

# Combinatorial semi-bandit with known covariance

**Authored by:**

Romy Degenne  
Vianney Perchet

## **Abstract**

The combinatorial stochastic semi-bandit problem is an extension of the classical multi-armed bandit problem in which an algorithm pulls more than one arm at each stage and the rewards of all pulled arms are revealed. One difference with the single arm variant is that the dependency structure of the arms is crucial. Previous works on this setting either used a worst-case approach or imposed independence of the arms. We introduce a way to quantify the dependency structure of the problem and design an algorithm that adapts to it. The algorithm is based on linear regression and the analysis uses techniques from the linear bandit literature. By comparing its performance to a new lower bound, we prove that it is optimal, up to a poly-logarithmic factor in the number of arms pulled.

## **1 Paper Body**

and setting

The multi-armed bandit problem (MAB) is a sequential learning task in which an algorithm takes at each stage a decision (or, ?pulls an arm?). It then gets a reward from this choice, with the goal of maximizing the cumulative reward [Robbins, 1985]. We consider here its stochastic combinatorial extension, in which the algorithm chooses at each stage a subset of arms [Audibert et al., 2013, Cesa-Bianchi and Lugosi, 2012, Chen et al., 2013, Gai et al., 2012]. These arms could form, for example, the path from an origin to a destination in a network. In the combinatorial setting, contrary to the the classical MAB, the inter-dependencies between the arms can play a role (we consider that the distribution of rewards is invariant with time). We investigate here how the covariance structure of the arms affects the difficulty of the learning task and whether it is possible to design a unique algorithm capable of performing optimally in all cases from the simple scenario with independent rewards to the more challenging scenario of general correlated rewards. Formally, at each stage  $t \in \{1, \dots, T\}$ , an algorithm pulls  $m \geq 1$  arms among  $d \geq m$ . Such a set of  $m$

arms is called an "action" and will be denoted by  $A_t \in \{0, 1\}^d$ , a vector with exactly  $m$  non-zero entries. The possible actions are restricted to an arbitrary fixed subset  $\mathcal{A} \subset \{0, 1\}^d$ . After choosing an action  $A_t$ , the algorithm receives the reward  $A_t^\top X_t$ , where  $X_t \in \mathbb{R}^d$  is the vector encapsulating the reward of the  $d$  arms at stage  $t$ . The successive reward vectors  $(X_t)_{t=1}^T$  are i.i.d with unknown mean  $\mu \in \mathbb{R}^d$ . We consider a semi-bandit feedback system: after choosing the action  $A_t$ , the algorithm observes the reward of each of the arms in that action, but not the other rewards. Other possible feedbacks previously studied include bandit (only  $A_t^\top X_t$  is revealed) and full information ( $X_t$  is revealed). The goal of the algorithm is to maximize the cumulated reward up to stage  $T \geq 1$  or equivalently to minimize the expected regret, which is the difference of the reward that would have been gained by choosing the best action in hindsight  $A^*$  and what was actually gained:  $\text{ERT} = \mathbb{E} \sum_{t=1}^T (A^{*\top} X_t - A_t^\top X_t)$ .

30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

For an action  $A \in \mathcal{A}$ , the difference  $\Delta A = (A^\top X_t - A_t^\top X_t)$  is called gap of  $A$ . We denote by  $\Delta_t$  the PT gap of  $A_t$ , so that regret rewrites as  $\text{ERT} = \mathbb{E} \sum_{t=1}^T \Delta_t$ . We also define the minimal gap of an arm,  $\Delta_{i,\min} = \min_{A \in \mathcal{A}: A_i=1} \Delta A$ . This setting was already studied Cesa-Bianchi and Lugosi [2012], most recently in Combes et al. [2015], Kveton et al. [2015], where two different algorithms are used to tackle on one hand the case where the arms have independent rewards and on the other hand the general bounded case. The regret guaranties of the two algorithms are different and reflect that the independent case is easier. Another algorithm for the independent arms case based on Thompson Sampling was introduced in Komiyama et al. [2015]. One of the main objectives of this paper is to design a unique algorithm that can adapt to the covariance structure of the problem when prior information is available. The following notations will be used throughout the paper: given a matrix  $M$  (resp. vector  $v$ ), its  $(i, j)$ th (resp.  $i$ th) coefficient is denoted by  $M(i, j)$  (resp.  $v(i)$ ). For a matrix  $M$ , the diagonal matrix with same diagonal as  $M$  is denoted by  $\text{diag}(M)$ . We denote by  $\epsilon_t$  the noise in the reward, i.e.  $X_t := \mu + \epsilon_t$ . We consider a subgaussian setting, in which we suppose that there is a positive semi-definite matrix  $C$  such that for all  $t \geq 1$ ,  $\epsilon_t \in \mathbb{R}^d$ ,  $\mathbb{E}[\epsilon_t] = 0$ ,  $\mathbb{E}[\epsilon_t \epsilon_t^\top] \leq C$ .

