

## Lecture #4: Continuous and linear bandits

In many applications, the agent can observe a context or just (eg user information in the case of online recommendation)

### Bandits with continuum of arms

Before that, let us consider the problem of continuous bandits

#### Setting (continuous bandits)

At each time step  $t=1, \dots, T$ :

- the agent pulls an arm  $a \in A \subset [0,1]^d$

- the agent observes and receives the reward  $Y_t = \mu(a_t) + \eta_t$  where  $\eta_t$  is

1 std-Gaussian  
0 mean  
independent noise

Goal: minimize regret  $R_T = T\mu^* - \sum_{t=1}^T \mu(a_t)$  where  $\mu^* = \sup_{a \in A} \mu(a)$

Without any assumption on  $\mu$  or  $A$ , we cannot do anything.

Otherwise, for  $A = [0,1]$  and any algorithm, we can choose  $x$  s.t.

$$\forall t \in \mathbb{N}, \underset{\mu=0}{\mathbb{P}}(a_t = x) = 0$$

Then for  $\mu(x) = 1I_{a=x}$ , the same algorithm would behave as if  $\mu=0$  and never pull  $x$ ,

so that its regret is  $R_T = T$ .

Continuity is actually enough to control the regret (i.e. having  $\mathbb{E}[R_T] = o(T)$ )

To get precise bounds, we will consider a stronger assumption on  $\mu$ .

$\mu$  is  $\beta$ -Hölder if there exists  $c > 0$  s.t.  $\forall a, a' \in A, |\mu(a) - \mu(a')| \leq c \|a - a'\|^\beta$ .

## (meta) Algo: Binning UCB

Input:  $\epsilon > 0$

Let  $K$  be an  $\epsilon$ -covering of minimal cardinal of  $A$ .

Pull each arm in  $K$  once

Then for any  $t \geq |K|$ :

$$\text{pull at } \arg\max_{a \in K} \hat{\mu}_k(t-1) + \sqrt{\frac{c \ln(t)}{N_k(t-1)}} \quad )$$

run UCB on the set of arms  $K$

(different constant because we have 1-sub Gaussian rewards)

## Theorem

Let  $\beta > 0$  and  $\epsilon > 0$ . Assume that  $\mu$  is  $\beta$ -Hölder and  $A \subset [0,1]^d$ .

Then running binning UCB with  $\mu^*$  yields the regret:

$$\mathbb{E}[R_T] \leq c \left( T \epsilon^\beta + \sqrt{\frac{T \ln T}{\epsilon^d}} \right) \quad \text{for some universal constant } c > 0.$$

---

Taking  $\epsilon$  of order  $\left(\frac{\ln T}{T}\right)^{\frac{1}{2\beta+1}}$  yields a regret  $\mathbb{E}[R_T] = O\left(T^{\frac{\beta+d}{2\beta+1}} \ln(T)^{\frac{\beta}{2\beta+1}}\right)$

## Proof:

$$R_T = \sum_{t=1}^T \mu^* - \mu_{\text{act}}$$

$$= \sum_{t=1}^T (\mu^* - \max_{a \in K} \mu(a)) + \sum_{t=1}^T \max_{a \in K} \mu(a) - \mu_{\text{act}}$$

$\leq c \varepsilon^\beta$  by  
Höldor assumption

regret of VCB run on  $K$ .

$$\text{so for } K = \text{card}(K), \quad \mathbb{E}[R_T] \leq cT\varepsilon^\beta + \underbrace{c\sqrt{K \ln(T)}}_{\text{distribution free bound of VCB}}$$

By minimisation of cardinal, a classical covering number bound yields  $K \leq \varepsilon^{-d}$ , which concludes

## Contextual bandits

Motivation

### Setting 1 (contextual bandits)

For each round  $t=1, \dots, T$ :

- agent observes context  $c_t$  (Conditionally chosen by nature)
- agent chooses action  $a_t \in [K]$ , depending on  $c_t$  and past observations.

