), in order to get  $\kappa$  arbitrary close to  $1 + 1/(e - 1)$ , but with a much worst time complexity. Actually, Sviridenko et al. (2013) proposed such an approach, with an approximation factor for  $L$  arbitrary close to 1, but not equal, so we would get back the undesirable term, which would require a complexity polynomial in  $t$  to control.

---

**Algorithm 3** GREEDY for maximizing  $L + F$  on  $\mathcal{B}$ .

**Input:**  $L, F, \mathcal{I}, m$ .

**Initialization:**  $S \leftarrow \emptyset$ .

**for**  $i \in [k]$  **do**
 $x \in \arg \max_{x \notin S, S \cup \{x\} \in \mathcal{I}} (L + F)(S \cup \{x\})$ .

 $S \leftarrow S \cup \{x\}$ .

**end for**
**Output:**  $S$ .

---

Notice that another advantage is that we do not need to assume  $F(A) > 0$  for  $A \neq \emptyset$  here.

**Theorem 3.** *Algorithm 3 outputs  $S \in \mathcal{B}$  such that*

$$L(S) + 2F(S) \geq L(O) + F(O), \quad \forall O \in \mathcal{B}.$$

*Its complexity is  $\mathcal{O}(mn)$ .*

Combining the results before, we get the following theorem.

**Theorem 4.** *With approximation techniques, the cumulative regret for the combinatorial semi-bandits is bounded as*

$$R_T \leq \mathcal{O}\left(\sqrt{n \log^2(m) T \log T}\right)$$

*with per-round time complexity of order  $\mathcal{O}(\log(mt)m^2n)$  (resp.,  $\mathcal{O}(mn)$ ) for  $\mathcal{A} = \mathcal{I}$  (resp.,  $\mathcal{A} = \mathcal{B}$ ).*

Notice that this new bound is better by a factor  $\sqrt{m}/\log m$  than the one of Kveton et al. (2014) in the case  $\mathcal{A} = \mathcal{B}$ .

### 4.3. Budgeted Matroid Semi-Bandits

In this subsection, we extend results of the two previous subsections to budgeted matroid semi-bandits. In budgeted bandits with single resource and infinite horizon (Ding et al., 2013; Xia et al., 2016a), each arm is associated with both a reward and a cost. The agent aims at maximizing the cumulative reward under a budget constraint for the cumulative costs. Xia et al. (2016b) studied budgeted bandits with multiple play, where a  $m$ -subset  $A$  of arms is selected at each round. An optimal (up to a constant term) offline algorithm chooses the same action  $A^*$  within each round, where  $A^*$  is the minimizer of the ratio “expected cost paid choosing  $A^*$ ” over “expected reward gained choosing  $A^*$ ”. In the setting of Xia et al. (2016b), the agent observes the *partial* random cost and reward of each arm in  $A$  (i.e., semi-bandit feedback), pays the sum of partial costs of  $A$  and gains the sum of partial rewards of  $A$ .  $A^*$  can be computed efficiently, and a Xia et al. (2016b) give an algorithm based on CUCB. It minimizes the ratio where the averages are replaced by UCBS. We extend this setting to matroid constraints. We assume that total costs/rewards are non-negative linear set functions of the chosen action  $A$ . The objective is to minimize a ratio of linear set functions. As previously, two kinds

of constraints can be considered for the minimization: either  $\mathcal{A} = \mathcal{I}$  or  $\mathcal{A} = \mathcal{B}$ . Theorem 1 implies that an optimistic estimation of this ratio is of the form  $\frac{L_1 - F_1}{L_2 + F_2}$ , where for  $i \in \{1, 2\}$ ,  $F_i$  are positive (except for  $\emptyset$ ), normalized, non-decreasing, submodular; and  $L_i$  are non-negative and linear.  $L_1 - F_1$  is a high-probability lower bound on the expected cost paid, and  $L_2 + F_2$  is a high-probability upper bound on the expected reward gained. Notice that the numerator,  $L_1 - F_1$ , can be negative, which can be an incitement to take arms with a high cost/low rewards. Therefore, we consider the minimization of the surrogate  $\left(\frac{L_1 - F_1}{L_2 + F_2}\right)^+$ . Indeed,  $(L_1 - F_1)/(L_2 + F_2)$  is a high probability lower bound on the ratio of expectation, so by monotonicity of  $x \mapsto x^+$  on  $\mathbb{R}$ ,  $(L_1 - F_1)^+/(L_2 + F_2)$  is also a high-probability lower bound. We assume  $L_2$  is normalized, but not necessarily  $L_1$ .  $L_1(\emptyset)$  can be seen as an entry price. When  $L_1$  is normalized, we assume that  $\emptyset$  is not feasible.