This is equivalent to the usual setting for bandits where we suppose that the individual arms are (i) subgaussian. Indeed if we have such a matrix  $C$  then each  $\epsilon_t$  is  $C$  (ii) -subgaussian. And under a subgaussian arms assumption, such a matrix always exists. This setting encompasses the case of bounded rewards. We call  $C$  a subgaussian covariance matrix of the noise (see appendix A of the

supplementary material). A good knowledge of  $C$  can simplify the problem greatly, as we will show. In the case of 1-subgaussian independent rewards, in which  $C$  can be chosen diagonal, a known lower bound  $d$  on the regret appearing in Combes et al. [2015] is  $\sqrt{d} \log T$ , while Kveton et al. [2015] proves a  $d \sqrt{m} \log T$  lower bound in general. Our goal here is to investigate the spectrum of intermediate cases  $\rho$  between these two settings, from the uninformed general case to the independent case in which one has much information on the relations between the arm rewards. We characterize the difficulty of the problem as a function of the subgaussian covariance matrix  $C$ . We suppose that we know a positive semi-definite matrix  $\Sigma$  such that for all vectors  $v$  with positive coordinates,  $v \succeq C v \succeq \Sigma v$ , property that we denote by  $C \succeq \Sigma$ .  $\Sigma$  reflects the prior information available about the possible degree of independence of the arms. We will study algorithms that enjoy regret bounds as functions of  $\Sigma$ . The matrix  $\Sigma$  can be chosen such that all its coefficients are non-negative and verify for all  $i, j$ ,  $\Sigma(ij) \leq \Sigma(ii) \Sigma(jj)$ . From now on, we suppose that it is the case. In the following, we will use  $t$  such that  $\Sigma t = C 1/2 t$  and write for the reward:  $X_t = \Sigma + C 1/2 t$ .

## 2

### Lower bound

We first prove a lower bound on the regret of any algorithm, demonstrating the link between the subgaussian covariance matrix and the difficulty of the problem. It depends on the maximal off-diagonal  $(ij)$  correlation coefficient of the covariance matrix. This coefficient is  $\rho = \max\{\Sigma(i,j)/\sqrt{\Sigma(i,i)\Sigma(j,j)} \mid i \neq j\} \leq C$ . (ii)  $C \leq \rho$  The bound is valid for consistent algorithms [Lai and Robbins, 1985], for which the regret on any problem verifies  $ER_t = o(t^\alpha)$  as  $t \rightarrow +\infty$  for all  $\alpha > 0$ . Theorem 1. Suppose to simplify that  $d$  is a multiple of  $m$ . Then, for any  $\rho > 0$ , for any consistent algorithm, there is a problem with gaps  $\rho$ ,  $\rho$ -subgaussian arms and correlation coefficients smaller than  $\rho$  on  $[0, 1]$  on which the regret is such that  $\liminf_{t \rightarrow \infty} \frac{ER_t}{\sqrt{t}} \geq \rho$ .

$$ER_t \geq \frac{1}{2} (d - m) \rho (1 + \rho(m - 1)) \log t$$

This bound is a consequence of the classical result of Lai and Robbins [1985] for multi-armed bandits, applied to the problem of choosing one among  $d/m$  paths, each of which has  $m$  different successive edges (Figure 1). The rewards in the same path are correlated but the paths are independent. A complete proof can be found in appendix B.1 of the supplementary material. 2

Figure 1: Left: parallel paths problem. Right: regret of OLS-UCB as a function of  $m$  and  $\rho$  in the parallel paths problem with 5 paths (average over 1000 runs).

## 3

### OLS-UCB Algorithm and analysis

Faced with the combinatorial semi-bandit at stage  $t \geq 1$ , the observations from  $t \geq 1$  stages form as many linear equations and the goal of an algorithm is to choose the best action. To find the action with the highest mean, we estimate the mean of all arms. This can be viewed as a regression problem. The design of our algorithm stems from this observation and is inspired by linear regression in the fixed design setting, similarly to what was done in the stochastic linear

bandit literature [Rusmevichientong and Tsitsiklis, 2010, Filippi et al., 2010]. There are many estimators for linear regression and we focus on the one that is simple enough and adaptive: Ordinary Least Squares (OLS). 3.1

Fixed design linear regression and OLS-UCB algorithm

For an action  $A \in \mathcal{A}$ , let  $I_A$  be the diagonal matrix with a 1 at line  $i$  if  $A(i) = 1$  and 0 otherwise. For a matrix  $M$ , we also denote by  $MA$  the matrix  $I_A M$ . At stage  $t$ , if all actions  $A_1, \dots, A_t$  were independent of the rewards, we would have observed a set of linear equations  $I_{A_1} X_1 = I_{A_1} \theta + I_{A_1} \epsilon_1 \dots I_{A_t} X_t = I_{A_t} \theta + I_{A_t} \epsilon_t$  and we could use the OLS estimator to estimate  $\theta$ , which is unbiased and has a known subgaussian constant controlling its variance. This is however not true in our online setting since the successive actions are not independent. At stage  $t$ , we define (i)

