

## Clustering with Bandit feedback Problem (CBP)

### Bandit learning protocol

Consider a **multi-armed bandit** with  $N$  arms. Each arm  $a \in \{1, \dots, N\}$  is associated with a multidimensional mean-vector  $\mu_a \in \mathbb{R}^d$  (with  $d$  possibly large).

For each time step  $t = 1, \dots, T$ ,

- chooses an arm  $A_t \in \{1, \dots, N\}$  (based on the passed observations)
- receives  $X_t$  with mean  $\mu_{A_t}$  and  $\sigma$ -subGaussian noise (e.g.,  $X_t \sim \mathcal{N}(\mu_a, \sigma^2 I_d)$ )

### Hidden partition assumption

We assume that there exists a **hidden partition**  $G^*$  of  $[N]$  into exactly  $K$  non-empty groups, such that **all arms in the group  $G_k^*$  share the same mean-vector  $\Lambda(k)$** .

### Objective: clustering in the PAC-setting

Given a prescribed probability  $\delta \in (0, 1)$ , the objective of the learner is to **recover exactly the unknown partition of the arms**. She collects observation until some time  $T$ , at which she is confidence enough to construct a partition  $\hat{G}$  equal to  $G^*$  with high probability (up to permutation of the groups).

An algorithm  $\mathcal{A}$  is  $\delta$ -PAC if for any environment  $\nu$ ,  $\mathbb{P}_{\mathcal{A}, \nu}(\hat{G} \sim G^* \text{ up to permutation }) \geq 1 - \delta$ .

### Objective: minimizing the budget spent

The performance of a  $\delta$ -PAC algorithm is measured by its budget  $T$  (by  $\mathbb{E}[T]$  or  $\|T\|_\infty$ ) – as the number of samples collected to construct  $\hat{G}$ .

For an environment  $\nu$ , we define two quantities, the minimal gap  $\Delta_*(\nu) = \min_{k \neq k'} \|\Lambda(k) - \Lambda(k')\|$ , and the balancedness  $\theta_*(\nu) = \min_k \frac{|G_k^*|}{N}$ . We denote as  $\mathcal{E}(\Delta, \theta)$  as the family of environment such that  $\Delta_* \geq \Delta$  and  $\theta_* \geq \theta$ .

Our main contribution is in showing that the complexity of the problem is characterizing by the following quantity:

$$T^* = N + \frac{\sigma^2}{\Delta^2} N \log \left( \frac{N}{\delta} \right) + \frac{\sigma^2}{\Delta^2} \sqrt{dKN \log \left( \frac{N}{\delta} \right)}.$$

