

Hello!

CS292F StatRL Lecture7

Exploration in Bandits

Instructor: Yu-Xiang Wang

Spring 2021

UC Santa Barbara

Notes / reminders

- Project proposal due today
 - Please submit on Gradescope.
- Start HW1 quickly.
 - It will be more time-consuming than HW0.
 - It will help you with the rest of the class.
- HW2 is to be released this week (hopefully by tomorrow)

Recap: Lecture 6

- Policy gradient methods
 - Policy gradient theorem
 - Unbiased Monte Carlo estimate of the gradient (REINFORCE)
 - Bootstrapping in policy gradient estimates
 - Function approximation and Actor-Critic
- Bandits problem setup
 - Regret definition
 - The need for exploration

k-armed bandits

Recap: Multi-arm bandits: Problem setup

- No state. k-actions $a \in \mathcal{A} = \{1, 2, \dots, k\}$
- You decide which arm to pull in every iteration

$$A_1, A_2, \dots, A_T$$

- You collect a cumulative payoff of $\sum_{t=1}^T R_t$

- For MAB, the regret is defined as follow

$$\underbrace{T \max_{a \in [k]} \mathbb{E}[R_t | a]}_{T \cdot \mu_{a^*}} - \sum_{t=1}^T \mathbb{E}_{a \sim \pi} [\mathbb{E}[R_t | a]]$$

$\sum_{a \in [k]} \mu_{A_t}$

$R_t \sim P(\cdot | A_t = a)$
 $\mathbb{E}[R_t | A_t = a] = \mu_a$

$R_t = \mu_a + \text{Noise}$
 $\mathbb{E}[\text{Noise}] = 0$

“No regret” means sublinear scaling in T . “Linear regret” is very bad.

- “No regret online learning”

$$\text{Regret}_T = o(T)$$

$$\lim_{T \rightarrow \infty} \frac{1}{T} \text{Regret}_T = 0$$

- A regret (upper) bound needs to apply to all problem instances

$$\text{Var}(y_t) \leq \sigma^2$$
$$0 \leq \mu_a \leq 1 \quad \forall a$$
$$|A| = k$$

- It suffices to identify one example to get a regret lower bound for a given algorithm.

- E.g., “Greedy strategy” has linear regret in MAB.

- Minimax lower bounds are information-theoretical

- They apply to all algorithms.

Recap: “Exploration first” strategy

- Let’s spend the first N step exploring.
 - Play each action for N / k times.

$$\hat{Q}_t(a) = \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}$$

- For $t = N + 1, N+2, \dots, T$:

$$A_t \doteq \arg \max_a Q_t(a),$$

This lecture

- Regret analysis for multi-armed bandits
 - Exploration first
 - epsilon-greedy
 - Upper Confidence Bound algorithm (AJKS 5.1)
- Linear bandits. (AJKS 5.2 – 5.3)
 - LinUCB algorithm
 - Regret analysis

Recap: Concentration inequalities --- finite-sample bounds of LLN and CLT

- **Hoeffding's inequality:** Assume X_1, \dots, X_n are independent and their support bounded:

$$S_n = X_1 + \dots + X_n$$

$$P(S_n - \mathbb{E}[S_n] \geq t) \leq \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right), = \frac{\delta}{2}$$

$$P(a_i \leq X_i \leq b_i) = 1$$

- Easy version, if $0 < X_i < B$, **with probability 1- δ** :

$$|\bar{X} - \mathbb{E}[\bar{X}]| \leq \sqrt{\frac{B^2}{2n} \log(2/\delta)}$$

Regret analysis of Exploration First

$$\frac{N}{K} \text{ times, w.p. } \geq 1 - \frac{\delta}{K} \quad |\hat{Q}(a) - Q(a)| \leq \sqrt{\frac{K}{2N} \log \frac{2K}{\delta}} \quad Q(a) =: \mu_a$$

for all $a \in A$, union bound $\sup_{a \in A} |\hat{Q}(a) - Q(a)| \leq \sqrt{\frac{K}{2N} \log \frac{2K}{\delta}} = \epsilon$