**Remark 2.** *Notice, If  $\mathcal{A} = \mathcal{I}$ , and  $L_1$  is normalized, then there is an optimal solution of the form  $\{s\} \in \mathcal{I}$ : If  $L_1 - F_1$  is negative for some  $S = \{s\} \subset \mathcal{I}$ , then such  $S$  is a minimizer. Otherwise, by submodularity (and thus subadditivity, since we consider normalized functions),  $L_1 - F_1$  is non-negative, and we have*

$$\begin{aligned} \frac{L_1(S) - F_1(S)}{L_2(S) + F_2(S)} &\geq \frac{\sum_{s \in S} L_1(\{s\}) - F_1(\{s\})}{\sum_{s \in S} L_2(\{s\}) + F_2(\{s\})} \\ &\geq \min_{s \in S} \frac{L_1(\{s\}) - F_1(\{s\})}{L_2(\{s\}) + F_2(\{s\})}. \end{aligned}$$

This minimization problem is at least as difficult as previous submodular maximization problems, taking  $L_1 = 1$  and  $F_1 = 0$ . In order to use our approximation algorithms, we consider the *Lagrangian function* associated to the problem (see e.g., Fujishige, 2005),

$$\mathcal{L}(\lambda, S) \triangleq L_1(S) - F_1(S) - \lambda(L_2(S) + F_2(S)),$$

for  $\lambda \geq 0$  and  $S \subset [n]$ . For a fixed  $\lambda \geq 0$ ,  $-\mathcal{L}(\lambda, \cdot)$  is a sum of a linear and a submodular function, and both Algorithms 2 and 3 can be used. However, the first step is to find  $\lambda$  sufficiently close to the optimal ratio

$$\lambda^* = \min_{A \in \mathcal{A}} \left( \frac{L_1(A) - F_1(A)}{L_2(A) + F_2(A)} \right)^+.$$

**Remark 3.** *For some  $\lambda \geq 0$ ,*

$$\min_{A \in \mathcal{A}} \mathcal{L}(\lambda, A) \geq 0 \Rightarrow \lambda \leq \lambda^*,$$

$$\min_{A \in \mathcal{A}} \mathcal{L}(\lambda, A) \leq 0 \Rightarrow \begin{cases} \lambda \geq \lambda^*, \text{ or} \\ \min_{A \in \mathcal{A}} L_1(A) - F_1(A) \leq 0, \\ \text{which further gives } \lambda^* = 0. \end{cases}$$

From Remark 3, if it was possible to compute  $\min_{A \in \mathcal{A}} \mathcal{L}(\lambda, A)$  exactly, then a binary search algorithm

**Algorithm 4** Binary search for minimizing the ratio  $(L_1 - F_1)^+ / (L_2 + F_2)$ .

**Input:**  $L_1, L_2, F_1, F_2, \text{ALGO}_\kappa, \eta > 0$ .

$\delta \leftarrow \frac{\eta \min_{\{s\} \in \mathcal{A}} F_1(\{s\})}{L_2(B) + \kappa^2 F_2(B)}$  with  $B = \text{ALGO}_\kappa(L_2 + \kappa F_2)$ .

$A \leftarrow A_0 \in \mathcal{A} \setminus \{\emptyset\}$  arbitrary.

**if**  $\mathcal{L}_\kappa(0, A) > 0$  **then**

$\lambda_1 \leftarrow 0, \quad \lambda_2 \leftarrow \frac{L_1(A) - F_1(A)}{L_2(A) + F_2(A)}$ .

**while**  $\lambda_2 - \lambda_1 \geq \delta$  **do**

$\lambda \leftarrow \frac{\lambda_1 + \lambda_2}{2}$ .

$S \leftarrow \text{ALGO}_\kappa(-\mathcal{L}(\lambda, \cdot))$ .

**if**  $\mathcal{L}_\kappa(\lambda, S) \geq 0$  **then**

$\lambda_1 \leftarrow \lambda$ .

**else**

$\lambda_2 \leftarrow \lambda$ .

$A \leftarrow S$ .

