

## Lecture #8: contextual/linear bandits

### Setting 1 (contextual bandits)

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

- agent observes context  $c_t \in \mathcal{C}$  (arbitrarily chosen by nature)

- agent chooses action  $a_t \in [K]$   $a_t$  is measurable w.r.t.

$$F_{t-1} = \sigma(U_0, C_1, Y_1, U_1, \dots, Y_{t-2}, U_{t-1}, c_t)$$

- agent observes and gets reward  $y_t$

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

with  $\eta_t | F_{t-1}$  is 1 sub-Gaussian  
0 mean

$$\begin{cases} \mathbb{E}[\eta_t | F_{t-1}] = 0 \\ \forall \lambda \in \mathbb{R}, \mathbb{E}[e^{\lambda \eta_t} | F_{t-1}] \leq \exp\left(\frac{\lambda^2}{2}\right) \end{cases}$$

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

← object to estimate

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

Without any assumption on  $r$ , independent bandit games for each context  $c$

- First possibility,  $r$  is "regular" (e.g. Lipschitz or Hölder) (see exercise sheet)
- A common assumption is that  $r$  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

$$r(k, c) = \langle \theta^*, \Psi(k, c) \rangle \quad \forall k, c$$

$$r(k, c) = \langle \theta^*, \Psi(k, c) \rangle \quad \forall k, c$$

This is equivalent to the following setting, with  $A_t = \{\psi(b, c_t) \mid b \in [k]\}$ :

## Setting 2 (linear bandits)

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

- agent observes decision set  $A_t \subset \mathbb{R}^d$
- agent chooses action  $a_t \in A_t$   $a_t$  is measurable w.r.t.  $F_{t-1} = \sigma(V_0, C_1, Y_1, U_1, \dots, Y_{t-1}, V_{t-1}, C_t)$
- agent observes and gets reward  $y_t$   
where  $y_t = \langle \theta^*, a_t \rangle + \gamma_t$   
with  $\gamma_t | F_{t-1}$  is  $\mathcal{N}$  sub-Gaussian  
mean

Particular cases:

- $A_t = \{e_1, \dots, e_d\} \rightarrow$  classical multi-armed bandits with  $d$  arms and  $\mu_a = \theta_k^*$
- $A_t \subset \{0,1\}^d \rightarrow$  combinatorial bandits.

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

LinUCB

The idea is to construct confidence sets  $C_t$  such that  $\theta^* \in C_t$  with high probability and pick at each round

$$a_t \in \operatorname{argmax}_{a \in A_t} \max_{\theta \in C_t} \langle a, \theta \rangle$$

(with  $C_t$  as small as possible)

UCB score of arm a.

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  $\lambda > 0$ :

$$P\left(\|\hat{\theta}_t - \theta^*\|_{V_t} \geq \sqrt{\lambda \|\theta^*\|_2^2 + \sqrt{2 \ln\left(\frac{1}{\delta}\right) + \ln\left(\frac{d \cdot t(V_t)}{\lambda}\right)}}\right) \leq \delta$$

The proof relies on the following concentration lemma

### Lemma

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

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

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

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

$$\begin{aligned} \text{Note that } \hat{\theta}_t &= V_t^{-1} \left( S_t + \sum_{s=1}^t a_s a_s^\top \theta^* \right) \\ &= V_t^{-1} S_t + V_t^{-1} (V_t - \lambda \text{Id}) \theta^* \end{aligned}$$

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

□

$$\leq \lambda^{-\frac{1}{2}} \|\theta^*\|_2$$

# Proof of the lemma

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

1) We show by induction that  $M_t(x)$  is a martingale, so that

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

$t \rightarrow t+1$

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

$$V_{t+1} = V_t + a_{t+1} a_{t+1}^T$$

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

$$\mathbb{E}[M_{t+1}(x) | \mathcal{F}_t] \leq M_t(x)$$

( $\eta_{t+1} | \mathcal{F}_t$  is 1 sub-Gaussian)

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

$$\bar{M}_t = \int M_t(x) d\nu(x)$$

is also a martingale  
by Tonelli and

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

$$S = x^T s_t - \frac{1}{2} x^T V_t x$$

$$= -\frac{1}{2} (x \cdot V_t^{-1} s_t)^T V_t (x - V_t^{-1} s_t) + \frac{1}{2} s_t^T V_t^{-1} s_t$$

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

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

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 UCB

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

## Theorem:

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

$$\mathbb{E}[R_T] \leq c_1 \sqrt{T m^2 \lambda + \ln(t) T + d \ln(1 + \frac{T}{\lambda})} \sqrt{d \ln(1 + \frac{T}{\lambda})} + c_2 m L \quad \text{for univ constants } c_1, c_2$$

Corollary: Taking  $\lambda = 1$  and considering the main factor in  $T$  we have:

$$\mathbb{E}[R_T] = O(d \sqrt{T} \ln T)$$

## Comments:

- distribution free bound
- if  $A_t$  is finite, and the same for every  $t$ , we can get a  $\log(t)$  instance dependent bound.

