

## Lecture #5: Lower bound

Last lecture, we proposed algorithms with (pseudo)regrets bounded as

$$R_T \leq c \sum_{k, \Delta_k > 0} \frac{\ln T}{\Delta_k}$$

(instance dependent regret)

Is it possible to do better?

This lecture focuses on lower bounding the achievable regret by any algorithm

For that we consider a model where the rewards distributions belong to some **known** distribution set  $\mathcal{D}$ .

i.e.  $\forall k \in [K], \nu_k \in \mathcal{D}$

unknown  $\xrightarrow{\hspace{1cm}}$  known

One can show matching upper and lower bounds (with associated strategies):

$R_T$  is at best of order  $\left( \sum_{k, \Delta_k > 0} \frac{\Delta_k}{\text{Kinf}(\nu_k, \mu^*, \mathcal{D})} \right) \ln T$

where

$$\text{Kinf}(\nu_k, \mu^*, \mathcal{D}) = \inf \left\{ \text{KL}(\nu_k, \nu') \mid \begin{array}{l} \nu' \in \mathcal{D} \\ \mathbb{E}[\nu'] > \mu^* \end{array} \right\}$$

Kullback-Leibler divergence

We will only prove the lower bound part

- Case 1:  $\mathcal{D} = \left\{ N(\mu, \sigma^2) \mid \mu \in \mathbb{R} \right\}$

Then

$$K_{\text{inf}}(r_k, \mu^*, \mathcal{D}) = \frac{\Delta_k^2}{2\sigma^2}$$

Best possible regret of order  $2\sigma^2 \sum_{k, \Delta_k > 0} \frac{\ln T}{\Delta_k}$

UCB has regret  $\leq 32\sigma^2 \sum_{k, \Delta_k > 0} \frac{\ln T}{\Delta_k}$

↪ optimal up to constant factor  
can be made optimal with finer version

- Case 2:  $\mathcal{D} = \left\{ \text{Ber}(p) \mid p \in [0, 1] \right\}$

Then

$$K_{\text{inf}}(r_k, \mu^*, \mathcal{D}) = \mu_k \ln \frac{\mu_k}{\mu^*} + (1 - \mu_k) \ln \frac{1 - \mu_k}{1 - \mu^*}$$

**But** before proving the lower bound, I guess that some reminder of basic and non-basic results about KL divergences would be needed!

For sake of time, I will only give these *key properties*, without any proof.

## Definition

let  $P, Q$  be two probability measures over  $(\Omega, \mathcal{F})$

$$KL(P, Q) = \begin{cases} +\infty & \text{if } P \text{ is not absolutely continuous wrt } Q \\ \int_{\Omega} \left( \frac{dP}{dQ} \ln \left( \frac{dP}{dQ} \right) \right) dQ = \int_{\Omega} \ln \left( \frac{dP}{dQ} \right) dP & \text{if } P \ll Q \\ Q(A) = 0 \Rightarrow P(A) = 0 \end{cases}$$

## Basic Facts

- existence of the defining integral when  $P \ll Q$ , because  $\Psi: x \mapsto x \ln x$  is bounded from below on  $[0, +\infty)$
- $KL(P, Q) \geq 0$  and  $KL(P, Q) = 0$  if and only if  $P = Q$ .  
Indeed,  $\Psi$  is strictly convex. Jensen's inequality indicates that

$$KL(P, Q) = \int_{\Omega} \psi\left(\frac{dP}{dQ}\right) dQ \geq \psi\left(\int_{\Omega} \frac{dP}{dQ} dQ\right) = \psi(1) = 0, \text{ with}$$

equality if and only if  $\frac{dP}{dQ}$  is  $Q$ -almost surely constant, i.e.  $P = Q$ .

A useful rewriting:

Assume  $P \ll Q$  and let  $\nu$  be any probability measure over  $(\Omega, \mathcal{F})$  with  $P \ll \nu$ ,  $Q \ll \nu$ . Denote  $f = \frac{dP}{d\nu}$ ,  $g = \frac{dQ}{d\nu}$ , then  $KL(P, Q) = \int_{\Omega} \ln\left(\frac{f}{g}\right) f d\nu$ .

useful when  $P$  and  $Q$  both admit densities over a classical reference measure (e.g. Lebesgue).