$$\begin{aligned} n_t &= \sum_{s=1}^t \sum_{i \in \mathcal{A}_s} I_{\{i \in \mathcal{A}_s\}} \\ &= \sum_{s=1}^t \sum_{i \in \mathcal{A}_s} I_{\{i \in \mathcal{A}_s\}} I_{\{j \in \mathcal{A}_s\}} \text{ and } D_t = \sum_{s=1}^t \sum_{i \in \mathcal{A}_s} I_{\{i \in \mathcal{A}_s\}} \end{aligned}$$

where  $n_t$  is the number of times arm  $i$  has been pulled before stage  $t$  and  $D_t$  is a diagonal matrix of these numbers. The OLS estimator is, for an arm  $i \in \mathcal{A}$ , (i)

$$\begin{aligned} \hat{\theta}_t(i) &= \frac{1}{n_t} \sum_{s=1}^t \sum_{i \in \mathcal{A}_s} X_s(i) \\ &= \frac{1}{n_t} \sum_{s=1}^t \sum_{i \in \mathcal{A}_s} X_s(i) \\ &= \frac{1}{n_t} \sum_{s=1}^t \sum_{i \in \mathcal{A}_s} X_s(i) \end{aligned}$$

Then for all  $A \in \mathcal{A}$ ,  $A_t$  ( $\theta_t$ ) in the fixed design setting has a subgaussian matrix equal to  $P_t D_t^{-1} ( \sum_{s=1}^t C A_s ) D_t^{-1}$ . We get confidence intervals for the estimates and can use an upper confidence bound strategy [Lai and Robbins, 1985, Auer et al., 2002]. In the online learning setting the actions are not independent but we will show that using this estimator still leads to estimates that are well concentrated around  $\theta$ , with confidence intervals given by the same subgaussian matrix. The algorithm OLS-UCB (Algorithm 1) results from an application of an upper confidence bound strategy with this estimator. We

now turn to an analysis of the regret of OLS-UCB. At any stage  $t \geq 1$  of the algorithm, let  $\rho_{(ij)}(t) = \max\{\rho_{(i,j)}(A_t), \rho_{(j,i)}(A_t)\}$  (ii) be the maximal off-diagonal correlation coefficient of  $A_t$  and  $\rho_{(jj)}(t)$  let  $\rho = \max\{\rho_{(ij)}(t) \mid t \leq T\}$  be the maximum up to stage  $T$ .

**Algorithm 1** OLS-UCB. Require: Positive semi-definite matrix  $\Sigma$ , real parameter  $\gamma \in (0, 1]$ . 1: Choose actions such that each arm is pulled at least one time. 2: loop: at stage  $t$ , 3:  $A_t = \arg \max_{A \in \mathcal{A}} \sum_{s=1}^t \langle A, p_s \rangle - \frac{\gamma}{2} \sum_{s=1}^t \langle A, A_s \rangle^2$  with  $E_t(A) = \frac{1}{2} \sum_{s=1}^t \langle A, A_s \rangle^2$ . 4: Choose action  $A_t$ , observe  $I_{A_t}$ . 5: Update  $\Sigma_t, D_t$ . 6: end loop

**Theorem 2.** The OLS-UCB algorithm with parameter  $\gamma \in (0, 1]$  and  $f(t) = \log t + (m+2) \log \log t + m e^{2 \log(1+\gamma)}$  enjoys for all times  $T \geq 1$  the regret bound

$$2 \sum_{i=1}^m \log m E[RT] \leq 16f(T) \left( 5(\gamma + 1) \rho + 45 \sum_{i=1}^m \rho_{i,\min} 1.6 i \rho[d] \right)$$

$8dm^2 \max_i \{C_{(ii)}\} \gamma_{\max} + 4 \gamma_{\max}^2 \rho_{\min}$  where  $d_{\min}$  stands for the smallest positive integer bigger than or equal to  $x$ . In particular,  $d_{\min} = 1$ . This bound shows the transition between a general case with a

$dm \log T$

regime and an independent

$d \log m \log T$

$T$  case with a upper bound (we recall that the lower bound is of the order of  $d \log T$ ). The weight of each case is given by the maximum correlation parameter  $\rho$ . The parameter  $\gamma$  seems to be an artefact of the analysis and can in practice be taken very small or even equal to 0.

Figure 1 illustrates the regret of OLS-UCB on the parallel paths problem used to derive the lower bound. It shows a linear dependency in  $\rho$  and supports the hypothesis that the true upper bound matches the lower bound with a dependency in  $m$  and  $\rho$  of the form  $(1 + \rho(m-1))$ . **Corollary 1.** The OLS-UCB algorithm with matrix  $\Sigma$  and parameter  $\gamma \in (0, 1]$  has a regret bounded as

$\sum_{i=1}^m \log m t_{(ii)} + 45 \sum_{i=1}^m \rho_{i,\min} \cdot E[RT] \leq O(dT \log T \max\{\rho\} (5(\gamma + 1) \rho + 1.6 i \rho[d]))$  Proof. We write that the regret up to stage  $T$  is bounded by  $\sum_{i=1}^m t_{(ii)}$  for actions with gap smaller than some  $\gamma$  and bounded using theorem 2 for other actions (with  $\rho_{\min}$ ). Maximizing over  $\gamma$  then gives the result.

