We first complete the <u>proof of the distribution-dependent lower bound.</u>

<u>Protocol</u>: with randomized strategies $(\Psi_t)_{t \geq 0}$, where

$$\Psi_t : \quad H_t = (U_0, Y_1, U_1, \ldots, Y_t, U_t) \longmapsto \Psi_t(H_t) = I_{t+1}$$

We have : $\begin{cases} Y_t \mid I_t \sim \nu_{I_t} & \forall t \\ U_0, U_1, \ldots & \text{iid} \sim \mathcal{U}_{[0,1]} \end{cases}$

<u>Regret</u>: $\qquad \overline{R}_T = T\mu^* - \mathbb{E}\left[\sum_{t=1}^{T} Y_t\right] = \sum_{a=1}^{K} \Delta_a \, \mathbb{E}_\nu[N_a(T)]$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ by the tower rule

where $\quad \mu_a = \mathbb{E}(\nu_a), \quad \mu^* = \max\limits_{a=1\ldots K} \mu_a,$

$\qquad\qquad \Delta_a = \mu^* - \mu_a \qquad\qquad\qquad\qquad$ and $\quad \mathbb{E}_\nu$ denotes the expectation when the underlying bandit problem is $\nu = (\nu_a)_{a \in \{1 \ldots K\}}$

$\hookrightarrow$ To control the regret (upper or lower bound it), it is thus necessary and sufficient to control the $\mathbb{E}_\nu[N_a(T)]$.

<u>Definitions</u>:

- Bandit model $\mathcal{D}$ : A subset of the set $\mathcal{M}_1(\mathbb{R})$ of all probability distributions over $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ with an expectation.

$\qquad\qquad\qquad\qquad\quad$ The decision-maker knows $\mathcal{D}$ but does not know the specific bandit problem $\nu = (\nu_a)_{a \in \{1 \ldots K\}} \in \mathcal{D}^K$ at hand.

- Strategy $\Psi$ consistent wrt. a model $\mathcal{D}$ : when $\forall \nu \in \mathcal{D}^K, \forall \alpha \in (0,1], \forall a \text{ s.t. } \Delta_a > 0,$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \mathbb{E}_\nu[N_a(T)] = o(T^\alpha)$

<u>Lemma</u> (Fundamental inequality for stochastic bandits) : For all strategies $\Psi$,

For all bandit problems $\nu = (\nu_a)_{a \in \{1 \ldots K\}}$ and $\nu' = (\nu'_a)_{a \in \{1 \ldots K\}}$ in $\mathcal{D}^K$ with $\nu_a \ll \nu'_a$ for all $a$,

For all random variables $Z$ taking values in $[0,1]$ and that are $\sigma(H_T)$-measurable,

$$\sum_{a=1}^{K} \mathbb{E}_\nu[N_a(T)] \, \mathrm{KL}(\nu_a, \nu'_a) = \mathrm{KL}(\mathbb{P}_\nu^{H_T}, \mathbb{P}_{\nu'}^{H_T})$$

$$\geq \mathrm{KL}\big(\mathrm{Ber}(\mathbb{E}_\nu[Z]), \mathrm{Ber}(\mathbb{E}_{\nu'}[Z])\big)$$

where $\mathbb{P}_\nu^{H_T}$ and $\mathbb{P}_{\nu'}^{H_T}$ denote the laws of $H_T$ when the strategy is $\Psi$ and when the underlying bandit problems are respectively $\nu$ and $\nu'$.

Proof:   • The inequality $\geqslant$ is a direct application of the data-processing inequality with expectations, see the previous lecture for its statement.

• For the equality:          We will explain how $\mathbb{P}_{\gamma}^{H_T}$ is constructed.

With no loss of generality, we can consider that

— the underlying probability space is $\Omega = [0,1] \times (\mathbb{R} \times [0,1])^T$

— $H_T$ is the identity over $\Omega$, i.e., that the $U_0, Y_1, U_1, \ldots Y_T, U_T$ are the projections on each component,

— $\mathbb{P}_{\gamma}$ is given by

$$\forall B \in \mathcal{B}([0,1]), \qquad \mathbb{P}_{\gamma}(U_0 \in B) = \eta(B)$$

$$\forall t \in \{0, \ldots T-1\}, \quad \forall B' \in \mathcal{B}(\mathbb{R}), \quad \forall B \in \mathcal{B}([0,1]), \qquad \mathbb{P}_{\gamma}(Y_{t+1} \in B' \text{ and } U_{t+1} \in B \mid H_t) = \gamma_{\psi_t(H_t)}(B') \, \eta(B)$$

where $\mathcal{B}(S)$ is the Borel-$\sigma$-algebra of a set $S \subseteq \mathbb{R}$
$\eta$ is the Lebesgue measure over $[0,1]$

In particular: $\mathbb{P}_{\gamma}^{H_t}$ refers to the first $1+2t$ marginals of $\mathbb{P}_{\gamma}^{H_T} = \mathbb{P}_{\gamma}$.
A similar construction can be done for the bandit problem $\gamma'$.