Lemma (data processing inequality)

Let  $P, Q$  be two probability measures over  $(\Omega, \mathcal{F})$ .

Let  $X: (\Omega, \mathcal{F}) \rightarrow (\Omega', \mathcal{F}')$  be any random variable.

Denote by  $P^X$  and  $Q^X$  the laws of  $X$  under  $P$  and  $Q$ .

Then  $KL(P^X, Q^X) \leq KL(P, Q)$ .

## Proposition (KL for product measures, independent case)

Let  $(\Omega, \mathcal{F})$  and  $(\Omega', \mathcal{F}')$  be two measurable spaces.

Let  $P, Q$  be two probability measures over  $(\Omega, \mathcal{F})$

$P', Q'$  over  $(\Omega', \mathcal{F}')$

and denote by  $P \otimes P'$  and  $Q \otimes Q'$  the product distributions over  $(\Omega \times \Omega', \mathcal{F} \otimes \mathcal{F}')$ . Then

$$KL(P \otimes P', Q \otimes Q') = KL(P, Q) + KL(P', Q').$$

## Consequence (Ganivier, Merad, Stoltz, 2016)

Data-processing inequality with expectations of random variables.

Let  $X : (\Omega, \mathcal{F}) \rightarrow ([0, 1], \mathcal{B}([0, 1]))$  be any  $[0, 1]$ -valued random variable

Then, denoting by  $E_P[X]$  and  $E_Q[X]$  the respective expectations of  $X$  under  $P$  and  $Q$ , we have:

$$E_P[X] \ln \frac{E_P[X]}{E_Q[X]} + (1 - E_P[X]) \ln \frac{1 - E_P[X]}{1 - E_Q[X]} = KL(\text{Ber}(E_P[X]), \text{Ber}(E_Q[X])) \leq KL(P, Q).$$

Proof by upper-bounding  $KL((P \otimes \mu)^{\text{Ber}(E)}, (Q \otimes \mu)^{\text{Ber}(E)})$

The chain rule - A generalization of the decomposition of the KL between product-distributions.

We will need it in a special case only, when the joint distributions follow from one of the marginal distributions via a stochastic kernel.

Definition Let  $(\Omega, \mathcal{F})$  and  $(\Omega', \mathcal{F}')$  be two measurable spaces; we denote by  $\mathcal{P}(\Omega', \mathcal{F}')$  the set of probability measures over  $(\Omega', \mathcal{F}')$ .

A (regular) stochastic kernel  $K$  is a mapping  $(\Omega, \mathcal{F}) \rightarrow \mathcal{P}(\Omega', \mathcal{F}')$   
 $w \mapsto K(w, \cdot)$

such that  $\forall B \in \mathcal{F}'$ ,  $w \mapsto K(w, B)$  is  $\mathcal{F}$ -measurable

Now consider two such kernels  $K$  and  $L$ , and two probability measures  $P$  and  $Q$  over  $(\Omega, \mathcal{F})$ . Then  $KP$  and  $LP$  defined below are probability measures over  $(\Omega \times \Omega', \mathcal{F} \otimes \mathcal{F}')$ , by some extension theorem (Carathéodory)

$$\forall A \in \mathcal{F}, \forall B \in \mathcal{F}', \quad K(P(A \times B)) = \int_{\Omega} \underbrace{\mathbb{1}_A(w)}_{\text{is indeed measurable}} K(w, B) dP(w)$$

$$L(Q)(A \times B) = \int_{\Omega} \mathbb{1}_A(\omega) L(\omega, B) dQ(\omega)$$

↓  
actually with no loss of generality.

Theorem (chain rule for KL): Assume  $P \ll Q$

As soon as  $(*) K(w, \cdot) \ll L(w, \cdot)$  for  $Q$ -almost all  $w \in \Omega$

with  $(**)$  the existence of a function  $g: (\omega, \omega') \mapsto \frac{dK(\omega, \cdot)}{dL(\omega, \cdot)}(\omega')$   
being  $\mathcal{F}_\Omega \otimes \mathcal{F}'$ -measurable, Up to a  $LQ$ -null set.

Then

$$KL(KP, LQ) = KL(P, Q) + \int_{\Omega} KL(K(w, \cdot), L(w, \cdot)) dP(w)$$

where  $w \mapsto KL(K(w, \cdot), L(w, \cdot))$  is indeed  $\mathcal{F}$ -measurable and  $> 0$  so  
that the integral in the right-hand side is well defined.

