

# Problèmes de Bandits – Identification du Meilleur Bras

Le cas du budget fixé

**Antoine Barrier**<sup>1,2</sup>, Aurélien Garivier<sup>1</sup>, Gilles Stoltz<sup>2</sup>

1. UMPA, ÉNS de Lyon      2. LMO, Université Paris-Saclay

[antoine.barrier@ens-lyon.fr](mailto:antoine.barrier@ens-lyon.fr) – <http://perso.ens-lyon.fr/antoine.barrier/fr/>

**UMPA**

**ENS**  
ENS DE LYON

 **Mathématiques**  
Orsay

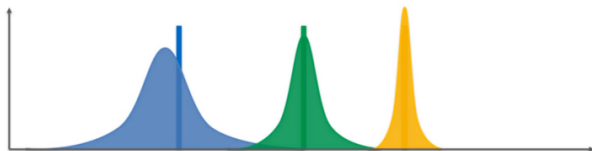
**université**  
PARIS-SACLAY

16 avril 2022, JDS, Lyon

- 1 Identification du meilleur bras
- 2 La difficulté de l'identification du meilleur bras à budget fixé
- 3 Vers l'obtention de bornes inférieures non paramétriques ?

# Problèmes de bandits

- $K$  groupes (*bras*). À chaque bras  $a \in \{1, \dots, K\}$  est associée une **distribution de probabilité**  $\nu_a$  (inconnue!) d'espérance  $\mu_a = \mathbb{E}[\nu_a]$ . On considère des bandits ayant un **unique bras optimal** noté  $a^* = a^*(\nu)$  tel que  $\underline{\mu^* = \mu_{a^*} > \max_{a \neq a^*} \mu_a}$



- À chaque tour  $t = 1, 2, \dots$ , un joueur/agent
    - choisit un bras  $A_t \in \{1, \dots, K\}$  en fonction des observations passées  $(A_1, X_1, A_2, X_2, \dots, A_{t-1}, X_{t-1})$ ,
    - reçoit une **récompense**  $X_t \sim \nu_{A_t}$  indépendante du passé
- cadre de l'**apprentissage séquentiel** : on reçoit les observations une par une et la stratégie du joueur influe sur celles-ci

**Applications** : systèmes de recommandation, essais cliniques, A/B testing, ...

# Plusieurs objectifs

## Le regret

On fixe un **nombre de tirages**  $T$  et on veut maximiser nos gains. En cumulant sur tous les tours on a une perte relative – un **regret** – par rapport à un oracle de

$$R(T) = \sum_{t=1}^T \mu^* - \mathbb{E}[X_t] = T\mu^* - \sum_{t=1}^T \mathbb{E}[\mu_{A_t}] = \sum_{a \neq a^*} (\mu^* - \mu_a) \mathbb{E}[N_a(T)]$$

où  $N_a(t) = \sum_{s=1}^t \mathbb{1}_{A_s=a}$  est le nombre de tirages du bras  $a$  après  $t$  tours.

**PREMIER OBJECTIF**     **Minimiser le regret**

→ **exploration vs exploitation**

- ✓ critère d'apprentissage :  $R(T) = o(T)$  est sous-linéaire en  $T$  : on tire de moins en moins les mauvais bras !
- ✓ objectif très étudié dans la littérature, nombreux algorithmes vérifiant des critères d'optimalité dans différents cadres

# Plusieurs objectifs

L'identification du meilleur bras (BAI pour Best Arm Identification)

## DEUXIÈME OBJECTIF Identifier rapidement le meilleur bras

→ exploration pure

- **BAI à confiance fixée** : à niveau de confiance  $\delta > 0$  fixé, un algorithme doit s'arrêter après un nombre d'observations *aléatoire*  $\tau$  et retourner un estimateur  $\hat{a}_\tau$  du meilleur bras de manière à garantir que  $\mathbb{P}_\nu(\hat{a}_\tau \neq a^*) \leq \delta$  (l'algorithme est dit  $\delta$ -correct).

→ Objectif : **minimiser le nombre moyen de tirages**  $\mathbb{E}_\nu[\tau]$  ✓

- **BAI à budget fixé** : à nombre total de tirages  $T$  fixé, un algorithme doit après  $T$  observations retourner un estimateur  $\hat{a}_T$  du meilleur bras.

→ Objectif : **minimiser la probabilité d'erreur**  $\mathbb{P}_\nu(\hat{a}_T \neq a^*)$  ✗

1 Identification du meilleur bras

2 La difficulté de l'identification du meilleur bras à budget fixé

3 Vers l'obtention de bornes inférieures non paramétriques ?

# Des résultats d'optimalité à confiance fixée

*On suppose que les distributions appartiennent à un même modèle exponentiel.*

On pose  $\text{Alt}(\nu)$  l'ensemble des problèmes de bandits ayant un bras optimal différent de celui de  $\nu$  et  $\Sigma_K = \{\mathbf{w} \in [0, 1]^K : \sum_{a=1}^K w_a = 1\}$ .