Now, the equality $\mathbb{P}_{\gamma}(Y_{t+1} \in B' \text{ and } U_{t+1} \in B \mid H_t) = \gamma_{\psi_t(H_t)}(B') \, \eta(B)$

indicates that $\mathbb{P}_{\gamma}^{H_{t+1}} = K_t \, \mathbb{P}_{\gamma}^{H_t}$ for the regular transition kernel

$$K_t(h, \cdot) = \gamma_{\psi_t(h)} \otimes \eta$$

regularity is:
for $E \in \mathcal{B}(\mathbb{R}) \otimes \mathcal{B}([0,1])$
$h \mapsto K_t(h, E)$ is measurable

Similarly, $\mathbb{P}_{\gamma'}^{H_{t+1}} = K_t' \, \mathbb{P}_{\gamma'}^{H_t}$ for $K_t'(h, \cdot) = \gamma'_{\psi_t(h)} \otimes \eta$

Let us check the assumptions of our chain rule:

$(*)$      $\forall h, \qquad K_t(h, \cdot) \ll K_t'(h, \cdot)$          as $\forall a, \gamma_a \ll \gamma_a'$ by assumption

$(**)$     $(h, (y,u)) \longmapsto \dfrac{dK_t(h, \cdot)}{dK_t'(h, \cdot)}(y,u) = \sum\limits_{a=1}^{K} \mathbb{1}_{\{\psi_t(h) = a\}} \dfrac{d\gamma_a}{d\gamma_a'}(y)$
is indeed bi-measurable.

Therefore, for $t \in \{0, \ldots T-1\}$,

$$KL\left(\mathbb{P}_{\gamma}^{H_{t+1}}, \mathbb{P}_{\gamma'}^{H_{t+1}}\right) = KL\left(\mathbb{P}_{\gamma}^{H_t}, \mathbb{P}_{\gamma'}^{H_t}\right) + \int KL\left(\gamma_{\psi_t(h)} \otimes \eta, \gamma'_{\psi_t(h)} \otimes \eta\right) d\mathbb{P}_{\gamma}^{H_t}(h)$$

$$= KL\left(\mathbb{P}_{\gamma}^{H_t}, \mathbb{P}_{\gamma'}^{H_t}\right) + \sum\limits_{a=1}^{K} KL\left(\gamma_a, \gamma_a'\right) \, \mathbb{P}_{\gamma}^{H_t}\{\psi_t(h) = a\}$$

Now, $\quad I_{t+1} = \Psi_t(H_t) \quad$ so that

$$\mathbb{P}_{\vec{\gamma}}^{H_t}\{\Psi_t(h) = a\} = \mathbb{P}_{\vec{\gamma}}\{\Psi_t(H_t) = a\} = \mathbb{P}_{\vec{\gamma}}\{I_{t+1} = a\}$$

$$= \mathbb{E}_{\vec{\gamma}}\left[\mathbb{1}_{\{I_{t+1} = a\}}\right]$$

Summing up:

- $\quad KL\left(\mathbb{P}_{\vec{\gamma}}^{H_0}, \mathbb{P}_{\vec{\gamma}'}^{H_0}\right) = KL\left(\mathbb{P}_{\vec{\gamma}}^{U_0}, \mathbb{P}_{\vec{\gamma}'}^{U_0}\right) = KL(\eta, \eta) = 0$

- $\forall t \in \{0, \dots T-1\}, \quad KL\left(\mathbb{P}_{\vec{\gamma}}^{H_{t+1}}, \mathbb{P}_{\vec{\gamma}'}^{H_{t+1}}\right) = KL\left(\mathbb{P}_{\vec{\gamma}}^{H_t}, \mathbb{P}_{\vec{\gamma}'}^{H_t}\right)$

$$+ \sum_{a=1}^{K} KL(\vec{\gamma}_a, \vec{\gamma}_a') \; \mathbb{E}_{\vec{\gamma}}\left[\mathbb{1}_{\{I_{t+1} = a\}}\right]$$

so that the stated result follows by induction.

<u>Back to the theorem</u> giving the lower bound:

<u>Theorem :</u>  For all bandit models $\mathcal{D} \subset \mathcal{M}_1(\mathbb{R})$,

For all (possibly randomized) strategies $\Psi$ consistent wrt $\mathcal{D}$,

For all bandit problems $\vec{\gamma} = (\vec{\gamma}_a)_{a \in \{1..K\}} \in \mathcal{D}^K$,

For all suboptimal arms $a$ (ie, arms $a$ with $\Delta_a > 0$),

$$\liminf_{T \to +\infty} \frac{\mathbb{E}_{\vec{\gamma}}[N_a(T)]}{\ln T} \geq \frac{1}{K_{inf}(\vec{\gamma}_a, \mu^*, \mathcal{D})}$$

where $\quad K_{inf}(\vec{\gamma}_a, \mu^*, \mathcal{D}) = \inf\left\{KL(\vec{\gamma}_a, \vec{\gamma}_a') : \vec{\gamma}_a' \in \mathcal{D} \text{ with } \mathbb{E}(\vec{\gamma}_a') > \mu^*\right\}$

with the convention $\inf \emptyset = +\infty$.

<u>Correction for Exercise #3 :</u>     The proof of this theorem based on the lemma above was your Exercise #2.

See next page.

