

## - Multi-armed bandits (MAB)

So far in online learning : agent receive instructive feedback

MAB: A basic model of agent making sequential decisions and receiving evaluative feedback.

A slot machines



a slot machine

"one armed bandit".

① Setup: action set  $A = \{1, \dots, A\}$ . (arm set)

For  $t=1, 2, \dots, T$ :

agent takes action  $a_t$  (arm)

agent receives reward  $r_t = f^*(a_t) + \xi_t$

$\Rightarrow$  (zero mean,  
+SG noise)

$f^*: A \rightarrow [0, 1]$ ,

agent's expected reward  $= \mathbb{E} \left[ \sum_{t=1}^T r_t \right] = \mathbb{E} \left[ \sum_{t=1}^T f^*(a_t) \right]$

optimal strategy if  $f^*$  is known: take  $a^* = \operatorname{argmax}_{a \in A} f^*(a)$

optimal expected reward  $= T \cdot f^*(a^*)$ .

Performance measure: regret:

$$\text{Reg}(T) = T f^*(a^*) - \mathbb{E} \left[ \sum_{t=1}^T f^*(a_t) \right] = \mathbb{E} \left[ \sum_{t=1}^T (f^*(a^*) - f^*(a_t)) \right]$$

How to design a strategy for the agent?

\* need to learn the  $f^*$  fn thru noisy observations

(exploration)  $\rightarrow$  Take diverse actions

\* need to take actions  $a_t$  that has large  $f^*$  value

(exploitation)  $\rightarrow$  Take actions that maximizes  $f^*$ .



Conflicting!

## ② Algorithms for MAB.

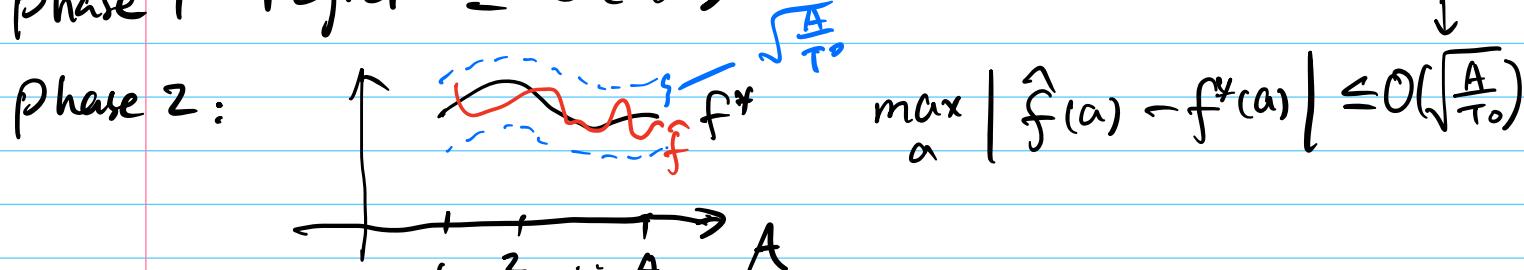
First idea: explore-then-exploit (commit)

Phase 1. use first  $T_0$  rounds to estimate  $f^*$  by taking actions in a round robin fashion  $\rightarrow \hat{f}$

Phase 2. from  $T_0+1$  round on. take action  $\hat{a} = \arg \max_{a \in A} \hat{f}(a)$

Analysis sketch:

Phase 1 regret  $\leq O(T_0)$



$$\Rightarrow f^*(a^*) - f^*(\hat{a}) \leq O(\sqrt{\frac{A}{T_0}})$$

$$\Rightarrow \text{Phase 2 regret} \leq O((T-T_0)\sqrt{\frac{A}{T_0}})$$

$$\Rightarrow \text{Reg}(T) \leq O(T_0 + T\sqrt{\frac{A}{T_0}}) = O(A^{\frac{1}{3}} T^{\frac{2}{3}}).$$

choosing  $T_0$  optimally

Second idea: intersperse explore & exploit

$\epsilon$ -greedy.

For  $t=1, 2, \dots, T$ :

Flip a coin  $Z$  w/ head prob.  $\varepsilon$ .