- agent observes and gets reward  $y_t$

where  $y_t = r(a_t, c_t) + \eta_t$

with

$\eta_t$

1 sub-Gaussian

0 mean  
independent noise

$$\forall t \in \mathbb{R}, \quad \mathbb{E}[e^{\lambda \eta_t}] \leq \exp\left(\frac{\lambda^2}{2}\right)$$

$r: [K] \times \mathcal{C} \rightarrow \mathbb{R}$  is called the reward function

$\leftarrow$  object to estimate

(pseudo)-regret defined as:  $R_T = \sum_{t=1}^T \left\{ \max_{k \in [K]} r(k, c_t) - r(a_t, c_t) \right\}$

Without any assumption on  $\pi$ , independent bandit games for each context  $c$ .

- First possibility,  $\pi$  is "regular" (e.g. Lipschitz or Hölder)

In that case, we can again run a binning version of UCB to discretize the **context set** (instead of action set)

## Binning UCB (contextual)

Input:  $\epsilon > 0$

Let  $X$  be an  $\epsilon$ -covering of minimal cardinal of  $\mathcal{C}$

For each time step  $t \geq 1$ :

- observe the context  $c_t$

- let  $B \in X$  a ball containing  $c_t$ :  $c_t \in B$

- pull the arm  $a_t$  following the UCB algorithm on the bin  $B$ , i.e. with

$$N_k^B(t) = \sum_{a=1}^r \mathbb{1}_{a \in k} \mathbb{1}_{c_a \in B}$$

$$\hat{\mu}_k^B(t) = \frac{1}{N_k^B(t)} \sum_{a=1}^r \mathbb{1}_{a \in k} \mathbb{1}_{c_a \in B} Y_a$$

## Theorem:

Let  $\beta > 0$  and  $\epsilon > 0$ . Assume that  $x \mapsto \mu(b, x)$  is  $\beta$ -Hölder for any  $b \in [K]$  and  $\mathcal{C} \subseteq [0, 1]^d$

The regret of (contextual) binning UCB is then bounded as:

$$\mathbb{E}[R_T] \leq C \left( T\epsilon^\beta + \sqrt{\frac{KT \ln T}{\epsilon^\beta}} \right)$$

Choosing  $\epsilon$  of order  $(\frac{K \log T}{T})^{\frac{1}{2\beta+d}}$ , we get  $\mathbb{E}[R_T] = O\left(T^{\frac{\beta+1}{2\beta+d}} (K \ln T)^{\frac{\beta}{2\beta+d}}\right)$

Proof is not as direct as the continuous bandits case: UCB does not act with iid variables here, but random variables of the type:  $X_a(t) = \mu_a + \eta_a(t)$  with  $\eta_a(t) \stackrel{\text{iid Gaussian}}{\sim} \mathcal{N}(\mathbb{E}[\eta_a(t)], \sigma^2)$

(left as exercise)

## Remarks:

algo  
with an optional distribution free UCB bound  
↓

- In all the previous algs, we can replace VCB by MOSS to get rid of the  $\log T$  terms.

- Instance dependent bounds? In the above results, we used the distribution free bound of VCB. Can we get something better, depending on the regularity level of  $\mu$ ?

Yes, with the  $\alpha$ -margin assumption for iid contexts  $c_r$  and all  $\xi \in [0, 1]$

The larger the  $\alpha$ , the easier the problem

$$\Pr \left( \min_{k, \Delta(k, c_r) > 0} \Delta(k, c_r) < \xi \right) \leq c \xi^\alpha ; \quad \text{where } \Delta(k, c) = \max_a \mu(k, c) - \mu(k, c_r)$$

## Theorem

Let  $\beta > 0$  and  $\alpha \in (0, 1)$ . Assume that  $x \mapsto \mu(k, x)$  is  $\beta$ -Hölder for any  $k \in [K]$ ,  $x \in [0, 1]^d$  and the  $\alpha$ -margin assumption.