**Comparison with other algorithms**

Previous works supposed that the rewards of the individual arms are in  $[0, 1]$ , which gives them a property. Hence we suppose  $\rho_{(ij)} \in [0, 1]$ ,  $C_{(ii)} = 1/2$  for our comparison.

1/2-subgaussian

In the independent case, our algorithm is the same as ESCB-2 from Combes et al. [2015], up to the  $d \log T$  parameter  $\rho$ . That paper shows that ESCB-2 enjoys an  $O(d \log T)$  upper bound but our analysis tightens it to  $O(d \log T)$

$m \log T$

$\rho$

$\log T$  In the general (worst) case, Kveton et al. [2015] prove an  $O(dm \log T)$  upper bound (which is tight) using CombUCB1, a UCB based algorithm

introduced in Chen et al. [2013] which at stage  $t$  uses  $q \in P(i)$  the exploration term  $1.5 \log t \sqrt{A_i / nt}$ . Our exploration term always verifies  $E_t(A) \leq q \leq P(i) \sqrt{f(t) \sqrt{A_i / nt}}$  with  $f(t) \leq \log t$  (see section 3.6). Their exploration term is a worst-case confidence interval for the means. Their broader confidence intervals however have the desirable property that one can find the action that realizes the maximum index by solving a linear optimization problem, making their algorithm computationally efficient, quality that both ESCB and OLS-UCB are lacking.

4

None of the two former algorithms benefits from guaranties in the other regime. The regret of ESCB in the general possibly correlated case is unknown and the regret bound for CombUCB1 is not improved in the independent case. In contrast, OLS-UCB is adaptive in the sense that its performance gets better when more information is available on the independence of the arms. 3.3

Regret Decomposition (i)

Let  $H_{i,t} = \{ \mu_i - \mu_i^* \leq \sqrt{2 \log t / nt} \}$  and  $H_t = \bigcup_{i=1}^m H_{i,t}$ .  $H_t$  is the event that at least one coordinate of  $\mu$  is far from the true mean. Let  $G_t = \{ A_i^* \leq \mu_i^* + \sqrt{2 \log t / nt} \}$  be the event that the estimate of the optimal action is below its true mean by a big margin. We decompose the regret according to these events:

$$\begin{aligned} & \sum_{t=1}^T \sum_{i=1}^m \mathbb{I}\{G_t, H_t\} + \\ & \sum_{t=1}^T \sum_{i=1}^m \mathbb{I}\{G_t\} + \\ & \sum_{t=1}^T \sum_{i=1}^m \mathbb{I}\{H_t\} \end{aligned}$$

Events  $G_t$  and  $H_t$  are rare and lead to a finite regret (see below). We first simplify the regret due to  $G_t \cap H_t$  and show that it is bounded by the "variance" term of the algorithm. Lemma 1. With the algorithm choosing at stage  $t$  the action  $A_t = \arg \max_A (A_i \leq \mu_i^* + \sqrt{2 \log t / nt})$ , we have  $\sum_{t=1}^T \mathbb{I}\{G_t, H_t\} \leq 2 \sum_{t=1}^T \mathbb{I}\{ \mu_{A_t} \leq \mu_{A_t}^* + \sqrt{2 \log t / nt} \}$ . Proof in appendix B.2 of the supplementary material. Then the regret is cut into three terms,  $RT \leq 2$

$$\begin{aligned} & \sum_{t=1}^T \sum_{i=1}^m \mathbb{I}\{ \mu_{A_t} \leq \mu_{A_t}^* + \sqrt{2 \log t / nt} \} + \\ & \sum_{t=1}^T \sum_{i=1}^m \mathbb{I}\{G_t\} + \\ & \sum_{t=1}^T \sum_{i=1}^m \mathbb{I}\{H_t\} . \end{aligned}$$

The three terms will be bounded as follows: The  $H_t$  term leads to a finite regret from a simple application of Hoeffding's inequality. The  $G_t$  term leads

to a finite regret for a good choice of  $f(t)$ . This is where we need to show that the exploration term of the algorithm gives a high probability upper confidence bound of the reward. The  $\text{Et}(\text{At})$  term, or variance term, is the main source of the regret and is bounded using ideas similar to the ones used in existing works on semi-bandits. 3.4

Expected regret from Ht

PT Lemma 2. The expected regret due to the event Ht is  $E[\sum_{t=1}^T I\{H_t\}]$

8dm2 maxi {C (ii) }?max ?2min

The proof uses Hoeffding’s inequality on the arm mean estimates and can be found in appendix B.2 of the supplementary material. 35

