Lecture 47: VKT distribution pressound and bandits with a continuum of erms

We have shown. (in exercise services #4) a minimax lower bound of order VET for atochestice bandits · distribution free upper bounds of order UKTlAT for UCB and SE.

Can we get a VKT upper bound?

1055 algaithm (Minimox Optimal Strategy in the Stochastic are of land tycklams)

Index policy rulying on
$$U_{\mathcal{L}}(h) = \hat{\mu}_{\mathcal{L}}(h) + \sqrt{\frac{1}{2N_{\mathcal{L}}(h)}} \ln_{+} \left(\frac{1}{KN_{\mathcal{L}}(h)}\right)$$

where $ln_{+} = max(ln_{/}0)$

te algo is defined as

Difference with UCB.

Vs (r/kNaO)

ZNe(h)

> no exploration after & was pulled to times (still exploration)

Therem MOSS artisfier for bondit model D= P(CO.1) RT (Moss, v) < K-1 + 45 VKT s ledk (the 45 constant can stillbe improved) - minimax ophinal up to constant factor Proof: First step In to K+1, $U_{gr}(r-1) \leq U_{ar}(r-1)$ by left of algorithm

thus Rr (K-1 + E [M-Vg-(t-1)] + E [Uar (t-1) · Mar]

at most K-1

subsphired public

first Kategor

TET + E [M (t-1) · Mar] TRT + Z E (Var (r.1)-Mr- (+))

Second step : control of each E[pt-Varl) tem by 201 (fat 3K)

forthot: E[p. Va. (b)] < E[(p. Va. (b)),] where 21= p-1 + from ford \$>1

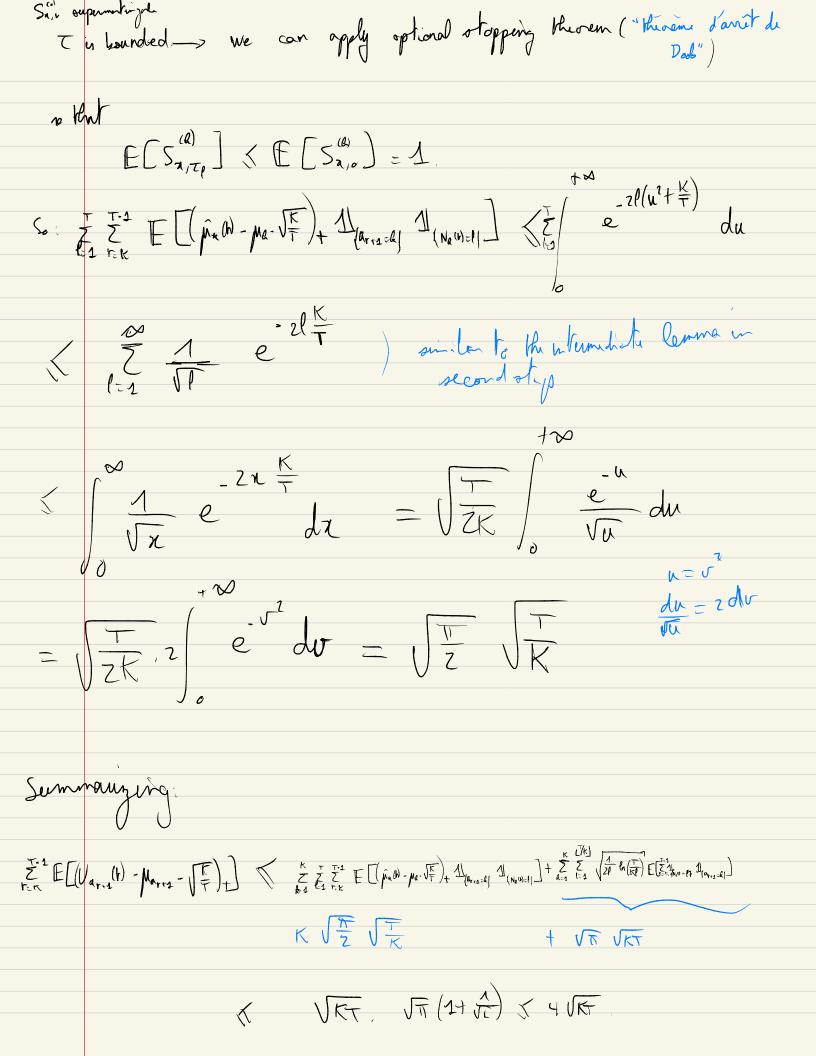
+ E[(\mu^* - V_e^* (r))_+ 1 {N_{k*}(r) >, re_o}]

third step: E (Ver (r-1)-phr · JE) in S 4 JKT = E E (Variation - Marina - (F)) E(Variali) - Marra - (F)+) = K T T-1 E (URG) - Ma - (F) + 1 [arrack] 1 (NeWell) we now use $(V_{\alpha}(r) - \mu_{\alpha} - \overline{V_{\tau}})_{+} \times (\hat{\mu}_{\alpha}(t) - \mu_{\alpha} - \overline{V_{\tau}})_{+} + \begin{cases} 6 & \text{if } N_{\alpha}(t) \geqslant \frac{t}{K} \\ \sqrt{2n_{\alpha}(r)} \ln(\frac{t}{KN_{\alpha}(r)}) & \text{if } N_{\alpha}(r) < \frac{t}{K} \end{cases}$ alo smalle blom