Regret for the Exploration Phase: $\frac{N}{K} \sum_a \max_{a'} Q(a') - Q(a) \leq N$

$$|0 \leq Q(a) \leq 1|$$

Regret for the Exploitation Phase: $\left[\hat{a}^* = \arg \max_a \hat{Q}(a) \right]$

$$(T-N) \cdot (Q(a^*) - Q(\hat{a}^*))$$

$$= (T-N) \cdot \left[\underbrace{Q(a^*) - \hat{Q}(a^*)}_{\leq \epsilon} + \underbrace{\hat{Q}(a^*) - \hat{Q}(\hat{a}^*)}_{\leq 0} + \underbrace{\hat{Q}(\hat{a}^*) - Q(\hat{a}^*)}_{\leq \epsilon} \right]$$

$$\leq (T-N) \cdot 2\epsilon \leq 2T \sqrt{\frac{K}{2N} \log \frac{2K}{\delta}}$$

Regret $= N + 2T \sqrt{\frac{K}{2N} \log \frac{2K}{\delta}} = O\left(T^{\frac{2}{3}} K^{\frac{1}{3}} \left(\log \frac{2K}{\delta}\right)^{\frac{1}{3}}\right)$

Choose $N = T^{\frac{2}{3}} K^{\frac{1}{3}} \left(\log \frac{2K}{\delta}\right)^{\frac{1}{3}}$

Regret analysis of Exploration First

ϵ -Greedy strategy: one way to balance exploration and exploitation

- You choose with probability $1 - \epsilon$

$$A_t \doteq \operatorname{argmax}_a Q_t(a),$$

- With probability ϵ , choose an action **uniformly at random!**
 - Including the argmax.
- Carefully choose ϵ parameter.

A sketch of the analysis for ϵ -greedy

- In expectation, each arm is chosen for at least ϵt times.
At time t, by Hoeffding's $N_t(a) \geq \frac{\epsilon t}{k} - \sqrt{\frac{k}{\epsilon t}} \geq \frac{\epsilon t}{2k}$

~~By~~

- Condition on the number of times, apply Hoeffding's inequality / union bound for all t and a

By Hoeffding's again $\left[\sup_a \left| \hat{Q}_t(a) - Q(a) \right| \leq \sqrt{\frac{k}{\epsilon t}} \right]$ *u.b.p*

estimator

- Regret bound is

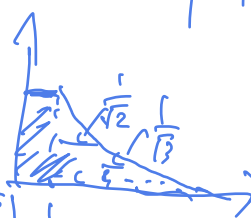
$$\epsilon T + \sum_{t=1}^T C \sqrt{\frac{k}{\epsilon t}}$$

Exploration

Exploitation

Minimise the regret by choosing $\epsilon_t \Rightarrow$

$\frac{1}{(3k)^{1/3}}$



$1 + \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{3}} + \dots + \frac{1}{\sqrt{T}} = ?$

$1 + \int_1^T \frac{1}{\sqrt{x}} dx = 1 + [2\sqrt{x}]_1^T = 1 + 2\sqrt{T} - 2 = 2\sqrt{T} - 1$

Optimism-in-the-face of uncertainty: Upper Confidence Bound algorithm

Alg (UCB)

① Play each action $a \in A$ once. k steps

② for $t = k+1, \dots, T$

• $A_t = \operatorname{argmax}_a \hat{Q}_t(a) + \sqrt{\frac{\log(\frac{t}{\delta})}{2N_t(a)}}$

$$N_t(a) = \sum_{i=1}^{t-1} \mathbb{1}(A_i = a)$$

③ • ~~where~~ where $\hat{Q}_t(a) = \frac{1}{N_t(a)} \left(R_a + \sum_{i=k+1}^{t-1} \mathbb{1}(A_i = a) R_i \right)$



Martingale

$$\underline{X_{ot} - X_{t-1}} = \mathbb{1}_{\{\text{toss } t \text{ is head}\}}$$