Expected regret from Gt

We want to bound the probability that the estimated reward for the optimal action is far from its mean. We show that it is sufficient to control a self-normalized sum and do it using arguments from Pe?a et al. [2008], or Abbasi-Yadkori et al. [2011] who applied them to linear bandits. The analysis also involves a peeling argument, as was done in one dimension by Garivier [2013] to bound a similar quantity. e Lemma 3. Let  $\epsilon_t \geq 0$ . With  $f(\epsilon_t) = \log(1/\epsilon_t) + m \log \log t + m^2 \log(1 + \epsilon_t)$  and an algorithm  $q \in \mathcal{P}^{t-1}$  given by the exploration term  $E_t(A) = 2f(\epsilon_t) A_i / (D_t^{1/2} \sum_{s=1}^t A_s) D_t^{1/2} A$ , then the event  $G_t = \{A_i \geq \epsilon_t A_i \text{ for } A_i \in E_t(A)\}$  verifies  $P\{G_t\} \geq 1 - \epsilon_t$ . With  $\epsilon_1 = 1$  and  $\epsilon_t = t^{-\log 2 / t}$  for  $t \geq 2$ , such that  $f(\epsilon_t) = f(t)$ , the regret due to  $G_t$  is finite in expectation, bounded by  $4 \max_i \mu_i$ .

Proof. We use a peeling argument: let  $\epsilon > 0$  and for  $\mathbf{a} = (a_1, \dots, a_m) \in \mathbb{N}^m$ , let  $D_{\mathbf{a}} \subseteq [T]$  be a  $(1-\epsilon)$  subset of indices defined by  $(t \in D_{\mathbf{a}} \iff A_t(\mathbf{a}) \geq (1-\epsilon)a_i \text{ for all } i \in [m])$ . For any  $B \subseteq [T]$ ,  $R_t$

$X^?$

$\text{PA} \vdash (\exists x \exists y (x \neq y \wedge \neg \text{Bt } x \wedge \text{Bt } y)) \rightarrow \text{Da} . a$

The number of possible sets  $D_a$  for  $t$  is bounded by  $(\log t / \log(1 + \epsilon))m$ , since each number of pulls

(i) nt for i ? A? is bounded by t. We now search a bound of the form P A?i (? ? ? ?t ) ? Bt —t ? Da . Suppose t ? Da and let D be a positive definite diagonal matrix (that depends on a). Pt?1 Pt?1 2 Let St = s=1 IAs ?A? C 1/2 s , Vt = s=1 CAs ?A? and IVt +D ( ) = 12 kSt k(Vt +D)?1 . Lemma 4. Let ?t 0 and let f?(?t ) be a function of ?t . With a choice of D such that IA? D ?IA? ?C Dt for all t in Da ,

q

no

$P(A_i | \{x_j\}_{j \neq i}) = \frac{1}{Z} \exp(-\sum_{j \neq i} \lambda_{ij} x_j)$

than 1. To control  $\mathbb{E}V_t(\cdot)$ , we are however interested in the expectation of this maximum and cannot interchange  $\max$  and  $\mathbb{E}$ . The method of mixtures circumvents this difficulty: it provides an approximation of the maximum by integrating the exponential against a multivariate normal centered at the point  $V_t \cdot 1 \cdot S_t$ , where the maximum is attained. The integrals over  $u$  and can then be swapped by Fubini's theorem to get an approximation of the expectation of the maximum using an integral of the expectations. Doing so leads to the following lemma, extracted from the proof of Theorem 1 of Abbasi-Yadkori et al. [2011].

Lemma 5.4 Let  $D$  be a positive definite matrix that does not depend on  $t$  and  $\det D \leq M_t(D) = \det(V \exp(IV_t + D))$ . Then  $\mathbb{E}[M_t(D)] \leq 1 + \frac{1}{n} \log t + D$ . We rewrite  $P(IV_t + D) \leq f_t(t)$  to introduce  $M_t(D)$ , ( $n \circ P(IV_t + D) \leq f_t(t) - tDa = P(M_t(D)) \leq p$

$$\frac{1}{\det(\text{Id} + D)^{1/2}} \int \exp(f_t(t)) dt Da$$

The peeling lets us bound  $V_t$ . Let  $Da$  be the diagonal matrix with entry  $(i, i)$  equal to  $(1 + \frac{1}{n})a_i$  for  $i \in A_t$  and 0 elsewhere. Lemma 6. With  $D = \frac{1}{n}C Da + I[d]A_t$ ,  $\det(\text{Id} + D)^{1/2} \int \exp(f_t(t)) dt Da \leq (1 + \frac{1}{n} \log t + \frac{1}{n} \log m)$

The union bound on the sets  $Da$  and Markov's inequality give

$$\mathbb{P} \left( \bigcup_i \left( \frac{1}{n} \log t + \frac{1}{n} \log m \right) \exp(f_t(t)) - t Da \leq P(M_t(D)) \leq (1 + \frac{1}{n} \log t + \frac{1}{n} \log m) \exp(f_t(t)) - t Da \right) \leq \frac{1}{n} \log t + \frac{1}{n} \log m + \exp(-\frac{1}{n} \log t) \log(1 + \frac{1}{n})$$