\[
\sum_{ZN_{\mathbb{L}}(\beta)} \frac{\tau}{\tau_{\mathbb{A}(\beta)}} \] and get therefore the upper bound Also $\sum_{l=1}^{T} \sqrt{\frac{1}{2l} \ln(\frac{1}{Kl})} < \sqrt{\frac{1}{2n} \ln(\frac{1}{Kn})} dn$ where $\sum_{l=1}^{T} \ln(\frac{1}{Kl}) < \sqrt{\frac{1}{2n} \ln(\frac{1}{Kn})} dn$ 12 = - K 2" du $\sqrt[3]{\frac{1}{2k}} \int_{1}^{1} \frac{1}{u^{-3/2}} \sqrt{\ln(u)} du$ = K W du = \langle \frac{1}{2k} \langle \frac{1}{2} \delta \ we will show that summarizing, we showed so for (in third step). < To the for each a E1 E[(Variath) - Maria - (F)+) < Z Z Z Z E [(px 0) - Ma - (F)+ 1 [(arraid) 1 (NeWell)]

we resort og ain to $7k_{1}t = Na(\hat{\mu}_{e}(t), \mu_{e})$ menturgle $S_{n/t} = e$ $S_{n/t} = e$ Maria ($\hat{\mu}_{e}(t)$) $p_{e}(t)$ menturgle

Solution $s_{n/t} = e$ where x = 4() + w) $\begin{array}{c} T & T-1 \\ \overline{L} & \overline{L} & \overline{L} \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} & \overline{L} & \overline{L} & \overline{L} \\ \overline{L} & \overline{L} &$ $\sum_{k=1}^{T}\sum_{r\in K} \left\{ \left(p_{k}(N) - p_{k} \cdot \nabla_{r} \right) + \left(p_{k} \cdot \nabla_{r} \right) + \left(p_{k} \cdot \nabla_{r} \right) \right\} \times \left(\left(p_{k}(N) - p_{k} \cdot \nabla_{r} \right) \right) \times \left(p_{k}(N) - p_{k} \cdot \nabla_{r} \right) \times \left(p_{k}(N) - p_{k} \cdot$ issu! depends on t...but can be replaced in some sense, by $S_{n,o} = 1$. $\begin{array}{c|c}
 & -2l(u^2 + \frac{K}{T}) \\
\hline
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\$ $\begin{array}{c|c}
 & & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & & \\
 & &$ When $T_{f} = \inf \left\{ f \in T : \underset{N_{k}(f)}{\text{ar+1}} = k \text{ and } \right\} \wedge T$



General conclusion Summarizing all steps, we bound the regul by

$$K.1 + \left(\sum_{r=k_1}^{T} 20 \sqrt{\frac{\kappa}{F\cdot 1}}\right) + \sqrt{\kappa\tau} + 4\sqrt{\kappa\tau} \langle \kappa-1 + 5\sqrt{\kappa\tau} + 20 \sqrt{\frac{\kappa}{\delta}} do$$

Bandits with continuum of arms indexed by a continuem?

Setting 1 Arms indexed by n EA, when A is some possebly large oct. With each arm ZEA is associated a probability distribution ve our R st. E(M) exists

At each round, the decision maker picks ar EA, gets areword of drawn atrandom according to var (given or); and this is the only feedback she gets.

Definition $f: x \in A \longrightarrow E(x_x)$ is the mean-poyoff function. (pundo)-Regret

RT = Toup f(n) = E[Z] Yr]

Setting 2 Copició case > noisy optimisation of a function

we fix $f: A \rightarrow IR$ The noise is given by a sequence of id random variables $\mathcal{E}_{1}, \mathcal{E}_{2}, ...$

when at Exis picked, Y= f(ar) + Ex

Li special case of setting #1 where who is the distribution of $f(x) + E_1$ (all these distributions have the same shape, given by the common distribution of the E_j)

We als now Definition let I be a set of possible bandit poblems is = (vn) rea. A. Foods the regul can be controlled (in non-uniform way) against I if: there exists astrategy p.k. VNEF R= o(T) Ex: A= (1,..., K) and F= P((0,1)) -> UCB does the jd. Counter-example: A=[0,1] and F= P(0,1]

all Landit publies (Vn) ac(0,1)

with distributions of horting upper [0,1] Inde! Consider (50) ac (0,1) the bandet publim in which each own x is associated with the Dirac mass on O. Since probability distributions can only have at most countably many atoms, S= {x & (0.1);] + | IP(ar=x) >0 under (50) ne(0.1) is countable. In particular, me can consider to E [0,1] 19. The strategy then behaves the same under ble publem $(\sqrt{x})_{x \in (0,1)}$ in which $\begin{cases} \sqrt{x} = \delta_0 & \forall x \neq n_0 \\ \sqrt{n_0} = \delta_1 \end{cases}$ With polar I, the strokegy never pulls to. Therefore, 1/20 as for any t and Ry = T. Actually, continuity is sufficient for the regret to be controlled as long as A is not too

we of course need conditions for the regret to be minimised

Theorem let A be a metric opace and let First be the oct of bonds problems (V2)xcA with . Va , Vx is a distribution over [0,1]