Then the regret of (contextual) binning VCB is then bounded as:

$$E[R_T] \leq C T^{\frac{\beta(1-\alpha)+1}{2\beta+d}} (K \ln T)^{\frac{\beta(1-\alpha)}{2\beta+d}}$$

for an optimised  $C > 0$

---

the proof is intricate

## Linear bandits

Another possible assumption is that  $\pi$  is linear with respect to a known feature

map  $\Psi: [K] \times \mathcal{C} \rightarrow \mathbb{R}^d$  and a parameter  $\Theta^* \in \mathbb{R}^d$  such that  
to estimate

$$\pi(k, c) = \langle \Theta^*, \Psi(k, c) \rangle \quad \forall k, c.$$

This is equivalent to the following setting, with  $A_T = \{\Psi(k, a_t) \mid k \in [K]\}$ :

## Setting 2 (linear bandits)

For each round  $t = 1, \dots, T$ :

- agent observes decision set  $A_t \subset \mathbb{R}^d$

- agent chooses action  $a_t \in A_t$

- agent observes and gets reward  $y_t$

$$\text{where } y_t = \langle \Theta^*, a_t \rangle + \gamma_t$$

with

$$\gamma_t$$

1 sub-Gaussian  
0 mean  
independent noise

Some defn of regret

Particular cases:

- $A_t = \{e_1, \dots, e_d\} \rightarrow$  classical multi-armed bandits with  $d$  arms and  $\mu_a = \Theta_a^*$

- $A_t \subset \{0, 1\}^d \rightarrow$  combinatorial bandits.

We want to build an adaptation of UCB for linear bandits, called

LinUCB.

- build confidence regions  $C_t$  such that  $\Theta^* \in C_t$  with high probability

- build confidence bounds for the arm means  $U_a(t) = \max_{\Theta \in C_t} \langle a, \Theta \rangle$

$$\underline{\Theta \in C_t}$$

- be optimistic: pull  $a_t \in \operatorname{argmax}_{a \in A_t} U_a(t)$

UCB score of arm  $a$ .

main question (

Before the confidence set, what is the estimate of  $\Theta^*$ ? (ie "empirical mean")

Regularised least-squares estimator:

$$\hat{\Theta}_T = \underset{\Theta \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{s=1}^T (Y_s - \langle \Theta, a_s \rangle)^2 + \lambda \|\Theta\|_2^2$$

$\lambda > 0$  is the penalty factor (or regularization parameter)

$\lambda > 0$  ensures uniqueness of the minimiser

We can indeed easily check that:

$$\hat{\Theta}_T = V_T^{-1} \sum_{s=1}^T a_s Y_s \quad \text{where } V_T = \lambda I_d + \sum_{s=1}^T a_s a_s^T$$



For any symmetric, positive definite matrix  $M \in \mathbb{R}^{d \times d}$  and vector  $u \in \mathbb{R}^d$ , we denote

$$\|u\|_M^2 := (u^T M u)$$

Theorem (linear bandits concentration)

For any  $\delta \in (0, 1)$ ,  $t \in \mathbb{N}$  and  $\lambda > 0$ , if for all  $s$ ,  $\max_{a \in \mathcal{A}} \|a\|_2 \leq 1$ , then with probability at least  $1 - \delta$ ,

$$\|\hat{\Theta}_T - \Theta^*\|_{V_T} \leq \sqrt{\lambda \|\Theta^*\|_2^2 + \sqrt{2 \ln(\frac{1}{\delta}) + d \ln(1 + \frac{T}{\lambda d})}}$$

The proof relies on the following concentration lemma

### Lemma

Let  $S_r = \sum_{s=1}^t Y_s$  as

For any  $\lambda > 0, r \in \mathbb{N}$  and  $\delta \in (0, 1)$ ,

$$\mathbb{P}\left(\|S_r\|_{V_r^{-1}}^2 \geq 2\ln\left(\frac{1}{\delta}\right) + \ln\left(\frac{\det(V_r)}{\lambda^d}\right)\right) \leq \delta.$$

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

$$\begin{aligned} \text{Note that } \hat{\theta}_r &= V_r^{-1} \left( S_r + \sum_{s=1}^r a_s a_s^\top \theta^* \right) \\ &= V_r^{-1} S_r + V_r^{-1} (V_r - \lambda \text{Id}) \theta^* \end{aligned}$$

$$\begin{aligned} \text{So } \|\hat{\theta}_r - \theta^*\|_{V_r} &= \|V_r^{-1} S_r - \lambda V_r^{-1} \theta^*\|_{V_r} \\ &\leq \|V_r^{-1} S_r\|_{V_r} + \lambda \|V_r^{-1} \theta^*\|_{V_r} \\ &= \|S_r\|_{V_r^{-1}} + \underbrace{\lambda \|\theta^*\|_{V_r}}_{\sqrt{\theta^{*\top} V_r^{-1} \theta^*}} \leq \|V_r^{-1}\|_{op}^{\frac{1}{2}} \|\theta^*\|_2 \\ &\leq \|S_r\|_{V_r^{-1}} + \sqrt{\lambda} \|\theta^*\|_2. \quad \leq \lambda_{\min}(V_r)^{-\frac{1}{2}} \|\theta^*\|_2 \\ &\leq \lambda^{-\frac{1}{2}} \|\theta^*\|_2 \end{aligned}$$

