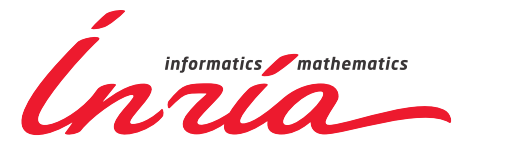


# IMPROVED LEARNING COMPLEXITY IN COMBINATORIAL PURE EXPLORATION BANDITS

VICTOR GABILLON, ALESSANDRO LAZARIC, MOHAMMAD GHAVAMZADEH, RONALD ORTNER & PETER BARTLETT



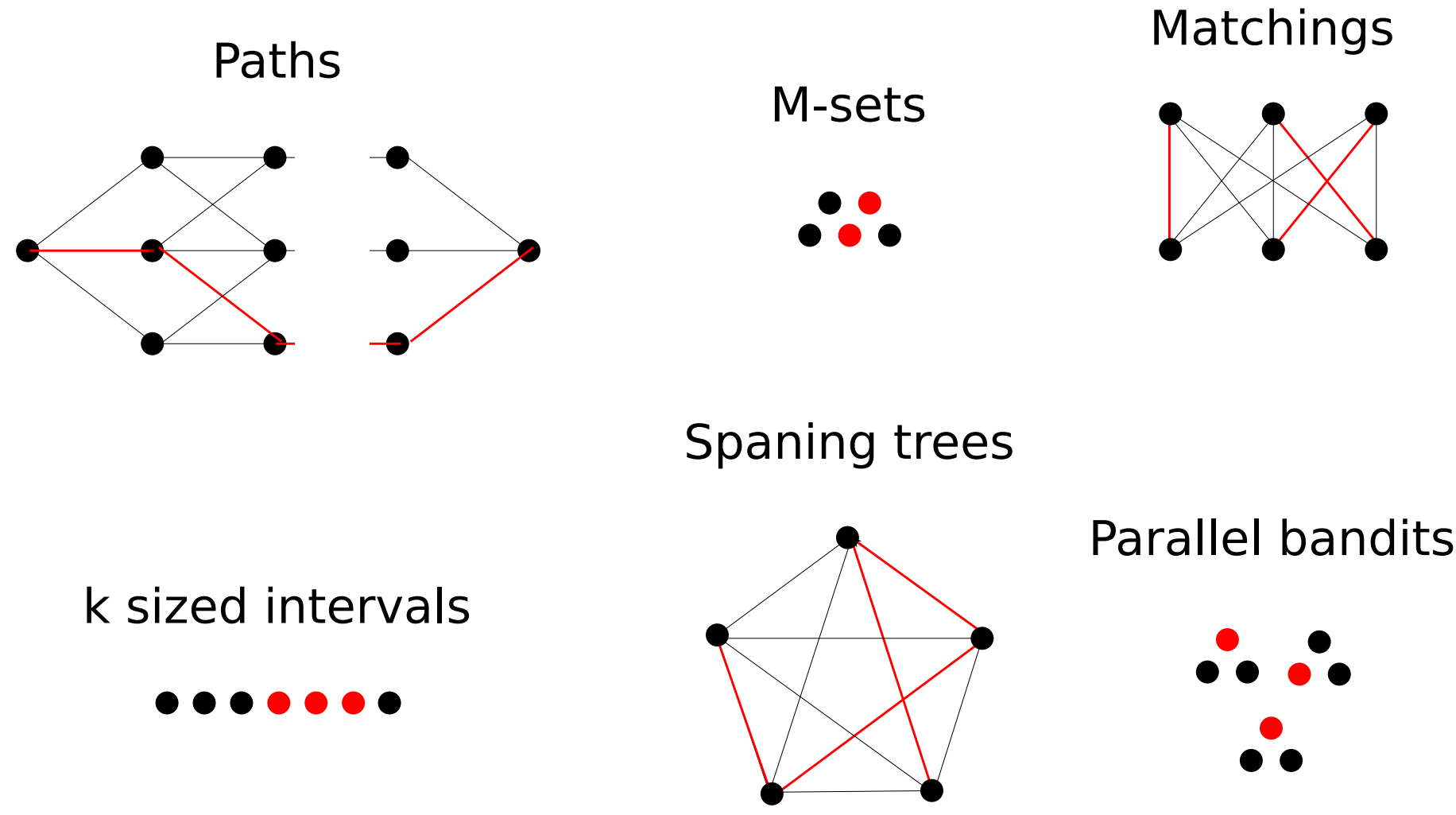
## 1) ABSTRACT

- We construct a new measure of **complexity** that provably characterizes the learning performance of the algorithms we propose for the fixed confidence and fixed budget settings.
- We show that this complexity is never bigger than the one in the existing work and illustrate a number of configurations in which it can be significantly smaller.
- While in general this improvement comes at the cost of increased computation, we provide a series of examples, including a planning problem, where the extra computation is not significant.

## 2) PROBLEM FORMULATION



- $\mathcal{K}$  a set of  $K$  arms with initially unknown distributions  $\nu_i \in [0, 1]$  and an expected value  $\mu_i$  that the learner aims to learn.
- The (combinatorial) decision space  $\mathcal{C} \subseteq 2^K$  contains decision sets (sets of arms)  $U \subseteq \mathcal{K}$



- Gap:  $\Delta_{U,V} = \mu_U - \mu_V$ , and  $U^* = \arg \max_{U \in \mathcal{C}} \mu_U$
- $U \oplus V = (U \setminus V) \cup (V \setminus U)$  the set of arms either in  $U$  or in  $V$ , but not in both.
- (A)symmetric  $\bar{d}_{U,V} = |U \oplus V|$  ( $d_{U,V} = |U \setminus V|$ )  
 $\mu_1, \dots, \mu_m$   
  