(Explore) If  $Z = \text{head}$ , take  $a_t \sim \text{Unif}(\{1 \dots A\})$

(Exploit) If  $Z = \text{tail}$ , take  $a_t = \hat{a}_t = \underset{a}{\operatorname{argmax}} \hat{f}_t(a)$

(Notation:  $\hat{f}_t(a) = \text{avg reward of } a \text{ at } t$

$$\frac{\text{total reward of } a}{\# \text{ times } a \text{ taken}} = \frac{\sum_{i=1}^t I(a_i=a) r_i}{N_{t-f}(a)} \longrightarrow \sum_{i=1}^t I(a_i=a)$$

### Analysis sketch

- at step  $t$ , each  $a$  has  $\approx t \frac{\varepsilon}{A}$  reward samples

$$\Rightarrow \max_a |\hat{f}_t(a) - f^*(a)| \leq O\left(\sqrt{\frac{A}{t\varepsilon}}\right).$$

$$\Rightarrow f^*(a^*) - f^*(\hat{a}_t) \leq O\left(\sqrt{\frac{A}{t\varepsilon}}\right)$$

$\Rightarrow$  "on average",  $a_t$  is  $(\varepsilon + \sqrt{\frac{A}{t\varepsilon}})$ -optimal

$$-\text{Reg}(T) \leq O\left(\sum_{t=1}^T \left(\varepsilon + \sqrt{\frac{A}{t\varepsilon}}\right)\right)$$

$$= O\left(\varepsilon T + \sqrt{\frac{AT}{\varepsilon}}\right).$$

setting  $\varepsilon$  optimally

$$= O\left(A^{\frac{1}{3}} T^{\frac{2}{3}}\right).$$

Can we do better?

A better idea: optimism in face of uncertainty  
(abbrev. optimism principle)

In words:  
Act according to the best plausible world.  
"optimistic world model"

Why it works:

If the optimistic world model is

correct ↘

no regret

✓

wrong ↙

learn new things

avoid making the same  
mistake in the future

✓

Real world agents are using this (perhaps w/o realizing it)  
the optimism bias (Sharot, 2011).

How to define the "best plausible world" in MAB?

For every action  $a$ , what's the highest plausible  $f^*(a)$   
given data?

upper confidence bound for  $f^*(a)$ !  
(UCB)

The UCB Algorithm:

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

- Define  $UCB_t(a) = \hat{f}_t(a) + b_t(a)$

define to be  $0$  if  $N_{t-1}(a) = 0$ .

confidence width.  $\Rightarrow$  "bonus" of action  $a$ .

$$b_t(a) = \sqrt{\frac{l}{N_{t-1}(a)+1}}, \quad l = 8 \ln(2AT).$$

- choose  $a_t = \underset{a \in A}{\operatorname{argmax}} UCB_t(a)$ .

(optimism property)

Validity of the UCB's:

Lemma: there exists an event  $E$ .  $P(E) \geq 1 - \frac{2}{T}$ .

and when  $E$  happens:

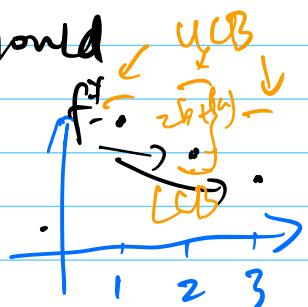
$$(*) \quad \forall a, \forall t, \quad |\hat{f}_t(a) - f^*(a)| \leq b_t(a) = \sqrt{\frac{l}{N_{t-1}(a)+1}}$$

$\hookrightarrow UCB_t(a) \geq f^*(a)$

Intuition: arm  $a$  reward estimate should

have precision  $\propto$

$$\frac{1}{\sqrt{\# \text{samples in arm } a}}$$



Pf idea:

$$\text{First idea: } E_{t,a} = \left\{ |\hat{f}_t(a) - f^*(a)| \leq b_t(a) \right\}$$

Hoeffding  $\Rightarrow \Pr[f_t(a) \geq 1 - \delta]$  ?

Problem:  $\hat{f}_t(a)$  is an average over  $N_{t+1}(a)$  samples,  
 $N_{t+1}(a)$  itself is random!

A fix:  $E_{t,a,n} = \left\{ \left| \hat{f}_t(a) - f^*(a) \right| \leq \sqrt{\frac{1}{n+1}} \beta(a), \right. \left. N_{t+1}(a) = n \right\}$

Hoeffding  $\Rightarrow \Pr(E_{t,a,n}) \geq 1 - \frac{2}{AT^3}$

Define  $E = \bigcap_{t=1}^T \bigcap_{a=1}^A \bigcap_{n=1}^T E_{t,a,n}$ .

$\Rightarrow \Pr(E) \geq 1 - \frac{2}{T}$ .

can show that on  $E$ , (\*) happens.