**Proof:**    We have

$$K_{inf}(\nu_a, \mu^*) = \inf\{KL(\nu_a, \nu'_a) : \nu'_a \in \mathcal{D} \text{ and } E(\nu'_a) > \mu^*\}$$

$$= \inf\{KL(\nu_a, \nu'_a) : \nu'_a \in \mathcal{D}, \ \nu_a \ll \nu'_a \text{ and } E(\nu'_a) > \mu^*\}$$

( cf. convention: $\inf \emptyset = +\infty$ and the fact that $KL(\nu_a, \nu'_a) = +\infty$
when $\nu_a \not\ll \nu'_a$ )

This is why we will

- Fix $\mathcal{D}, \Psi, \nu$ and $a$ s.t. $\Delta_a > 0$

- Fix an alternative model $\nu''$ of the form

$$\begin{cases} \nu''_k = \nu_k & \forall k \neq a \\ \nu''_a & \text{s.t.} \ \nu''_a \in \mathcal{D}_j, \ \nu_a \ll \nu''_a \text{ and } E(\nu''_a) > \mu^* \end{cases}$$

That is, $\nu$ and $\nu''$ only differ at $a$; $a$ is the unique optimal arm in $\nu''$

- Take $Z = N_a(T)/T$ which is indeed $[0,1]$ — valued
$\sigma(H_T)$ — measurable

Our fundamental inequality yields, since $\nu$ and $\nu''$ only differ at $a$:

$$E_\nu[N_a(T)] \ KL(\nu_a, \nu''_a) \geq KL\left(Ber\left(E_\nu\left[N_a(T)/T\right]\right), Ber\left(E_{\nu''}\left[\frac{N_a(T)}{T}\right]\right)\right)$$

$$\geq -\ln 2 + \left(1 - E_\nu[N_a(T)/T]\right) \ln \frac{1}{1 - E_{\nu''}[N_a(T)/T]}$$

indeed: $KL(Ber(p), Ber(q))$

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

$$= \underbrace{p \ln \frac{1}{q}}_{> 0} + (1-p) \ln \frac{1}{1-q} + \underbrace{\left(p \ln p + (1-p) \ln(1-p)\right)}_{\geq -\ln 2}$$

$$\geq -\ln 2 + (1-p) \ln \frac{1}{1-q}$$

for all $p, q \in (0,1)$ and even for all $p, q \in [0,1]$ ( study the cases
$q=0$ and $q=1$ separately )

Now, the considered strategy $\Psi$ is consistent and:

- in the problem $\nu$, $a$ is suboptimal: $E_\nu[N_a(T)/T] \longrightarrow 0$

— in the problem $\nu''$, all arms $k \neq a$ are suboptimal;

for all $\alpha \in (0,1]$, $\qquad T - E_{\nu''}[N_a(T)] = \sum_{k \neq a} E_{\nu''}[N_k(T)] = o(T^\alpha)$

$\hookrightarrow$ in particular, for $T$ large enough,

$$\frac{1}{1 - E_{\nu''}[N_a(T)/T]} = \frac{T}{T - E_{\nu''}[N_a(T)]} \geqslant \frac{T}{T^\alpha} = T^{1-\alpha}$$

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

$$\frac{E_\nu[N_a(T)]}{\ln T} KL(\bar{\nu}_a, \bar{\nu}'_a) \geqslant -\frac{\ln 2}{\ln T} + \underbrace{\left(1 - E_\nu\left[\frac{N_a(T)}{T}\right]\right)}_{\to 0} \underbrace{\frac{\ln T^{1-\alpha}}{\ln T}}_{= 1-\alpha}$$

thus

$$\liminf_{T \to +\infty} \frac{E_\nu[N_a(T)] KL(\bar{\nu}_a, \bar{\nu}'_a)}{\ln T} \geqslant 1 - \alpha$$

Letting $\alpha \to 0$,

$$\liminf_{T \to +\infty} \frac{E_\nu[N_a(T)] KL(\bar{\nu}_a, \bar{\nu}'_a)}{\ln T} \geqslant 1$$

Whether $KL(\bar{\nu}_a, \bar{\nu}'_a) < +\infty$ or $= +\infty$, we thus get

$$\liminf_{T \to +\infty} \frac{E_\nu[N_a(T)]}{\ln T} \geqslant \frac{1}{KL(\bar{\nu}_a, \bar{\nu}'_a)}$$

The left-hand side is independent of $\nu''_a \in \mathcal{D}$ s.t. $\bar{\nu}'_a \gg \bar{\nu}_a$ and $E(\nu''_a) > \mu^*$, so that taking the supremum of the right-hand side over these $\nu''_a$, we get the desired $1/K_{inf}(\bar{\nu}_a, \mu^*)$ lower bound.

Remark: Non-asymptotic lower bounds are possible ( for « super-consistent » strategies and « well-behaved » models), but are heavily technical.

Exercise 1.          $1/K_{inf}(\bar{\nu}_a, \mu^*, \mathcal{D})$          vs.          $8/\Delta_a^2$     for UCB

Recall that in the model $\mathcal{D} = \mathcal{P}([0,1])$, the UCB algorithm enjoys the following performance bound:

$$\forall \vec{\nu} \in \mathcal{P}([0,1])^K, \qquad \forall a \text{ s.t. } \Delta_a > 0,$$

$$\mathbb{E}_{\vec{\nu}}[N_a(T)] \leq \frac{8}{\Delta_a^2} \ln T + 2.$$

Actually, there are refinements of UCB that get the distribution-dependent constant $8/\Delta_a^2$ arbitrarily close to $2/\Delta_a^2$.

