CS 5 9000-RL
Linear bandit
Agenda
- Linear regression
- Regret bound

---

*Theorem: For $\delta \subset (0,1)$, with probability at least $1-\delta$, for any $t \in [n]$, we have:
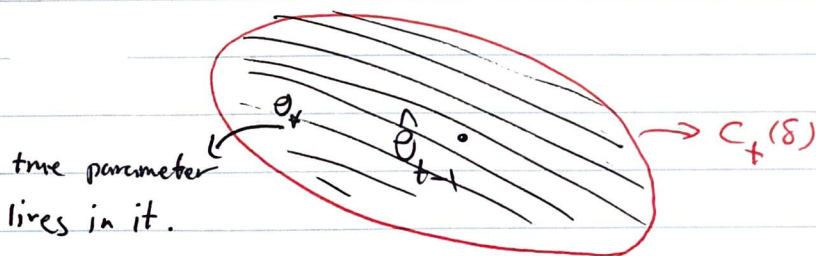
$$\|\hat{\theta}_t - \theta_*\|_{V_t(v)} \leq \sqrt{\beta_t(\delta)}: \sqrt{\lambda}\|\theta_*\| + \sqrt{2\log(\tfrac{1}{\delta}) + \log\left(\frac{\det V_t(v)}{\det(v)}\right)}$$

for $V = \lambda I$.

Furthermore, if $\|\theta_*\| \leq S$, define confidence interval/set

$$C_t(\delta) = \left\{ \theta \in \mathbb{R}^d : \|\hat{\theta}_{t-1} - \theta\|_{V_{t-1}(v)} \leq \sqrt{\lambda} S + \sqrt{2\log(\tfrac{1}{\delta}) + \log\left(\frac{\det(V_{t-1}(v))}{\det v}\right)} \right\}$$

Then, $\mathbb{P}\left(\text{exist } t \in [n] : \theta_* \notin C_t(\delta)\right) \leq \delta$



true parameter lives in it.

$$\|\theta - \hat{\theta}_{t-1}\|^2_{V_{t-1}(v)} \leq \beta_t(\delta)$$

Note: Setting $V = \lambda I \Rightarrow \det(v) = \lambda^d$

Can we simplify $\det(V_t(\lambda I))$?

Remember that $V_t(\lambda) = \lambda I + \sum_{s=1}^{t} A_s A_s^T$ for $1 \leq s \leq t$

Assume that all $a \in D_s$; $\|a\| \leq L$ for all $1 \leq s \leq t$.
Also note that $V_t(\lambda)$ is positive definite matrix (why?)

Let $S_1, \ldots, S_d$ denote the eigenvalues of $V_t(\lambda)$.

Therefore, $\det(V_t(\lambda)) = \prod_{i=1}^{d} S_i$ and the trace $(V_t \lambda)$

$= \sum_{i=1}^{d} S_i$ . By inequality of arthmatic and

geometric mean of positive numbers we have:

$$\sqrt[d]{\prod_{i=1}^{d} S_i} \leq \frac{\sum_{i=1}^{d} S_i}{d}$$

Therefore,

$$\det(V_t(\lambda)) \leq \left( \frac{\text{trace}(V_t \lambda)}{d} \right)^d$$

Let's simplify the trace:

we know that, $\text{trace}(V_t(\lambda)) = \text{trace}(\lambda I) + \sum_{s=1}^{t} \text{trace}(A_s A_s^T)$

$$= d\lambda + \sum_{s=1}^{t} \|A_s\|_2^2 \leq d\lambda + t L^2$$

$\otimes$ $\log(\det(V_t(\lambda))) \leq d \log\left(\lambda + \frac{t L^2}{d}\right)$

Now using these simplifications:

$$\beta_t(\delta) \leq \sqrt{\lambda} S + \sqrt{2 \log\left(\frac{1}{\delta}\right) + d \log\left(\lambda + \frac{t L^2}{d}\right) - d \log(\lambda)}$$

$$\leq \sqrt{\lambda} S + \sqrt{2 \log\left(\frac{1}{\delta}\right) + d \log\left(\lambda + \frac{n L^2}{d}\right) - d \log(\lambda)}$$