**end if**

**end while**

**end if**

**Output:**  $A$ .

would find  $\lambda^*$ . This dichotomy method can be extended to  $\kappa$ -approximation algorithms by defining the *approximation Lagrangian* as

$$\mathcal{L}_\kappa(\lambda, S) \triangleq L_1(S) - \kappa F_1(S) - \lambda(L_2(S) + \kappa F_2(S)),$$

for  $\lambda \geq 0$  and  $S \subset [n]$ . The idea is to use the following approximation guarantee for a  $\kappa$ -approximation algorithms outputting  $S$  (with objective function  $-\mathcal{L}$ ),

$$\min_{A \in \mathcal{A}} \mathcal{L}_\kappa(\lambda, A) \leq \mathcal{L}_\kappa(\lambda, S) \leq \min_{A \in \mathcal{A}} \mathcal{L}(\lambda, A).$$

Thus, for a given  $\lambda$ , either the l.h.s is strictly negative or the r.h.s is non-negative, depending on the sign of  $\mathcal{L}_\kappa(\lambda, S)$ .

Therefore, from Remark 3, a lower bound  $\lambda_1$  on  $\lambda^*$ , and an upper bound  $\lambda_2$  on  $\min_{A \in \mathcal{A}} \left( \frac{L_1(A) - \kappa F_1(A)}{L_2(A) + \kappa F_2(A)} \right)^+$  can be computed. The detailed method is given in Algorithm 4. Notice that it takes as input some  $\text{ALGO}_\kappa$ , that can be either Algorithm 2 or Algorithm 3, depending on the assumption on the constraint (either  $\mathcal{A} = \mathcal{I}$  or  $\mathcal{A} = \mathcal{B}$ ). We denote the output as  $\text{ALGO}_\kappa(L + F)$ , for some linear set function  $L$  and some submodular set function  $F$ , for maximizing the objective  $L + F$  on  $\mathcal{A}$ , so that  $S = \text{ALGO}_\kappa(L + F)$  satisfies  $L(S) + \kappa F(S) \geq \max_{A \in \mathcal{A}} L(A) + F(A)$ , i.e.,  $\kappa = 2(1 + \varepsilon)$  if  $\text{ALGO}_\kappa = \text{Algorithm 2}$ ,  $\mathcal{A} = \mathcal{I}$  and  $\kappa = 2$  if  $\text{ALGO}_\kappa = \text{Algorithm 3}$ ,  $\mathcal{A} = \mathcal{B}$ . In Theorem 5, we state the result for the output of Algorithm 4 and prove it in Appendix E.

**Theorem 5.** *Algorithm 4 outputs  $A$  such that*

$$\left( \frac{L_1(A) - (\kappa + \eta)F_1(A)}{L_2(A) + \kappa F_2(A)} \right)^+ \leq \lambda^*,$$

where  $\lambda^*$  is the minimum of  $\left( \frac{L_1 - F_1}{L_2 + F_2} \right)^+$  over  $\mathcal{I}$  if  $\text{ALGO}_\kappa = \text{Algorithm 2}$ , and over  $\mathcal{B}$  if  $\text{ALGO}_\kappa = \text{Algorithm 3}$ . For  $\mathcal{C}_t$  given by any algorithm of Table 2,  $F(A) = \max_{\delta \in \mathcal{C}_t^+ - \bar{\mu}_{t-1}} \mathbf{e}_A^\top \delta$ , the complexity is of order  $\log(mt/\eta)$  times the complexity of  $\text{ALGO}_\kappa$ .

## 5. Experiments

We provide experiments for a *graphic matroid*, on a five nodes complete graph, as did Combes et al. (2015). We thus have  $n = 10, m = 4$ . We consider two experiments. In the first one we use  $\mathcal{A} = \mathcal{B}$ ,  $\mu_i^* = 1 + \Delta \mathbb{I}\{i \leq m\}$ , for all  $i \in [n]$ , and in the second,  $\mathcal{A} = \mathcal{I}$ , where we set  $\mu_i^* = \Delta(2\mathbb{I}\{i \leq m - 1\} - 1), \forall i \in [n]$ . We take  $\Delta = 0.1$ , with rewards drawn from independent unit variance Gaussian distributions. Figure 2 illustrates the comparison between CUCB and our implementations of ESCB (Combes et al., 2015) using Algorithm 3 (left) and 2 (right, with  $\varepsilon = 0.1$ ), showing the behavior of the regret vs. time horizon. We observe that although we are approximating the confidence region within a factor at least 2 (and thus force the exploration), our efficient implementation outperforms CUCB. Indeed, we take advantage (gaining a factor  $\sqrt{m}/\log m$ ) of the previously inefficient algorithm (that we made efficient) to use reward independence, which the more conservative CUCB is not able to. The latter algorithm has still a better per round-time complexity of  $\mathcal{O}(n \log m)$  and may be more practical on larger instances.