For  $\frac{1}{n} = e^{-1}$  and  $f_t(t)$  as in lemma 3, this is bounded by  $\frac{1}{n} \log t$ . The result with  $\frac{1}{n}$  instead of  $C$  is a consequence of  $C + \frac{1}{n}$ . With  $\frac{1}{n} = 1$  and  $\frac{1}{n} \log t = 1/(t \log^2 t)$  for  $t \geq 2$ , the regret due to  $G_t$  is  $\mathbb{E}[\sum_{t=1}^T X_t] \leq \frac{1}{n} \log t + \frac{1}{n} \log m + \frac{1}{n} \log t \log^2 t$

$$\sum_{t=1}^T X_t \leq \frac{1}{n} \log t + \frac{1}{n} \log m + \frac{1}{n} \log t \log^2 t$$

Bounding the variance term

The goal of this section is to bound  $\mathbb{E}t(A_t)$  under the event  $\{\frac{1}{n} \log t \leq \mathbb{E}t(A_t)\}$ . Let  $\frac{1}{n} \in [0, 1]$  such that for  $q$  all  $i, j \in A_t$  with  $i \neq j$ ,  $\frac{1}{n}(ij) \leq \frac{1}{n}(ii) \frac{1}{n}(jj)$ . From the Cauchy-Schwartz inequality,  $(ij)$

$$(i)(j) \leq \frac{1}{n}$$

Using these two inequalities,  $\sum_{t=1}^T X_t$

$$\sum_{t=1}^T A_t \leq \sum_{t=1}^T \left( \frac{1}{n} \log t + \frac{1}{n} \log m + \frac{1}{n} \log t \log^2 t \right)$$



$s=1$   
 $i,j?At$   
 $(i) (j) nt nt$   
 $? (1 ? ?t )$   
 $X ?(ii) i?At$   
 $(i) nt$   
 $s X$   
 $?(ii)$   
 $i?At$   
 $(i) nt$   
 $+ ?t ($   
 $)2 .$

We recognize here the forms of the indexes used in Combes et al. [2015] for independent arms (left term) and Kveton et al. [2015] for general arms (right term). Using  $?t ? Et (At )$  we get  $s X ?(ii) X ?(ii) ?2t ? ( ? + 1 ? ?t ) + ?t ( )2 . (1) (i) (i) 8f (t) nt nt i?At$

$i?At$

The strategy from here is to find events that must happen when (1) holds and to show that these events cannot happen very often. For positive integers  $j$  and  $t$  and for  $e ? \{1, 2\}$ , we define the set of arms (ii)

(i)

$j$  in  $At$  that were pulled less than a given threshold:  $St,e = \{i ? At , nt ? ?j,e 8f (t)?$

$ge (m,?t ) \}$ ,  $?2t j At . (St,e )j?0$

0 with  $ge (m, ?t )$  to be stated later and  $(?i,e )i?1$  a decreasing sequence. Let also  $St,e =$  is decreasing for the inclusion of sets and we impose  $\lim j?+? ?j,e = 0$ , such that there is an index  $j? j?$  with  $St,e = ?$ . We introduce another positive sequence  $(?j,e )j?0$  and consider the events that  $j$  at least  $m?j,e$  arms in  $At$  are in the set  $St,e$  and that the same is false for  $k \text{ } i j$ , i.e. for  $t ? 1, j j k$   $At,e = \{ -St,e - ? m?j,e ; ?k \text{ } j, -St,e - i m?k,e \}$ . To avoid having some of these events being 0 impossible we choose  $(?j,e )j?0$  decreasing. We also impose  $?0,e = 1$ , such that  $-St,e - = m?0,e . j$  Let then  $At,e = ?+? j=1 At,e$  and  $At = At,1 ? At,2$ . We will show that  $At$  must happen for (1) to be true. First, remark that under a condition on  $(?j,e )j?0$ ,  $At$  is a finite union of events, 0 Lemma 7. For  $e ? \{1, 2\}$ , if there exists  $j0,e$  such that  $?j0,e ,e ? 1/m$ , then  $At,e = ?jj=1 Ajt,e$ . We now show that  $At$  is impossible by proving a contradiction in (1).

