

# The True Shape of Regret in Bandit problems

Aurélien Garivier, **Pierre Ménard**, Gilles Stoltz

July 5, 2016

## Environment and strategy

- K arms bandit problem,  $\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_K))$  with  $\mu_i \in (0, 1)$ .  
Game, for each round  $1 \leq t \leq T$ :
  1. Player pulls arm  $A_t \in \{1, \dots, K\}$ .
  2. He gets a reward  $Y_t \sim \mathcal{B}(\mu_{A_t})$ .
- Information available at time  $t$ :  $Y_{1:t} = (Y_1, \dots, Y_t)$ .

## Environment and strategy

- K arms bandit problem,  $\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_K))$  with  $\mu_i \in (0, 1)$ .  
Game, for each round  $1 \leq t \leq T$ :
  1. Player pulls arm  $A_t \in \{1, \dots, K\}$ .
  2. He gets a reward  $Y_t \sim \mathcal{B}(\mu_{A_t})$ .
- Information available at time  $t$ :  $Y_{1:t} = (Y_1, \dots, Y_t)$ .
- Goal of the player, minimize the expected regret :

$$R_{\nu, T} = T\mu^* - \mathbb{E}_{\nu} \left[ \sum_{t=1}^T Y_t \right] = \sum_{a=1}^K (\mu^* - \mu_a) \mathbb{E}_{\nu} [N_a(T)].$$

where  $\mu^* = \max_{a=1, \dots, K} \mu_a$  and  $N_a(T) = \sum_{t=1}^T \mathbb{I}_{\{A_t=a\}}$ .

## Theorem (Asymptotic lower bound from Lai & Robbins)

*For all reasonable strategies (consistent), for all bandits problems  $\nu$ , for all suboptimal arms  $a$ ,*

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_{\nu}[N_a(T)]}{\ln T} \geq \frac{1}{\text{kl}(\mu_a, \mu^*)}.$$

where  $\text{kl}$  the Kullback-Leibler divergence for Bernoulli distributions :

$$\forall p, q \in [0, 1]^2, \quad \text{kl}(p, q) := p \ln \frac{p}{q} + (1 - p) \ln \frac{1 - p}{1 - q} \geq 2(p - q)^2.$$

## Theorem (Asymptotic lower bound from Lai & Robbins)

*For all reasonable strategies (consistent), for all bandits problems  $\nu$ , for all suboptimal arms  $a$ ,*

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_\nu[N_a(T)]}{\ln T} \geq \frac{1}{\text{kl}(\mu_a, \mu^*)}.$$

where  $\text{kl}$  the Kullback-Leibler divergence for Bernoulli distributions :

$$\forall p, q \in [0, 1]^2, \quad \text{kl}(p, q) := p \ln \frac{p}{q} + (1 - p) \ln \frac{1 - p}{1 - q} \geq 2(p - q)^2.$$

## Theorem (UCB algorithm from Auer & Cesa-Bianchi & Fischer)

*Algorithm UCB, for all bandits problems  $\nu$ , for all suboptimal arm  $a$ :*

$$\mathbb{E}_\nu[N_a(T)] \leq \frac{8 \ln(T)}{(\mu^* - \mu_a)^2} + 2$$

*(right constant with KL-UCB algorithm from Cappe & al).*

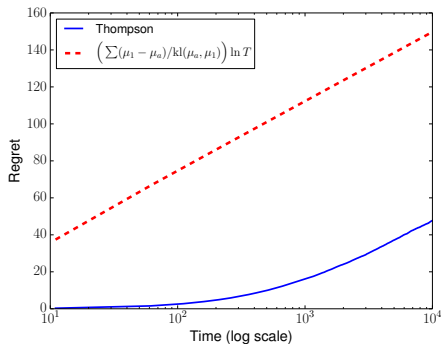
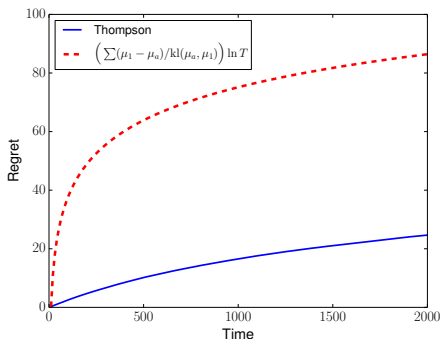


Figure : Bernoulli bandit problem with parameters :  
 $(\mu_a)_{1 \leq a \leq 6} = (0.05, 0.04, 0.02, 0.015, 0.01, 0.005)$

- Logarithmic regret for large  $T$  (**asymptotic** lower bound).
- Transition phase between.
- Linear regret for  $T$  small.