- We say that a sequence of r.v. X_1, \dots, X_n, \dots is a Martingale if for any n

$$\underline{\mathbf{E}(|X_n|)} < \underline{\infty}$$

$$\mathbf{E}(X_{n+1} \mid \underline{X_1, \dots, X_n}) = X_n.$$

- Example:
 - Random-walk: Total number of heads minus tails in n coin tosses

Azuma-Hoeffding's inequality

- **Azuma-Hoeffding's inequality:** Assume X_1, \dots, X_n are Martingale differences

$$S_n = X_1 + \dots + X_n$$

S_n is a martingale

$$\mathbb{P}[S_n \geq \epsilon] \leq e^{-\frac{2\epsilon^2}{\sum_{i=1}^n \underbrace{(b_i - a_i)^2}_{N_t(a)}}$$

$E(S_n | S_1, \dots, S_{n-1}) = S_{n-1}$

- Apply Azuma-Hoeffding's inequality to our problem

$$Q_t(a) = \frac{1}{N_t(a)} \left(R_t + \sum_{i=1}^{t-1} \mathbb{1}(A_i=a) \cdot R_i \right)$$

$$S_t(a) = R_t + \sum_{i=1}^{t-1} \mathbb{1}(A_i=a) \cdot R_i$$

~~$S_t(a)$~~

$$R_i = \mu_a + \sum_{i=1}^{t-1} \mathbb{1}(A_i=a) R_i - E[\mathbb{1}(A_i=a) R_i | \text{Hist}_{i-1}]$$

we know from UCB, $\mathbb{1}(A_i=a)$ is fixed
 $X_i =: \mathbb{1}(A_i=a) \cdot R_i$ - conditional on the X_1, \dots, X_{i-1}

$b_i = a_i = 0$, $b_i = 1$, $a_i = 0$ for those where $A_i = a$

$E[X_n | S_1, \dots, S_{n-1}] = 0$

$$R_i \sim P(R_i | A_i = a_i)$$

$$\underline{R_1, R_2, \dots, R_T}$$

Regret analysis of UCB

for $A_t = a$, $b_i = 0, a_i = 1$
 for $A_t \neq a$, $b_i = 0, a_i = 0$

$$R_t - Q(a) + \sum_{i=k_t}^{t-1} \mathbb{1}(A_i = a) (R_i - Q(a)) \text{ is martingale}$$

$$| \cdot | \leq \sqrt{2N_t(a) \log \frac{KT}{\delta}} \text{ u.p. } 1 - \frac{\delta}{KT}$$

union bound over all $a \in A$, all $k_t \leq t \leq T$ u.p. $1 - \delta$

$$\bar{Q}_t(a) = Q(a) + \sqrt{\frac{2 \log \frac{KT}{\delta}}{N_t(a)}}$$

UCB

$$\sup_{t,a} \frac{1}{N_t(a)} \left| R_t - Q(a) + \sum_{i=k_t}^{t-1} \mathbb{1}(A_i = a) (R_i - Q(a)) \right| \leq \sqrt{2 \log \frac{KT}{\delta}}$$

$$Q(a^*) - Q(A_t) = \underbrace{Q(a^*) - \bar{Q}_t(a^*)}_{\leq 0} + \bar{Q}_t(a^*)$$

≤ 0
UCB rule

$$N_t(a) \cdot (\bar{Q}_t(a) - Q(a))$$

$$|\bar{Q}_t(a) - Q(a)| \leq \sqrt{\frac{2 \log \frac{KT}{\delta}}{N_t(a)}} \text{ for all } a, t$$

simultaneously u.p. $1 - \delta$

$$\bar{Q}(A_t) + \bar{Q}(A_t) - Q(A_t) \leq 2 \cdot \epsilon$$

gaps

Regret analysis of UCB

$$\Delta_a = Q(a^*) - Q(a)$$

from last slide

$$\text{Regret} = \sum_{a=1}^k \Delta_a + \sum_{t=K+1}^T \underbrace{Q(a^*) - Q(A_t)}$$