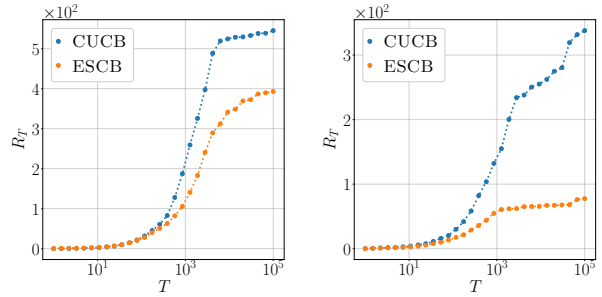


Figure 1. Cumulative regret for the minimum spanning tree setting in up to  $10^5$  rounds, averaged over 100 independent simulations. **Left:** for  $\mathcal{A} = \mathcal{B}$ . **Right:** for  $\mathcal{A} = \mathcal{I}$ .

## 6. Discussion

In this paper, we gave several approximation schemes for the confidence regions and applied them to combinatorial semi-bandits with matroid constraints and their budgeted version. We believe our approximation methods can be extended to approximation regret for non-linear objective functions (e.g., for influence maximization, Wang & Chen, 2018), if the maximization algorithm keeps the same approximation factor for the objective, either with or without the bonus.



**Acknowledgments** Vianney Perchet has benefited from the support of the ANR (grant n.ANR-13-JS01-0004-01), of the FMJH Program Gaspard Monge in optimization and operations research (supported in part by EDF), from the Labex LMH and from the CNRS through the PEPS program. The research presented was also supported by European CHIST-ERA project DELTA, French Ministry of Higher Education and Research, Nord-Pas-de-Calais Regional Council, Inria and Otto-von-Guericke-Universität Magdeburg associated-team north-European project Allocate, and French National Research Agency project BoB (grant n.ANR-16-CE23-0003), FMJH Program PGMO with the support to this program from Criteo.

## References

- Audibert, J.-Y., Bubeck, S., and Lugosi, G. Regret in online combinatorial optimization. *Mathematics of Operations Research*, 39:31–45, 2014.
- Auer, P., Cesa-Bianchi, N., and Fischer, P. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256, 2002.
- Bach, F. Learning with Submodular Functions: A Convex Optimization Perspective. 2011.
- Brualdi, R. A. Comments on bases in dependence structures. *Bulletin of the Australian Mathematical Society*, 1(2):161–167, 1969.
- Cesa-Bianchi, N. and Lugosi, G. Combinatorial bandits. In *Journal of Computer and System Sciences*, volume 78, pp. 1404–1422, 2012.
- Chen, W., Wang, Y., and Yuan, Y. Combinatorial Multi-Armed Bandit: General Framework and Applications. In *International Conference on Machine Learning*, pp. 151–159, 2013.
- Chen, W., Wang, Y., Yuan, Y., and Wang, Q. Combinatorial Multi-Armed Bandit and Its Extension to Probabilistically Triggered Arms. jul 2014.
- Combes, R., Shahi, M. S. T. M., Proutiere, A., and Others. Combinatorial bandits revisited. In *Advances in Neural Information Processing Systems*, pp. 2116–2124, 2015.
- Degenne, R. and Perchet, V. Combinatorial semi-bandit with known covariance. dec 2016.
- Ding, W., Qin, T., Zhang, X.-d., and Liu, T.-y. Multi-Armed Bandit with Budget Constraint and Variable Costs. In *Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence*, 2013.
- Edmonds, J. Matroids and the Greedy Algorithm. *Mathematical Programming*, 1(1):127–136, 1971.
- Edmonds, J. and Fulkerson, D. Transversals and matroid partition. *Journal of Research of the National Bureau of Standards Section B Mathematics and Mathematical Physics*, 69B(3):147–153, 1965.
- Filmus, Y. and Ward, J. A tight combinatorial algorithm for submodular maximization subject to a matroid constraint. In *Proceedings - Annual IEEE Symposium on Foundations of Computer Science, FOCS*, pp. 659–668, 2012.
- Fujishige, S. *Submodular functions and optimization*. Annals of discrete mathematics. 2005.
- Gai, Y., Krishnamachari, B., and Jain, R. Combinatorial network optimization with unknown variables: Multi-armed bandits with linear rewards and individual observations. *Transactions on Networking*, 20(5):1466–1478, 2012.
- He, S., Zhang, J., and Zhang, S. Polymatroid Optimization, Submodularity, and Joint Replenishment Games. *Operations Research*, 60(1):128–137, 2012.
- Jun, K.-S., Bhargava, A., Nowak, R., and Willett, R. Scalable Generalized Linear Bandits: Online Computation and Hashing. may 2017.
- Koolen, W., Warmuth, M. K., and Kivinen, J. Hedging structured concepts. In *COLT*, volume 118653, pp. 93–105, 2010.
- Krause, A. and Golovin, D. Submodular function maximization. In *Tractability*, volume 9781107025, pp. 71–104, 2011.
- Kveton, B., Wen, Z., Ashkan, A., Eydgahi, H., and Eriksson, B. Matroid bandits: Fast combinatorial optimization with learning. In *Uncertainty in Artificial Intelligence*, 2014.
- Kveton, B., Wen, Z., Ashkan, A., and Szepesvari, C. Tight regret bounds for stochastic combinatorial semi-bandits. In *International Conference on Artificial Intelligence and Statistics*, 2015.
- Lai, T. L. and Robbins, H. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1): 4–22, 1985.
- Lee, J., Mirrokni, V. S., Nagarajan, V., and Sviridenko, M. Maximizing Nonmonotone Submodular Functions under Matroid or Knapsack Constraints. *SIAM Journal on Discrete Mathematics*, 23(4):2053–2078, 2010.
- Li, L., Chu, W., Langford, J., and Schapire, R. E. A contextual-bandit approach to personalized news article recommendation. *International World Wide Web Conference*, 2010.

