

Lecture 4.6: \sqrt{KT} distribution free bound

and bandits with a continuum of arms

We have shown. (excise last lecture)

a minimax lower bound of order \sqrt{KT} for stochastic bandits

distribution free upper bounds of order $\sqrt{KT \ln T}$ for UCB and SE.

Can we get a \sqrt{KT} upper bound?

Ross algorithm (Minimax Optimal strategy in the Stochastic case of bandit problems)

Index policy relying on $U_k(t) = \hat{\mu}_k(t) + \sqrt{\frac{1}{2N_k(t)} \ln_+ \left(\frac{T}{KN_k(t)} \right)}$

where $\ln_+ = \max(\ln, 0)$.

The algo is defined as

- | For $t=1, \dots, K$: pull $a_t \in t$
- | For $t \geq K+1$: pull $a_t \in \operatorname{argmax}_{a \in [K]} U_a(t-1)$

Difference with UCB:

bonus

$$\sqrt{\frac{2\ln t}{N_k(t)}} \quad \text{vs} \quad \sqrt{\frac{\ln_+ \left(\frac{T}{KN_k(t)} \right)}{2N_k(t)}}$$

→ no exploration after t was pulled
 T times (still exploitation)

Theorem MOSS satisfies for bandit model $\mathcal{D} = \mathbb{P}(0,1)$

$$\sup_{\nu \in \Delta^K} \mathbb{E}[R_T(\text{MOSS}, \nu)] \leq K-1 + 45\sqrt{KT}$$

→ minimax optimal, up to constant factor

(the 45 constant can still be improved)

Proof: First step for $t \geq K+1$, $U_{\mu^*}(t-1) \leq U_{\alpha_t}(t-1)$ by defn of algorithm

$$\text{Thus } R_T \leq K-1 + \underbrace{\sum_{t=K+1}^T \mathbb{E}[\mu^* - U_{\mu^*}(t-1)]}_{\substack{\text{at most } K-1 \\ \text{suboptimal pull}}}_{\substack{\text{first } K \text{ steps}}} + \underbrace{\sum_{t=K+1}^T \mathbb{E}[U_{\alpha_t}(t-1) - \mu_{\alpha_t}]}_{\sqrt{KT}} + \sum_{t=K+1}^T \mathbb{E}[(U_{\alpha_t}(t-1) - \mu_{\alpha_t})_+ \sqrt{\frac{K}{T}}]$$

Second step: control of each $\mathbb{E}[\mu^* - U_{\mu^*}(t)]$ term by $20\sqrt{T}$ ($t \geq K$)

for that:

$$\mathbb{E}[\mu^* - U_{\mu^*}(t)] \leq \mathbb{E}[(\mu^* - U_{\mu^*}(t))_+]$$

$$\leq \sum_{l=0}^{+\infty} \mathbb{E}[(\mu^* - U_{\mu^*}(t))_+ \mathbf{1}_{\{N_{\mu^*}(t) \in [x_{l+1}, x_l]\}}]$$

where
 $x_1 = \beta^{-1} \frac{t}{K}$
 for some fixed $\beta > 1$

$$+ \mathbb{E}[(\mu^* - U_{\mu^*}(t))_+ \mathbf{1}_{\{N_{\mu^*}(t) > x_0\}}]$$

$$\text{Now, } V_{\hat{\mu}^*(t)} = \hat{\mu}_{\hat{\mu}^*(t)} + \begin{cases} 0 & \text{if } N_{\hat{\mu}^*(t)} \geq \frac{T}{K} = t_0 \\ \sqrt{\frac{1}{2N_{\hat{\mu}^*(t)}} \ln\left(\frac{T}{KN_{\hat{\mu}^*(t)}}\right)} & \text{if } N_{\hat{\mu}^*(t)} < t_0 \end{cases}$$

$$\geq \underbrace{\sqrt{\frac{1}{2N_0} \ln\left(\frac{T}{KN_0}\right)}}_{:= \epsilon} \quad \text{if } N_{\hat{\mu}^*(t)} \in [t_0, T]$$

$$\text{So } \mathbb{E}[\mu^* \cdot V_{\hat{\mu}^*(t)}] \leq \mathbb{E}\left[(\mu^* \cdot \hat{\mu}_{\hat{\mu}^*(t)}) + \mathbf{1}_{\{N_{\hat{\mu}^*(t)} > \frac{T}{K}\}}\right] + \sum_{l=0}^{+\infty} \mathbb{E}\left[(\mu^* \cdot \hat{\mu}_{\hat{\mu}^*(t)} - \epsilon_l)_+ \mathbf{1}_{\{N_{\hat{\mu}^*(t)} \in [t_0, t_0 + \Delta]\}}\right]$$