$$\leq K + \sum_{t=K+1}^T 2\sqrt{\frac{2 \log \frac{2T}{\delta}}{N_t(A_t)}}$$

$$= K + 2\sqrt{\frac{2 \log \frac{2T}{\delta}}{\delta}} \sum_{a=1}^k \sum_{i=1}^{N_t(a)} \frac{1}{j_i}$$

$$\leq K + 4\sqrt{\frac{2 \log \frac{2T}{\delta}}{\delta}} \sum_{a=1}^k \sqrt{N_t(a)}$$

$$\left(\sum_{i=1}^{N_t(a)} \frac{1}{j_i} \leq 2\sqrt{N_t(a)} \right)$$

$$\leq K + 4\sqrt{\frac{2 \log \frac{2T}{\delta}}{\delta}} \sqrt{k \cdot \sum_{a=1}^k (N_t(a))^2}$$

$$= K + c \cdot \sqrt{kT \log \frac{2T}{\delta}}$$

Gap dependent analysis

$$\text{Claim: } N_t(a) \leq \frac{2\sqrt{2} \log \left(\frac{2T}{\delta} \right)}{\Delta_a^2}$$

Substitute into (*)

if $\bar{Q}_t(a^*) > \bar{Q}_t(a)$ for each a then $N_t(a)$ will no longer get larger.

$$\begin{aligned} \bar{Q}_t(a) &\leq Q(a) + \frac{\sqrt{2 \log \frac{2T}{\delta}}}{N_t(a)} \\ &= Q(a^*) - \Delta_a + \frac{\sqrt{2 \log \frac{2T}{\delta}}}{N_t(a)} \\ &\leq \bar{Q}_t(a^*) - \Delta_a + \frac{\sqrt{2 \log \frac{2T}{\delta}}}{N_t(a)} \end{aligned}$$

$$\Leftrightarrow \Delta_a \leq \frac{\sqrt{2 \log \frac{2T}{\delta}}}{N_t(a)} \Leftrightarrow N_t(a) \leq \frac{2 \log \frac{2T}{\delta}}{\Delta_a^2}$$

Summary of Exploration in Multi-Armed Bandits

- Explore-First

$$O\left(T^{\frac{2}{3}} k^{\frac{1}{3}}\right)$$

- eps-greedy

$$O\left(T^{\frac{2}{3}} k^{\frac{1}{3}}\right)$$

- UCB

$$\min_{\text{UCB}} \left\{ \sqrt{\sum_{\Delta_a < \tau} \log \frac{kT}{\delta}} + \sum_{\Delta_a > \tau} \frac{\log \frac{kT}{\delta}}{\Delta_a} \right\} = O\left(\sqrt{TK}\right)$$

(Note: The first term is labeled "Candy-Schwarz" and the second term is labeled "gap dependent bound")

$O\left(\sqrt{TK}\right)$
 t-hrde
 Constant
 log factors

Notes on MAB

UCB chooses $\hat{Q}(a) + \sqrt{\frac{1}{N_t(a)}}$ ← Exploration Bonus

$$P(0 \leq R_t \leq 1) = 1 \quad \forall t$$

- We considered “stochastic setting”
 - Adversarial setting (“a rigged casino”)
 - Reward sequence is arbitrary / no expectation in the regret.
- Exponential weight algorithm for Explore-Exploit. (Exp3) achieves the same regret.
 - Read Auer et al. (2001) The Nonstochastic Multiarmed Bandit Problem