 $\bar{d}_{U,V} = m + 1$  and  $d_{V,U} = m$  and  $d_{U,V} = 1$

**The fixed budget setting:** Given a budget  $n$ , the learner minimizes the probability  $\tilde{\delta}$  of not identifying the best decision set, i.e.,  $\tilde{\delta} = \mathbb{P}[\hat{U}^*(n) \neq U^*]$ . See paper for the fixed confidence setting.

## 4) FIXED BUDGET ALGORITHM

The algorithm is an extension of Audibert et. al. [1]. It is composed of  $K$  phases. At the end of each phase, we compute the empirical  $\hat{\mu}_i(k)$ ,  $\hat{\Delta}_{U,V}(k)$ ,  $\hat{U}^*(k)$ ,  $\hat{G}_{U,V}(k) = \hat{\Delta}_{U,V}(k)/\bar{d}_{U,V}$ ,  $\hat{C}_U(k) = \arg \max_{V \in \mathcal{C}: \hat{\mu}_V(k) > \hat{\mu}_U(k)} \hat{G}_{V,U}(k)$ ,  $\hat{G}_i(k) = \min_{U \in \mathcal{C}: i \in U \oplus \hat{C}_U(k)} \hat{G}_{\hat{C}_U(k), U}(k)$ . We set  $n_k = \left\lceil \frac{n-K}{\log(K)(K+1-k)} \right\rceil$ ,  $k \in \mathcal{K}$ .

**Parameters:** number of rounds  $n$ , set of arms  $\mathcal{K}$ , decision set  $\mathcal{C}$  and cumulative pulls schedule  $n_0, n_1, \dots, n_K$ .

Let  $\mathcal{K}_1 = \mathcal{K}$ ,  $k = 1$

**while**  $|\mathcal{K}_k| \geq 1$  **do**

    Pull each arm  $i \in \mathcal{K}_k$  for  $n_k - n_{k-1}$  rounds.

    Compute  $\hat{U}^*(k) = \arg \max_{U \in \mathcal{C}} \hat{\mu}_U(k)$ .

    Find  $j_k = \max_{i \in \mathcal{K}_k} \hat{G}_i(k)$ .

    Deactivate arms  $j_k$ , i.e., set  $\mathcal{K}_{k+1} = \mathcal{K}_k \setminus j_k$ .