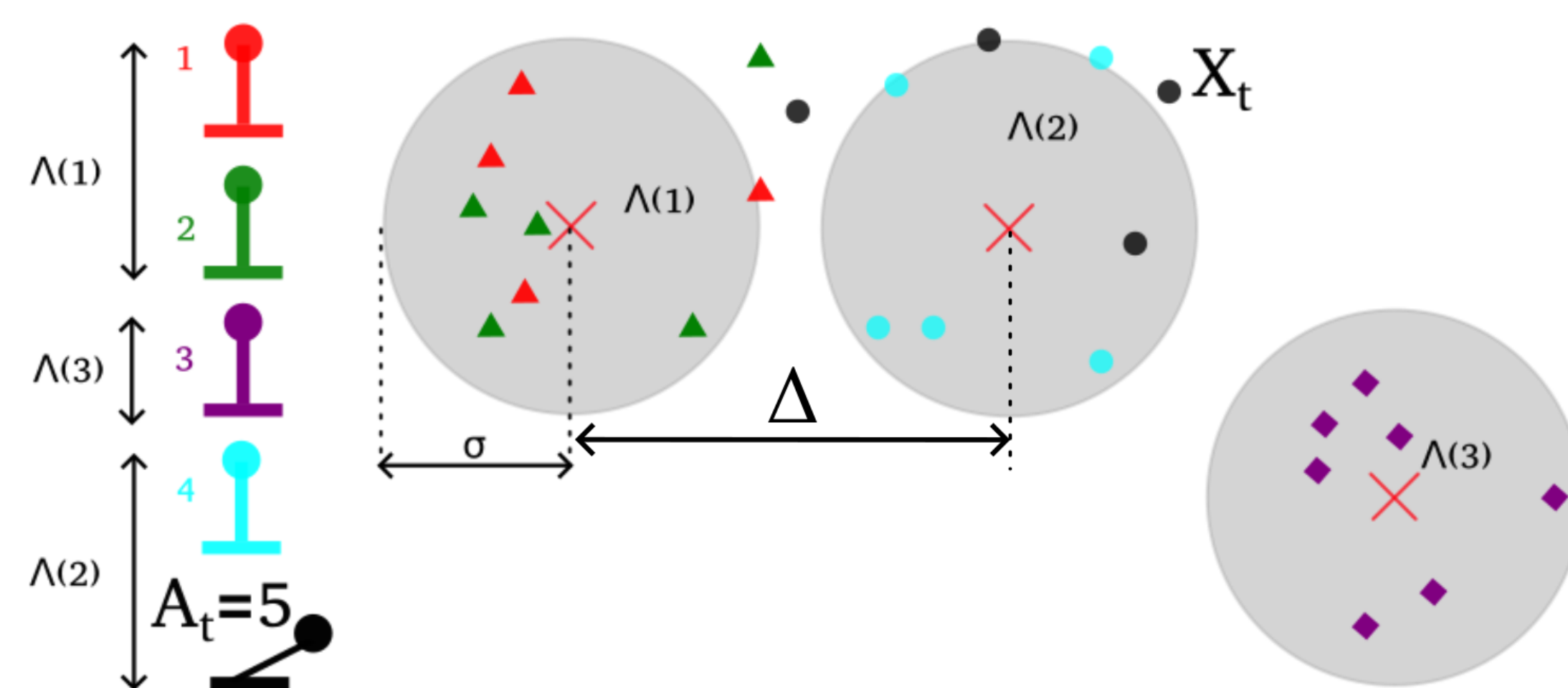


Figure 1. In this illustration,  $N = 5$ ,  $K = 3$ ,  $d = 2$ ,  $\Delta_* = \|\Lambda(1) - \Lambda(3)\|$  and  $\theta_* = 1/5$ . Based on  $X_1, \dots, X_{t-1}$ , the algorithm chooses  $A_t = 5$  and observes  $X_t$  centred on  $\mu_5 = \Lambda(2)$ .

## Algorithms

**Algorithm 1:** Sequential Representative Identification (SRI)

**Input:**  $\delta, \Delta, \theta$

**Result:**  $S$  a set of arms

Pick randomly  $a_0 \in [N]$  ;

Set  $S = \{a_0\}$

$\hat{\mu}_{a_0}, \hat{\mu}'_{a_0} \leftarrow \text{empirical\_mean}(a_0, n_{\max})$  ;

*/\* Estimate  $\mu_{a_0}$  \*/*

**for**  $u = 1, \dots, U$  **do**

Sample uniformly at random  $a_u \in [N]$

**for**  $s = s_0, \dots, r$  **do**

$\hat{\mu}_{a_u}, \hat{\mu}'_{a_u} \leftarrow \text{empirical\_mean}(a_u, n_0 2^s)$  ;

*/\* Estimate  $\mu_{a_u}$  \*/*

**if**  $\min_{b \in S} \langle \hat{\mu}_a - \hat{\mu}_b, \hat{\mu}'_a - \hat{\mu}'_b \rangle \leq \frac{\Delta^2}{2}$  **then**

Break ; */\* reject  $a_u$  \*/*

**if**  $s = r$  **then** */\* if  $a_u$  passed all tests \*/*

$S \leftarrow S \cup \{a_u\}$  */\* Add  $a_u$  to  $S$  \*/*

$\hat{\mu}_{a_u}, \hat{\mu}'_{a_u} \leftarrow \text{empirical\_mean}(a_u, n_{\max})$