It remains to show that  $\frac{\det(V_r)}{\lambda^d} \leq \left(1 + \frac{\tau}{\lambda}\right)^d$ , i.e.  $\det(V_r) \leq (\lambda + \tau)^d$ .

Indeed, we have:  $\det(V_r) \leq \left(\frac{\text{tr}(V_r)}{d}\right)^d \leq \left(\lambda + \frac{1}{d} \sum_{s=1}^r \text{tr}(a_s a_s^\top)\right)^d \leq \left(\lambda + \frac{\tau}{d}\right)^d$

□

# Proof of the lemma

(optional)

For any  $x \in \mathbb{R}^d$ , define  $M_r(x) = \exp\left(\langle x, s_r \rangle - \frac{1}{2} \|x\|_{V_r + \lambda I}^2\right)$

1) We show by induction that  $M_r(x)$  is a supermartingale, so that

$$\mathbb{E}[M_r(x)] \leq M_0(x) = 1$$

$r \rightarrow r+1$

$$M_{r+1}(x) = \exp\left(\langle x, s_{r+1} \rangle - \frac{1}{2} (x^T (V_{r+1} - \lambda I) x)\right)$$

$$V_{r+1} = V_r + a_{r+1} a_{r+1}^T$$

$$= M_r(x) \cdot \exp\left(\langle x, a_{r+1} \rangle \eta_{r+1} - \frac{1}{2} \langle x, a_{r+1} \rangle^2\right).$$

$$\mathbb{E}[M_{r+1}(x) | \mathcal{F}_r] \leq M_r(x)$$

( $\eta_{r+1}$  is 1 sub-Gaussian)

2) Let  $v = \mathcal{N}(0, \lambda^{-1} I_d)$ .

$$\bar{M}_r = \int M_r(x) d\nu(x)$$

is also a supermartingale  
by Tonelli and

$$\bar{M}_r = \frac{1}{\sqrt{(2\pi)^d \lambda^d}} \int_{\mathbb{R}^d} \exp\left(\langle x, s_r \rangle - \frac{1}{2} \|x\|_{V_r + \lambda I}^2 - \frac{1}{2} \|x\|_{\lambda I}^2\right) d\nu$$

$$S = x^T s_r - \frac{1}{2} x^T V_r x$$

$$= -\frac{1}{2} (x - V_r^{-1} s_r)^T V_r (x - V_r^{-1} s_r) + \frac{1}{2} s_r^T V_r^{-1} s_r$$

$$= -\frac{1}{2} \left\| \alpha \cdot V_t^{-1} S_t \right\|_{V_t^{-1}}^2 + \frac{1}{2} \| S_t \|_{V_t^{-1}}^2$$

$$\bar{M}_t = \exp\left(\frac{1}{2} \| S_t \|_{V_t^{-1}}^2\right) \cdot \left(\frac{\lambda}{2\pi}\right)^{d/2} \int_{\mathbb{R}^d} \exp\left(-\frac{1}{2} \left\| \alpha \cdot V_t^{-1} S_t \right\|_{V_t^{-1}}^2\right) d\alpha$$

upto scaling,  
pdf of  $N(V_t^{-1} S_t, V_t)$

$$= \exp\left(\frac{1}{2} \| S_t \|_{V_t^{-1}}^2\right) \overline{\frac{\lambda^{d/2}}{\det(V_t)}}$$

$$\| S_t \|_{V_t^{-1}}^2 = 2 \ln(\bar{M}_t) - \ln\left(\frac{\lambda^d}{\det(V_t)}\right)$$