Linear bandits: MAB with an infinite number of actions

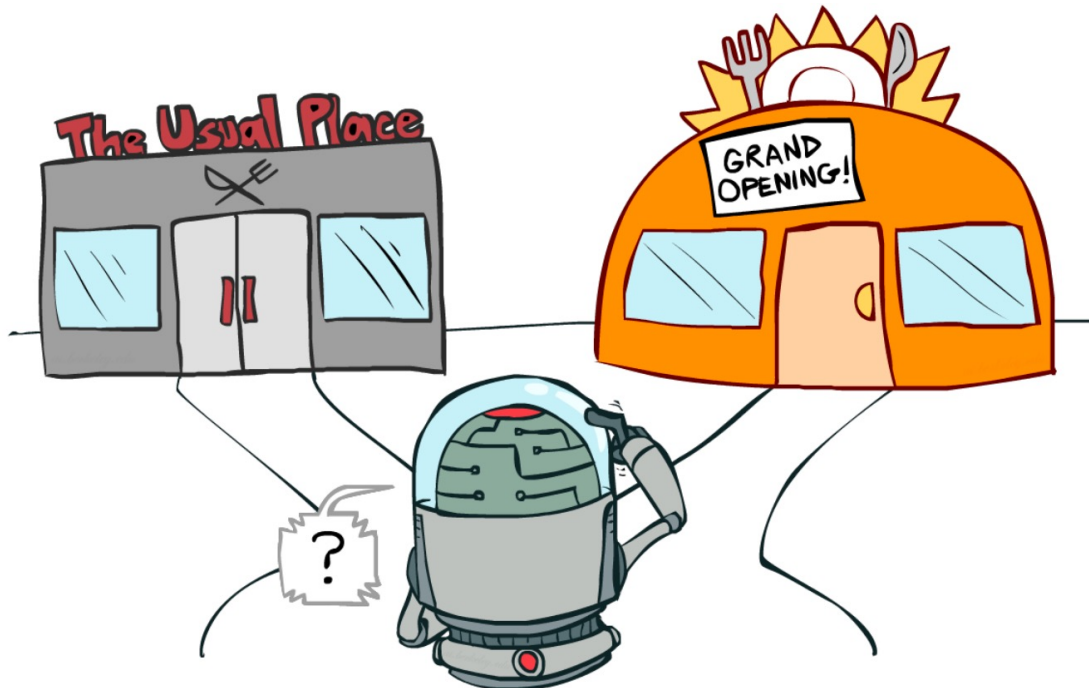
- Each action is determined by a “feature vector”

Features of action 1:

[Noodles, Tom Yum Soup, Poor service]

Features of action 2:

[Burger, Fries, Onion Ring, Fried Chicken]



Linear bandits: problem setup

- Action space is a compact set

$$A \subset \mathbb{R}^d, A \text{ is compact}$$

- Reward is linear + noise.

$$R_t = \langle A_t, \mu_k \rangle + \eta_t \quad \eta_t \text{ independent}$$

when $\mu_k \in \mathbb{R}^d$ σ^2 subgaussian

- Agent chooses a sequence of actions

$$A_1, \dots, A_T$$

- The regret is defined similarly

$$\text{Regret}_T = T \cdot \langle A^*, \mu_k \rangle - \sum_{t=1}^T \langle A_t, \mu_k \rangle$$

Assumption: $|A| \leq B \Leftrightarrow \forall a \in A, \|a\|_2 \leq B$
 $\|\mu_k\|_2 \leq W$

New notation

$$A: \mathbb{C} \rightarrow \mathbb{D}$$

$$a \mapsto x$$

$$\max_{a \in \mathbb{D}} \langle a, \mu_k \rangle - \sum_{t=1}^T \langle x_t, \mu_k \rangle$$

The LinUCB algorithm: Optimism in the Face of Uncertainty.

- Consider the ridge regression at each time t .

$$\hat{\mu}_t = \underset{\mu \in \mathcal{N}}{\operatorname{argmin}} \sum_{i=1}^t (r_{i,t} - \mu^T x_{i,t})^2 + \lambda \|\mu\|^2 \quad \left| \quad \hat{\mu}_t = \sum_{i=1}^t \sum_{k=1}^K x_{i,t} x_{i,t}^T x_{i,t} \right.$$

- Construct high probability confidence set of the parameter vector

$$\mathcal{B}_{t,\beta} = \left\{ \mu \mid (\mu - \hat{\mu}_t)^T \Sigma_t (\mu - \hat{\mu}_t) \leq \beta_t \right\}$$

where $\Sigma_t = \sum_{i=1}^t x_i x_i^T + \lambda I_d$

← ellipsoid



- Choose actions that maximize the UCB.

$$X_t = \underset{x \in \mathcal{D}}{\operatorname{argmax}} \max_{\mu \in \mathcal{B}_{t,\beta}} \langle x, \mu \rangle \quad \text{UCB-rule}$$

Regret bound of LinUCB

Sublinear regret: $R_T \leq O^*(d\sqrt{T})$

poly dependence on d , no dependence on the cardinality $|D|$.