More details: Slivkins, "Introduction to Multi-armed bandits", 2019. See. 1.3.1. ↗

Regret analysis of UCB:

big-O ignoring  
log factors

Thm: Assume that

UCB guarantees that  $\text{Reg}(T) \leq \tilde{O}(\sqrt{AT})$

Remark: — regret bound much better than  $AT^{2/3}$  in ETC.  
↳ E-greedy

— in general unimprovable if there is a (lower)

- can also show UCB has instance dep<sup>bound</sup> regret:  $\text{Reg}(T) = \tilde{\mathcal{O}}\left(\sum_{a \neq a^*} \frac{\ln T}{\Delta a}\right)$ .  $\Delta a = f^*(a^*) - f^*(a)$
- $\Delta a$  longer, arm  $a$  less config.*

PF:

Recall:  $\text{Reg}(T) = \mathbb{E}\left[\sum_{t=1}^T (f^*(a^*) - f^*(a_t))\right]$

$$= \mathbb{E}\left[\left(\sum_{t=1}^T \text{reg}(t)\right)(I(E) + I(E^c))\right]$$

Note:  $\mathbb{E}\left[\left(\sum_{t=1}^T \text{reg}(t)\right) I(E^c)\right] \leq T \cdot \mathbb{E}[I(E^c)]$

$$= T \cdot P(E^c)$$

$$\leq T \cdot \frac{2}{T} = 2. \quad \textcircled{1}$$

We just need to control  $\sum_{t=1}^T \text{reg}(t)$  when  $E$  happens.

Insight 1: can bound  $\text{reg}(t)$  by  $b_t(a_t)$ : bonus on the action taken

$$\text{reg}(t) = f^*(a^*) - f^*(a_t)$$

validity of UCB  $\leq UCB_t(a^*) - f^*(a_t)$

choice of  $a_t$   $\leq UCB_t(a_t) - f^*(a_t)$

$$= \underbrace{\hat{f}_t(a_t) + b_t(a_t)}_{\text{event } E} - \underline{f^*(a_t)}$$

$$\leq 2b_t(a_t).$$

Insight 2:  $\sum_{t=1}^T b_t(a_t)$  is controlled.

$$n_t(a) = \sqrt{\frac{l}{N_{t-1}(a)+1}}$$

(every time  $a$  is chosen as  $a_t$ ,  
its bonus  $b_t(a)$  will shrink a bit)

$$\sum_{t=1}^T \text{reg}(t)$$

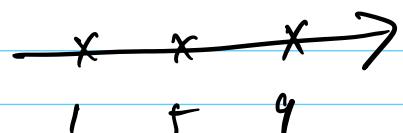
$$\leq 2 \sum_{t=1}^T b_t(a_t)$$

$$= 2 \sum_{a=1}^A \sum_{t: a_t=a} b_t(a)$$

(group according  
to  $a_t$ )

a decreasing sequence say.  $a=1$ .

$$= 2 \sum_{a=1}^A \sum_{t: a_t=a} \sqrt{\frac{l}{N_{t-1}(a)+1}}$$



$$= 2 \sum_{a=1}^A \sum_{n=1}^{N_T(a)} \sqrt{\frac{l}{n}}$$

$$\sqrt{\frac{l}{0+1}} \sqrt{\frac{l}{1+1}} \sqrt{\frac{l}{2+1}}$$

$$= 4\sqrt{l} \sum_{a=1}^A \sqrt{N_T(a)}$$

$$\left( \sum_{n=1}^N \sqrt{\frac{l}{n}} \leq 2\sqrt{N} \right)$$

$$\leq 4\sqrt{l} \sqrt{AT}$$

$$\left( \frac{1}{A} \sum_{a=1}^A \sqrt{N_T(a)} \right)$$

$$\leq \sqrt{\frac{1}{A} \sum_{a=1}^A N_T(a)}$$

$$= \sqrt{\frac{T}{A}}$$

by Jensen & concavity  
of  $x \mapsto \sqrt{x}$ )

Therefore,

$$\mathbb{E}\left[\left(\sum_{t=1}^T \text{reg}(t)\right) I(E)\right] \leq 4\sqrt{LAT} \quad (2)$$

Summing up (1) (2)

$$\Rightarrow \text{Reg}(T) \leq 4\sqrt{LAT} + 2 = \tilde{O}(\sqrt{AT}).$$

Remarks on designing optimism-based algs:

(1) Define bonus so UCB's are valid wh.p. (so  $\text{reg}_t \leq b_t(c_{\text{UB}})$ )

(2). subject to (1). design bonus to be as tight as

possible. (so  $\sum_t b_t(c_{\text{UB}})$  is small).

# Stochastic linear (contextual) bandits

## ① Setup :

Application: recommendation with personalization

MAB: recommend a single product so that it has maximum overall satisfaction over population

Contextual Bandits: for different users. recommend different products that fit their respective needs.

Protocol:

For  $t=1, 2, \dots, T$ :

observes context  $x_t \in \mathcal{X}$

takes action  $a_t \in A$

receive reward  $r_t = f^*(x_t, a_t) + \xi_t$  zero-mean, i.i.d.

unknown

Goal: maximize  $\mathbb{E} \left[ \sum_{t=1}^T r_t \right] = \mathbb{E} \left[ \sum_{t=1}^T f^*(x_t, a_t) \right]$

Assumption (Linear Realizability)

This lecture:-  $f^* \in F = \left\{ f(x, a) = \langle \theta, \phi(x, a) \rangle : \|\theta\|_2 \leq 1 \right\}$  with  $\phi$  known

- assumption:  $\|\phi(x, a)\|_2 \leq 1$  for all  $x, a$ .