Lemma 8. Under the event  $At,1$ , if there exists  $j0$  such that  $?j0 ,1 ? 1/m$ , then  $? ? j0 X ?(ii) X m?2t ? ? ? ? j?1,1 j,1 j ,1 ? i + 0 ?$ . (i)  $8f (t)g1 (m, ?t ) j=1 ?j,1 ?j0 ,1 nt i?At$

$P+? ?$  Under the event  $At,2$ , if  $\lim j?+? ?j,2 / ?j,2 = 0$  and  $j=1 s X$

$?(ii)$

$i?At$

$nt$

(i)

$?p$

$\exists j_{1,2} \exists j_{2,2} \exists j_{j,2}$

exists, then

$\forall t \sum_{j=1}^m X_{j_{1,2},2} \exists j_{j,2} \cdot \sum_{j=1}^m X_{j_{2,2},2} f(t) g_2(m, t) \leq m^2$

A proof can be found in appendix B.2 of the supplementary material. To ensure that the conditions of these lemmas are fulfilled, we impose that  $(i_{1,1})_{i \geq 0}$  and  $(i_{2,1})_{i \geq 0}$  have limit 0 and  $\lim_{j \rightarrow \infty} j_{j,2} / j_{j,2} = 0$ . Let  $j_{0,1}$  be the smallest integer such that  $j_{0,1} \geq 1/m$ . Let  $P_{j_{0,1}} \exists j_{1,1} \exists j_{2,1} P_{j_{1,1}} \exists j_{2,1} \exists j_{j,2} \exists l_1 = j_{0,1} + j = 1$  and  $l_2 = j = 1 \exists j_{1,2}$ . Using the two last lemmas with (1),  $\exists j_{1,1} \exists j_{2,1} \exists j_{j,2}$

we get that if  $A_t$  is true,  $\sum_{j=1}^m \sum_{t=1}^T f(t) g_2(m, t) \leq m^2$

$(\sum_{j=1}^m + 1) \sum_{t=1}^T$

$m l_1 m^2 l_2 + \sum_{t=1}^T g_1(m, t) g_2(m, t)$

Taking  $g_1(m, t) = 2(\sum_{j=1}^m + 1) \sum_{t=1}^T m l_1$  and  $g_2(m, t) = 2 \sum_{t=1}^T m^2 l_2$ , we get a contradiction. Hence with these choices  $A_t$  must happen. The regret bound will be obtained by a union bound on the events that form  $A_t$ . First suppose that all gaps are equal to the same  $\gamma$ . 7

Lemma 9. Let  $\gamma$  (ii)

$d_{j,e} f(T) \max_i \{ \sum_{t=1}^T m_{j,e}^2$

$=$

$\max_{t \geq 1} \sum_{t=1}^T$

$\}_{ge(m, \gamma)}$

For  $j$

$N_{\gamma}$ , the event  $A_{j,t,e}$  happens at most

$\gamma$

times. (i)

Proof. Each time that  $A_{j,t,e}$  happens, the counter of plays  $n_t$  of at least  $m_{j,e}$  arms is incremented. (ii)

After

$\sum_{j,e} f(T) \max_i \{ \sum_{t=1}^T$

$\}_{ge(m, \gamma)}$

(i)

increments, an arm cannot verify the condition on  $n_t$  any more.

There are  $d$  arms, so the event can happen at most  $d \sum_{j,e} 1$

$e$

$\sum_{j,e} f(T) \max_i \{ \sum_{t=1}^T \}_{ge(m, \gamma)} \sum_{t=1}^T$

times.

If all gaps are equal to  $\gamma$ , an union bound for  $A_t$  gives  $\sum_{j=1}^m \sum_{t=1}^T \sum_{e=1}^E X_{j,1} X_{j,2} \sum_{t=1}^T f(T) \gamma \cdot E[ d (\sum_{j=1}^m + 1) \sum_{t=1}^T m l_2 + \sum_{t=1}^T m l_2^2 \sum_{i=1}^I \{ H_t - G_t \} ] \leq 16 \max\{ \sum_{j=1}^m \}_{ge(m, \gamma)} \sum_{t=1}^T \sum_{i=1}^I \{ d \} \sum_{j=1}^m \sum_{t=1}^T \sum_{e=1}^E j_{j,2}$  The general case requires more involved manipulations but the result is similar and no new important idea is used. The following lemma is proved in appendix B.2 of the supplementary material: Lemma 10. Let  $\gamma$  (i)  $= \max\{ t, i \mid A_t \} \sum_{t=1}^T$ . The regret from the event  $H_t - G_t$  is such that  $\sum_{j=1}^m \sum_{t=1}^T \sum_{e=1}^E X_{j,1} X_{j,2} \sum_{t=1}^T \sum_{i=1}^I \{ (\sum_{j=1}^m + 1) \sum_{t=1}^T m l_2 \sum_{t=1}^T I \{ H_t - G_t \} \} \leq 16 f(T) E[ \cdot + \sum_{t=1}^T m l_2^2 \sum_{t=1}^T \sum_{i=1}^I \{ \min_{j=1}^m \sum_{t=1}^T \sum_{e=1}^E j_{j,2} \} \sum_{t=1}^T \sum_{i=1}^I \{ d \} ]$

Finally we can find sequences  $(j_1)_{j \geq 1}$ ,  $(j_2)_{j \geq 1}$ ,  $(j_1)_{j \geq 0}$  and  $(j_2)_{j \geq 0}$  such that !

$$2 \sum_{t=1}^T \sum_{i=1}^m \log m \cdot (i) \cdot 5(\epsilon + 1 \cdot \epsilon) + 45 \epsilon \cdot m \cdot \mathbb{I}\{H_t \neq G_t\} \leq 16f(T) \cdot \mathbb{E}[\sum_{i=1}^m \min_{t=1}^{1.6} i \cdot d]$$

See appendix C of the supplementary material. In Combes et al. [2015],  $\epsilon_{i,1}$  and  $\epsilon_{i,1}$  were such  $\epsilon$  that the  $\log 2 \cdot m$  term was replaced by  $m$ . Our choice is also applicable to their ESCB algorithm. Our use of geometric sequences is only optimal among sequences such that  $\epsilon_{i,1} = \epsilon_{i,1}$  for all  $i \geq 1$ . It is unknown to us if one can do better. With this control of the variance term, we finally proved Theorem 2.

