# Linear Bandits

## Sham M. Kakade (and Wen Sun)

Hi!

# Outline

- (distribution free) Agnostic learning is not possible in RL:
  we showed that to get $O(\log |\Pi|)$ sample complexity we need either:
  - $\text{poly}(|\mathcal{S}|)$ samples OR
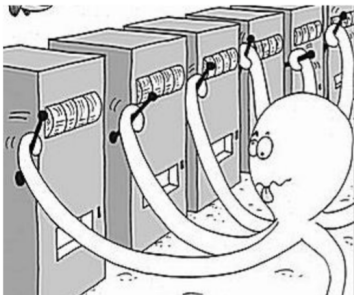  - $\text{poly}(H)$ samples.

    *exp(H)*

  in order to learn the best policy in some policy class.
- upshot: we need stronger assumptions for RL analysis.

# Multi-Armed-Bandits: High-level picture

## Setting

- Set of alternatives (arms)
- Each arm has a reward distribution

- Learner adaptively selects arms
- Challenge: Distributions not known

# Upper Confidence Bound (UCB)

Pick arm with highest Upper Confidence Bound

By Hoeffding and union bound, with probability $\geq 1 - \delta$, it holds $\forall a \in [k], t \in [T]$:

$$\mu(a) \in [LCB^t(a), UCB^t(a)]$$

$$UCB^t(a) = \tilde{\mu}^t(a) + \sqrt{\frac{\log(2kT/\delta)}{2n^t(a)}}$$

$$LCB^t(a) = \tilde{\mu}^t(a) - \sqrt{\frac{\log(2kT/\delta)}{2n^t(a)}}$$

*Actual reward means*

**Claim :** In the event that all confidence intervals confidence intervals hold, the regret is at most

$$\sum_t (UCB^t(a^t) - LCB^t(a^t)) + \delta \cdot T$$

<u>Proof:</u> $Reg^t = \mu(a^\star) - \mu(a^t)$

$$\leq UCB^t(a^\star) - LCB^t(a^t)$$
$$\leq UCB^t(a^t) - LCB^t(a^t)$$

# Upper Confidence Bound (UCB)

$$UCB^t(a) = \widetilde{\mu}^t(a) + \sqrt{\frac{\log(2kT/\delta)}{2n^t(a)}} \qquad\qquad LCB^t(a) = \widetilde{\mu}^t(a) - \sqrt{\frac{\log(2kT/\delta)}{2n^t(a)}}$$

**Claim :** In the event that all confidence intervals confidence intervals hold, the regret is at most $\sum_t(UCB^t(a^t) - LCB^t(a^t)) + + \delta \cdot T$

**Regret bound by confidence sum**

$$\sum_t(UCB^t(a^t) - LCB^t(a^t)) \leq 2 \cdot \sum_t \sqrt{\frac{\log\left(\frac{2kT}{\delta}\right)}{2n^t(a^t)}} = \sum_a \sum_{j=1}^{N(a)} \sqrt{\frac{\log\left(\frac{2kT}{\delta}\right)}{2 \cdot j}}$$

$$\leq \sum_a \sum_{j=1}^{\frac{T}{k}} \sqrt{\frac{\log\left(\frac{2kT}{\delta}\right)}{2 \cdot j}} \leq k \cdot \sqrt{\log\left(\frac{2kT}{\delta}\right) \cdot \frac{T}{k}} = O\left(\sqrt{T \cdot k \cdot \log\left(\frac{kT}{\delta}\right)}\right)$$

regret $\leq O(\sqrt{T} \cdot \# \; arms)$

# Outline

Lin Bandit Model

Abe & Long '99

On each round, we must choose a decision $x_t \in D \subset R^d$.