- Notation:

$f^* = \langle \theta^*, \phi(x, a) \rangle$ ,  $\theta^*$  needs to be learned.

What's the best strategy had  $f^*$  been known?

$\forall x$ , take action  $a = \pi_{f^*}(x) := \arg \max_{a \in A} f^*(x, a)$ .

Regret notion:

$P\text{Reg}(T)$ : pseudo regret  
optimal policy

$$\text{Reg}(T) = \mathbb{E} \left[ \sum_{t=1}^T \max_a f^*(x_t, a) - \sum_{t=1}^T f^*(x_t, a_t) \right].$$

reward by following  
optimal policy

reward of the  
agent

Practical relevance: "A Contextual Bandit approach  
to personalized News Article Recommendation".  
(Li, Chu, Langford, Schapire, 2010, WWW)  
Test of Time Award.

MAB as linear bandits:

Every MAB problem can be viewed as a  $K$ -dimensional  
Linear bandit problem by defining:

$$- x_t = z_0 \quad \forall t \quad (\text{dummy context})$$

$$- \phi(z_0, a) = e_a = \begin{pmatrix} 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{pmatrix} \leftarrow a\text{-th location} \in \mathbb{R}^K$$

$$- \Theta^* = (f^*(1) \dots f^*(K)) \text{, thus}$$

$$f^*(a) = \langle \Theta^*, \phi(z_0, a) \rangle . + a .$$

② Decision algorithms for linear bandits:

Designing algorithms

Idea: optimism principle, again.



Act according to the best plausible world.

- \* "world model" = reward model =  $\Theta$ .
- \* "plausible world" = set of  $\Theta$ 's that can plausibly be  $\Theta^*$ . i.e.  $\underset{\text{At } t}{\text{can explain observations}}$   
 $(x_s, a_s, r_s)_{s=1}^{t-1}$

Similar to constructing confidence intervals in MAB.

here we construct confidence set for  $\Theta^*$ , say.

$$P(\forall t, \Theta^* \in \mathbb{H}_t) \geq 1 - \frac{1}{T}.$$

Algorithm ( OFUL: optimism in the face of uncertainty for linear bandits )

For  $t = 1, 2, \dots, T$ :

- construct confidence set  $\mathbb{H}_t$  for  $\Theta^*$ .
  - observe  $x_t$ .
  - Take action  $a_t = \operatorname{argmax}_{a \in A} \underbrace{\langle \Theta, \phi(x_t, a) \rangle}_{\Theta \in \mathbb{H}_t}$
- The largest plausible reward action  $a$  can get

Q1: How to construct confidence set  $\mathbb{H}_t$ ?

Q2: How to analyze the OFUL alg?

Q1: Recall: in MAB, confidence interval for  $f^*(a)$

$$= \left[ \underbrace{f_t^*(a)}_{\text{center}} \pm \underbrace{b_t(a)}_{\text{width deduced by concentration ineq.}} \right]$$

center: some best guess of  $f^*(a)$   
(e.g. Hoeffding)

What's our "best guess" of  $\theta^*$  based on  $(x_s, a_s, r_s)_{s=1}^{t+1}$ ?

Recall:  $r_s = \langle \theta^*, \phi(x_s, a_s) \rangle + \varepsilon_s$

Estimating  $\theta^*$  from data is a linear regression problem

Standard solution: least squares + regularization

$$\hat{\theta}^t(\lambda) = \underset{\theta}{\operatorname{argmin}} \sum_{s=1}^{t+1} \underbrace{(\langle \theta, \phi(x_s, a_s) \rangle - r_s)^2}_{\phi_s} + \lambda \|\theta\|_2^2$$

How close is  $\hat{\theta}^t(\lambda)$  to  $\theta^*$ ?

Note: in general, cannot claim, say  $\|\hat{\theta}^t(\lambda) - \theta^*\|_2 \leq \frac{1}{\sqrt{t}}$

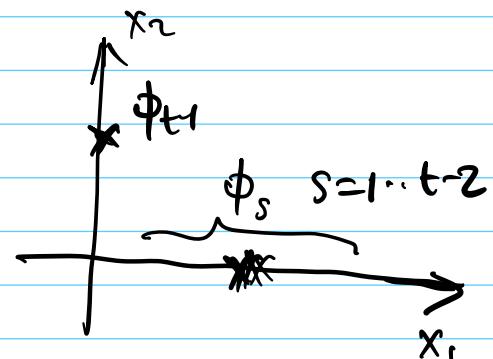
Problem: data may be highly "degenerate."