But how do those $8/\Delta_a^2$ and $2/\Delta_a^2$ constants compare to $1/K_{inf}(\bar{\nu}_a, \mu^*, \mathcal{P}([0,1]))$?

(1)     For $p, q \in [0,1]$, we denote

$$kl(p,q) = KL(Ber(p), Ber(q))$$

Show that          $\forall (p,q) \in [0,1]^2, \qquad kl(p,q) \geq 2(p-q)^2.$

(2)     Show Pinsker's inequality:          let $(\Omega, \mathcal{F})$ be a measurable space,

let $\mathbb{P}, \mathbb{Q}$ be two distributions over $(\Omega, \mathcal{F})$, then:

$$\|\mathbb{P} - \mathbb{Q}\|_{TV} \underset{\uparrow}{=} \sup_{A \in \mathcal{F}} |\mathbb{P}(A) - \mathbb{Q}(A)| \leq \sqrt{\frac{1}{2} KL(\mathbb{P}, \mathbb{Q})}$$

the total variation
distance between $\mathbb{P}$ and $\mathbb{Q}$

Even better,   show the   stronger form:

$$\sup_{\substack{Z, \mathcal{F}\text{-measurable} \\ \text{taking values} \\ \text{in } [0,1]}} |\mathbb{E}_{\mathbb{P}}[Z] - \mathbb{E}_{\mathbb{Q}}[Z]| \leq \sqrt{\frac{1}{2} KL(\mathbb{P}, \mathbb{Q})}$$

(3)     Exhibit a lower bound on $K_{inf}(\bar{\nu}_a, \mu^*, \mathcal{P}([0,1]))$ and conclude that some work is needed to get an upper bound matching our lower bound!

Lower bound on the regret for adversarial bandits

$\longleftrightarrow$ Minimax lower bound on the regret for stochastic bandits

We discussed so far $E_\nu[N_a(T)]$ in terms of distribution-dependent bounds like

$$C_a(\nu) \ln T$$

What about distribution-free / minimax bounds?

To that end we restrict our attention to the model $\mathcal{D} = \mathcal{P}([0,1])$, the set of all probability distributions over $[0,1]$.

| Stochastic bandits | Adversarial bandits |
|---|---|
| With each arm $a$ is associated $\nu_a \in \mathcal{P}([0,1])$ | An opponent selects the payoffs $g_{k,t}$ |
| For $t = 1, 2 \dots$ | For $t = 1, 2 \dots$ |
| — The decision maker picks $I_t \in \{1..K\}$ | — The opponent picks $(g_{1,t} \dots g_{K,t}) \in [0,1]^K$ while, simultaneously, |
| — Her reward $Y_t$, which is such that $Y_t \mid I_t \sim \nu_{I_t}$, is her only piece of information | — The decision-maker picks $I_t \in \{1..K\}$ |
| | — Her payoff is $g_{I_t,t}$ and this is the only piece of information she gets |
| Aim: control the regret | Aim: control the regret |
| $$\bar{R}_T = T \max_{a=1..K} E(\nu_a) - E\left[\sum_{t=1}^{T} Y_t\right]$$ | $$R_T = \max_{k=1..K} \sum_{t=1}^{T} g_{k,t} - \sum_{t=1}^{T} g_{I_t,t}$$ |

Typical adversarial results   ( Auer, Cesa-Bianchi, Freund, Schapire, 2002, later improved by Audibert and Bubeck, 2009):

Strategies such that for all opponents picking gains in $[0,1]$,   for all $T \geq 1$,

{ with probability at least $1-\delta$,            $R_T \leq C \sqrt{TK \ln(K/\delta)}$

{   $E[R_T] \leq C \sqrt{TK \ln K}$

where                                                         for some numerical constant $C$
the probability and $E$
are wrt decision-maker's internal
randomization

For "oblivious" opponents (ie when the $g_{k,t}$ do not "react" to        the $\sqrt{\ln K}$ can
the decision-maker's actions    ):        be dropped.

It is in particular
the case when $g_{k,t} \sim \nu_k$ $\forall t$ in an independent way

In this stochastic model:

$$\mathbb{E}[R_T] = \mathbb{E}\left[\max_{k=1..K} \sum_{t=1}^{T} g_{k,t}\right] - \mathbb{E}\left[\sum_{t=1}^{T} g_{I_t,t}\right] \quad \Big)\ g_{I_t,t} \text{ is } y_t$$

$$\geqslant T \max_{k=1..K} \mathbb{E}[g_{k,1}] - \mathbb{E}\left[\sum_{t=1}^{T} y_t\right]$$

$$= T \max_{k=1..K} \mathbb{E}(\gamma_k) - \mathbb{E}\left[\sum_{t=1}^{T} y_t\right] = \overline{R}_T$$

The adversarial results entail in particular that there exists a strategy of the decision-maker such that

$$\sup_{\gamma_1 \dots \gamma_K \in \mathcal{P}([0,1])} \overline{R}_T \leqslant \sup_{\substack{\text{opponents} \\ \text{picking } g_{k,t} \in [0,1]}} \mathbb{E}[R_T] \leqslant C\sqrt{TK} \qquad \text{for some numerical constant } C$$