Lemma: $\mathbb{E}\left[(\mu^* \cdot \hat{\mu}_{\hat{\mu}^*(t)} - \epsilon)_+ \mathbf{1}_{\{N_{\hat{\mu}^*(t)} > n_0\}}\right] \leq \frac{1}{\sqrt{n_0}} e^{-2n_0 \epsilon^2}$

Proof of the lemma:

$$\begin{aligned} \mathbb{E}\left[(\mu^* \cdot \hat{\mu}_{\hat{\mu}^*(t)} - \epsilon)_+ \mathbf{1}_{\{N_{\hat{\mu}^*(t)} > n_0\}}\right] &= \int_0^{+\infty} \mathbb{P}(\mu^* \cdot \hat{\mu}_{\hat{\mu}^*(t)} - \epsilon \geq u \text{ and } N_{\hat{\mu}^*(t)} > n_0) du \\ &= \int_0^{+\infty} \mathbb{P}(Z_t^* \geq (\epsilon + u)N_{\hat{\mu}^*(t)} \text{ and } N_{\hat{\mu}^*(t)} > n_0) du \end{aligned}$$

when $Z_t^* = N_{\hat{\mu}^*(t)} (\mu^* - \hat{\mu}_{\hat{\mu}^*(t)}) = \sum_{k=1}^T (\mu^* - X_{\hat{\mu}^*(t)}(k)) \mathbf{1}_{\{\hat{\mu}^*(t) = k\}}$ is a martingale.

see proof
UCB regret

and for all $x \in \mathbb{R}$, $S_{x,t} = e^{xZ_t^* - \frac{x^2}{8} N_{\hat{\mu}^*(t)}}$ is a supermartingale

Thus by Markov-Chernoff, we continue the bounding as, for $x > 0$

$$= \int_0^{+\infty} \mathbb{P}\left(e^{xZ_t^* - \frac{x^2}{8} N_{\hat{\mu}^*(t)}} \geq \exp(N_{\hat{\mu}^*(t)}(\epsilon + u) - \frac{x^2}{8}) \text{ and } N_{\hat{\mu}^*(t)} > n_0\right) du$$

$$x = 4(\epsilon + u)$$

so that

$$x(\epsilon + u) - \frac{x^2}{8} = 2(\epsilon + u)^2$$

$$\leq \int_0^{+\infty} \sum_{l=n_0}^{+\infty} e^{-2l(\epsilon + u)^2} \mathbb{E}\left[S_{4(\epsilon + u), l} \mathbf{1}_{\{N_{\hat{\mu}^*(t)} = l\}}\right] du$$

$$\leq \int_0^{+\infty} e^{-2n_0(c^2+u^2)} \mathbb{E}[S_{4(c+u), t} \mathbb{1}_{\{N_{c+u}(t) \geq n_0\}}] du$$

we know
 $\mathbb{E}[S_{4(c+u), t}] \leq 1$
 (and $s > d$)

so all in all:

$$\mathbb{E}[(\mu^* - \hat{\mu}_{\alpha^*}(t) - \varepsilon) \mathbb{1}_{\{N_{c+u}(t) \geq n_0\}}] \leq e^{-2n_0 c^2} \int_0^{+\infty} e^{-2n_0 u^2} du = e^{-2n_0 c^2} \cdot \sqrt{\frac{\pi}{8n_0}} \leq \frac{e^{-2n_0 c^2}}{\sqrt{n_0}}$$

integral of Gaussian density (upto norm)