### Remark:

- 1) The assumptions  $(*)$  and  $(**)$  will be satisfied for the applications we have in mind.
- 2) They can be relaxed: - it suffices to assume that  $\Omega'$  is a topological space with a countable base and it is the Borel  $\sigma$ -algebra.

i.e there exists some countable collection  $(O_m)_{m \geq 1}$  of open sets of  $\mathcal{S}'$  such that each open set  $V$  of  $\mathcal{S}'$  can be written

$$V = \bigcup_{i: O_i \subseteq V} O_i \text{, that is, as a countable union of elements of } (O_m)_{m \geq 1}.$$

$$(O_m)_{m \geq 1}.$$

E:  $\mathcal{S}'$  a separable metric space  $\rightarrow$  we will consider

$$\mathcal{S}' = [0, 1] \times (\mathbb{R} \times [0, 1])^{\mathbb{N}}$$

3) A typical kernel is given by  $K(w, B) = \Pr(Y \in B \mid X = X(w))$

The chain rule then rewrites:  $KL(P^{(x,y)}, Q^{(x,y)}) = KL(P^{(x,y)}, Q^{(x,y)}) + KL(P^{(y)}, Q^{(y)})$

Now we have stated the useful properties of the KL, let's get back to the lower bound.

Recall bandit setting

- to each arm  $k$  is associated a probability distribution

$$v_k \in \mathcal{D}$$

-  $\mathcal{D}$  is the bandit model ( $\mathcal{D} \subset P_1(\mathbb{R})$ )

- A bandit instance is denoted by  $\nu = (\nu_k)_{k \in [K]}$
- Goal, minimise the regret, which can be rewritten as:

$$R_T = \sum_{k=1}^K \Delta_k E[N_k(T)]$$

Bounding the regret  $\Leftrightarrow$  bounding  $E[N_k(T)]$

What are the best possible (by an algorithm) bounds?

- What is a randomised strategy  $\pi$ ?
- a sequence of measurable functions  $(\pi_t)_{t \geq 1}$  with

$$\pi_{t+1}: H_t = (U_0, X_{a_1}(1), U_1, \dots, X_{a_t}(t), U_t) \mapsto \underbrace{\pi_{t+1}(H_t)}_{\substack{\text{history of decisions + transformation for} \\ \text{first } t \text{ rounds}}} = a_{t+1} \quad \underbrace{\pi_{t+1}(H_t)}_{\substack{\text{arm picked at } t+1}}$$

Lemma: (Fundamental inequality for stochastic bandits)

For all bandit problems  $\nu = (\nu_k)_{k \in [K]}$  and  $\nu' = (\nu'_k)_{k \in [K]}$  in  $\mathcal{D}^K$  with  $\nu \leq \nu'$  for all  $k$ ,

for all strategies and random variables  $Z$  taking values in  $[0,1]$  that are  $\pi(H_t)$ -measurable,

law of  $H_t$  under  $\nu$  (and  $\pi$ )

$$\begin{aligned} \sum_{k=1}^K E[N_k(T)] & KL(\nu_k, \nu'_k) = KL(P_\nu^{H_T}, P_{\nu'}^{H_T}) \\ & \geq KL(Bn(E[Z]), Bn(E_{\nu'}[Z])) \end{aligned}$$

dependence in strategy  $\pi$  hidden everywhere here.

# Proof

The inequality 3 is a direct application of the data processing inequality with expectations.

For the equality:

(1) we show by induction that  $P_v^{H_T} = K_T(K_{T+1} \dots (K_1) \lambda_0)$

where  $K_T$  is the transition kernel:

$$h \in [0,1] \times (\mathbb{R} \times [0,1])^{T-1} \mapsto K_T(h, \cdot) = \nu_{H_T}(h) \otimes \lambda_0 \quad \text{with } \nu_T \sim \lambda_0$$

prob. measure on  $\mathbb{R} \times [0,1]$

$$T=0: H_0 = U_0 \sim \lambda_0 : P_v^{U_0} = \lambda_0$$

$$\underbrace{T, T+1}_{\vdots} \forall A \in \mathcal{B}([0,1] \times (\mathbb{R} \times [0,1]^T)), \forall B' \in \mathcal{B}(\mathbb{R}), \forall B \in \mathcal{B}([0,1]) :$$

