CS 59000 - RL
Linear bandit
  Agenda.
    - Regret bound
    - More structure?
    - Thompson Sampling

---

Theorem: The regret of Lin UCB satisfies,

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

with probability at least $1-\delta$.

---

Before proving this theorem, let us study the following lemma.

---

Lemma: Let $V$ be a positivedef finite matrix, $\underline{v} = \text{trace}(V)$ and $x_1, \dots, x_n \in \mathbb{R}^d$ with $\|x_i\|_2 \leqslant L < \infty$ a sequence of vectors. Let $V_t(V) = \sum_{s=1}^{t} x_s x_s^\top + V$, then:

$$\sum_{t=1}^{n} \left(1 \wedge \|x_t\|^2_{V_{t-1}^{-1}(V)}\right) \leqslant 2 \log\left(\frac{\det V_n(V)}{\det V}\right)$$

---

Proof: Using the fact that for any $u \geq 0$, $u \wedge 1 \leq 2 \log(1+u)$ _why?_
we have

$$\sum_{t=1}^{n} \left( 1 \wedge \|x_t\|^2_{V_t(V)^{-1}} \right) \leq 2 \sum_{t=1}^{n} \log\left( 1 + \|x_t\|^2_{V_{t-1}(V)^{-1}} \right)$$

on the other hand, $V_t(V) = V_{t-1}(V) + x_t x_t^T$

$$= V_{t-1}(V)^{1/2} \left( I + V_{t-1}(V)^{-\frac{1}{2}} x_t x_t^T V_{t-1}(V)^{-\frac{1}{2}} \right) V_{t-1}(V)^{1/2}$$

Therefore:

$$\det\left( V_t(V) \right) = \det\left( V_{t-1}(V) \right) \det\left( I + V_{t-1}(V)^{-\frac{1}{2}} x_t x_t^T V_{t-1}(V)^{-\frac{1}{2}} \right)$$

$$= \det\left( V_{t-1}(V) \right) \left( 1 + \|x_t\|^2_{V_{t-1}(V)^{-1}} \right)$$

Ergo

$$\det\left( V_t(V) \right) = \det(V) \prod_{t=1}^{n} \left( 1 + \|x_t\|_{V_{t-1}(V)^{-1}} \right)$$

Hence,

$$\sum_{t=1}^{n} \left( 1 \wedge \|x_t\|^2_{V_{t-1}(V)^{-1}} \right) \leq$$

$$\leq 2 \log\left( \frac{\det(V_n(V))}{\det(V)} \right)$$

which is the statement of the lemma.

Page3

Proof of the regret theorem :
Consider the per step reget of lin UCB.
what is it?  $\langle A_t^*, \theta^* \rangle - \langle A_t, \theta^* \rangle$

$A_t^* = \text{argmax } \langle a, \theta^* \rangle$
$\qquad a \in D_t$

↳ let's play a little bit.

$$\langle A_t^*, \theta^* \rangle - \langle A_t, \theta^* \rangle \leqslant \langle A_t, \tilde{\theta} \rangle - \langle A_t, \theta^* \rangle$$

$A_t = \text{argmax } \langle a, \tilde{\theta}_t \rangle$
$\qquad a \in D_t$

↳ $\langle A_t^*, \tilde{\theta} \rangle \leqslant \langle A_t, \tilde{\theta} \rangle$

— where optimism kicks in

$$= \langle A_t, \tilde{\theta}_t - \theta^* \rangle$$

$\langle A_t, \tilde{\theta}_t - \theta^* \rangle =$

$A_t^T V_{t-1}^{-\frac{1}{2}}(v) V_{t-1}(v)^{\frac{1}{2}} (\tilde{\theta} - \theta^*)$

$\longleftarrow \quad \leqslant \|A_t\|_{V_{t-1}^{-1}(v)} \|\tilde{\theta}_t - \theta^*\|_{V_t(v)}$

$$\leqslant \|A_t\|_{V_{t-1}^{-1}(v)} \left( \|\tilde{\theta}_t - \hat{\theta}_{t-1} + \hat{\theta}_{t-1} - \theta^*\| \right)$$

$$\leqslant \|A_t\|_{V_{t-1}^{-1}(v)} \left( 2 \sqrt{B_{t-1}(\delta)} \right)$$

on the other hand we know that

$$|\langle A_t^*, \theta^* \rangle - \langle A_t, \theta^* \rangle| \leqslant 2$$

Page 4:

$$\langle A_t^*, \hat{\theta^*} \rangle - \langle A_t, \theta^* \rangle \leq 2 \wedge \left( 2 \|A\|_{V_{t-1}(V)} \sqrt{\beta_t(\delta)} \right)$$

when $\beta_n(\delta) \geq 1 \vee \beta_t(\delta)$ we have

$$\langle A_t^*, \theta^* \rangle - \langle A, \theta^* \rangle \leq 2\sqrt{\beta_n(\delta)} \left( 1 \wedge \|A_t\|_{V_t(V)^{-1}} \right) \qquad \textcolor{red}{why?}$$