Going back to the main proof:

$$\begin{aligned} \mathbb{E}[\mu^* \cdot V_{\alpha^*}(t)] &\leq \sqrt{\frac{K}{F}} + \sum_{l=0}^{+\infty} \underbrace{\frac{1}{\sqrt{x_{l+1}}} e^{-2x_{l+1} c_l^2}}_{\frac{1}{\sqrt{x_{l+1}}} \exp\left(-2x_{l+1} \cdot \frac{1}{2\beta} \ln\left(\frac{t}{Kn}\right)\right)} \\ &= \frac{1}{\sqrt{x_{l+1}}} \exp\left(-\frac{1}{\beta} \cdot l \ln(\beta)\right) \\ &= \sqrt{\frac{K}{F}} \beta^{\frac{l+1}{2}} \exp\left(-\frac{l}{\beta} \ln(\beta)\right) \\ &= \sqrt{\frac{K}{F}} \beta^{\frac{1}{2} + l\left(\frac{1}{2} - \frac{1}{\beta}\right)} \xrightarrow{\text{we want } \beta \in (1, L)} \end{aligned}$$

$$\text{Taking } \beta = \frac{3}{2}: \mathbb{E}[\mu^* \cdot V_{\alpha^*}(t)] \leq \sqrt{\frac{K}{F}} + \sqrt{\frac{K}{F}} \beta^{\frac{1}{2}} \cdot \sum_{l=0}^{+\infty} \left(\beta^{\left(\frac{1}{2} - \frac{1}{\beta}\right)}\right)^l$$

$$= \sqrt{\frac{K}{F}} \left(1 + \underbrace{\sqrt{\frac{3}{2}} \cdot \frac{1}{1-\alpha}}_{\leq 19}\right) \quad \text{with } \alpha = \left(\frac{3}{2}\right)^{\left(\frac{1}{2} - \frac{2}{\beta}\right)} \in (0, 1)$$

$$\text{Third step: } \sum_{t=K+1}^T \mathbb{E}[(U_{a_{t+1}}(t) - \mu_{arr} - \sqrt{\frac{K}{T}})_+] \leq 4\sqrt{KT}$$

$$= \sum_{t=K}^{T-1} \mathbb{E}[(U_{a_{t+1}}(t) - \mu_{arr} - \sqrt{\frac{K}{T}})_+]$$

$$\sum_{t=K}^{T-1} \mathbb{E}[(U_{a_{t+1}}(t) - \mu_{arr} - \sqrt{\frac{K}{T}})_+] = \sum_{k=1}^K \sum_{t=1}^T \sum_{r=k}^{T-1} \mathbb{E}[(U_k(r) - \mu_a - \sqrt{\frac{K}{T}})_+ \mathbb{1}_{\{a_{t+1}=k\}} \mathbb{1}_{\{N_k(r)=1\}}]$$

we now use $(U_k(r) - \mu_a - \sqrt{\frac{K}{T}})_+ \leq (\hat{\mu}_k(t) - \mu_a - \sqrt{\frac{K}{T}})_+ + \begin{cases} 0 & \text{if } N_k(r) \geq \frac{T}{K} \\ \sqrt{\frac{1}{2N_k(r)} \ln(\frac{T}{K N_k(r)})} & \text{if } N_k(r) < \frac{T}{K} \end{cases}$

and get therefore the upper bound:

$$\leq \sum_{k=1}^K \sum_{t=1}^T \sum_{r=k}^{T-1} \mathbb{E}[(\hat{\mu}_k(r) - \mu_a - \sqrt{\frac{K}{T}})_+ \mathbb{1}_{\{a_{t+1}=k\}} \mathbb{1}_{\{N_k(r)=1\}}] + \sum_{k=1}^K \sum_{t=1}^{[T/K]} \sqrt{\frac{1}{2t} \ln(\frac{T}{Kt})} \mathbb{E}[\underbrace{\sum_{r=k}^{T-1} \mathbb{1}_{\{N_k(r)=1\}} \mathbb{1}_{\{a_{t+1}=k\}}}_{\leq 1 \text{ a.s.}}]$$