$$\begin{aligned}
 P_v^{H_{T+2}}(A \times B' \times B) &= \mathbb{P}_v(H_T \in A \text{ and } X_{\alpha_{T+2}}(T+1) \in B' \text{ and } U_{T+2} \in B) \\
 &= \mathbb{E}_v[\mathbb{1}_A(H_T) \mathbb{P}_v[X_{\alpha_{T+2}}(T+1) \in B' \text{ and } U_{T+2} \in B | H_T]] \\
 &= \mathbb{E}_v[\mathbb{1}_A(H_T) \cdot \nu_{H_{T+2}(H_T)}(B') \cdot \lambda_0(B)] \\
 &= \mathbb{E}_v[\mathbb{1}_A(H_T) K_{T+2}(H_T, B' \times B)] \\
 &= \int \mathbb{1}_A(h) K_{T+2}(h, B' \times B) dP_v^{H_T}(h) \\
 &= K_{T+2} P_v^{H_T}(A \times B' \times B)
 \end{aligned}$$

$\downarrow$  defn of  $K_{T+2}$

$\downarrow$  rewriting

$\downarrow$  defn of  $K_{T+2}, P_v^{H_T}$

$\rightarrow$  we've shown the induction

(2) we check that the assumptions of the chain rule are satisfied.

- the  $K_r$  are regular transition kernels:  $\Pi \in C(B(\mathbb{R}) \otimes B([0,1]))$ ,

$$h \mapsto K_r(h, E) = \sum_{k=1}^K \mathbb{1}_{\{\pi_r(k)=h\}} (\nu_k \otimes \lambda_r)(E) \quad \text{is measurable as}$$

$\pi_r$  is measurable (with respect to considered spaces)

- Assumption (1):  $\forall h, K_r(h, \cdot) \ll K'_r(h, \cdot)$  as  $\forall k, \nu_k \ll \nu'_k$  by ass.

- Assumption (2):  $(h, (y, u)) \mapsto \frac{dK_r(h, \cdot)}{dK'_r(h, \cdot)}(y, u)$

$$= \sum_{a=1}^K \mathbb{1}_{\{\pi_r(h)=a\}} \frac{d\nu_a(y)}{d\nu'_a(y)}$$

is indeed bi-measurable (product of measurable functions)

(3) We then may apply the chain rule and show by induction the desired result based on:

$$- KL(P_\nu^{H_0}, P_{\nu'}^{H_0}) = KL(\lambda_0, \lambda_0) = 0$$

$$\begin{aligned} - \text{For } t \geq 0, \quad & KL(P_\nu^{H_{t+1}}, P_{\nu'}^{H_{t+1}}) = KL(K_t P_\nu^{H_t}, K'_{t+1} P_{\nu'}^{H_t}) \quad \xrightarrow{\text{chain rule}} \\ & = KL(P_\nu^{H_t}, P_{\nu'}^{H_t}) + \int KL(K_{t+1}(h, \cdot), K'_{t+1}(h, \cdot)) dP_\nu^{H_t}(h) \end{aligned}$$

$$= KL(P_\nu^{H_t}, P_{\nu'}^{H_t}) + \int KL(\nu_{\pi_{t+1}(h)} \otimes \lambda_0, \nu'_{\pi_{t+1}(h)} \otimes \lambda_0) dP_\nu^{H_t}(h)$$

$$= KL(P_{\pi}^{H_T}, P_{\nu}^{H_T}) + \sum_{k=1}^K KL(\nu_k, \nu'_k) \cdot \underbrace{\int 1_{\tilde{a}_{T+1}(h)=k} dP_{\nu}^{H_T}(h)}_{\mathbb{E}[1_{\tilde{a}_{T+1}(H_T)=k}]} \\ = \mathbb{E}[1_{\tilde{a}_{T+1}=k}]$$

$$= KL(P_{\pi}^{H_T}, P_{\pi}^{H_T}) + \sum_{k=1}^K KL(\nu_k, \nu'_k) \mathbb{E}[1_{a_{T+1}=k}]$$

by induction

$$KL(P_{\pi}^{H_T}, P_{\nu}^{H_T}) = \sum_{t=1}^T \sum_{h=1}^K KL(\nu_a, \nu'_{a'}) \mathbb{E}[1_{a_t=k}]$$

$$= \sum_t KL(\nu_a, \nu'_{a'}) \mathbb{E}[N_a(t)]$$

□

We are now equipped to prove the lower bound.