Theorem 5.3 (AJKS)

Suppose: bounded noise $|\eta_t| \leq \sigma$, that $\|\mu^*\| \leq W$, and that $\|x\| \leq B$ for all $x \in D$. Set $\lambda = \sigma^2 / W^2$ and

$$\beta_t := \sigma^2 \left(2 + 4d \log \left(1 + \frac{TB^2 W^2}{d} \right) + 8 \log(4/\delta) \right).$$

With probability greater than $1 - \delta$, that for all $t \geq 0$,

$$R_T \leq c\sigma\sqrt{T} \left(d \log \left(1 + \frac{TB^2 W^2}{d\sigma^2} \right) + \log(4/\delta) \right)$$

where c is an absolute constant.

(Dani, Hayes & Kakde, 2009)

(From this slide onwards mostly taken from Sham Kakade)

Two components of the regret analysis

- Uniform (**over all t**) confidence bound

Proposition 5.5 (AJKS)

(Confidence) Let $\delta > 0$. We have that

$$\Pr(\forall t, \mu^* \in \text{BALL}_t) \geq 1 - \delta.$$

- Sum of Squares Regret bound

Proposition 5.6 (AJKS)

(Sum of Squares Regret Bound) Define:

$$\text{regret}_t = \mu^* \cdot x^* - \mu^* \cdot x_t$$

Suppose $\|x\| \leq B$ for $x \in D$. Suppose β_t is increasing and larger than 1. Suppose $\mu^ \in \text{BALL}_t$ for all t , then*

$$\sum_{t=0}^{T-1} \text{regret}_t^2 \leq 4\beta_T d \log \left(1 + \frac{TB^2}{d\lambda} \right)$$

Proof of the main regret bound

- By Cauchy-Schwarz

$$\sum_{t=0}^{T-1} \text{regret}_t \leq \sqrt{T \sum_{t=0}^{T-1} \text{regret}_t^2} \leq \sqrt{4T\beta_T d \log \left(1 + \frac{TB^2}{d\lambda} \right)}.$$

Plan of the proof

1. First prove the Proposition that bounds the sum of square regret
 - By bounding instantaneous regret
 - And then bounding the sum of squares with “Information Gain”
2. Prove the uniform confidence bound
 - Basically show that the choice of β_t “works”.

“Width” of Confidence Ball

Lemma

Let $x \in D$. If $\mu \in \text{BALL}_t$ and $x \in D$. Then

$$|(\mu - \hat{\mu}_t)^\top x| \leq \sqrt{\beta_t x^\top \Sigma_t^{-1} x}$$

Proof: By Cauchy-Schwarz, we have:

$$\begin{aligned} |(\mu - \hat{\mu}_t)^\top x| &= |(\mu - \hat{\mu}_t)^\top \Sigma_t^{1/2} \Sigma_t^{-1/2} x| = |(\Sigma_t^{1/2} (\mu - \hat{\mu}_t))^\top \Sigma_t^{-1/2} x| \\ &\leq \|\Sigma_t^{1/2} (\mu - \hat{\mu}_t)\| \|\Sigma_t^{-1/2} x\| = \|\Sigma_t^{1/2} (\mu - \hat{\mu}_t)\| \sqrt{x^\top \Sigma_t^{-1} x} \leq \sqrt{\beta_t x^\top \Sigma_t^{-1} x} \end{aligned}$$

where the last inequality holds since $\mu \in \text{BALL}_t$. ■

Instantaneous Regret is bounded by the width of the ellipsoid.

Define

$$w_t := \sqrt{x_t^\top \Sigma_t^{-1} x_t}$$

which is the “normalized width” at time t in the direction of our decision.