4

#### Conclusion

We defined a continuum of settings from the general to the independent arms cases which is suitable for the analysis of semi-bandit algorithms. We exhibited a lower bound scaling with a parameter that quantifies the particular setting in this continuum and proposed an algorithm inspired from linear regression with an upper bound that matches the lower bound up to a  $\log 2 \cdot m$  term. Finally we showed how to use tools from the linear bandits literature to analyse algorithms for the combinatorial bandit case that are based on linear regression. It would be interesting to estimate the subgaussian covariance matrix online to attain good regret bounds without prior knowledge. Also, our algorithm is not computationally efficient since it requires the computation of an argmax over the actions at each stage. It may be possible to compute this argmax less often and still keep the regret guaranty, as was done in Abbasi-Yadkori et al. [2011] and Combes et al. [2015]. On a broader scope, the inspiration from linear regression could lead to algorithms using different estimators, adapted to the structure of the problem. For example, the weighted least-square estimator is also unbiased and has smaller variance than OLS. Or one could take advantage of a sparse covariance matrix by using sparse estimators, as was done in the linear bandit case in Carpentier and Munos [2012].

**Acknowledgements** The authors would like to acknowledge funding from the ANR under grant number ANR-13-JS010004 as well as the Fondation Mathématiques Jacques Hadamard and EDF through the Program Gaspard Monge for Optimization and the Irsdi project Tecolere. 8

## 2 References

Yasin Abbasi-Yadkori, David Pal, and Csaba Szepesvari. Improved Algorithms for Linear Stochastic Bandits. *Neural Information Processing Systems*, pages 1719, 2011. Jean-Yves Audibert, Sébastien Bubeck, and Gábor Lugosi. Regret in online combinatorial optimization. *Mathematics of Operations Research*, 39(1):31745, 2013. Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235256, 2002. Alexandra Carpentier and Rémi Munos. Bandit Theory meets Compressed Sensing for high dimensional Stochastic Linear Bandit.

Advances in Neural Information Processing Systems (NIPS), pages 251?259, 2012. Nicolo Cesa-Bianchi and G?bor Lugosi. Combinatorial bandits. *Journal of Computer and System Sciences*, 78 (5):1404?1422, 2012. Wei Chen, Yajun Wang, and Yang Yuan. Combinatorial multi-armed bandit: General framework and applications. *Proceedings of the 30th International Conference on Machine Learning (ICML)*, pages 151?159, 2013. Richard Combes, M. Sadegh Talebi, Alexandre Proutiere, and Marc Lelarge. Combinatorial Bandits Revisited. *Neural Information Processing Systems*, pages 1?9, 2015. Sarah Filippi, Olivier Capp?, Aur?lien Garivier, and Csaba Szepesv?ri. Parametric Bandits: The Generalized Linear Case. *Neural Information Processing Systems*, pages 1?9, 2010. Yi Gai, Bhaskar Krishnamachari, and Rahul Jain. Combinatorial network optimization with unknown variables: Multi-armed bandits with linear rewards and individual observations. *IEEE/ACM Transactions on Networking*, 20(5):1466?1478, 2012. Aur?lien Garivier. Informational confidence bounds for self-normalized averages and applications. *2013 IEEE Information Theory Workshop, ITW 2013*, 2013. Junpei Komiyama, Junya Honda, and Hiroshi Nakagawa. Optimal Regret Analysis of Thompson Sampling in Stochastic Multi-armed Bandit Problem with Multiple Plays. *Proceedings of the 32nd International Conference on Machine Learning*, 2015. Branislav Kveton, Zheng Wen, Azin Ashkan, and Csaba Szepesvari. Tight regret bounds for stochastic combinatorial semi-bandits. *Proceedings of the 18th International Conference on Artificial Intelligence and Statistics*, 2015. Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4?22, 1985. Victor H Pe?a, Tze Leung Lai, and Qi-Man Shao. Self-normalized processes: Limit theory and Statistical Applications. Springer Science & Business Media, 2008. Herbert Robbins. Some aspects of the sequential design of experiments. In *Herbert Robbins Selected Papers*, pages 169?177. Springer, 1985. Paat Rusmevichientong and John N. Tsitsiklis. Linearly Parameterized Bandits. *Mathematics of Operations Research*, (1985):1?40, 2010.