## Consistent strategy

Strategy which always pulls the same arm  $\rightarrow$  assumptions on the strategy.

### Definition

A strategy is consistent if for all bandit problems  $\nu$ , for all suboptimal arms  $a$ , i.e., for all arms  $a$  such that  $\Delta_a > 0$ , it satisfies  $\mathbb{E}_\nu [N_a(T)] = o(T^\alpha)$  for all  $0 < \alpha \leq 1$ .

$$\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_K)) \quad \nu' = (\mathcal{B}(\mu'_1), \dots, \mathcal{B}(\mu'_K))$$

$$\sum_{a=1}^K \mathbb{E}_{\nu}[N_a(T)] \text{kl}(\mu_a, \mu'_a) \geq \text{kl}(\mathbb{E}_{\nu}[Z], \mathbb{E}_{\nu'}[Z]), \quad (\text{M})$$

where  $Z$  is a  $\sigma(Y_{1:T})$ -measurable random variable with values in  $[0, 1]$ .  
Typically  $Z = N_a(T)/T$ .



$$\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_K)) \quad \nu' = (\mathcal{B}(\mu'_1), \dots, \mathcal{B}(\mu'_K))$$

$$\boxed{\sum_{a=1}^K \mathbb{E}_{\nu}[N_a(T)] \text{kl}(\mu_a, \mu'_a) \geq \text{kl}(\mathbb{E}_{\nu}[Z], \mathbb{E}_{\nu'}[Z])}, \quad (\text{M})$$

where  $Z$  is a  $\sigma(Y_{1:T})$ -measurable random variable with values in  $[0, 1]$ .  
Typically  $Z = N_a(T)/T$ .

Sketch of proof :

$$\sum_{a=1}^K \mathbb{E}_{\nu}[N_a(T)] \text{kl}(\mu_a, \mu'_a) = \text{KL}(\mathbb{P}_{\nu}^{Y_{1:T}}, \mathbb{P}_{\nu'}^{Y_{1:T}})$$
$$\text{KL}(\mathbb{P}_{\nu}^{Y_{1:T}}, \mathbb{P}_{\nu'}^{Y_{1:T}}) \geq \text{kl}(\mathbb{E}_{\nu}[Z], \mathbb{E}_{\nu'}[Z]),$$

where :

- $\mathbb{P}_{\nu}^{Y_{1:T}}$  and  $\mathbb{P}_{\nu'}^{Y_{1:T}}$  respective distributions of  $Y_{1:T}$  under  $\mathbb{P}_{\nu}$  and  $\mathbb{P}_{\nu'}$
- chain rule for Kullback-Leibler divergences
- contraction of entropy

## Proof.

- Contraction of entropy :

Let  $V \sim \mathcal{U}[0, 1]$  independent of  $Y_{1:T}$ , and the event  $E = \{Z \geq V\}$  then

$$\begin{aligned} \text{KL}(\mathbb{P}_\nu^{Y_{1:T}}, \mathbb{P}_{\nu'}^{Y_{1:T}}) &= \text{KL}(\mathbb{P}_\nu^{Y_{1:T}} \otimes \mathcal{U}[0, 1], \mathbb{P}_{\nu'}^{Y_{1:T}} \otimes \mathcal{U}[0, 1]) \\ &\geq \text{KL}\left((\mathbb{P}_\nu^{Y_{1:T}} \otimes \mathcal{U}[0, 1])^{\mathbb{1}_E}, (\mathbb{P}_{\nu'}^{Y_{1:T}} \otimes \mathcal{U}[0, 1])^{\mathbb{1}_E}\right) \\ &= \text{kl}\left((\mathbb{P}_\nu^{Y_{1:T}} \otimes \mathcal{U}[0, 1])(E), (\mathbb{P}_{\nu'}^{Y_{1:T}} \otimes \mathcal{U}[0, 1])(E)\right). \end{aligned}$$

To conclude, for  $\alpha = \nu$  or  $\nu'$  (Fubini theorem):

$$(\mathbb{P}_\alpha^{Y_{1:T}} \otimes \mathcal{U}[0, 1])(E) = \mathbb{E}_\alpha[Z].$$