Lemma

Fix $t \leq T$. If $\mu^* \in \text{BALL}_t$, then

$$\text{regret}_t \leq 2 \min(\sqrt{\beta_t} w_t, 1) \leq 2\sqrt{\beta_T} \min(w_t, 1)$$

Proof: Let $\tilde{\mu} \in \text{BALL}_t$ denote the vector which minimizes the dot product $\tilde{\mu}^\top x_t$. By choice of x_t , we have

$$\tilde{\mu}^\top x_t = \max_{\mu \in \text{BALL}_t} \max_{x \in D} \mu^\top x \geq (\mu^*)^\top x^*,$$

where the inequality used the hypothesis $\mu^* \in \text{BALL}_t$. Hence,

$$\begin{aligned} \text{regret}_t &= (\mu^*)^\top x^* - (\mu^*)^\top x_t \leq (\tilde{\mu} - \mu^*)^\top x_t \\ &= (\tilde{\mu} - \hat{\mu}_t)^\top x_t + (\hat{\mu}_t - \mu^*)^\top x_t \leq 2\sqrt{\beta_t} w_t \end{aligned}$$

“Geometric potential” argument: Converting summation to product

Lemma 5.9 (AJKS)

We have:

$$\det \Sigma_T = \det \Sigma_0 \prod_{t=0}^{T-1} (1 + w_t^2).$$

Proof: By the definition of Σ_{t+1} , we have

$$\begin{aligned} \det \Sigma_{t+1} &= \det(\Sigma_t + x_t x_t^\top) = \det(\Sigma_t^{1/2} (I + \Sigma_t^{-1/2} x_t x_t^\top \Sigma_t^{-1/2}) \Sigma_t^{1/2}) \\ &= \det(\Sigma_t) \det(I + \Sigma_t^{-1/2} x_t (\Sigma_t^{-1/2} x_t)^\top) = \det(\Sigma_t) \det(I + v_t v_t^\top), \end{aligned}$$

where $v_t := \Sigma_t^{-1/2} x_t$. Now observe that $v_t^\top v_t = w_t^2$ and ... ■

Taking logarithm (get information gain), then bounding it with data-independent terms.

Lemma

For any sequence x_0, \dots, x_{T-1} such that, for $t < T$, $\|x_t\|_2 \leq B$, we have:

$$\log \left(\det \Sigma_{T-1} / \det \Sigma_0 \right) = \log \det \left(I + \frac{1}{\lambda} \sum_{t=0}^{T-1} x_t x_t^\top \right) \leq d \log \left(1 + \frac{TB^2}{d\lambda} \right).$$

Proof: Denote the eigenvalues of $\sum_{t=0}^{T-1} x_t x_t^\top$ as $\sigma_1, \dots, \sigma_d$, and note:

$$\sum_{i=1}^d \sigma_i = \text{Trace} \left(\sum_{t=0}^{T-1} x_t x_t^\top \right) = \sum_{t=0}^{T-1} \|x_t\|^2 \leq TB^2.$$

Using the AM-GM inequality,

$$\begin{aligned} \log \det \left(I + \frac{1}{\lambda} \sum_{t=0}^{T-1} x_t x_t^\top \right) &= \log \left(\prod_{i=1}^d (1 + \sigma_i / \lambda) \right) \\ &= d \log \left(\prod_{i=1}^d (1 + \sigma_i / \lambda) \right)^{1/d} \leq d \log \left(\frac{1}{d} \sum_{i=1}^d (1 + \sigma_i / \lambda) \right) \leq d \log \left(1 + \frac{TB^2}{d\lambda} \right) \end{aligned}$$

Bounding the Sum of Square Instantaneous Regret

$$\sum_{t=0}^{T-1} \text{regret}_t^2 \leq \sum_{t=0}^{T-1} 4\beta_t \min(w_t^2, 1) \leq 4\beta_T \sum_{t=0}^{T-1} \min(w_t^2, 1)$$

Plan of the proof

1. First prove the Proposition that bounds the sum of square regret
 - By bounding instantaneous regret
 - And then bounding the sum of squares with “Information Gain”
2. Prove the uniform confidence bound
 - Basically show that the choice of β_t “works”.