Finally, using Jenson's inequality we have

$$\hat{R}_n = \sum_{t=1}^{n} \langle A_t^*, \theta^* \rangle - \langle A_t, \theta^* \rangle$$

$$\leq \sqrt{n \sum_{t=1}^{n} \left( \langle A_t^*, \theta^* \rangle - \langle A_t, \theta^* \rangle \right)^2}$$

$$\leq 2\sqrt{n \, \beta_n \sum_{t=1}^{n} \left( 1 \wedge \|A_t\|_{V_t(V)^{-1}} \right)}$$

$$\leq 2\sqrt{2n \, \beta_n \, \log \left( \frac{\det(V_n(V))}{\det V} \right)}$$

Remark, setting $V = \lambda I$, where $\det(V) = \lambda^d$ and the fact that $\det(V_n(\lambda)) \leq \left( \frac{\lambda + nL^2}{d} \right)^d$

we have $\hat{R}_n \leq 2\sqrt{2n\beta_n \left( d \, \log \left( \frac{\lambda + nL^2/d}{\lambda} \right) \right)}$

Remark:

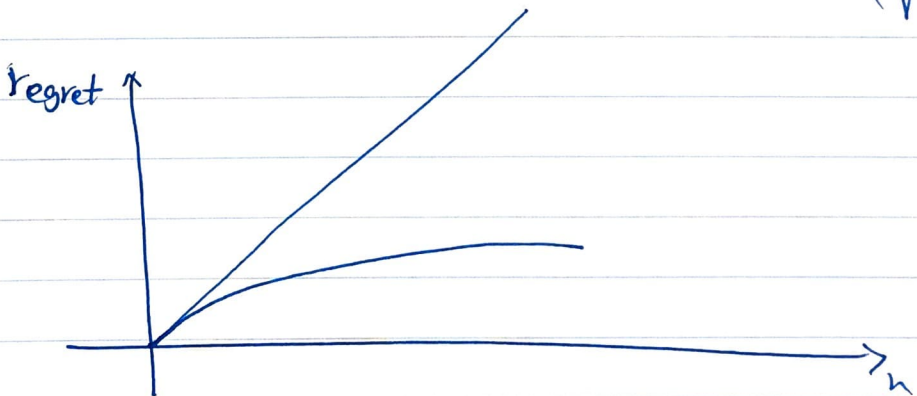$$\hat{R}_n \leqslant 2\sqrt{2nd \cdot \beta_n \log\left(\frac{\lambda + \frac{nL^2}{d}}{\lambda}\right)}$$

$$\leqslant 2\sqrt{2nd \log\left(\frac{\lambda + \frac{nL^2}{d}}{\lambda}\right)}\left(\sqrt{\lambda} S + \sqrt{d \log\left(\frac{1 + \frac{nL^2}{d\lambda}}{\delta^2/d}\right)}\right)$$

where we used $\sqrt{\beta_n(\delta)} = \sqrt{\lambda} S + \sqrt{2\log\left(\frac{1}{\delta}\right) + d \log\left(\frac{\lambda + \frac{nL^2}{d}}{\lambda}\right)}$

Remark

$$\hat{R}_n = \tilde{O}\left(d\sqrt{n}\right)$$

your per step regret vanishes with rate $\tilde{O}\left(\frac{d}{\sqrt{n}}\right)$



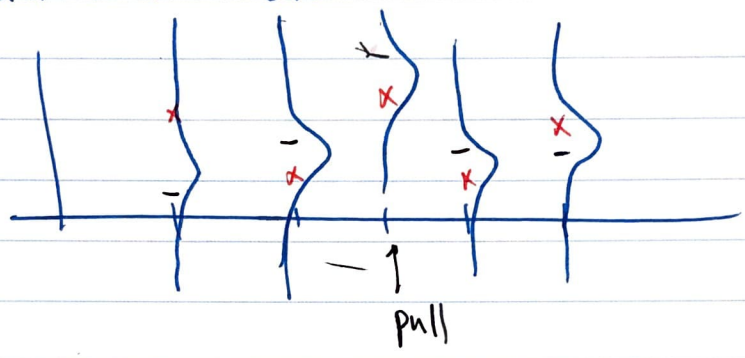$\Omega\left(d\sqrt{n}\right) \rightarrow$ lower bound is $\sqrt{n}$

$\Theta\left(d\sqrt{n}\right)$ or $\theta\left(d\sqrt{n}\right)$, Algorithm matches the lower bound.
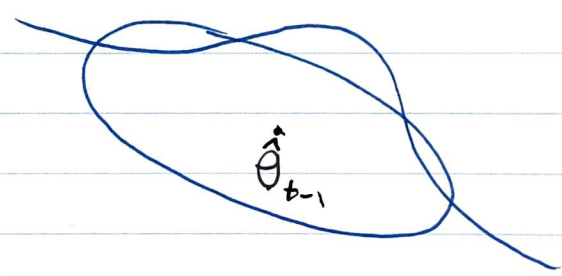
Page 6

Is ther any other approach than optimism?

Thompson Sampling. [ Thompson Sampling for Contextual Bandit with Linear Pay off.
— Linear Thompson Sampling Revisibed.

For multi-armed bandit:



pull

For linear bandit



$\hat{\theta}_{t-1}$

I draw a $\theta_{TS}$ then

$A_t = \text{argmax} \langle a, \theta_{TS} \rangle$

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

Saprse linear bandit: We assume $\theta^*$ is an spare vector
means $\| \theta^* \|_0 \leq S$.

Low dimensional bandit: Arms are close to a linear subspace with small dimension

Adversarial linear bandit: $\theta_1 \cdots \theta_n$, $x_t = \langle A_t, \theta_t \rangle$