Proof:

Let us bound the instantaneous regret first.

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

Define the good event

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

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

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

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

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

$$\leq \max_{\theta \in C_{r+1}} \langle \theta, a_r \rangle \quad \text{by defn of } a_r.$$

$$= \langle \tilde{\theta}_r, a_r \rangle \text{ for some } \tilde{\theta}_r \in C_{r+1}.$$

Cauchy-Schwarz gives:

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

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

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

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

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

$$r_r \leq 2\alpha_r (1 + \|a_r\|_{V_{r+1}}) \quad (\because E_r \text{ holds})$$

overall,

$$R_T \leq \sum_{r=2}^T \mathbb{E}[r_r \mathbf{1}_{C_r}] + mL \left( \frac{1}{(r-1)^2} \mathbf{1} \right) + mL$$

$$\leq 2 \sum_{r=1}^T \alpha_r \left( 1 + \|a_r\|_{V_{r+1}}^{-1} \right) + m L \left( 1 + \frac{\pi^2}{6} \right)$$

$$\leq 2 \sqrt{\sum_{r=1}^T \alpha_r^2} \sqrt{\sum_{r=1}^T \left( 1 + \|a_r\|_{V_{r+1}}^{-1} \right)^2} + m L \left( 1 + \frac{\pi^2}{6} \right)$$

$$\leq 2 \sqrt{2 \sum_{r=1}^T \lambda_m^2 + 4 \ln(1) + \ln\left(\frac{\det(V_r)}{\lambda^d}\right)} \sqrt{\sum_{r=1}^T \left( 1 + \|a_r\|_{V_{r+1}}^{-1} \right)^2} + m L \left( 2 + \frac{\pi^2}{6} \right)$$

Bound on:  $\ln\left(\frac{\det(V_r)}{\lambda^d}\right)$

$$\frac{1}{\lambda^d} \lambda_i \leq \left(\frac{\sum_{i=1}^d \lambda_i}{d}\right)^d \quad (\text{arithmetic vs geometric mean})$$

$$\frac{\det(V_r)}{\lambda^d} = \det\left(\frac{V_r}{\lambda}\right) \leq \left(\frac{\ln\left(\frac{V_r}{\lambda}\right)}{d}\right)^d = \left(\frac{\ln(V_r)}{\lambda d}\right)^d$$

$$\ln(V_r) = \ln\left(\lambda I + \sum_{i=1}^r a_i a_i^\top\right) = \lambda d + \sum_{i=1}^r \ln(a_i a_i^\top)$$

$$\leq \lambda d + r L^2.$$

$$\approx \ln\left(\frac{\det(V_r)}{\lambda^d}\right) \leq d \ln\left(1 + \frac{r L^2}{\lambda d}\right)$$

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

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

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

$$= \ln\left(\det\left(\frac{V_r}{V_0}\right)\right)$$

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

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

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

$I + \alpha_r \alpha_r^\top$  has eigenvalues:  $(1 + \|\alpha_r\|^2, 1, \dots, 1)$

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

eigenvalue  $\alpha_r$

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

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

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

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

Thanks to previous bound.

In conclusion, gathering every thing we get.

$$R_T \leq 4 \sqrt{2 T m^2 \lambda + 4 \ln(T) T + d \ln \left( 1 + \frac{T L^2}{\lambda d} \right)} \sqrt{d \ln \left( 1 + \frac{T L^2}{\lambda d} \right)} + mL \left( 2 + \frac{L^2}{\delta} \right)$$

D

LinUCB has regret  $R_T = O(d \sqrt{T} \ln T)$ .

Can we do better?

**Theorem** (minimax lower bound, linear bandits)

Let  $A_T = [-1, 1]^d$  and  $\Theta = \left\{ -\frac{1}{\sqrt{T}}, \frac{1}{\sqrt{T}} \right\}^d$ . Then for any algorithm, there exists  $\theta^* \in \Theta$  s.t.:

$$\mathbb{E}[R_T(\theta^*)] \geq \frac{e^{-2}}{8} d \sqrt{T}$$

Here  $m = \sqrt{\frac{d}{T}}$  and  $L = \sqrt{d}$ . So for  $T \geq d$ , LinUCB is optimal, up to a  $\ln T$  term.