We need to prove that the true parameter is in the version space w.h.p.

- Recall the version space is:

Proof: Since $r_\tau = \mathbf{x}_\tau \cdot \mu^* + \eta_\tau$, we have:

$$\begin{aligned}\hat{\mu}_t - \mu^* &= \Sigma_t^{-1} \sum_{\tau=0}^{t-1} r_\tau \mathbf{x}_\tau - \mu^* = \Sigma_t^{-1} \sum_{\tau=0}^{t-1} \mathbf{x}_\tau (\mathbf{x}_\tau \cdot \mu^* + \eta_\tau) - \mu^* \\ &= \Sigma_t^{-1} \left(\sum_{\tau=0}^{t-1} \mathbf{x}_\tau (\mathbf{x}_\tau)^\top \right) \mu^* - \mu^* + \Sigma_t^{-1} \sum_{\tau=0}^{t-1} \eta_\tau \mathbf{x}_\tau \\ &= \lambda \Sigma_t^{-1} \mu^* + \Sigma_t^{-1} \sum_{\tau=0}^{t-1} \eta_\tau \mathbf{x}_\tau\end{aligned}$$

By the triangle inequality,

$$\begin{aligned}\sqrt{(\hat{\mu}_t - \mu^*)^\top \Sigma_t (\hat{\mu}_t - \mu^*)} &\leq \left\| \lambda \Sigma_t^{-1/2} \mu^* \right\| + \left\| \Sigma_t^{-1/2} \sum_{\tau=0}^{t-1} \eta_\tau \mathbf{x}_\tau \right\| \\ &\leq \sqrt{\lambda} \|\mu^*\| + ??.\end{aligned}$$

How can we bound “??” To be continued...



Self-normalized Martingale concentration bound.

Lemma (Self-Normalized Bound for Vector-Valued Martingales)

(Abassi et. al '11) Suppose $\{\varepsilon_i\}_{i=1}^{\infty}$ are mean zero random variables (can be generalized to martingales), and ε_i is bounded by σ . Let $\{X_i\}_{i=1}^{\infty}$ be a stochastic process. Define $\Sigma_t = \Sigma_0 + \sum_{i=1}^t X_i X_i^{\top}$. With probability at least $1 - \delta$, we have for all $t \geq 1$:

$$\left\| \sum_{i=1}^t X_i \varepsilon_i \right\|_{\Sigma_t^{-1}}^2 \leq \sigma^2 \log \left(\frac{\det(\Sigma_t) \det(\Sigma_0)^{-1}}{\delta^2} \right).$$

Continue the proof by applying concentration, and the bound for information-gain

$$\begin{aligned}
 \sqrt{(\hat{\mu}_t - \mu^*)^\top \Sigma_t (\hat{\mu}_t - \mu^*)} &= \|(\Sigma_t)^{1/2} (\hat{\mu}_t - \mu^*)\| \\
 &\leq \left\| \lambda \Sigma_t^{-1/2} \mu^* \right\| + \left\| \Sigma_t^{-1/2} \sum_{\tau=0}^{t-1} \eta_\tau x_\tau \right\| \\
 &\leq \sqrt{\lambda} \|\mu^*\| + \sqrt{2\sigma^2 \log(\det(\Sigma_t) \det(\Sigma^0)^{-1} / \delta_t)}.
 \end{aligned}$$

$$\delta_t = (3/\pi^2)/t^2$$

$$1 - \Pr(\forall t, \mu^* \in \text{BALL}_t) = \Pr(\exists t, \mu^* \notin \text{BALL}_t) \leq \sum_{t=1}^{\infty} \Pr(\mu^* \notin \text{BALL}_t) < \sum_{t=1}^{\infty} (1/t^2)(3/\pi^2) = 1/2.$$

Final remarks on Linear Bandits

- The regret of LinUCB is optimal up to
- Strong assumption on realizability.
 - Agnostic linear bandits?
- Contextual version: a finite list of available actions are given at each t .