Example  
 $d=2$ :

$$\phi_1 = \phi_2 = \dots = \phi_{t-2} = e_1 \quad \phi_{t-1} = e_2$$

$$r_1 = \theta_1^* + \varepsilon_1$$

$$r_2 = \theta_1^* + \varepsilon_2$$



$$r_{t-2} = \theta_1^* + \varepsilon_{t-2}$$

$$r_{t-1} = \theta_2^* + \varepsilon_t$$

(only one data pt is related).

The data reveals little information about  $\theta_2^*$ , so cannot hope to show  $|\hat{\theta}_2^t - \theta_2^*| = O(\frac{1}{\sqrt{t}})$ .

In general, we will show the following estimation guarantee:

Lemma:  $\exists \text{ event } E, P(E) \geq 1 - \frac{1}{t}, \forall t$  for all  $t$ ,

$$\theta^* \in \mathbb{H}_t(\lambda) = \dots \quad \|\hat{\theta}^t(\lambda) - \theta^*\|_{V_{t-1}(\lambda)} \leq \beta_t(\lambda) = \tilde{O}(\sqrt{\lambda} + d)$$

Specifically,  $\lambda = 1 \Rightarrow$

$$\theta^* \in \mathbb{H}_t(1) = \dots \quad \|\hat{\theta}^t(1) - \theta^*\|_{V_{t-1}(1)} \leq \beta_t(1) = \tilde{O}(d)$$

This, although not a pointwise guarantee, turns out to be still useful.

Here:  $\|x\|_M := \sqrt{x^T M x}$  For  $M \succ 0$   
 $= \|M^{\frac{1}{2}} x\|_2$  psd  $M = U(\lambda_1, \dots, \lambda_d) U^T$  This is a norm:  
 $- V_{t-1}(\lambda) = \sum_{s=1}^{t-1} \phi_s \phi_s^T + \lambda I$   $\|ax\|_M = |a| \|x\|_M$   
 $\|x+y\|_M \leq \|x\|_M + \|y\|_M$   
 $\langle x, y \rangle \leq \|x\|_M \|y\|_M$   
 generalized Cauchy-Schwarz.

Intuition check on previous example w/  $\lambda = 0$ :

$$* \quad \hat{\theta}^t = \arg \min_{\theta} \sum_{s=1}^{t-2} (\theta_1 - (\theta_1^* + \varepsilon_s))^2 + (\theta_2 - (\theta_2^* + \varepsilon_{t-1}))^2$$

$$\text{let } z = \hat{\theta}^t - \theta^* \quad \Rightarrow |z_1| \leq \sqrt{\frac{1}{t}}, |z_2| \leq O(1).$$

$$* \quad \text{Note: } V_{t-1} = (t-1) e_1 e_1^T + 1 \cdot e_1 e_2^T = \begin{pmatrix} t-1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\Rightarrow \|z\|_V^2 = (z_1, z_2) \begin{pmatrix} t-1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$$

$$= (t-1) \mathbf{z}_1^2 + \mathbf{z}_2^2$$

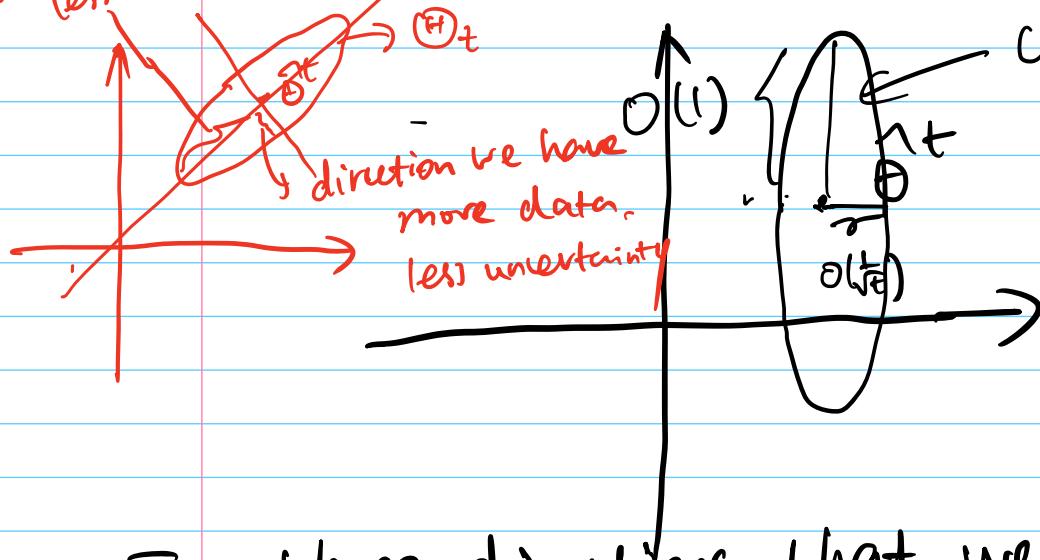
$$\leq O(t \cdot \frac{1}{\epsilon} + 1 \cdot \frac{1}{\epsilon}) = O(1)$$