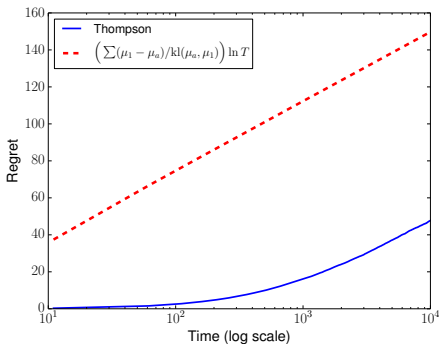
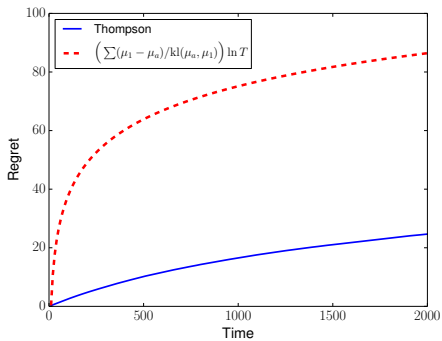


Figure : Bernoulli bandit problem with parameters :  
 $(\mu_a)_{1 \leq a \leq 6} = (0.05, 0.04, 0.02, 0.015, 0.01, 0.005)$

- Linear regret for  $T$  small.
- Logarithmic regret for large  $T$  (**asymptotic** lower bound).
- Transition phase between.

## Absolute lower bound for small $T$

In the remainder of this section  $\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_K))$  with an unique optimal arm  $i^*$ .

## Absolute lower bound for small $T$

In the remainder of this section  $\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_K))$  with an unique optimal arm  $i^*$ .

Uniform strategy : pull an arm uniformly at random at each round.

### Definition

A strategy is smarter than the uniform strategy if for all bandit problems  $\nu$ , for all  $T \geq 1$ ,

$$\mathbb{E}_\nu [N_{i^*}(T)] \geq \frac{T}{K}$$
$$\mathbb{E}_\nu [N_a(T)] \leq \frac{T}{K} \quad \text{if } a \text{ suboptimal.}$$

## Theorem

*For all strategies that are smarter than the uniform strategy, for all bandit problems  $\nu$ , for all suboptimal arms  $a$ , for all  $T \geq 1$ ,*

$$\mathbb{E}_\nu [N_a(T)] \geq \frac{T}{K} \left(1 - \sqrt{2T \text{kl}(\mu_a, \mu^*)}\right).$$

*In particular,*

$$\forall T \leq \frac{1}{8 \text{kl}(\mu_a, \mu^*)}, \quad \mathbb{E}_\nu [N_a(T)] \geq \frac{T}{2K}.$$

Linear regret

a suboptimal arm.

Modified bandit problem with  $\mu'_a > \mu^*$ :

$$\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

$$\nu' = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))$$

a suboptimal arm.

Modified bandit problem with  $\mu'_a > \mu^*$ :

$$\begin{aligned}\nu &= (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K)) \\ \nu' &= (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))\end{aligned}$$

Main inequality (M),

$$\begin{aligned}\mathbb{E}_\nu[N_a(T)] \text{kl}(\mu_a, \mu'_a) &\geq \text{kl}\left(\mathbb{E}_\nu[N_a(T)]/T, \mathbb{E}_{\nu'}[N_a(T)]/T\right) \\ \left(\mathbb{E}_\nu[N_a(T)]/T \leq 1/K \leq \mathbb{E}_{\nu'}[N_a(T)]/T\right) &\geq \text{kl}\left(\mathbb{E}_\nu[N_a(T)]/T, 1/K\right) \\ \text{( Pinsker inequality )} &\geq \frac{K}{2} \left(\mathbb{E}_\nu[N_a(T)]/T - 1/K\right)^2\end{aligned}$$



a suboptimal arm.

