

# Stochastic bandit lower bound

Pierre Ménard

May 16, 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  $t \geq 1$  :

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  $t$  :

$$I_t = (Y_1, \dots, Y_t).$$

## Regret

- Optimal arm and gap :

$$\mu^* = \max_{a=1,\dots,K} \mu_a \quad \text{and} \quad \Delta_a = \mu^* - \mu_a.$$

- Number of time arm  $a$  is pulled :

$$N_a(T) = \sum_{t=1}^T \mathbb{I}_{\{A_t=a\}}.$$

- Goal of the player, minimize the expected regret :

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

(tower rule)

Tow blocks inequality :

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

where

- $\mathbb{P}_{\nu}^{I_T}$  and  $\mathbb{P}_{\nu'}^{I_T}$  respective distributions of  $I_T$  under  $\mathbb{P}_{\nu}$  and  $\mathbb{P}_{\nu'}$
- $\text{kl}$  the Kullback-Leibler divergence for Bernoulli distributions :

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

- $Z$  a  $\sigma(I_T)$ -measurable random variable with values in  $[0, 1]$ .

Tow blocks inequality :

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

where

- $\mathbb{P}_{\nu}^{I_T}$  and  $\mathbb{P}_{\nu'}^{I_T}$  respective distributions of  $I_T$  under  $\mathbb{P}_{\nu}$  and  $\mathbb{P}_{\nu'}$
- $\text{kl}$  the Kullback-Leibler divergence for Bernoulli distributions :

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

- $Z$  a  $\sigma(I_T)$ -measurable random variable with values in  $[0, 1]$ .

Typically  $Z = N_a(T)/T$ .

## Proof.

- Equality in F, an application of chain rule for Kullback-Leibler divergences

:

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



## Proof.

- Equality in F, an application of chain rule for Kullback-Leibler divergences :

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

- Inequality in F, use contraction of entropy : Let  $V \sim \mathcal{U}[0, 1]$  independent of  $I_T$ , and the event  $E = \{Z \geq V\}$  then

$$\begin{aligned} \text{KL}(\mathbb{P}_{\nu}^{I_T}, \mathbb{P}_{\nu'}^{I_T}) &= \text{KL}(\mathbb{P}_{\nu}^{I_T} \otimes \mathcal{U}, \mathbb{P}_{\nu'}^{I_T} \otimes \mathcal{U}) \geq \text{KL}((\mathbb{P}_{\nu}^{I_T} \otimes \mathcal{U})^{\mathbb{I}_E}, (\mathbb{P}_{\nu'}^{I_T} \otimes \mathcal{U})^{\mathbb{I}_E}) \\ &= \text{kl}((\mathbb{P}_{\nu}^{I_T} \otimes \mathcal{U})(E), (\mathbb{P}_{\nu'}^{I_T} \otimes \mathcal{U})(E)). \end{aligned}$$

The proof is concluded by noting that for all  $\alpha = \nu$  or  $\nu'$ ,

$$(\mathbb{P}_{\alpha}^{I_T} \otimes \mathcal{U})(E) = \mathbb{E}_{\alpha}[Z].$$



## Definition

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

Lower bound from Lai & Robbins :

## Theorem (asymptotic distribution-dependent lower bounds)

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

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



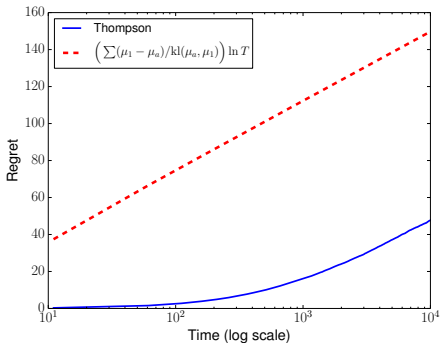
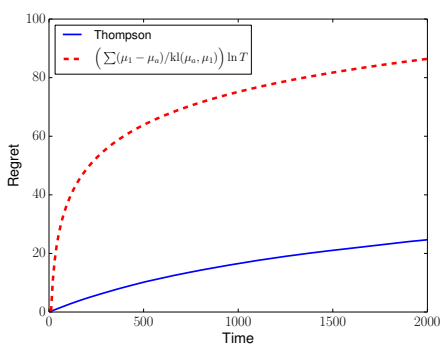


Figure : Bernoulli bandit problem with parameters :  
 $(\mu_a)_{1 \leq a \leq 6} = (0.05, 0.04, 0.02, 0.015, 0.01, 0.005)$

- Linear regret for T small.
- Logarithmic regret for large T (**asymptotic** lower bound).

## Absolute lower bound for a suboptimal arm

In what follows  $\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_K))$  with an unique optimal arm  $i^*$ .

## Absolute lower bound for a suboptimal arm

In what follows  $\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_K))$  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  $T \geq 1$ ,

$$\begin{aligned} \mathbb{E}_\nu [N_{i^*}(T)] &\geq \frac{T}{K} \\ \mathbb{E}_\nu [N_a(T)] &\leq \frac{T}{K} \quad a \text{ supotimal.} \end{aligned}$$

## Theorem

*For all strategies that are smarter than the uniform strategy, for the bandit problems  $\nu$ , for all 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  $\nu' = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu'_a), \dots, \mathcal{B}(\mu_K))$  with  $\mu'_a > \mu^*$ .

Tow blocks inequality,

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

smarter than the uniform :  $\mathbb{E}_{\nu}[N_a(T)]/T \leq 1/K \leq \mathbb{E}_{\nu'}[N_a(T)]/T$  and  $q \mapsto \text{kl}(p, q)$  is increasing on  $[p, 1]$ ,

$$\geq \text{kl}\left(\mathbb{E}_{\nu}[N_a(T)]/T, 1/K\right)$$

Pinsker inequality,

$$\geq \frac{K}{2} \left(\mathbb{E}_{\nu}[N_a(T)]/T - 1/K\right)^2$$

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

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