- Defn |**
- a strategy is consistent w.r.t a model  $\mathcal{D}$  if,  
for all bandit instances  $\nu \in \mathcal{D}^K$ ,  $\forall \alpha \in (0, 1]$ ,  $\forall k$  s.t.  $D_k > 0$ ,
- $$\mathbb{E}[N_k(T)] = o(T^\alpha)$$

for well behaved models, there exist consistent strategies  
e.g. UCB with  $\mathcal{D} = \mathcal{P}(0, 1)$ .

- (asymptotic)
- typical bounds for good strategies
- $\forall \nu \in \mathcal{D}^K, \exists \delta, \Delta_\nu > 0, \limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_k(T)]}{\ln T} \leq C_k(\nu)$
- optimal such term:  $C_k(\nu) = \frac{1}{\text{Kinf}(\nu_k, \mu^*, \mathcal{D})}$  problem dependent term

$$\text{when } \text{Kinf}(\nu_k, \mu^*, \mathcal{D}) = \inf \left\{ \text{KL}(\nu_k \| \nu'_k) \mid \begin{array}{l} \nu'_k \in \mathcal{D} \\ \mathbb{E}(\nu'_k) > \mu^* \end{array} \right\}$$

we will now prove one part of this optimality: a lower bound on  $C_k(\nu)$ .

Theorem (Lai and Robbins, 1985,  
Biandas and Katselis, 1996)

For all bandit models  $\mathcal{D} \in \mathcal{P}_s(\mathbb{R})$ ,

for any consistent strategy wrt  $\mathcal{D}$ ,

for any bandit instance  $\nu \in \mathcal{D}^K$ ,

for all suboptimal arms  $k$  ( $i.e. \Delta_k > 0$ ),  $\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_k(T)]}{\ln T} \geq \frac{1}{\text{Kinf}(\nu_k, \mu^*, \mathcal{D})}$

### Corollary

for all bandit models  $\mathcal{D}$ , any consistent strategy wrt  $\mathcal{D}$ , all bandit instances  $\nu \in \mathcal{D}^K$ :

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[R_T]}{\ln T} \geq \sum_{\substack{k \\ \Delta_k > 0}} \frac{\Delta_k}{\text{Kinf}(\nu_k, \mu^*, \mathcal{D})}$$

## Proof of the Theorem (based on the Lemma)

$$K_{\text{Inf}}(v_a, D, \mu^*) = \inf \left\{ KL(v_a, v'_a) \mid v'_a \in D, v_a \ll v'_a \text{ and } \mathbb{E}(v') > \mu^* \right\}.$$

Convention if  $\phi = +\infty$

This is why we will:

- fix  $D$ , strategy  $T$ ,  $v$  and  $\delta$  s.t.  $\Delta_\delta > 0$  ( $T$  is consistent w.r.t  $D$ )

- fix an alternative model  $v'$  with

$$\begin{cases} v'_i = v_i & \text{for all } i \neq k \\ v'_k \text{ s.t. } v'_k \in D, v_k \ll v'_k \text{ and } \mathbb{E}(v'_k) > \mu^* \end{cases}$$

That  $v$  and  $v'$  only differ at  $k$ , the unique optimal arm in  $v'$ .

- Take  $\tau = \frac{N_{a(T)}}{T}$  which is  $[0, 1]$ -valued  
 $\sigma(H_T)$ -measurable

Our fundamental inequality (Lemma) yields, since  $v$  and  $v'$  only differ at  $k$ :

$$\mathbb{E}_v[N_{a(T)}] \cdot KL(v_a, v'_a) \geq KL(Bn(\mathbb{E}_v[\frac{N_{a(T)}}{T}]), Bn(\mathbb{E}_{v'}[\frac{N_{a(T)}}{T}]))$$

$$\geq -\ln(2) + \left(1 - \mathbb{E}_v\left[\frac{N_{a(T)}}{T}\right]\right) \ln\left(\frac{1}{1 - \mathbb{E}_{v'}\left[\frac{N_{a(T)}}{T}\right]}\right)$$

$$\text{Indeed } KL(Bn(p), Bn(q)) = p \ln\left(\frac{p}{q}\right) + (1-p) \ln\left(\frac{1-p}{1-q}\right)$$

$$= p \ln\left(\frac{1}{q}\right) + (1-p) \ln\left(\frac{1}{1-q}\right) + (p \ln(p) + (1-p) \ln(1-p))$$