---

**Proof of theorem $\otimes$:**

$$\hat{\theta}_t = V_t(\lambda)^{-1} \sum_{s=1}^{t} A_s X_s$$

remember that $X_s = \langle A_s, \theta_* \rangle + \eta_s$, therefore

$$\hat{\theta}_t = V_t(\lambda)^{-1} \left( \sum_{s=1}^{t} A_s A_s^T \theta_* + \underbrace{\sum_{s=1}^{t} A_s \eta_s}_{\to S_t} \right)$$

$$= V_t(\lambda)^{-1} V_t \theta_* + V_t(\lambda)^{-1} S_t$$

Page 4)

Using this equality, we have:

$$\|\hat{\theta}_t - \theta^*\|_{V_t(\lambda)} = \left\| V_t(\lambda)^{-1} S_t + \left( V_t(\lambda)^{-1} V_t - I \right) \theta_* \right\|_{V_t(\lambda)}$$

$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underset{V_t(\lambda)^{-1} V_t(\lambda)}{}$$

Using the fact that $\|\circ\|_{V_t(\lambda)}$ is a norm, and triangle inequality

$$\|\hat{\theta}_t - \theta_*\|_{V_t(\lambda)} \leq \|S_t\|_{V_t(\lambda)^{-1}} + \left\| \left( V_t(\lambda)^{-1} V_t - I \right) \theta_* \right\|_{V_t(\lambda)}$$

$$= \|S_t\|_{V_t(\lambda)^{-1}} + \left\| \lambda V_t(\lambda)^{-1} \theta_* \right\|_{V_t(\lambda)}$$

$$\leq \|S_t\|_{V_t(\lambda)^{-1}} + \lambda^{\frac{1}{2}} \|\theta_*\| .$$

From the past we know that
$$P\left( t \in [n] ; \|S_t\|_{V_t(\lambda)^{-1}} > 2 \log\left(\frac{1}{\delta}\right) + \log\left( \frac{\det(V_t(\lambda))}{\det(V)} \right) \right) \leq \delta$$

Ergo $\|\hat{\theta}_t - \theta_*\|_{V_t(\lambda)} \leq 2\log\left(\frac{1}{\delta}\right) + \log\left( \frac{\det(V_t(\lambda))}{\det(\lambda I)} \right) + \lambda^{\frac{1}{2}} \|\theta_*\|$

$$\underset{\leq S}{\underbrace{\qquad\qquad}}$$

which the statement of the theorem.

We proved this theorem for almost any sequence $A_t$.

---

Stochastic linear bandit:

- At each time step $t$, the agent is given a decision set $D_t$, from which it needs to choose an action. The reward of choosing $A_t \in D_t$ is as follows:

$$X_t = \langle A_t, \theta_* \rangle + \eta_t$$

---

what is the oracle's expected reward at time $t$?

$$\max_{a \in D_t} \langle a, \theta_* \rangle.$$

Note that this can be random!

Since $D_t$ can be random, or adversarially chosen, let $\hat{R}_n$ denote random regret defined as follows:

$$\hat{R}_n = \sum_{t=1}^{n} \max_{a \in D_t} \langle a, \theta_* \rangle - \sum_{t=1}^{n} X_t$$

---

Consider a setting where $|\langle a, \theta_* \rangle| < 1$ and $\|a\| < L$ for all $a \in \cup_t D_t$.

# Lin UCB (LinRel, OFUL)

Pseudo code of Lin UCB

- At time step $t$, comput $\hat{\theta}_{t-1}$ = $\underset{\theta \in R^d}{\text{argmin}} \left( \sum_{s=1}^{t-1} \left( X_s - \langle A_s, \theta \rangle \right)^2 + \lambda \left( \|\theta\| \right)_2^2 \right)$

  the estimate

  we have at the

  besining of time $t$

  i.e. $\hat{\theta}_{t-1} = V_{t-1}(\lambda)^{-1} \sum_{s=1}^{t-1} A_s X_s$

- Construct the confidene st $C_t (\delta)$

$$C_t(\delta) = \{ \theta \in R^d ; \|\theta - \hat{\theta}_{t-1}\|^2_{V_{t-1}(\lambda)} \leq \beta_t(\delta) \}$$

- Optimism step: Choose an optimal arm of the most optimistic model

$$A_t = \underset{a \in D_t}{\text{arg max}} \quad \underset{\theta \in C_t}{\text{max}} \quad \langle a, \theta \rangle$$

and $\tilde{\theta}_t$ is the corresponding optimistc model

Theorem: The regret of Lin UCB satifies:

$$\hat{R}_n < \sqrt{8n \, \beta_n(\delta) \, \log \left( \frac{\det V_n(\lambda)}{\det V} \right)}$$

with probability at least $1-\delta$.