- Mohri, M. and Munoz, A. Optimal regret minimization in posted-price auctions with strategic buyers. In *Advances in Neural Information Processing Systems*, pp. 1871–1879, 2014.
- Nemhauser, G. L., Wolsey, L. A., and Fisher, M. L. An analysis of approximations for maximizing submodular set functions-I. *Mathematical Programming*, 14(1):265–294, 1978.
- Neu, G. and Bartók, G. An efficient algorithm for learning with semi-bandit feedback. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 8139 LNAI, pp. 234–248, 2013.
- Oliveira, C. A. and Pardalos, P. M. A survey of combinatorial optimization problems in multicast routing, 2005.
- Perfect, H. Applications of Menger’s graph theorem. *Journal of Mathematical Analysis and Applications*, 22(1): 96–111, 1968.
- Robbins, H. Some aspects of the sequential design of experiments. *Bulletin of the American Mathematics Society*, 58:527–535, 1952.
- Russo, D. and Roy, B. V. An information-theoretic analysis of thompson sampling. *Journal of Machine Learning Research*, 17(68):1–30, 2016.
- Russo, D. and Van Roy, B. Learning to Optimize Via Posterior Sampling. jan 2013.
- Sakaue, S., Ishihata, M., and Minato, S.-i. Efficient Bandit Combinatorial Optimization Algorithm with Zero-suppressed Binary Decision Diagrams. In *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84, pp. 585–594, 2018.
- Schrijver, A. *Combinatorial Optimization: Polyhedra and Efficiency*. 2008.
- Stobbe, P. and Krause, A. Efficient Minimization of Decomposable Submodular Functions. oct 2010.
- Sviridenko, M., Vondrák, J., and Ward, J. Optimal approximation for submodular and supermodular optimization with bounded curvature. nov 2013.
- Talebi, M. S. and Proutiere, A. An Optimal Algorithm for Stochastic Matroid Bandit Optimization. In *The 2016 International Conference on Autonomous Agents & Multiagent Systems*, pp. 548–556, 2016.
- Talebi, M. S., Zou, Z., Combes, R., Proutiere, A., and Johansson, M. Stochastic Online Shortest Path Routing: The Value of Feedback. sep 2013.
- Tran-Thanh, L., Stein, S., Rogers, A., and Jennings, N. R. Efficient crowdsourcing of unknown experts using bounded multi-armed bandits. *Artificial Intelligence*, 214: 89–111, 2014.
- Wang, S. and Chen, W. Thompson Sampling for Combinatorial Semi-Bandits. mar 2018.
- Whitney, H. On the abstract properties of linear dependence. *American Journal of Mathematics*, 57(3):509–533, 1935.
- Xia, Y., Ding, W., Zhang, X.-D., Yu, N., and Qin, T. Budgeted bandit problems with continuous random costs. In *Asian Conference on Machine Learning*, 2016a.
- Xia, Y., Qin, T., Ma, W., Yu, N., and Liu, T.-Y. Budgeted multi-armed bandits with multiple plays. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, pp. 2210–2216. AAAI Press, 2016b.