*/\* Estimate  $\mu_{a_u}$  \*/*

**if**  $|S| = K$  or budget  $> T_{\max}$  **then**

Break */\* Terminate  $u$  loop \*/*

**return**  $S$  */\* Return a set of representatives \*/*

**Algorithm 2:** Active Distance-based Classifier (ADC)

**Input:**  $\delta, \Delta$  and  $S = \{b_1, \dots, b_K\}$

**Result:**  $\hat{G}$  a partition of the arms

Compute  $I = T^*/N$  and  $J = T^*/K$

**for**  $j \in [K]$  **do**

$\hat{\mu}(j), \hat{\mu}'(j) \leftarrow \text{empirical\_mean}(b_j, J)$  ;

*/\* Estimate the centers \*/*

$\hat{g}(b_j) \leftarrow j$

**for**  $a \in [n] \setminus S$  **do**

$\hat{\mu}_a, \hat{\mu}'_a \leftarrow \text{empirical\_mean}(a, I)$  ;

Compute

$\hat{g}(a) \in \text{argmin}_{j=1, \dots, K} \langle \hat{\mu}_a - \hat{\mu}(j), \hat{\mu}'_a - \hat{\mu}'(j) \rangle$

; */\* Classify arm  $a$  \*/*

**return**  $\{\hat{g}(1), \dots, \hat{g}(n)\}$  */\* Return a clustering \*/*

**Algorithm 3:** Active Clustering Bandits (ACB)

**Input:**  $\delta, \Delta, \theta$

$\hat{S} \leftarrow \text{SRI}(\delta/2, \Delta, \theta)$  ; */\* Alg 1 \*/*

**return**  $\hat{G} = \text{ADC}(\delta/2, \Delta, \hat{S})$  ; */\* Alg 2 \*/*

## Lower bound

We derive a lower bound, combining methods from information theory and high-dimensional statistics:

### Theorem 1

For any algorithm  $\mathcal{A}$ , any  $\Delta > 0$ ,  $\theta > 2/N$ , there exists an environment  $\nu \in \mathcal{E}(\Delta, \theta)$ , such that

$$\mathbb{E}_{\mathcal{A}, \nu}[T] \geq cT^*.$$

## Upper bound

We introduce ACB, an algorithm which works as a two step procedure (describe in the above column in pseudocode):

1. (SRI): identifying  $S$ , a set of arms with exactly one arm from each cluster
2. (ADC): estimate the common means of the clusters and classify the arms with a distance-based classifier,