3)

$$\mathbb{P}\left(\| S_t \|_{V_t^{-1}}^2 \geq 2 \ln\left(\frac{1}{\delta}\right) + \ln\left(\frac{\det(V_t)}{\lambda^d}\right)\right) = \mathbb{P}\left(\ln(\bar{M}_t) \geq \ln\left(\frac{1}{\delta}\right)\right)$$

$$= \mathbb{P}\left(\bar{M}_t \geq \frac{1}{\delta}\right) \leq \mathbb{E}[\bar{M}_t] \leq \delta.$$

Alg Lin VCB

For each  $t \in \mathbb{N}$

$$\text{Play } a_t \in \arg\max_{a \in A_t} \max_{\theta \in \Theta_{t-1}} \langle \theta, a_t \rangle$$

suppose we know  $m$  with  $\|\theta\|_2 \leq m$

can be computed efficiently for our specific form of  $\theta_t$  and nice  $A_t$ .

with  $\hat{\theta}_t = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \sum_{s=1}^t (y_s - \langle \theta, a_s \rangle)^2 + \lambda \|\theta\|_2^2$

$$V_t = \lambda I + \sum_{s=1}^t a_s a_s^\top$$

and  $\mathcal{C}_t = \left\{ \theta \in \mathbb{R}^d \mid \|\hat{\theta}_t - \theta\|_{V_t} \leq \sqrt{\lambda m + \sqrt{4 \ln(t) + d \ln(4 + \frac{t}{\lambda})}} \right\}$

## Theorem:

If  $\|\theta^*\|_2 \leq 1$  and for any  $t$ ,  $\max_{a \in \mathcal{A}^t} \|a\|_2 \leq 1$ , then the regret of LinUCB satisfies for any  $\lambda > 0$ ,

$$\mathbb{E}[R_T] \leq c_d \sqrt{T \ln T}$$

where  $c_d$  is a constant that only depends on  $d$ .

## Comments:

- distribution free bound.
- if  $\mathcal{A}_t$  is finite, and the same for every  $t$ , we can get a  $\log(t)$  instance dependent bound:  $\mathbb{E}[R_t] \leq c \sqrt{T_d \ln(T_K)}$
- another possible improvement when  $d \gg 1$  is to assume that  $\theta^*$  is  $m_\theta$ -sparse  
Then, we can get a regret of order  $\tilde{O}(\sqrt{d m_\theta T})$

Proof:

Let us bound the instantaneous regret first.

$$r_t = \langle \theta^*, A_t^* - a_t \rangle \quad \text{when } A_t^* \in \underset{a \in A_t}{\arg\max} \langle \theta^*, a \rangle$$

Define the  
good event

$$\mathcal{E}_t = \left\{ \theta^* \in \mathcal{C}_{t-1} \right\}$$

Thanks to our concentration theorem,  $P(\neg \mathcal{E}_t) \leq \frac{1}{(t-1)^2}$

$$\begin{aligned} \mathbb{E}[r_t] &\leq 2 \cdot P(\neg \mathcal{E}_t) + \mathbb{E}[r_t \mathbf{1}_{\mathcal{E}_t}] \\ &\leq \frac{2}{(t-1)^2} + \mathbb{E}[r_t \mathbf{1}_{\mathcal{E}_t}]. \end{aligned}$$

If  $\mathcal{E}_t$ ,  $\theta^* \in \mathcal{C}_{t-1}$  so:

$$\langle \theta^*, A_t^* \rangle \leq \max_{\theta \in \mathcal{C}_{t-1}} \langle \theta, A_t^* \rangle$$

$$\leq \max_{\theta \in \mathcal{C}_{t-1}} \langle \theta, a_t \rangle \quad \text{by defn of } a_t$$

$$= \langle \tilde{\theta}_t, a_t \rangle \text{ for some } \tilde{\theta}_t \in \mathcal{C}_{t-1}.$$