**Théorème (borne inférieure, [GK16])**

*Pour toute stratégie  $\delta$ -correcte, on a*

$$\forall \nu, \quad \mathbb{E}_\nu[\tau_\delta] \geq T(\nu) \text{kl}(\delta, 1 - \delta) \underset{\delta \rightarrow 0}{\sim} T(\nu) \log(1/\delta)$$

$$\text{où} \quad T(\nu)^{-1} = \sup_{\mathbf{w} \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\nu)} \sum_{a=1}^K w_a \text{KL}(\nu_a, \lambda_a)$$

[GK16] Garivier, A. et Kaufmann, E. (2016), **Optimal Best Arm Identification with Fixed Confidence**, In *29th Conference on Learning Theory (COLT)*

# Des résultats d'optimalité à confiance fixée

*Schéma de preuve.*

Soit  $a^* = a^*(\nu)$ . Pour tout  $\lambda \in \text{Alt}(\nu)$ , on obtient en utilisant la chain rule pour KL, l'inégalité de data-processing et l'hypothèse de stratégie  $\delta$ -correcte que

$$\sum_{a=1}^K \mathbb{E}_{\nu} [N_a(\tau_{\delta})] \text{KL}(\nu_a, \lambda_a) \geq \text{kl}(\mathbb{P}_{\nu}(\hat{a}_{\tau_{\delta}} \neq a^*), \mathbb{P}_{\lambda}(\hat{a}_{\tau_{\delta}} \neq a^*)) \geq \text{kl}(\delta, 1 - \delta)$$

et donc en passant à l'inf sur  $\lambda$

$$\begin{aligned} \text{kl}(\delta, 1 - \delta) &\leq \mathbb{E}_{\nu} [\tau_{\delta}] \inf_{\lambda \in \text{Alt}(\nu)} \sum_{a=1}^K \underbrace{\frac{\mathbb{E}_{\nu} [N_a(\tau_{\delta})]}{\mathbb{E}_{\nu} [\tau_{\delta}]}}_{=w_a(\nu)} \text{KL}(\nu_a, \lambda_a) \\ &\leq \mathbb{E}_{\nu} [\tau_{\delta}] \inf_{\lambda \in \text{Alt}(\nu)} \sum_{a=1}^K w_a(\nu) \text{KL}(\nu_a, \lambda_a) \\ &\leq \mathbb{E}_{\nu} [\tau_{\delta}] \sup_{w \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\nu)} \sum_{a=1}^K w_a \text{KL}(\nu_a, \lambda_a) \quad (\text{car } w(\nu) \in \Sigma_K) \end{aligned}$$



# Une preuve similaire à budget fixé ?

Essayons de faire de même en budget fixé. Pour cela on a besoin de l'inégalité

$$\text{kl}(p, q) \geq p \log(1/q) - \ln(2) \iff q \geq \exp\left(-\frac{\text{kl}(p, q) + \ln(2)}{p}\right)$$

ce qui donne

$$\begin{aligned} \mathbb{P}_{\nu}(\hat{a}_T \neq a^*) &\geq \exp\left(-\frac{\text{kl}(\mathbb{P}_{\lambda}(\hat{a}_T \neq a^*), \mathbb{P}_{\nu}(\hat{a}_T \neq a^*)) + \ln(2)}{\mathbb{P}_{\lambda}(\hat{a}_T \neq a^*)}\right) \\ &\geq \exp\left(-\frac{\sum_{a=1}^K \mathbb{E}_{\lambda}[N_a(T)] \text{KL}(\lambda_a, \nu_a) + \ln(2)}{\mathbb{P}_{\lambda}(\hat{a}_T \neq a^*)}\right) \end{aligned}$$

**Problème** : on a été contraint d'échanger les rôles de  $\nu$  et  $\lambda$  par rapport à la preuve précédente (en cause notamment la non symétrie de la kl).

→ Le fait d'avoir  $\mathbb{E}_{\lambda}[N_a(T)]$  et non  $\mathbb{E}_{\nu}[N_a(T)]$  est l'une de raisons qui rend compliquée la recherche d'une borne inférieure à budget fixé.

- 1 Identification du meilleur bras
- 2 La difficulté de l'identification du meilleur bras à budget fixé
- 3 Vers l'obtention de bornes inférieures non paramétriques ?