Lower bound:

Theorem: For all (randomized) strategies of the decision-maker, for all $K \geqslant 2$ and $T \geqslant K \ln 2 / 2$,

$$\sup_{\substack{\text{opponents} \\ \text{picking } g_{j,t} \in [0,1]}} \mathbb{E}[R_T] \geqslant \sup_{\substack{\text{individual} \\ \text{sequences } g_{k,t} \in [0,1]}} \mathbb{E}[R_T] \geqslant \sup_{\gamma_1, \dots \gamma_K \in \mathcal{P}([0,1])} \overline{R}_T \geqslant \frac{1}{20}\sqrt{TK}$$

the only inequality we need to prove!

Proof: (i) Bandit problem $\gamma^{(0)} = (\text{Ber}(\tfrac{1}{2}), \dots \text{Ber}(\tfrac{1}{2}))$

versus bandit problems $\gamma^{(i)} = (\text{Ber}(\tfrac{1}{2}), \dots \text{Ber}(\tfrac{1}{2}), \text{Ber}(\tfrac{1}{2}+\varepsilon), \text{Ber}(\tfrac{1}{2}), \dots)$

$\qquad\qquad$ for $i \in \{1..K\}$ $\qquad\qquad$ in $i$-th position
$\qquad\qquad$ and $\quad \varepsilon \in (0, \tfrac{1}{2})$

There exists $k \in \{1..K\}$ such that $\mathbb{E}_{\gamma^{(0)}}[N_k(T)] \leqslant T/K$. For this $k$:

The fundamental inequality indicates that

$$\sum_{a=1}^{K} \mathbb{E}_{\gamma^{(0)}}[N_a(T)] \, KL\left(\gamma_a^{(0)}, \gamma_a^{(k)}\right) \geqslant kl\left(\mathbb{E}_{\gamma^{(0)}}\left[\frac{N_k(T)}{T}\right], \mathbb{E}_{\gamma^{(k)}}\left[\frac{N_k(T)}{T}\right]\right)$$

$$= \mathbb{E}_{\gamma^{(0)}}[N_k(T)] \, KL\left(\gamma_k^{(0)}, \gamma_k^{(k)}\right) \leqslant \frac{T}{K} \, kl\left(\tfrac{1}{2}, \tfrac{1}{2}+\varepsilon\right)$$

Using Pinsker's inequality:

$$kl\left(\mathbb{E}_{\gamma^{(0)}}\left[\frac{N_k(T)}{T}\right], \mathbb{E}_{\gamma^{(k)}}\left[\frac{N_k(T)}{T}\right]\right)$$

$$\geqslant 2\left(\mathbb{E}_{\gamma^{(0)}}\left[\frac{N_k(T)}{T}\right] - \mathbb{E}_{\gamma^{(k)}}\left[\frac{N_k(T)}{T}\right]\right)^2$$

Solving for $\mathbb{E}_{\gamma^{(k)}}\left[\frac{N_k(T)}{T}\right]$:

$$\mathbb{E}_{\gamma^{(k)}}\left[\frac{N_k(T)}{T}\right] \leq \mathbb{E}_{\gamma^{(0)}}\left[\frac{N_k(T)}{T}\right] + \sqrt{\frac{T}{2K}\, kl\left(\tfrac{1}{2}, \tfrac{1}{2}+\varepsilon\right)}$$

$$\leq \frac{1}{K} + \sqrt{\frac{T}{2K}\, kl\left(\tfrac{1}{2}, \tfrac{1}{2}+\varepsilon\right)}$$

$$\leq \frac{1}{K} + \sqrt{(2\ln 2)\,\frac{T\varepsilon^2}{K}} \qquad \text{for } \varepsilon \in \left(0, \tfrac{1}{2\sqrt{2}}\right]$$

Indeed: $\varepsilon \in \left(0, \tfrac{1}{2\sqrt{2}}\right]$ entails $4\varepsilon^2 \leq \tfrac{1}{2}$, as $\ln\left(\tfrac{1}{1-u}\right) \leq (2\ln 2)\, u \quad \forall u \in \left(0, \tfrac{1}{2}\right]$

and $kl\left(\tfrac{1}{2}, \tfrac{1}{2}-\varepsilon\right) = \frac{1}{2}\ln\left(\frac{1}{1-2\varepsilon}\right) + \frac{1}{2}\ln\left(\frac{1}{1+2\varepsilon}\right) = \frac{1}{2}\ln\frac{1}{1-4\varepsilon^2}$

$$\leq (4\ln 2)\,\varepsilon^2$$

(2) Let's go back to the regret.