Cauchy-Schwarz gives:

$$r_r = \langle \theta^*, A_r^* \cdot a_r \rangle \leq \langle \tilde{\theta}_r - \theta^*, a_r \rangle \leq \|\tilde{\theta}_r - \theta^*\|_{V_{r-1}} \|a_r\|_{V_{r-1}}$$

$$\leq \|a_r\|_{V_{r-1}} \left( \|\tilde{\theta}_r - \hat{\theta}_{r-1}\|_{V_{r-1}} + \|\theta^* - \hat{\theta}_{r-1}\|_{V_{r-1}} \right)$$

$$\leq 2 \|a_r\|_{V_{r-1}} \cdot \left( \sqrt{\lambda} + \sqrt{4 \ln(t) + t \ln\left(1 + \frac{1}{\lambda d}\right)} \right)$$

define  $\alpha_r = \max(\cdot, 1)$

also by assumption)  $r_r \leq 2$ , so

$$r_r \leq 2 \alpha_r \left( 1 \wedge \|a_r\|_{V_{r-1}} \right) \quad (\text{if } E_r \text{ holds})$$

$\alpha_y = \min(\alpha_y)$

overall:

$$R_T \leq \sum_{T=2}^T \mathbb{E}[s_T \mathbf{1}_{E_T}] + \sum_{t=1}^T \left( \frac{1}{(t+1)^2} \wedge 1 \right)$$

$$\leq 2 \sum_{r=2}^T \alpha_r \left( 1 \wedge \|a_r\|_{V_{r-1}}^{-1} \right) + c$$

$$\leq 2 \sqrt{\sum_{r=2}^T \alpha_r^2} \sqrt{\sum_{r=2}^T \left( 1 \wedge \|a_r\|_{V_{r-1}}^{-1} \right)^2} + c$$

$$\leq c_1 \sqrt{\sum_{r=2}^T d \ln(T)} \sqrt{\sum_{r=2}^T \left( 1 \wedge \|a_r\|_{V_{r-1}}^{-1} \right)^2} + c$$

$$\leq c_1 \sqrt{d \ln(T)} \sqrt{\sum_{r=2}^T \left( 1 \wedge \|a_r\|_{V_{r-1}}^{-1} \right)^2} + c$$

Bound on  $\sum_{r=2}^T \left( 1 \wedge \|a_r\|_{V_{r-1}}^{-1} \right)$

$$u \wedge 1 \leq 2 \ln(1+u)$$

$$\stackrel{10}{\leq} 2 \sum_{r=2}^T \left( 1 \wedge \|a_r\|_{V_{r-1}}^{-1} \right) \leq 2 \sum_{r=1}^T \ln \underbrace{\left( 1 + \|a_r\|_{V_{r-1}}^{-1} \right)^2}_{= \ln \left( \det \left( \frac{V_r}{V_0} \right) \right)}$$

$$\text{Indeed, } V_r = V_{r-1} + \alpha_r \alpha_r^\top = V_{r-1}^{1/2} \left( I + V_{r-1}^{-1/2} \alpha_r \alpha_r^\top V_{r-1}^{-1/2} \right) V_{r-1}^{1/2}$$

$$\Rightarrow \det(V_r) = \det(V_{r-1}) \cdot \det(I + V_{r-1}^{-1/2} \alpha_r \alpha_r^\top V_{r-1}^{-1/2})$$

$\alpha_r \alpha_r^\top$  is a rank one matrix.

If  $y^\top$  has eigenvalues:  $(\underbrace{1 + \|y\|^2, 1, \dots, 1}_{\text{eigenvalues } y})$

$$\det(V_r) = \det(V_{r-1}) \cdot (1 + \|V_{r-1}^{-1/2} \alpha_r\|_2^2)$$

$$= \det(V_{r-1}) (1 + \|\alpha_r\|_{V_{r-1}^{-1}}^2)$$

$$\text{So by induction } \ln(\det(V_r)) = \ln(\det(V_0)) + \sum_{t=1}^r \ln(1 + \|\alpha_t\|_{V_{t-1}^{-1}}^2)$$

$$\text{so } \sum_{t=1}^T (1 + \|\alpha_t\|_{V_{t-1}^{-1}}^2) \leq 2 \ln \left( \frac{\det(V_T)}{\det(V_0)} \right)_d$$

$$\leq 2d \ln \left( 1 + \frac{T}{\lambda_d} \right)$$

$$\leq c_d d \ln(T)$$

Thanks to previous bound.

In conclusion, gathering every thing we get.

$$R_T \leq c_d d \ln T \sqrt{T} + c.$$

D