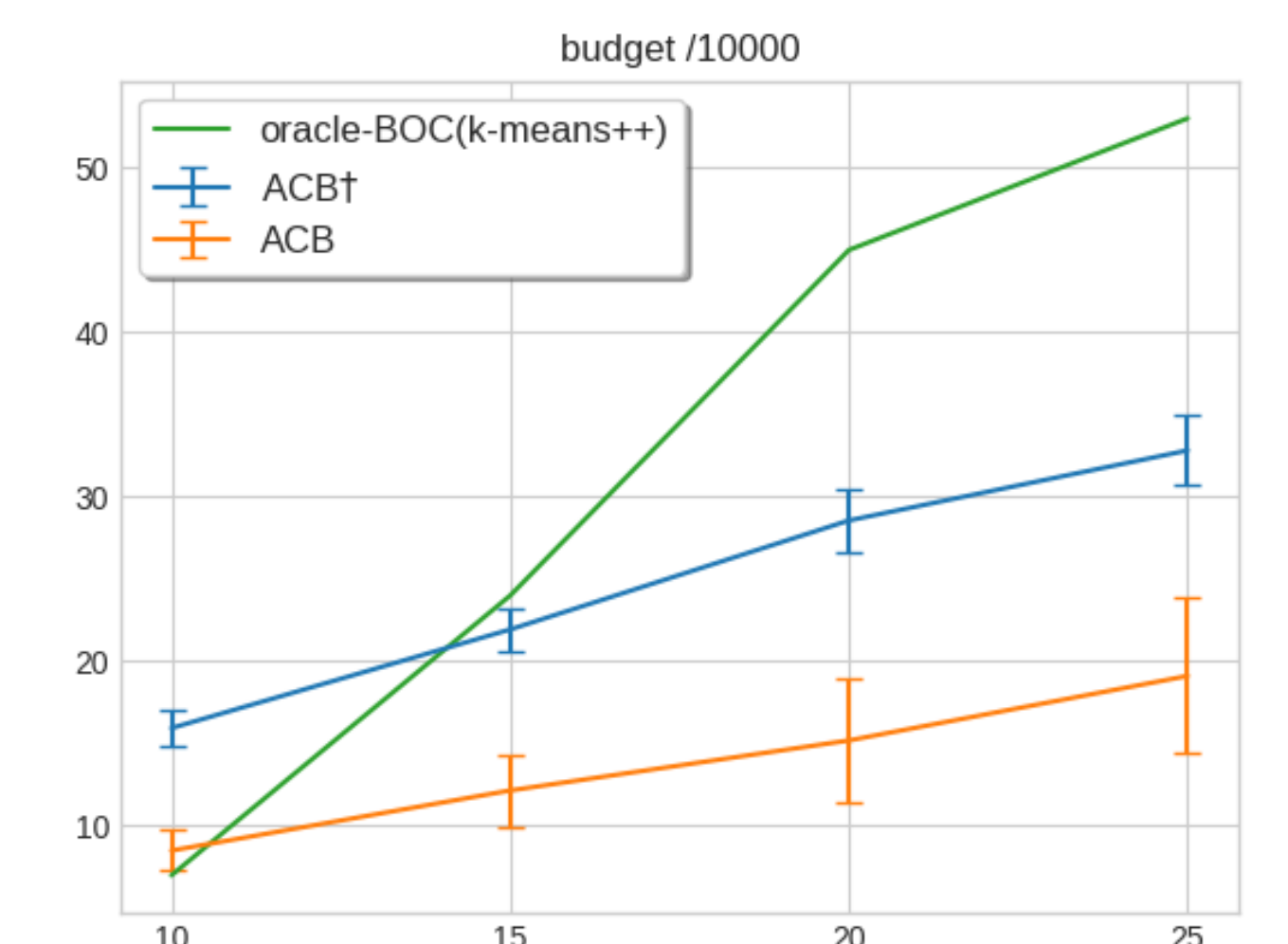
## Contributions

We answer the following questions:

1. **Can we improve the budget of a simple uniform sampling strategy ?**  
**Yes**, we provide the ACB Algorithm, a polynomial-time algorithm which outperforms the uniform sampling strategy.
2. **Can we achieve optimality ?**  
**Yes**, ACB is  $\delta$ -PAC, and we bound its budget, which matches the lower bound  $T^*$  in most regimes (for  $\theta$  not too small, e.g., with balanced groups).
3. **Is there an information-computation gap for ACP?**  
**No**, there is no computational gap (contrary to the batch setting), ACB is optimal and computationally efficient.

## Numerical experiments

Figure 2. Comparison of the necessary budget for ACB and oracle-BOC with varying number of clusters. In blue (resp. orange) the (empirical) budget of ACB† (resp. ACB) computed with 100 simulations. Algorithm ACB knows  $\Delta, \theta$ , while ACB does not know  $\Delta$  (we use a doubling trick). In green, the smallest budget for which oracle-BOC (uniform sampling followed by kmeans++) makes less than 10% of error out of 100 experiments.



## References

- [1] V. Thuot, A. Carpentier, C. Giraud, and N. Verzelen. Clustering with bandit feedback: breaking down the computation/information gap. In G. Kamath and P.-L. Loh, editors, *Proceedings of The 36th International Conference on Algorithmic Learning Theory*, volume 272 of *Proceedings of Machine Learning Research*, pages 1221–1284. PMLR, 24–27 Feb 2025.