In $\gamma^{(k)}$: $\quad \overline{R}_T = \sum_{a \neq k} \underbrace{\varepsilon}_{\substack{\text{gap} \\ \text{of arm } a}} \underbrace{\mathbb{E}_{\gamma^{(k)}}[N_a(T)]}_{\substack{\text{number of} \\ \text{times } a \text{ is pulled}}} = \varepsilon\left(T - \mathbb{E}_{\gamma^{(k)}}[N_k(T)]\right)$

$$\geq T\varepsilon\left(1 - \underbrace{\frac{1}{K}}_{\geq \frac{1}{2}} - \sqrt{(2\ln 2)\,T\varepsilon^2/K}\right)$$

to be optimized over $\varepsilon \in \left(0, \tfrac{1}{2\sqrt{2}}\right]$

Optimal $\varepsilon$ s.t. $\quad \frac{1}{2} - 2\varepsilon\sqrt{(2\ln 2)\,T/K} = 0, \quad$ ie, $\quad \varepsilon = \frac{1}{4} \times \frac{1}{\sqrt{(2\ln 2)\,T/K}}$

which is $\leq \tfrac{1}{2\sqrt{2}}$ as soon as $T\ln 2/K \geq \tfrac{1}{2}$,

For this $\varepsilon$, the lower bound is $\qquad\qquad\qquad\qquad$ ie $T \geq K\ln 2/2$

$$T\varepsilon\underbrace{\left(\tfrac{1}{2} - \varepsilon\sqrt{(2\ln 2)\,T/K}\right)}_{=\,\frac{1}{4}} = T\varepsilon/4 = \frac{1}{16\sqrt{2\ln 2}}\sqrt{TK} \geq \frac{1}{20}\sqrt{TK}$$

Note about the constants: By further restricting the considered $T$ (ie, by imposing $T \geq \gamma K$ with a larger $\gamma$), we can improve the $\frac{1}{20}\sqrt{TK}$ bound (ie, get $c'\sqrt{TK}$ with $c' \geq \tfrac{1}{20}$).

This is done in the upper bound on $\ln\left(\tfrac{1}{1-u}\right)$: eg, $\leq 4\ln\left(\tfrac{4}{3}\right)u$ for $u \in [0, \tfrac{1}{4}]$

↳ Entails: $\sup \overline{R}_T \geq \frac{\sqrt{2}-1}{\sqrt{32\ln(4/3)}}\sqrt{TK} \geq 0.136\sqrt{TK}$ for $T \geq \frac{1}{4\ln(4/3)}K \approx 0.87\,K$

## Open question:

So, the minimax rates for the regret for stochastic bandits or oblivious opponents are $\sqrt{TK}$.

What is the minimax rate against general, reactive opponents?

→ Should the upper bound be improved?

→ Should the proof technique be improved? (in particular, look for sequences of payoffs with real and strong correlations/ dependencies in ght the past).

Stochastic bandits :                    What about arms indexed by a continuum?

Setting 1 :    Arms indexed by $x \in A$, where $A$ is some possibly large set

With each arm $x \in A$ is associated a probability distribution $\nu_x$ over $\mathbb{R}$ s.t. $E(\nu_x)$ exists

At each round, the decision-maker picks $I_t \in A$, gets a reward $Y_t$ drawn at random according to $\nu_{I_t}$ (given $I_t$); and this is the only feedback she gets.

Definition :    $f : x \in A \longmapsto E(\nu_x)$ is the mean-payoff function

Regret :    $$\overline{R}_T = T \sup_{x \in A} f(x) - E\left[\sum_{t=1}^{T} Y_t\right]$$

Setting 2 :    [ special case ] $\longrightarrow$ Noisy optimization of a function.

We fix $f : A \to \mathbb{R}$

The noise is given by a sequence of iid random variables $\varepsilon_1, \varepsilon_2, \dots$

When $I_t \in A$ is picked, $Y_t = f(I_t) + \varepsilon_t$

$\hookrightarrow$ Special case of setting #1 where $\nu_x$ is the distribution of $f(x) + \varepsilon_1$ (all these distributions have the same shape given by the common distribution of the $\varepsilon_j$)

We of course need conditions for the regret to be minimized.

Definition :    Let $\mathcal{F}$ be a set of possible bandit problems $\nu = (\nu_x)_{x \in A}$

$\longrightarrow$ The regret can be controlled (in a non-uniform way) against $\mathcal{F}$ if :

we also say that $(A, \mathcal{F})$ is tractable

there exists a strategy s.t. $\forall \nu \in \mathcal{F}$, $\overline{R}_T = o(T)$.

Ex: $A = \{1, \ldots K\}$ and $\mathcal{F} = \left(\mathcal{P}([0,1])\right)^K$, the set of all $K$-tuples of probability distributions over $[0,1]$

$\underbrace{\qquad}$ the case of finitely many arms with bounded distributions

$\rightarrow$ UCB does the job.

Counter-example: $A = [0,1]$ and $\mathcal{F} = \left(\mathcal{P}([0,1])\right)^{[0,1]}$

$\uparrow$ illustrating that continuity is a minimal requirement.

all bandit problems $(\nu_x)_{x \in [0,1]}$ with distributions $\nu_x$ having support $[0,1]$.

Indeed: Consider $(\delta_0)_{x \in [0,1]}$ the bandit problem in which each arm $x$ is associated with the Dirac mass on $0$

Fix any strategy: it gets $Y_t = 0$ $\forall t$ and

uses a sequence of (possibly) random choices $I_t, t \geq 1$

Since probability distributions can only have at most countably many atoms,

$$\mathcal{Y} = \left\{ x \in [0,1] : \exists t \mid \mathbb{P}\{I_t = x\} > 0 \text{ under } (\delta_0)_{x \in [0,1]} \right\}$$

is countable. In particular, $[0,1] \setminus \mathcal{Y}$ is non empty.