$k \leftarrow k + 1$

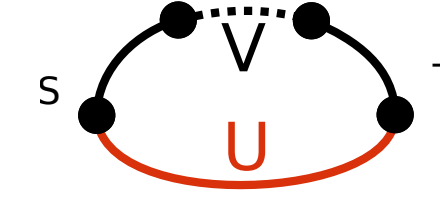
**end while**

Return  $J_n = \arg \max_{U \in \mathcal{C}} \hat{\mu}_U(n)$

## 3) SAMPLE COMPLEXITIES

- **Discriminate two decision sets  $U$  and  $V$**   
Simply uniformly allocate pulls over the arms in  $U \oplus V$  and stop at the first step  $t$  when the lower-bound on the gap is positive, i.e.,

$$\hat{\Delta}_{V,U}(t) - \sum_{i \in U \oplus V} \beta_i(t) > 0.$$

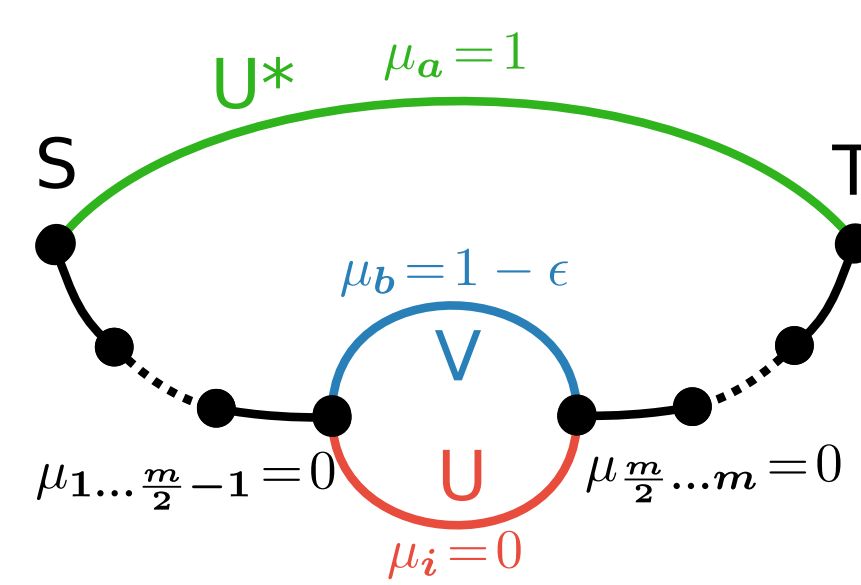


This stopping time suggests the idea to define the **sample complexity** for  $U, V$  as

$$H_{U,V} = \frac{\bar{d}_{U,V}^2}{\Delta_{U,V}^2} \quad \begin{array}{l} \bullet \text{ inverse dependency on } \Delta_{U,V} \\ \bullet \bar{d}_{U,V} \text{ is like a variance term} \end{array}$$

- $C_U$ : the complement of  $U$

When trying to discard a suboptimal set  $U$  from  $\mathcal{C}$  it may be easier to compare it to a set  $V \neq U^*$ ,  $\mu_V > \mu_U$  with a smaller complexity  $H_{U,V}$ .



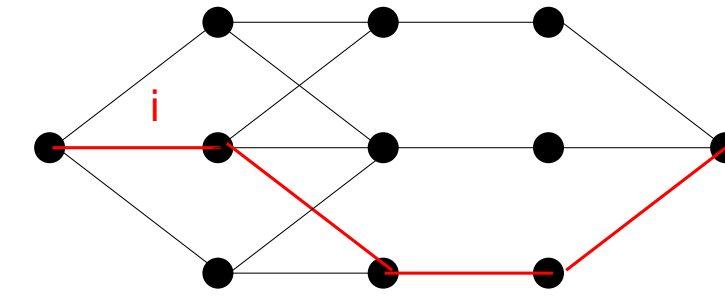
For any decision  $U \neq U^*$ , its complement is

$$C_U = \arg \min_{V \in \mathcal{C}: \mu_V > \mu_U} H_{U,V}.$$

For any arm  $i \in \mathcal{K}$  different sets  $U$  may force different requirements on the number of times  $i$  should be pulled.