Also

$$\begin{aligned} \sum_{t=1}^{[T/K]} \sqrt{\frac{1}{2t} \ln(\frac{T}{Kt})} &\leq \sqrt{\int_0^{[T/K]} \sqrt{\frac{1}{2x} \ln(\frac{T}{Kx})} dx} \\ &\leq \sqrt{\frac{T}{2K}} \int_1^{+\infty} u^{-3/2} \sqrt{\ln(u)} du \\ &= \sqrt{\frac{T}{2K}} \int_0^{+\infty} 2v^2 e^{-\frac{v^2}{2}} dv \end{aligned}$$

variance of standard Gaussian, up to $\sqrt{2\pi}$ renormalization

$$= \sqrt{\pi} \sqrt{\frac{T}{K}}$$

summarizing, we showed so far (in third step): we will show that

$$< \sqrt{\frac{T}{K}} \sqrt{\frac{T}{K}} \text{ for each } k$$

$$\sum_{t=K}^{T-1} \mathbb{E}[(U_{a_{t+1}}(t) - \mu_{arr} - \sqrt{\frac{K}{T}})_+] \leq \sum_{k=1}^K \sum_{t=1}^T \sum_{r=k}^{T-1} \mathbb{E}[(\hat{\mu}_k(r) - \mu_a - \sqrt{\frac{K}{T}})_+ \mathbb{1}_{\{a_{t+1}=k\}} \mathbb{1}_{\{N_k(r)=1\}}] + K \cdot \sqrt{\pi} \sqrt{\frac{T}{K}}$$

We resort again to $Z_{k,t} = N_a(\hat{\mu}_k(t) \cdot \mu_k)$ martingale

$$S_{a,t}^{(a)} = e^{\alpha Z_{k,t} + \frac{\sigma^2}{8} N_a(t)}$$

$$\text{where } \alpha = 4/\sqrt{K} + u$$

For each k ,

$$\sum_{t=1}^T \sum_{r=k}^{T-1} \mathbb{E} \left[\left(\hat{\mu}_k(t) - \mu_k - \frac{\sqrt{K}}{T} \right) + \mathbf{1}_{\{a_{r+1}=k\}} \mathbf{1}_{\{N_k(t)=l\}} \right] =$$

$$\sum_{l=1}^T \sum_{r=k}^{T-1} \int_0^{+\infty} P \left(\alpha Z_{k,r} - \frac{\sigma^2}{8} N_k(t) \geq N_k(t) \left(\alpha(u + \sqrt{\frac{K}{T}}) - \frac{u^2}{8} \right) \text{ and } a_{r+1}=k \text{ and } N_k(t)=l \right) du$$

$$\leq \sum_{l=1}^T \sum_{r=k}^{T-1} \int_0^{+\infty} e^{-2l(u^2 + \frac{K}{T})} \underbrace{\mathbb{E} \left[S_{a,r}^{(a)} \mathbf{1}_{\{a_{r+1}=k\}} \mathbf{1}_{\{N_k(t)=l\}} \right]}_T$$

The sum over t of them will be ≤ 1

Issue: depends on t ...

but can be replaced in some sense, by $S_{n,0}^{(a)} = 1$.

$$\leq \sum_{l=1}^T \int_0^{+\infty} e^{-2l(u^2 + \frac{K}{T})} \mathbb{E} \left[\sum_{r=k}^{T-1} S_{a,r}^{(a)} \mathbf{1}_{\{a_{r+1}=k\}} \mathbf{1}_{\{N_k(t)=l\}} \right] du$$

$$\leq \sum_{l=1}^T \int_0^{+\infty} e^{-2l(u^2 + \frac{K}{T})} \mathbb{E} \left[S_{a,\tau_l}^{(a)} \right] du$$

$$\text{when } \tau_l = \inf \left\{ t \in [T] : a_{r+1}=k \text{ and } N_k(t)=l \right\} \wedge T$$