But the strategy behaves the same under the problem $(\nu'_x)_{x \in [0,1]}$ in which $\begin{cases} \nu'_x = \delta_0 & \forall x \neq x_0 \\ \nu'_{x_0} = \delta_1 & \text{for one fixed } x_0 \in [0,1] \setminus \mathcal{Y} \end{cases}$

With probability $1$, it thus never hits $x_0$.

Therefore, $Y_t = 0$ a.s. $\forall t$ and $\bar{R}_T = T - \mathbb{E}\left[\sum_{t=1}^{T} Y_t\right] = T.$

Actually, continuity is necessary and sufficient for the regret to be controlled, as long as $A$ is not too large.

Theorem: Let $A$ be a cont metric space and let $\mathcal{F}$ be the set of bandit problems $(\nu_x)_{x \in A}$ with:

$\rightarrow$ $\forall x$, $\nu_x$ is a distribution over $[0,1]$

$\rightarrow$ a continuous mean-payoff function $f : x \mapsto \mathbb{E}(\nu_x)$

The regret can be controlled against $\mathcal{F}^{cont}$ if and only if $A$ is separable.

Corollary. Let $\mathcal{F}^{all}$ be the family of all bandit models $(\nu_x)_{x \in A}$ with distributions $\nu_x$ over $[0,1]$. Then the regret against $\mathcal{F}^{all}$ can be controlled if and only if $A$ is at most countable.

(let $A$ be any set.)

Before we prove these facts, consider the following more concrete example, in which, by strengthening the regularity requirement on the mean-payoff function, we can even get rats.

Exercise 2.   Let $A = [0,1]$ and let $\mathcal{F}^{Lip}$ be the family of bandit models $(\nu_x)_{x \in [0,1]}$ with distributions $\nu_x$ over $[0,1]$ and with mean-payoff functions that are Lipschitz.

Exhibit a strategy based on UCB + a sequence of discretizations of $[0,1]$ into $K$ bins (to be refined over time) such that:

$$\forall \nu \in \mathcal{F}^{Lip}, \qquad \overline{R}_T \leq \left(3L + 6\sqrt{8\ln T + 2}\right)T^{2/3} + 2$$

where $L$ is the Lipschitz constant of the mean-payoff function of $\nu$

Hint:

First, prove a performance bound by splitting $[0,1]$ into $[(i-1)/K, i/K]$ with $i = 1, \dots K$ for a fixed $K$, where each bin $[(i-1)/K, i/K]$ plays the role of an $i$ in a bandit problem with finitely many arms. Then discuss how to pick $K$ over the time, as we do in the next proof.

Proof of the Corollary:

We endow $A$ with the discrete topology, ie, choose the distance

$$d(x,y) = \mathbb{1}_{\{x \neq y\}}.$$   Then:

1. All applications $f: A \to \mathbb{R}$ are continuous
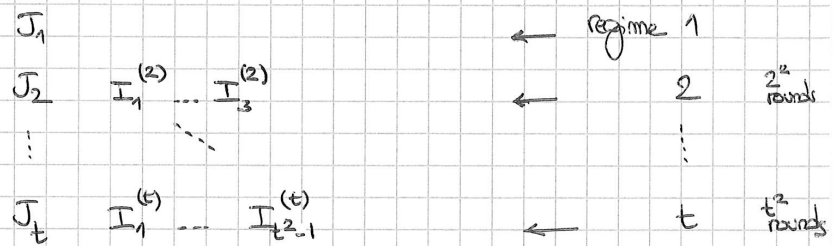2. $A$ is separable if and only if $A$ is at most countable.

Proof of the Theorem:    It relies on the possibility or impossibility of uniform exploration of the arms.

1) **If A is separable:**    let $(x_n)_{n \in \mathbb{N}}$ be a collection of points in A that is dense

In particular, the probability distribution $\quad d = \sum_{n \geq 0} \frac{1}{2^{n+1}} \delta_{x_n}$ is such that $\quad d(V) > 0 \quad$ for all open sets $V \subset A$.

We pick elements $\quad J_1 \quad J_2 \quad I_1^{(2)} \quad \dots \quad J_t \quad I_1^{(t)} \quad \dots \quad I_{t^2-1}^{(t)} \quad$ as follows :

$$J_1 \qquad\qquad\qquad\qquad\qquad \leftarrow \text{regime } 1$$
$$J_2 \quad I_1^{(2)} \dots I_3^{(2)} \qquad\qquad \leftarrow \quad 2 \qquad 2^2 \text{ rounds}$$
$$\vdots \qquad\qquad\qquad\qquad\qquad\qquad \vdots$$
$$J_t \quad I_1^{(t)} \dots I_{t^2-1}^{(t)} \qquad\qquad \leftarrow \quad t \qquad t^2 \text{ rounds}$$

where $\begin{cases} \text{the } J_s \text{ are drawn at random according to } d \\ \text{the } I_s^{(r)}, \ 1 \leq s \leq r^2-1, \text{ follow from the UCB strategy with arms } J_1 \dots J_r \end{cases}$