Thus the **complexity of an arm  $i \in \mathcal{K}$**  is

$$H_i = \max_{U \in \mathcal{C}: i \in U \oplus C_U} H_{U,C_U}.$$



The **global complexity** is  $H = \sum_{i \in \mathcal{K}} H_i$ .

The **simplicity** is  $G_{U,V} = \frac{\Delta_{U,V}}{\bar{d}_{U,V}}$ .

## 7) DISCUSSION

- Is our sample complexity optimal? Solution of an optimization problem?
- Asymmetric distance  $d_{U,V}$  instead of  $\bar{d}_{U,V}$ ?  $\sigma_i^2/\Delta_i^2$  instead of the standard  $(\sigma_i^2 + \sigma_{i^*}^2)/\Delta_i^2$  in the single bandit fixed budget setting?

$$\mathbb{P}[\hat{U}^*(n) \neq U^*] \leq 2K^2 \exp\left(-\frac{n-K}{32\log(K)\bar{H}}\right),$$

with  $\bar{H} = \max_{i \in \mathcal{K}} iH_i$ ,  $H_{(1)} \geq H_{(2)} \geq \dots \geq H_{(K)}$ . As noted in [1], it holds that  $\bar{H} \leq H \leq \bar{H}\log(K)$ .

*Sketch of the proof.* By **double** induction through the phases of the algorithm:

- (i) The arms that are rejected are well classified as belonging or not to  $U^*$
- (ii) The deactivated arm has been sufficiently sampled so that if will be used later as a "good" reference arm (as part of some  $C_U$ ).

## 5) COMPARISON WITH CHEN ET AL.

Chen et al. (2014) [2] proposed a different measure of sample complexity of arm  $i \in \mathcal{K}$ :

$$\Delta_i^\odot = \begin{cases} \mu^* - \max_{U \in \mathcal{C}: i \in U} \mu_U & \text{if } i \notin U^* \\ \mu^* - \max_{U \in \mathcal{C}: i \notin U} \mu_U & \text{if } i \in U^* \end{cases}$$

The complexity is  $H_i^\odot = \frac{\text{width}(\mathcal{C})^2}{(\Delta_i^\odot)^2}$ .

- Our complexity uses the **complement  $C_U$**  instead of  $U^*$ .
- $\text{width}(\mathcal{C})$  defines a **distance  $\bar{d}$**  directly for  $\mathcal{C}$  while we define it for every pair  $U, V$ :

$$\begin{array}{l} H_i^\odot = \frac{\text{width}(\mathcal{C})^2}{(\Delta_{U^*,V}^\odot)^2} = \frac{m^2}{1/m^2} = m^4 \\ H_i = \frac{\bar{d}_{U^*,V}^2}{\Delta_{U^*,V}^2} = \frac{2^2}{(1/m)^2} = 4m^2 \leq H_i^\odot \\ \text{width}(\mathcal{C}) \approx m \end{array}$$