$S_{n,i}^{(k)}$ supermartingale
 T is bounded \rightarrow we can apply optional stopping theorem ("théorème d'arrêt de Doob")

so that

$$\mathbb{E}[S_{\tau, T_p}^{(k)}] \leq \mathbb{E}[S_{\tau, 0}^{(k)}] = 1.$$

$$\text{So: } \sum_{t=1}^T \sum_{n=k}^{T-1} \mathbb{E}\left[\left(\hat{\mu}_n(t) - \mu_n - \sqrt{\frac{K}{T}}\right)_+ \mathbf{1}_{\{\text{Arr}_n=t\}} \mathbf{1}_{\{N_n(t)=1\}}\right] \leq \int_0^{\infty} e^{-2t(u^2 + \frac{K}{T})} du$$

$$\leq \sum_{l=1}^{\infty} \frac{1}{\sqrt{l}} e^{-2l\frac{K}{T}} \quad) \quad \text{similar to the Intermediate Lemma in second step}$$

$$\leq \int_0^{\infty} \frac{1}{\sqrt{x}} e^{-2x\frac{K}{T}} dx = \sqrt{\frac{T}{2K}} \int_0^{\infty} \frac{e^{-u}}{\sqrt{u}} du$$

$$= \sqrt{\frac{T}{2K}} \cdot 2 \int_0^{+\infty} e^{-v^2} dv = \sqrt{\frac{\pi}{2}} \sqrt{\frac{T}{K}}$$

$$u = v^2 \\ \frac{du}{dv} = 2v$$

Summarizing:

$$\sum_{n=k}^{T-1} \mathbb{E}\left[\left(V_{\text{Arr}_n}(t) - \mu_{\text{Arr}_n} - \sqrt{\frac{K}{T}}\right)_+\right] \leq \sum_{k=1}^K \sum_{t=1}^T \sum_{n=k}^{T-1} \mathbb{E}\left[\left(\hat{\mu}_n(t) - \mu_n - \sqrt{\frac{K}{T}}\right)_+ \mathbf{1}_{\{\text{Arr}_n=t\}} \mathbf{1}_{\{N_n(t)=1\}}\right] + \underbrace{\sum_{k=1}^K \sum_{t=1}^{T-k} \sqrt{\frac{1}{2t} \ln\left(\frac{T}{K}\right)} \mathbb{E}\left[\sum_{n=t+1}^{T-1} \mathbf{1}_{\{\text{Arr}_n=t\}} \mathbf{1}_{\{N_n(t)=1\}}\right]}_{K \sqrt{\frac{\pi}{2}} \sqrt{\frac{T}{K}}} + \underbrace{\sqrt{\pi} \sqrt{KT}}$$

$$\propto \sqrt{KT}, \quad \sqrt{\pi} \left(1 + \frac{1}{\sqrt{2}}\right) \leq 4\sqrt{KT}$$

General conclusion

Summarizing all steps, we bound the regret by

$$K \cdot 1 + \left(\sum_{t=K+1}^T 20 \sqrt{\frac{K}{T-1}} \right) + \sqrt{KT} + 4\sqrt{KT} \leq K \cdot 1 + 5\sqrt{KT} + 20 \int_0^T \sqrt{\frac{K}{s}} ds \\ = K \cdot 1 + 45\sqrt{KT} \quad \blacksquare$$

MOSSE is minimax optimal, but is it optimal w.r.t. instance dependent bounds?

No. Reaching simultaneous optimality on both ends is hard, but has been proven for some algorithms recently. (see Tsallis-TNF paper)

Bandits with continuum of arms

Stochastic bandits: what about arms indexed by a continuum?

Setting 1 Arms indexed by $a \in A$, where A is some possibly large set. With each arm $a \in A$ is associated a probability distribution ν_a over \mathbb{R} , s.t. $E(\nu_a)$ exists.

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

Definition $f: a \in A \mapsto E(Y_a)$ is the mean-payoff function.

(prob.)-Regret: $R_T = T \sup_{a \in A} f(a) - \sum_{t=1}^T f(a_t)$

Setting 2 [Special case] \rightarrow noisy optimization of a function

we fix $f: A \rightarrow \mathbb{R}$. The noise is given by a sequence of iid random variables $\epsilon_1, \epsilon_2, \dots$

when $a_t \in A$ is picked, $Y_t = f(a_t) + \epsilon_t$

\hookrightarrow special case of setting #1 where ν_a is the distribution of $f(a) + \epsilon_1$ (all these distributions have the same shape, given by the common distribution of the ϵ_j)

We of course need conditions for the regret to be minimised

We already
knew
 (A, F) is
measurable

Definition Let F be a set of possible bandit problems $\omega = (\omega_n)_{n \in \mathbb{N}}$.