In regime r :    $\qquad r^2 \max_{s \leq r} \mu_{J_s} - E\left[ \sum_{s=1}^{r^2} Y_{S_{r-1}+s} \right]$

regime r starts at time $S_{r-1}+1$ where $[1+2^2+\dots+(r-1)^2] = S_{r-1}$

$$\leq 1 + c\sqrt{r^3 \ln r}$$

↑ for $J_r$    ↑ well-chosen numerical constant    ↑ distribution-free regret bound for UCB on $r^2-1$ steps with r arms (we saw this bound as an exercise)

Let $\varepsilon > 0$,    let $\tilde{r}_\varepsilon$ the first (random) time when $\mu_{J_r} = f(J_r) \geq \sup_{x \in A} f(x) - \varepsilon$

We have $\tilde{r}_\varepsilon < +\infty$ a.s. because :

- by continuity of $f$, there exists an open set $V_\varepsilon$ with
$$\forall x' \in V_\varepsilon, \quad f(x') \geq \sup_A f - \varepsilon \ ;$$

— we have $\quad \tilde{r}_\varepsilon \leq \inf\{r \geq 1 ; \quad J_r \in V_\varepsilon\} < +\infty \quad$ a.s.

as this $\nearrow$ random variable follows a geometric distribution with parameter $\gamma(V_\varepsilon) \cup > 0$.

For $r \geq \tilde{r}_\varepsilon$, $\qquad \max_{s \leq r} \mu_{J_s} + \varepsilon \geq \sup_A f$

So that $\qquad \overline{R}_T = T \sup_A f - \mathbb{E}\left[\sum_{t=1}^{T} Y_t\right]$

the "complete" regimes

$$\leq \sum_{r=1}^{\tilde{r}_\varepsilon - 1} r^2 \quad + \quad T\varepsilon \quad + \quad \sum_{r=\tilde{r}_\varepsilon}^{r_T - 1} \left(1 + c\sqrt{r^3 \ln T}\right)$$

$$+ \; r_T^2$$

↑ lengths of regimes $r \leq \tilde{r}_\varepsilon - 1$ $< +\infty$ a.s.

↑ regime $r_T$ may be incomplete

where $r_T$ is such that time $T$ belongs to regime $r_T$: we have $r_T^3$ of order $T$ ie $r_T$ of order $T^{1/3}$

and

$$\sum_{r \leq r_T} \left(1 + c\, r^{3/2}\sqrt{\ln r}\right)$$

$$\leq \sum_{r \leq r_T} \left(1 + c\, r^{3/2}\sqrt{\ln r_T}\right) = O\left(r_T^{5/2}\sqrt{\ln r_T}\right)$$

$$= O\left(T^{5/6}\sqrt{\ln T}\right)$$

Thus,

$$\limsup_{T \to +\infty} \frac{\overline{R}_T}{T} \leq \varepsilon \qquad \text{a.s.}$$

but since $\overline{R}_T$ is a deterministic quantity and this is true $\forall \varepsilon > 0$, we have

$$\lim_{T \to +\infty} \frac{\overline{R}_T}{T} = 0 \qquad \text{as requested.}$$

2)    **If A is not separable:**

\* We use the following characterization of separability ( which relies on Zorn's lemma ) :

> A metric space $X$ is separable if and only if it contains no uncountable subset $\mathfrak{D}$ s.t.
> $$\rho = \inf\{ d(x,y) : x,y \in \mathfrak{D} \} > 0.$$

In particular, if $A$ is not separable, there exist an uncountable subset $\mathfrak{D} \subset A$ and $\rho > 0$ such that the balls $B(a, \rho/2)$, with $a \in \mathfrak{D}$, are all disjoint.

↳ No probability distribution over $A$ can give a positive mass to all these balls.

\* We consider the bandit models $\nu^{(a)}$ inducing mean-payoff functions
$$f^{(a)} : x \in A \mapsto \left(1 - \frac{d(x,a)}{\rho/2}\right)^+ \ ; \quad \text{in particular,} \quad \nu_x^{(a)} = \delta_0 \text{ for } x \notin B(a, \rho/2).$$
↑ $f^{(a)}$ is indeed continuous.

We proceed as in the example showing the necessity of continuity and consider the bandit model $(\delta_0)_{x \in A}$, as well as any strategy and the laws induced by the $I_t$ under this model : let $d_t$ be the law of $I_t$ under $(\delta_0)_{x \in A}$ and let $d = \sum_{t \geq 1} \frac{1}{2^t} d_t$.

There exists $a \in A$ s.t. $d(B(a, \rho/2)) = 0$, that is, s.t., $\forall t \geq 1$, $\mathbb{P}\big( I_t \in B(a, \rho/2) \text{ under } (\delta_0)_{x \in A}\big) = 0$.

The considered strategy is therefore such that the $I_t$ have the same distribution under $(\delta_0)_{x \in A}$ and $\nu^{(a)}$. In particular,
$$\mathbb{E}\left[ \sum_{t=1}^{T} y_t \right] = 0 \text{ in both cases,} \quad \text{but in the latter case,}$$
$\sup f^{(a)} = 1$, so that $\overline{R}_T = T$ against $\nu^{(a)}$. The regret is **not** controlled against $\nu^{(a)} \in \mathcal{F}^{cont}$.