Modified bandit problem with  $\mu'_a > \mu^*$ :

$$\begin{aligned}\nu &= (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K)) \\ \nu' &= (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))\end{aligned}$$

Main inequality (M),

$$\begin{aligned}\mathbb{E}_\nu[N_a(T)] \text{kl}(\mu_a, \mu'_a) &\geq \text{kl}\left(\mathbb{E}_\nu[N_a(T)]/T, \mathbb{E}_{\nu'}[N_a(T)]/T\right) \\ \left(\mathbb{E}_\nu[N_a(T)]/T \leq 1/K \leq \mathbb{E}_{\nu'}[N_a(T)]/T\right) &\geq \text{kl}\left(\mathbb{E}_\nu[N_a(T)]/T, 1/K\right) \\ \text{( Pinsker inequality )} &\geq \frac{K}{2} \left(\mathbb{E}_\nu[N_a(T)]/T - 1/K\right)^2\end{aligned}$$

Still with  $\mathbb{E}_\nu[N_a(T)]/T \leq 1/K$  :

$$\text{kl}(\mu_a, \mu'_a) T/K \geq \frac{K}{2} \left(\mathbb{E}_\nu[N_a(T)]/T - 1/K\right)^2$$

## Collective lower bound for small $T$

### Theorem

*For all strategies that are smarter than the uniform strategy, for all bandit problems  $\nu$ , for all suboptimal arms  $a$ ,*

$$\forall T \leq \frac{K?}{8\text{kl}(\mu_a, \mu^*)}, \quad \mathbb{E}_\nu [N_a(T)] \geq \frac{T}{2K}.$$

## Collective lower bound for small $T$

### Theorem

*For all strategies that are smarter than the uniform strategy, for all bandit problems  $\nu$ , for all suboptimal arms  $a$ ,*

$$\forall T \leq \frac{K?}{8 \text{kl}(\mu_a, \mu^*)}, \quad \mathbb{E}_\nu [N_a(T)] \geq \frac{T}{2K}.$$

### Theorem

*Under weak (symmetry, ...) assumptions on the strategy, for all bandit problems  $\nu$ ,*

$$\sum_{a \neq i^*} \mathbb{E}_\nu [N_a(T)] \geq T \left( 1 - \frac{1}{K} - \frac{\sqrt{2T \text{kl}(\mu_a, \mu^*)}}{K} - \frac{2T \text{kl}(\mu_a, \mu^*)}{K} \right).$$

# Non-asymptotic bounds for large $T$

## Theorem

*For all reasonable strategies (refinement of consistence), for all bandit problems  $\nu$ , for all suboptimal arms  $a$ ,*

$$\mathbb{E}_{\nu}[N_a(T)] \geq \frac{\ln T}{\text{kl}(\mu_a, \mu^*)} - O(\ln(\ln T)),$$

*with a closed-form expression for the last term.*

Where, for  $T$  large enough

$$O(\ln(\ln T)) = \frac{1}{(1 - \mu^*)\text{kl}(\mu_a, \mu^*)} (\ln T)^{-3} + C_{\psi, \mathcal{D}} H(\nu) \frac{\ln(T)^2}{T} + \ln(K C_{\psi, \mathcal{D}} (\ln T)^9) + \frac{\ln 2}{\text{kl}(\mu_a, \mu^*)}.$$

## General bandit problems

$$\nu = (\nu_1, \dots, \nu_K),$$

$\nu_a \in \mathcal{D}$  a probability distribution and a real number  $x$ , we introduce

$$\mathcal{K}_{\text{inf}}(\nu_a, x) = \inf \left\{ \text{KL}(\nu_a, \nu'_a) : \nu'_a \in \mathcal{D} \text{ and } E(\nu'_a) > x \right\};$$

by convention, the infimum of the empty set equals  $+\infty$ .

→ replace  $\text{kl}(\mu_a, \mu^*)$  by  $\mathcal{K}_{\text{inf}}(\nu_a, \mu^*)$

# References I

 P. Auer, N. Cesa-Bianchi, and P. Fischer.

Finite-time analysis of the multiarmed bandit problem.

*Machine Learning*, 47(2-3):235–256, 2002.

 A.N. Burnetas and M.N. Katehakis.

Optimal adaptive policies for sequential allocation problems.

*Advances in Applied Mathematics*, 17(2):122–142, 1996.

 O. Cappé, A. Garivier, O.-A. Maillard, R. Munos, and G. Stoltz.

Kullback-Leibler upper confidence bounds for optimal sequential allocation.

*Annals of Statistics*, 41(3):1516–1541, 2013.

 Aurélien Garivier, Pierre Ménard, and Gilles Stoltz.

Explore first, exploit next: The true shape of regret in bandit problems.

*arXiv preprint arXiv:1602.07182*, 2016.

 T. L. Lai and H. Robbins.

Asymptotically efficient adaptive allocation rules.

*Advances in Applied Mathematics*, 6:4–22, 1985.