The regret can be controlled (in a non-uniform way) against F if:

there exists a strategy s.t. $\forall \omega \in F, \mathbb{E}[R_T] = o(T)$.

E: $A = \{1, \dots, K\}$ and $F = P([0, 1])^K$ \rightarrow UCB does the job.

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

all bandit problems $(\omega_n)_{n \in [0, 1]}$
with distributions ω_n having support $[0, 1]$

Ende! Consider $(\delta_x)_{x \in [0, 1]}$ the bandit problem in which each arm x is associated with the Dirac mass on 0.

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

$\mathcal{S} = \{x \in [0, 1] : \exists t \mid \mathbb{P}(a_t = x) > 0 \text{ under } (\delta_x)_{x \in [0, 1]}\}$ is countable. In particular,

we can consider $x_0 \in [0, 1] \setminus \mathcal{S}$. The strategy then behaves the same under

the problem $(\omega'_x)_{x \in [0, 1]}$ in which $\begin{cases} \omega'_x = \delta_0 & \forall x \neq x_0 \\ \omega'_{x_0} = \delta_1 \end{cases}$

With prob 1, the strategy never pulls x_0 .

Therefore, $y_t = 0$ as for any t and $R_T = T$.

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

Theorem Let A be a metric space and let F^{cont} be the set of bandit problems

$(\nu_x)_{x \in A}$ with

- $\forall x, \nu_x$ is a distribution over $[0,1]$
- a continuous mean-payoff function $f: x \mapsto E(\nu_x)$

The expected regret can be controlled against F^{cont} if and only if A is separable

Corollary Let A be any set. Let F^{all} be the family of all bandit models $(\nu_x)_{x \in A}$ with

distributions ν_x over $[0,1]$. Then the regret against F^{all} can be controlled if and only if A is at most countable.

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 rates. (see exercise section #5)

Proof of the corollary: we endow A with the discrete topology, i.e., choose the

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

1. All applications $f: A \rightarrow \mathbb{R}$ 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.

We pick actions in a triangular fashion:

Regime 1: UCB based on x_1, x_2 : $a_1^{(1)}, \dots, a_4^{(1)}$
(fresh start)

Regime r : UCB based on x_1, \dots, x_r, x_{r+1} : $a_1^{(r)}, \dots, a_{(r+1)}^{(r)}$
(push start)

In regime r :

starts at time

$$S_r + 1 = 2^2 + 3^2 + \dots + r^2 + 1$$

$$(r+1)^2 \max_{s \leq r} f(x_s) - \mathbb{E} \left[\sum_{t=S_r+1}^{S_{r+1}} Y_t \right] \leq c \sqrt{r^3 \ln r}$$

distribution free bound
of VCB on $(r+1)$ steps with
 $(r+1)$ arms (see section #2)

Now, let $\varepsilon > 0$ and let $\tilde{r}_\varepsilon \in \mathbb{N}^*$ s.t. $f(\tilde{x}_{\tilde{r}_\varepsilon}) \geq \sup_A f - \varepsilon$

(\tilde{r}_ε exists by separability of A and continuity of f)

In particular, $\max_{s \leq \tilde{r}_\varepsilon} f(x_s) \geq \sup_A f - \varepsilon$

we denote by r_T the index of the regime where T lies:

we have that S_{r_T} is of order of r^3

so r_T is of the order of $T^{2/3}$, i.e. $r_T = O(T^{2/3})$.

The regret can be decomposed (for T large enough) as

$$\mathbb{E}[R_T] = T \sup_A f - \mathbb{E} \left[\sum_{t=1}^{r_T} Y_t \right] = \text{sum of the regrets of each regime}$$

$$\begin{aligned} &< \underbrace{\sum_{r=1}^{\tilde{r}_\varepsilon} (r+1)^2}_{\substack{\text{initial regimes,} \\ \text{regret bounded by} \\ \text{then lengths } = O(1)}} + \underbrace{\sum_{r=\tilde{r}_\varepsilon}^{r_T-1} ((r+1)^2 \varepsilon + c \sqrt{r^3 \ln r})}_{\substack{\text{cf bounds (1) and (4)} \\ \text{and (4*)}}} + (r_T+1)^2 \\ &\quad \underbrace{\text{regime } r_T}_{\substack{\text{may be incomplete } O(T^{2/3})}} \\ &\leq TE + \sum_{r=\tilde{r}_\varepsilon}^{r_T-1} r^{3/2} \sqrt{\ln r} \\ &\leq TE + O(r_T^{5/2} \sqrt{\ln r_T}) \\ &= TE + O(T^{5/6} \sqrt{\ln T}) \end{aligned}$$