\* Why is  $\|\mathbf{z}\|_{V_{t-1}}^2 \leq O(1)$  guarantee useful?

$$\Rightarrow \|\theta^* - \hat{\theta}^t\|_{V_{t-1}}^2 \leq O(1) \text{ gives a}$$

*direction we have more data, more uncertainty*

Constraint on  $\theta^*$ :



For those directions that we have more data, our estimate on  $\theta^*$  is more accurate.

Proof of the lemma: ( $V_{t-1}(\lambda)$  abbrev.  $V_{t-1}$  below)

$$\text{First, } \hat{\theta}_t(\lambda) = V_{t-1}^{-1} \left( \sum_{s=1}^{t-1} \phi_s r_s \right)$$

$$= V_{t-1}^{-1} \left( \underbrace{\sum_{s=1}^{t-1} \phi_s \phi_s^T \theta^*}_{V_{t-1}} + \underbrace{\sum_{s=1}^{t-1} \phi_s \varepsilon_s}_{m_{t-1}!!} \right)$$

$$= \theta^* - V_{t-1}^{-1} \lambda \cdot \theta^* + V_{t-1}^{-1} m_{t-1}$$

$$\Rightarrow \hat{\theta}_t(\lambda) - \theta^* = - V_{t-1}^{-1} \lambda \theta^* + V_{t-1}^{-1} m_{t-1}$$

$$\Rightarrow \|\hat{\theta}_t(\lambda) - \theta^*\|_{V_{t-1}} \leq \|V_{t-1}^{-1} \lambda \theta^*\|_{V_{t-1}} + \|V_{t-1}^{-1} m_{t-1}\|_{V_{t-1}}$$

defn of Mahalanobis norm

$$= \|\lambda \theta^*\|_{V_{t-1}^{-1}} + \|m_{t-1}\|_{V_{t-1}^{-1}}$$

$$\text{w.p. } 1 - \frac{1}{T^2} \quad \text{exercise} \quad \leq \sqrt{\lambda} + \sqrt{2 \ln(T^2) + d \ln\left(1 + \frac{t}{d\lambda}\right)}$$

union bound over all  $t \Rightarrow$  Lemma.

In the last step, to deal with  $\|(m_{t-1})\|_{V_{t-1}^{-1}}$ , we use

self-normalized Tail Inequality for vector-valued martingales:

Lemma (Abbasi-Yadkori, Dál, Szepesvári, 2011)

At.

$$P\left(\left\|\sum_{s=1}^t \phi_s \varepsilon_s\right\|_{V_t^{-1}(\lambda)} \geq \sqrt{2 \ln \frac{1}{\delta} + \ln \frac{\det(V_t(\lambda))}{\det(V_0(\lambda))}}\right) \geq 1 - \delta$$

$$\leq d \ln\left(1 + \frac{t}{d\lambda}\right).$$

Intuition: when all  $\varepsilon_s \sim N(0, 1)$  and  $\phi_s$ 's are chosen ahead of time

$$m_t = \sum_{s=1}^t \phi_s \varepsilon_s \stackrel{d}{\approx} N(0, V_t(\theta))$$

$$V_t^{-\frac{1}{2}}(\lambda) \cdot m_t \stackrel{d}{\approx} N(0, I_d)$$

$$\| V_t^{-\frac{1}{2}}(\lambda) m_t \|_2 \stackrel{\text{whp.}}{\leq} O(d + \ln \frac{1}{\delta}) .$$

problem with this intuition:  $\phi_s$ 's are not fixed:

$\phi_s$  depends on  $\phi, \varepsilon_1, \dots, \phi_s, \varepsilon_{s-1}$ .

Q1 summary:

At. define  $\Theta_t = \{ \theta : \| \theta - \hat{\theta}^t(\lambda) \|_{V_{t-1}(\lambda)} \leq \beta_t(\lambda) \}$

$$P(\text{At. } \theta^* \in \Theta_t) \geq 1 - \frac{1}{T} .$$

Algorithm induced by this  $\Theta_t$  construction:

$$a_t = \underset{a \in A}{\operatorname{argmax}} \left[ \max_{\theta \in \Theta_t} \langle \theta, \phi(x_t, a) \rangle \right] \rightarrow \begin{aligned} & \text{Interpretation:} \\ & = a \cdot UCB_t(x_t, a) \end{aligned}$$