# Avec une hypothèse de monotonie

Continuons tout de même nos calculs en supposant que les bras sont rangés de sorte que  $a^* = 1$  et  $\mu_1 > \mu_2 \geq \dots \geq \mu_K$ .

$$\mathbb{P}_\nu(\hat{a}_T \neq 1) \geq \exp\left(-\frac{\sum_{a=1}^K \mathbb{E}_\lambda[N_a(T)] \text{KL}(\lambda_a, \nu_a) + \ln(2)}{\mathbb{P}_\lambda(\hat{a}_T \neq 1)}\right)$$

Si la stratégie est performante, elle va identifier le meilleur bras donc sous le bandit alternatif  $\lambda$ , on aura  $\mathbb{P}_\lambda(\hat{a}_T \neq 1) \rightarrow_{T \rightarrow +\infty} 0$ .

Fixons  $2 \leq k \leq K$  et considérons le bandit alternatif  $\lambda = (\nu_k, \nu_2, \nu_3, \dots, \nu_K)$ . Alors on a obtenu

$$\mathbb{P}_\nu(\hat{a}_T \neq 1) \geq \exp\left(-\mathbb{E}_\lambda[N_1(T)] \underbrace{\text{KL}(\lambda_1, \nu_1)}_{\nu_k} (1 + o(1))\right)$$

Si de plus la stratégie est *monotone*, alors  $\mathbb{E}_\lambda[N_1(T)] \leq \frac{T}{k}$  et donc

$$\forall 2 \leq k \leq K, \quad \mathbb{P}_\nu(\hat{a}_T \neq 1) \geq \exp\left(-\frac{T}{\frac{k}{\text{KL}(\nu_k, \nu_1)}} (1 + o(1))\right)$$

## Théorème (Borne inférieure, cas d'une stratégie monotone)

Soit une stratégie

- *symétrique* : le choix des bras à tirer est indépendant de leur indice,
- *consistante* :  $\mathbb{P}_\nu(\hat{a}_T \neq a^*) = o_{T \rightarrow +\infty}(1)$  pour tout bandit  $\nu$ ,
- *monotone* : les espérances de tirages  $(\mathbb{E}_\nu[N_a(T)])_{1 \leq a \leq K}$  sont classées dans l'ordre inverse des espérances des bras  $(\mu_a)_{1 \leq a \leq K}$  pour tout bandit  $\nu$ .

Alors pour tout problème de bandit  $\nu$  on a

$$\mathbb{P}_\nu(\hat{a}_T \neq a^*) \geq C \exp\left(-\frac{2T}{\max_{2 \leq k \leq K} \frac{1}{\text{KL}(\nu_{(k)}, \nu_{(1)})}} (1 + o_{T \rightarrow +\infty}(1))\right)$$

pour une constante  $C > 0$ , où l'on définit un ordre tel que

$$\mu_{(1)} > \mu_{(2)} \geq \dots \geq \mu_{(K)}.$$

# Comparaison avec la borne de [ABM10]

Théorème (Borne inférieure, cas gaussien [ABM10])

Soit une stratégie symétrique et consistante. Alors pour tout problème de bandit  $\nu$  dont les bras sont gaussiens de variance  $\sigma^2$  on a

$$\mathbb{P}_\nu(\hat{a}_T \neq a^*) \geq C \exp\left(-\frac{2.5T}{\sigma^2 \max_{2 \leq k \leq K} \frac{k}{(\mu_{(1)} - \mu_{(k)})^2}} (1 + o_{T \rightarrow +\infty}(1))\right)$$

pour une constante  $C > 0$ .

Notre borne donne quant à elle

$$\mathbb{P}_\nu(\hat{a}_T \neq a^*) \geq C \exp\left(-\frac{T}{\sigma^2 \max_{2 \leq k \leq K} \frac{k}{(\mu_{(1)} - \mu_{(k)})^2}} (1 + o_{T \rightarrow +\infty}(1))\right)$$

[ABM10] Audibert, J.-Y., Bubeck, S. et Munos, R. (2010), **Best Arm Identification in Multi-Armed Bandits**, In *23th Conference on Learning Theory (COLT)*

- le cadre de BAI à budget fixé est aujourd'hui moins bien compris que le BAI à confiance fixée et semble plus difficile,
- avec des hypothèses sur les stratégies, on obtient tout de même des bornes inférieures à base de quantités informatives,
- ces hypothèses sont pour l'instant trop restrictives pour espérer des résultats avec des modèles non paramétriques.



Jean-Yves AUDIBERT, Sébastien BUBECK et Rémy MUNOS :  
Best Arm Identification in Multi-Armed Bandits.

*In Proceedings of the 23rd Annual Conference on Learning Theory*, janvier 2010.



Aurélien GARIVIER et Emilie KAUFMANN :  
Optimal Best Arm Identification with Fixed Confidence.

*In Conference on Learning Theory*, volume 49, pages 998–1027. PMLR, juin 2016.