All in all, $\limsup_T \frac{\mathbb{E}[R_T]}{T} \leq \varepsilon$ which is true for any $\varepsilon > 0$

that is $\lim \frac{\mathbb{E}[R_T]}{T} = 0$

2) If A is not separable

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

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

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

\Rightarrow No probability distribution over A can give a positive mass to all these balls.

* we consider the bandit models $v^{(a)}$ inducing mean-payoff function

$$f^{(a)}: x \in A \longrightarrow \left(1 - \frac{d(x, a)}{\rho/2}\right)_+$$

continuous \nearrow

$$\text{in particular, } v_x^{(a)} = \delta_0 \text{ for } x \notin B(a, \rho/2)$$

We proceed as in the example showing the necessity of continuity when $A = [0, 1]$ and consider the bandit model $(\delta_x)_{x \in A}$, as well as any strategy and the laws induced by the a_T under this model: Let λ_T be the law of a_T under $(\delta_x)_{x \in A}$ and let $\lambda = \sum_{t=1}^T \frac{1}{t} \lambda_t$

As only countably many balls can have a positive mass under λ , there exists a.s.t. $\lambda(B(a, \rho/2)) = 0$, i.e.

$$\forall t \geq T, \mathbb{P}(a_t \in B(a, \rho/2) \text{ under } (\delta_x)_{x \in A}) = 0$$

The considered strategy is therefore such that the a_T have the same distribution under $(\delta_x)_{x \in A}$ and $v^{(a)}$. In particular, $E\left[\sum_{t=1}^T Y_t\right] = 0$ in both cases, but in the latter case $\sup_A f^{(a)} = 1$, so that $R_T = T$ against $v^{(a)}$. The regret is thus not controlled against $v^{(a)} \in \mathcal{F}^{\text{cont}}$ □

We consider the setting of stochastic with a continuum of arms indexed by $\mathcal{A} = [0, 1]$, with a mean-payoff function f that is α -Hölder with $\alpha \in (0, 1]$, i.e., there is $L > 0$ such that

$$\forall x, x' \in [0, 1], \quad |f(x) - f(x')| \leq L|x - x'|^\alpha.$$

We now consider the following algorithm that discretizes the action space into K bins:

- discretize $[0, 1]$ into K bins, with $B_i = [\frac{i-1}{K}, \frac{i}{K}]$ for any $i = 1, \dots, K$,
- run MOSS algorithm over K arms, where picking the arm $I_t \in [K]$ corresponds to picking an action a_t (chosen arbitrarily) in the bin B_{I_t} .

1) Show that the regret of this “discretized” algorithm satisfies:

$$\mathbb{E}[R_T] \leq K + 45\sqrt{KT} + \frac{TL}{K^\alpha}.$$

2) Assume that T, α are known in advance. Show that for a good choice of K , the regret is of order $T^{\frac{\alpha+1}{2\alpha+1}}$.

Solution: 1) Define $f_i = \min_{x \in B_i} f(x)$ and $f_b^* = \max_{i \in [K]} f_i$. We have

$$\begin{aligned} \mathbb{E}[R_T] &\leq \mathbb{E}\left[\sum_{t=1}^T f^* - f_{I_t}\right] \\ &\leq T(f^* - f_b^*) + \mathbb{E}\left[\sum_{t=1}^T f_b^* - f_{I_t}\right]. \end{aligned}$$

The second term is the regret of MOSS on the discretized game, so it is bounded by $K + 45\sqrt{KT}$. The first term on the other hand is the approximation error. Each bin is of size $\frac{1}{K}$, so using the α -Hölder property, $f^* - f_b^* \leq \frac{L}{K^\alpha}$, which leads to the result.

2) This is a consequence of choosing $K = T^{\frac{1}{2\alpha+1}}$.