. a continuous mean payoff function $f: x \mapsto E(x_n)$ The regul can be controlled against F^{cont} if and only if A is separable

Corollary let Abe any set, let Fall be the family of all bandit models (2) xea with distributions 22 our [0,1]. Then the regular against Fall can be controlled if and only if A is at most countable.

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

(see exercise session # 5)

Proof of the corollary: we endow A with the diouete topology, i.e, choose the diotance $d(x,y) = 4|_{x+y}$. 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 whiles on the possibility of impossibility of impo

1) If Aisseparable: let (In) now be a collection of points in A that in dense.

We pick actions in a triangular fashion:

Regime 1: UCB based on as & 1 a1, ..., a4 (1)

Regime n: UCB based on ty ... x, re, 1: a1 ... a(r+1):

 $(r+1)^2 \max_{s \leqslant r} f(x_s) - \mathbb{E}\left[\sum_{k=s,+1}^{s+1} y_k\right] \leqslant c \sqrt{r^3 \ln r}$ In regime ni starts at time Sr+1 = 2+32+...+ r2+1 distribution free bound of UCB on (r-1) origin with (r-1) mms (me union #2) Now, let $\varepsilon > 0$ and let $\tilde{r}_{\varepsilon} \in \mathbb{N}^{\bullet}$ s.t. $f(x_{\widetilde{\varepsilon}}) \geqslant \sup_{A} \int_{-\varepsilon}^{\infty} -\varepsilon$ (Fe exists by reportably of A and continuity of f) In particular, $\max_{S \leqslant \widehat{r_{i}}} f(x_{i}) \geqslant \sup_{A} f - \mathcal{E}$ we denote by of the index of the regime where This we have that Sr is oftender of r3 Noting is of the order of $T^{2/3}$, i.e. $\Gamma_T = O(T^{2/3})$. The regul can be decomposed (for Tlarge enough) as R_ = T sup f - E[\frac{7}{r_{2}} \gamma_{r}] = sum of the regular of each regime $\left\langle \sum_{\lambda=1}^{n_{c}-1} (n+1)^{2} + \sum_{r \in r_{c}}^{r_{c}-1} ((r+1)^{2} \mathcal{E} + c \sqrt{r^{c} \ell_{n} r}) + (r_{r}+1)^{2} \right\rangle$ TE +0 (15 /6 15) = TE + O(T516 VAT)

All in all, tunoup $\frac{R_T}{T} \leqslant E$ which is true for any E>0that is $\lim_{T \to 0} \frac{R_T}{T} = 0$

2) If A is not reparable

We use the following characterisation of separability (which relies on Forn's lemma).

A metric open χ is separable if and only if it contains no uncountable subset D st. $p=\inf d(d(x,y):x,y\in D)>0$.

In particular, if A is not separable, there exists an uncountable subset $D \subset A$ and $Q \supset O$ such that the balls B(a, Vz) with $a \in D$ are all diggoint.

> No probability distribution over A congive a positive mass to all these balls

we consider the bandit models $v^{(n)}$ inducing mean-payoff function $\int_{0}^{(n)} : \kappa \in A \longrightarrow \left(1 - \frac{d(x, a)}{e^{rx}}\right) +$ in particular, $v^{(n)} = \delta_0$ for $x \notin B(x, e^{(x)})$

we proceed as in the example showing the necessity of continuity when A = [0,1] and consider the bandit model $(\mathfrak{F}_{\bullet})_{x \in A}$, as well as any studgy and the laws induced by the ap under the model: let λ_{F} be the low of ap under $(\mathfrak{F}_{\bullet})_{x \in A}$ and let $\lambda = \Sigma \frac{1}{2^{\mathsf{F}}} \lambda_{\mathsf{F}}$

As only countably many bolls can have a positive mass under λ , there exists a s.t. $\lambda(B(a, p/z)) = 0$, the $\forall + >1$, $P(a_1 \in B(a_1, p/z))$ under $(E_b)_{x \in A} = 0$.

The considered strategy is therefore such that the ap have the same distribution under $(\overline{J}_{\cdot})_{\times A}$ and $v^{(a)}$. In particular, $\mathbb{E}\left[\overline{L}_{\cdot 2}^{\times}\right]=0$ in both cases, but in the latter case $\sup_{A}\int_{0}^{(a)}=1$, so that $R_{T}=T$ against $v^{(a)}$. The regret is thus not controlled against $v^{(a)}\in \mathcal{F}$ and