$\geq 0$  $> -\ln 2$ 

$$\geq -\ln 2 + (1-p) \ln\left(\frac{1}{1-p}\right) \text{ for all } (p, q) \in [0, 1] \text{ (and even for } p \in [0, 1])$$

$\pi$  is consistent, so

- instance  $\rightarrow \pi$  is suboptimal  $E_{\pi} \left[ \frac{N_a(T)}{T} \right] \xrightarrow{T \rightarrow \infty} 0$

- instance  $v' \rightarrow$  all  $i \neq k$  are suboptimal

$$\text{for any } \alpha \in (0, 1], E_{v'}[N_a(T)] = o(T^\alpha)$$

$$\text{In particular: } T \cdot E_{v'}[N_a(T)] = \sum_{i \neq k} E_{v'}[N_k(T)] = o(T^\alpha)$$

$$\text{so: } \frac{1}{1 - E_{v'} \left[ \frac{N_a(T)}{T} \right]} = \frac{T}{T - E_{v'}[N_a(T)]} = \frac{T}{o(T^\alpha)}$$

$$\geq T^{1-\alpha} \text{ for } T \text{ large enough}$$

Substituting back and dividing by  $\ln T$ : for any  $\alpha \in (0, 1]$  and  $T$  large enough

$$\frac{E_{v'}[N_a(T)]}{\ln T} \xrightarrow{\text{KL}(v_e, v_e)} -\frac{\ln 2}{\ln T} + \left(1 - E_{v'} \left[ \frac{N_a(T)}{T} \right]\right) \frac{\ln T^{1-\alpha}}{\ln T}$$

$$\text{Thus } \liminf_{T \rightarrow +\infty} \frac{E_{v'}[N_a(T)]}{\ln T} \geq \frac{(1-\alpha)}{\text{KL}(v_e, v_e)} \quad \begin{array}{l} (\text{true whether the KL is } < \infty \\ \text{or } +\infty \\ (\text{true necessarily } > 0)) \end{array}$$

$$\text{for any } \alpha \in (0, 1], \liminf_{T \rightarrow +\infty} \frac{E_{v'}[N_a(T)]}{\ln T} \geq \frac{1}{\text{KL}(v_e, v_e)}$$

Holds for any  $v'_a \in \mathcal{D}$  s.t.  $v_a < v'_a$  and  $\mathbb{E}(v'_a) > \mu^*$ , so that taking the supremum of the right hand side on these  $v'_a$  yields the lower bound:

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\ln T} \geq \frac{1}{\text{Inf}_{v'_a \in \mathcal{D}}(\mu^*)}$$

### Comments on the lower bound:

- see a comparison with an upper bounds in exercise sheet
- algorithms with optimal instance dependent bounds are known (e.g. KL-UCB, Thompson sampling), but requires a long and technical analysis.
- This is an asymptotic lower bound for  $T \rightarrow \infty$ . What about small  $T$ ? → see exercises sheet

What if we fix  $T$  and choose arbitrarily the bandit instance  $v$ ?

### Theorem (minimax lower bound)

Let  $\mathcal{D} = \{N(\mu, 1) \mid \mu \in \mathbb{R}\}$ ,  $K \geq 2$  and  $T \geq K \cdot 1$ . Then, there exists a universal constant  $c > 0$  such that,

for any policy  $\pi$ , there exists  $v \in \mathcal{D}^K$  s.t.

$$\mathbb{E}[R_T(\pi, v)] \geq c \sqrt{KT}$$

Proof in exercise below

$$\text{minimax} \leftarrow \min_{\pi} \max_{v \in \mathcal{D}^K} R_T(\pi, v) \geq c \sqrt{KT}$$