$\max_{\theta \in \Theta_t} \langle \theta, \phi(x_t, a) \rangle$

$\forall x : \| x \|_{V_{t-1}(\lambda)} \leq \beta_t(\lambda) \Rightarrow$

$\max_{\theta \in \Theta_t} \langle \hat{\theta}^t(\lambda) + x, \phi(x_t, a) \rangle \geq \langle \theta^*, \phi(x_t, a) \rangle = f^*(x_t, a)$

|| Exercise

$$\langle \hat{\theta}^t(\lambda), \phi(x_t, a) \rangle + \beta_t(\lambda) \| \phi(x_t, a) \|_{V_{t-1}(\lambda)^{-1}}$$

$$\phi(z_a) = e_a$$

$$V_{t-1}(\lambda) = \sum_{s=1}^{t-1} \phi(z_{s,a}) \phi(z_{s,a})^\top + I$$

Estimated reward

For MAB:  $\sqrt{K} \sqrt{\frac{1}{N_{t-1}(a) + 1}}$

uncertainty / exploration bonus

$N_{t+1}(A) + 1$

## Q2 Regret analysis

Ihm: on  $E$  ( $R(E) \geq -\frac{1}{T}$ ),  $P\text{Reg}(T) \leq \tilde{\mathcal{O}}(d\sqrt{T})$ .

( $\Rightarrow \text{Reg}(T) \leq \tilde{\mathcal{O}}(d\sqrt{T})$ .)

(For MAB, what regret bound do we get?)

Pf: Recall  $P\text{Reg}(T) = \sum_{t=1}^T \text{reg}(t)$ , where

$$\text{reg}(t) = \max_a f^*(x_t, a) - f^*(x_t, a_t)$$

How to bound? Similar to MAB analysis

Step 1:  $\text{reg}(t) \leq 2 b_t(a_t)$

why:  $\text{reg}(t) = \max_a f^*(x_t, a) - f^*(x_t, a_t)$

$$\leq \max_a UCB_t(a) - f^*(x_t, a_t)$$

$$= UCB_t(a_t) - f^*(x_t, a_t)$$

$$= \langle \hat{\theta}_t^{(1)} - \theta^*, \phi(x_t, a_t) \rangle + b_t(a_t)$$

$$\leq b_t(a_t)$$

$$\leq \underbrace{\|\hat{\theta}_t^{(1)} - \theta^*\|_{V_{t-1}^{(1)}}}_{\leq \beta_t^{(1)}} \| \phi(x_t, a_t) \|_{V_{t-1}^{(1)}} + b_t(a_t)$$

$$= 2 b_t(a_t)$$

Step 2: bounding  $\sum_t b_t(a_t)$

$$\begin{aligned} \text{how: } \sum_t b_t(a_t) &= \sum_t \beta_{t(1)} \|\phi(x_t, a_t)\|_{V_{t(1)}^{-1}}^2 \\ &\leq \underbrace{\left(\max_t \beta_{t(1)}\right)}_{\tilde{\sigma}(\sqrt{d})} \underbrace{\left(\sum_{t=1}^T \|\phi_t\|_{V_{t(1)}^{-1}}\right)}_{(*)}. \end{aligned}$$

$$(*) = \sum_{t=1}^T \|\phi_t\|_{V_{t(1)}^{-1}}$$

$$\leq \sqrt{\sum_{t=1}^T \|\phi_t\|_{V_{t(1)}^{-1}}^2} \cdot \sqrt{T}$$

$$= \sqrt{\sum_{t=1}^T \phi_t^T V_{t(1)}^{-1} \phi_t} \cdot \sqrt{T}$$

$$\text{Note: } V_{t(1)} \leq 2V_{t(1)} \Rightarrow V_{t(1)}^{-1} \leq 2V_{t(1)}^{-1}$$

*special case:  $d=1$ .*  $A \leq B \Leftrightarrow A - B \geq 0$  . positive semi definite order of matrices

property:  $A \leq B \Rightarrow \forall x. x^T A x \leq x^T B x$

$$A \leq B \Rightarrow B^{-1} \leq A^{-1}$$

$$\leq \sqrt{\sum_{t=1}^T \phi_t^T V_{t(1)}^{-1} \phi_t} \cdot \sqrt{2T}$$

The elliptic potential lemma: (EPL)

suppose  $u_1, \dots, u_T \in \mathbb{R}^d$ .  $A_t = \mu I + \sum_{s=1}^t u_s u_s^T$ , then

$$\sum_{t=1}^T \|u_t\|_{A_t^{-1}}^2 \leq \ln \frac{\det(A_T)}{\det(A_0)} \quad (A_0 = \mu I)$$

$$\text{if } \forall t \|u_t\| \leq D \\ \leq d \ln \left( 1 + \frac{D^2 T}{d\mu} \right) = \tilde{\mathcal{O}}(d)$$

$$\leq \tilde{\mathcal{O}}(\sqrt{dT})$$

$$\Rightarrow \sum_{t=1}^T b_t(a_t) \leq \tilde{\mathcal{O}}(\sqrt{d}) \cdot \tilde{\mathcal{O}}(\sqrt{dT}) = \tilde{\mathcal{O}}(d\sqrt{T})$$