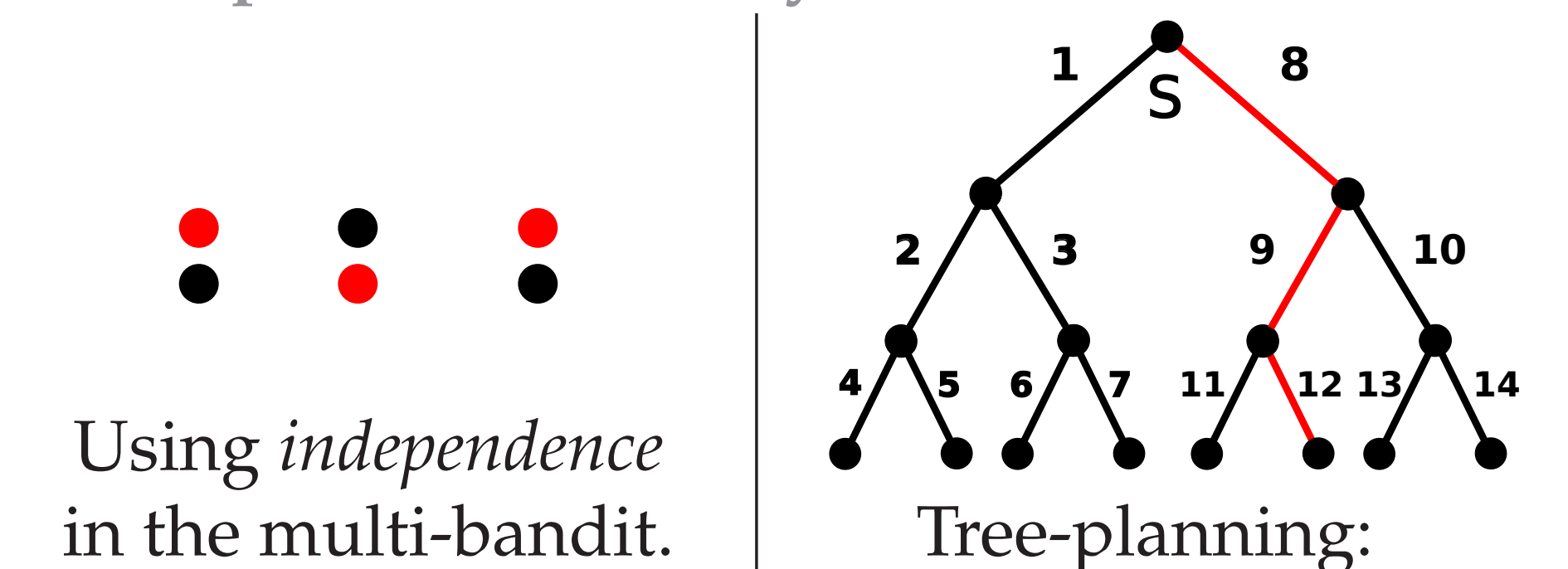
For all  $i \in \mathcal{K}$ ,  $H_i^\odot \geq H_i$

*Sketch of Proof.* In the first part we show that the definition  $H_i$  is indeed better when comparing sets  $U$  to their complements  $C_U$  instead of  $U^*$ . Consider the drawing on the left. There are three decision sets  $U, V$  and  $U^*$  containing the three edges/arms  $i, b$  and  $a$ .  $V = C_U, U^* = C_V$  and  $b \in V$ . As  $b \in V \oplus C_V$ , by definition the complexity of  $b$  should be at least  $H_{V,C_V}$ . As  $b \in U \oplus C_U$ , by definition the complexity of  $b$  should be at least  $H_{U,C_U}$ . This last requirement is an extra requirement compared to the complexity defined by Chen et al. Fortunately, it is shown that  $H_{U,C_U} \leq H_{V,C_V}$ .

## 6) COMPUTATIONAL COMPLEXITY

Comparing two by two every solution has a worst case *computational complexity* that can be exponential in  $K$ .

Examples of tractability:



**Tree-planning case:**  $K \approx |\mathcal{C}|$

- as maximizing the expected sum of rewards over  $m$  consecutive actions from a starting state when the state dynamics is deterministic and the reward distributions are unknown.
- the set of decisions (paths)  $V \neq U$  can be clustered into  $m$  groups of sets depending on which is the first node where they differ from  $U$  among the  $m$  possible ones.

Since the distance  $d_{U,V}$  is constant within these clusters, identifying the complement  $\arg \min H_{U,V}$  within each cluster corresponds to finding the set  $V$  with the largest value.

We can also show a planning tree example where the learning complexity is significantly improved over [2].

[1] J.-Y. Audibert, S. Bubeck, and R. Munos. Best Arm Identification in Multi-Armed Bandits. In COLT, 2010.

[2] S. Chen, T. Lin, I. King, M. Lyu, and W. Chen. Combinatorial pure exploration of multi-armed bandits. In NIPS, 2014.