Proof relies on the following lemma, which is the equivalent of the fundamental inequality for linear bandits.

**Lemma** (fundamental inequality, linear bandits)

For all linear bandit instances,  $\Theta$  and  $\Theta'$  in  $\mathbb{R}^d$

for all strategies and events  $A$  that are  $\sigma(H_T)$ -measurable,

$$\mathbb{E}_{\theta} \left[ \sum_{t=1}^T \text{KL}(N(a_t, \theta), N(a_t, \theta')) \right] = \text{KL}(\mathbb{P}_{\theta}^{H_T}, \mathbb{P}_{\theta'}^{H_T})$$

↑  
law of  $H_T$  under  $\theta$  (and  $\pi$ )

and  $\frac{1}{2} e^{-\text{KL}(\mathbb{P}_{\theta}^{H_T}, \mathbb{P}_{\theta'}^{H_T})} \leq \mathbb{P}_{\theta}^{H_T}(A) + \mathbb{P}_{\theta'}^{H_T}(A^c)$  (Bartognolle-Huber inequality)

## Proof of the Theorem

Consider any  $\theta \in \Theta$  and  $i \in [d]$ . Define  $\theta'$  as  $\theta'_j = \begin{cases} \theta_j & \text{if } i \neq j \\ -\theta_j & \text{if } i = j \end{cases}$

1) By the above equality:

$$\begin{aligned} \text{KL}(P_{\theta}^{H_T}, P_{\theta'}^{H_T}) &= \mathbb{E}_{\theta} \left[ \sum_{t=1}^T \text{KL}(N(a_t, \theta), N(a_t, \theta'), 1) \right] \\ &= \frac{1}{2} \sum_{t=1}^T \mathbb{E}[|a_t - \theta'|^2] \quad \text{KL between Gaussian of same variance} \\ &\leq \frac{1}{2} \sum_{t=1}^T \frac{4}{T} = 2 \end{aligned}$$

2) Moreover, define the event  $A(\theta, i) = \left\{ \sum_{t=1}^T \mathbb{1}_{\{\text{sgn}(a_{it}) \neq \text{sgn}(\theta_{it})\}} \geq \frac{T}{2} \right\}$ .

Then Bretagnolle-Huber inequality states:

$$P_{\theta}^{\circ}(A(\theta, i)) + P_{\theta'}^{\circ}(A(\theta', i)^c) \geq \frac{1}{2} e^{-\text{KL}(P_{\theta}^{H_T}, P_{\theta'}^{H_T})} \geq \frac{1}{2} e^{-2}$$

$\underbrace{\hspace{10em}}$

1

$$P_{\theta}^{\circ}(A(\theta, i)) + P_{\theta'}^{\circ}(A(\theta', i))$$

so that, when summing over all possible  $\theta$ :

(Catergorizing hammon argument)

$$\sum_{\theta \in \Theta} P_{\theta}^{\circ}(A(\theta, i)) + P_{\theta'}^{\circ}(A(\theta', i)) = \sum_{\theta \in \Theta} P_{\theta}^{\circ}(A(\theta), i) \geq \frac{|\Theta|}{2} e^{-2}$$

Now summing overall  $i$ :

$$\frac{1}{|\Theta|} \sum_{\theta \in \Theta} \sum_{i \in [d]} P_{\theta}^{\circ}(A(\theta), i) \geq \frac{d}{4} e^{-2}.$$

In particular, there exists  $\theta \in \Theta$  s.t.  $\sum_{i \in [d]} P_{\theta}^{\circ}(A(\theta), i) \geq \frac{d}{4} e^{-2}$ .

Consider such a  $\theta$  in the following.

Then we can lower bound the regret on  $\theta$ :

$$R_T(\theta) = \sum_{t=1}^T \max_{a \in A_t} \langle a - a_t, \theta \rangle$$

$$= \sum_{t=1}^T \sum_{i=1}^d (\text{sgn}(\theta_i) - \alpha_{it}) \cdot \theta_i$$

$$= \frac{1}{\sqrt{T}} \sum_{t=1}^T \sum_{i=1}^d |\text{sgn}(\theta_i) - \alpha_{it}|$$

$$\geq \frac{1}{\sqrt{T}} \sum_{t=1}^T \sum_{i=1}^d \mathbb{1}_{\{\text{sgn}(\theta_i) \neq \text{sgn}(\alpha_{it})\}}$$

$$\text{so that } E[R_T(\theta)] \geq \frac{1}{\sqrt{T}} \sum_{i=1}^d \sum_{t=1}^T P(A(\theta_i, t))$$

$$\geq \frac{\sqrt{T}}{3} d e^{-2} \quad \square$$