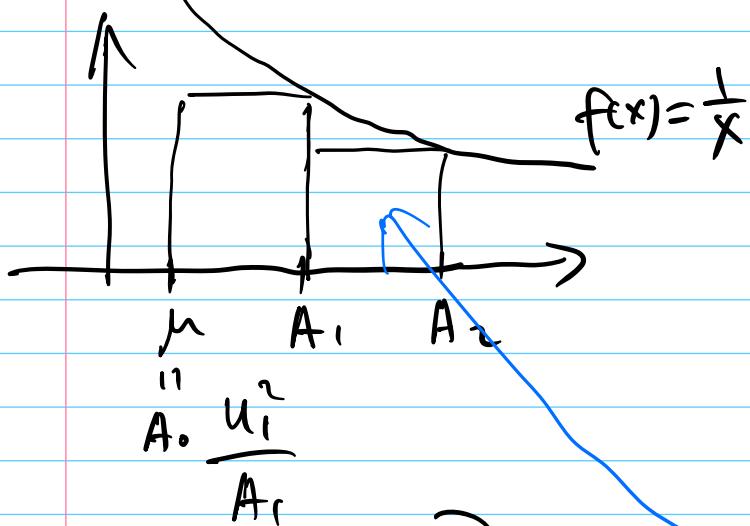
↗

Intuition of EPL

$\therefore$  when  $d=1$ .

$$A_t = \mu + \sum_{s=1}^t u_s^2$$

$$\sum_{t=1}^T \frac{u_t^2}{A_t} \leq \int_{\mu}^{A_T} \frac{1}{x} dx = \ln \frac{A_T}{\mu}$$



Pf idea of EPL:

$$\text{Note: } \ln \frac{\det(A_t)}{\det(A_{t-1})}$$

$$\textcircled{1} \quad \text{show that } \|u_t\|_{A_t^{-1}}^2 \leq \ln \frac{\det(A_t)}{\det(A_{t-1})}$$

$$(1-d: ) \quad \frac{u_t^2}{A_t} \leq \ln \frac{A_t}{A_{t-1}}$$

$$= \ln(\det(A_{t-1}^{-\frac{1}{2}} A_t A_t^T A_{t-1}^{-\frac{1}{2}}))$$

$$= \ln(\det(I + A_{t-1}^{-\frac{1}{2}} U_t U_t^T A_t^{-\frac{1}{2}}))$$

Exercise ( Hint:  $\det(M) = \prod_{i=1}^d \lambda_i(M)$ )

$$\stackrel{?}{=} \ln(1 + \|U_t\|_{A_{t-1}^{-\frac{1}{2}}}^2)$$

$$\geq \frac{\|U_t\|_{A_{t-1}^{-\frac{1}{2}}}^2}{1 + \|U_t\|_{A_{t-1}^{-\frac{1}{2}}}^2} \quad (x \geq 0, \ln(1+x) \geq \frac{x}{1+x})$$

$$= \|U_t\|_{A_t^{-\frac{1}{2}}}^2 \quad (A_t^{-\frac{1}{2}} = (A_t + U_t U_t^T)^{-\frac{1}{2}}$$

$$\text{Sherman-Morrison} \\ = A_{t-1}^{-\frac{1}{2}} - \frac{A_{t-1}^{-\frac{1}{2}} U_t U_t^T A_{t-1}^{-\frac{1}{2}}}{1 + U_t^T A_{t-1}^{-\frac{1}{2}} U_t}$$

⑦ Show  $\ln \frac{\det(A_T)}{\det(A_0)} = \tilde{\sigma}(d)$

Fact: for p.s.d matrices  $A = U \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{pmatrix} U^T$ ,  $\det(A) = \prod_{i=1}^d \lambda_i$

$$\text{For } A_0 = \begin{pmatrix} \mu & & \\ & \ddots & \\ & & \mu \end{pmatrix} \Rightarrow \det(A_0) = \mu^d.$$

$$\text{For } A_T = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{pmatrix} \Rightarrow \det(A_T) = \prod_{i=1}^d \lambda_i \leq \left(\frac{1}{d} \sum_{i=1}^d \lambda_i\right)^d$$

$$= \left(\frac{1}{d} \operatorname{tr}(A_T)\right)^d$$

$$\text{Ex} \leq \left(\frac{1}{d} (d\mu + T D^2)\right)^d \quad \times$$