This exercise aims at showing a minimax lower bound of the regret of the form  $R_T \geq c\sqrt{KT}$ . We restrict ourselves to the bandit model  $\mathcal{D} = \{\mathcal{N}(\mu, 1) \mid \mu \in \mathbb{R}\}$ , but similar arguments can be used for more general models (e.g. Bernoulli bandits). Fix in the following  $K \geq 2$  and  $T \geq \frac{K-1}{2}$ . The minimax regret is defined as

$$R_T^* = \inf_{\text{strategy } \pi} \sup_{\text{instance } \nu} \mathbb{E}[R_T(\pi, \nu)]$$

Let  $\varepsilon > 0$ . We consider the following  $K + 1$  bandit instances  $(\nu^j)_{j \in [K+1]}$ , where

$$\begin{aligned} \nu_k^j &= \mathcal{N}(0, 1) \quad \text{for any } k \in [K] \text{ such that } j \neq k \\ \nu_k^k &= \mathcal{N}(\varepsilon, 1) \quad \text{for any } k \in [K]. \end{aligned}$$

1) Justify that

$$R_T^* \geq \inf_{\pi} \sup_{\varepsilon \in (0, 1)} \max_{i \in [K]} \varepsilon(T - \mathbb{E}_{\nu^i}^{\pi}[N_i(T)]),$$

and that for any strategy  $\pi$ , there exists  $k_0$  such that  $\mathbb{E}_{\nu^{k_0}}[N_{k_0}(T)] \leq \frac{T}{K}$ .

2) Use the fundamental inequality and Pinsker's inequality to show that

$$\mathbb{E}_{\nu^{k_0}}[N_{k_0}(T)] \frac{\varepsilon^2}{2} \geq 2 \left( \mathbb{E}_{\nu^{k_0}}\left[\frac{N_{k_0}(T)}{T}\right] - \mathbb{E}_{\nu^{k_0}}\left[\frac{N_{k_0}(T)}{T}\right] \right)^2.$$

3) Combine the above results to derive

$$R_T^* \geq \sup_{\varepsilon \in (0, 1)} \varepsilon T \left( 1 - \frac{1}{K} - \varepsilon \sqrt{\frac{T}{2K}} \right)$$

and conclude that  $R_T^* \geq \frac{1}{8\sqrt{2}} \sqrt{KT}$ .

**Solution:** 1) The first point is just taking a subset over all the possible instances and rewriting the regret as  $\varepsilon(T - \mathbb{E}_{\nu^i}^{\pi}[N_i(T)])$ . The second point is because the sum over all  $k$  is equal to  $T$ , so at least one of them is below (or equal) to the average.

2) Direct use of fundamental inequality and Pinsker's inequality, with the fact that for  $i \neq 0$ ,  $\text{KL}(\nu^0, \nu^i) = \frac{\varepsilon^2}{2}$ .

3) We have

$$\mathbb{E}_{\nu^0}[N_{k_0}(T)] \frac{\varepsilon^2}{2} \geq 2 \left( \mathbb{E}_{\nu^0}\left[\frac{N_{k_0}(T)}{T}\right] - \mathbb{E}_{\nu^{k_0}}\left[\frac{N_{k_0}(T)}{T}\right] \right)^2.$$

So that  $\mathbb{E}_{\nu^{k_0}}\left[\frac{N_{k_0}(T)}{T}\right] \leq \mathbb{E}_{\nu^0}\left[\frac{N_{k_0}(T)}{T}\right] + \varepsilon \sqrt{\frac{\mathbb{E}_{\nu^0}[N_{k_0}(T)]}{2}}$ . Since  $\mathbb{E}_{\nu^0}[N_{k_0}(T)] \leq \frac{T}{K}$ , this implies

$$\mathbb{E}_{\nu^{k_0}}\left[\frac{N_{k_0}(T)}{T}\right] \leq \frac{1}{K} + \varepsilon \sqrt{\frac{T}{2K}}.$$

This then yields

$$R_T(\pi, \nu^{k_0}) \geq \varepsilon T \left( 1 - \frac{1}{K} - \varepsilon \sqrt{\frac{T}{2K}} \right).$$

Using question 1), this yields the first point, by noticing that the bound actually does not depend on the choice of the strategy  $\pi$ . We conclude by taking  $\varepsilon = \sqrt{\frac{K(K-1)^2}{2T}} \frac{1}{K^2}$ , which is smaller than 1 since  $T \geq \frac{K-1}{2}$ , and noticing that for  $K \geq 2$ ,  $\frac{K-1}{K} \geq \frac{1}{2}$ .