# Handling Large Actions Spaces

- On each round, we must choose a decision $x_t \in D \subset R^d$.
- Obtain a reward $r_t \in [-1, 1]$, where

$$\mathbb{E}[r_t | x_t = x] = \mu^\star \cdot x \in [-1, 1],$$

- On each round, we must choose a decision $x_t \in D \subset R^d$.
- Obtain a reward $r_t \in [-1, 1]$, where

$$\mathbb{E}[r_t | x_t = x] = \mu^\star \cdot x \in [-1, 1],$$

- so the the conditional expectation of $r_t$ is linear)
- Also, we have the *noise sequence*,

$$\eta_t = r_t - \mu^\star \cdot x_t$$

is i.i.d noise.

observe

$r_t = \mu^{\star} \cdot X + \eta$

Unknown

$$\sum_{t=0}^{T-1} r_k \sim \sum_{t=0}^{T-1} \mu^\star \cdot x_t$$

If $x_0, \ldots x_{T-1}$ are our decisions, then our cumulative regret is

$$R_T = \overbrace{\left(\mu^\star \cdot x^\star\right)} - \sum_{t=0}^{T-1} \mu^\star \cdot x_t$$

where $x^\star \in D$ is an optimal decision for $\mu^\star$, i.e.

$$x^\star \in \text{argmax}_{x \in D} \, \mu^\star \cdot x$$

online Lin. Programing if D was a polytope.

# Outline

After *t* rounds, define our uncertainty region $\text{BALL}_t$: with center, $\widehat{\mu}_t$, and shape, $\Sigma_t$, using the $\lambda$-regularized least squares solution:

*seen*

$x_1 \cdots x_{t-1}$

$r_1 \cdots r_{t-1}$

$$\widehat{\mu}_t = \arg\min_{\mu} \sum_{\tau=0}^{t-1} \|\mu \cdot x_t - r_t\|_2^2 + \lambda\|\mu\|_2^2$$

$$= \Sigma_t^{-1} \sum_{\tau=0}^{t-1} r_\tau x_\tau,$$

$$\Sigma_t = \lambda I + \sum_{\tau=0}^{t-1} x_t x_t^\top, \text{ with } \Sigma_0 = \lambda I.$$

# The "Confidence Ball"

After $t$ rounds, define our uncertainty region $\text{BALL}_t$: with center, $\widehat{\mu}_t$, and shape, $\Sigma_t$, using the $\lambda$-regularized least squares solution:

$$\widehat{\mu}_t = \arg \min_{\mu} \sum_{\tau=0}^{t-1} \|\mu \cdot x_t - r_t\|_2^2 + \lambda \|\mu\|_2^2$$

$$= \Sigma_t^{-1} \sum_{\tau=0}^{t-1} r_\tau x_\tau,$$

$$\Sigma_t = \lambda I + \sum_{\tau=0}^{t-1} x_t x_t^\top, \text{ with } \Sigma_0 = \lambda I.$$

Define the uncertainty region: $\{ \mu \mid (\widehat{\mu}_t - \mu) \Sigma_t (\widehat{\mu}_t - \mu) \leq \beta_t \}$

$$\text{BALL}_t = \left\{ (\widehat{\mu}_t - \mu^\star)^\top \Sigma_t^{-1} (\widehat{\mu}_t - \mu^\star) \leq \beta_t \right\},$$

where $\beta_t$ is a parameter of the algorithm.

1. Input: $\lambda$, $\beta_t$
2. For $t = 0, 1, \ldots$
   1. Execute
      $$x_t = \text{argmax}_{x \in D} \max_{w \in \text{BALL}_t} w \cdot x$$
      and observe the reward $r_t$.
   2. Update $\text{BALL}_{t+1}$.

using $x_t, r_t$

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

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

$$|D| \sim \left(\tfrac{1}{\varepsilon}\right)^d$$

# LinUCB Regret Bound

Sublinear regret: $R_T \leq O^\star(d\sqrt{T})$ *wlog* ∃ *scaling* $W \approx O(\sqrt{d})$
poly dependence on $d$ , no dependence on the cardinality $|D|$. $\|x\| \leq 1$

## Theorem

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

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

*t·t* (handwritten annotation over $T$)

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

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

*Tight in $d$ & $T$*

*where c is an absolute constant.*

(due to Dani, K., and Hayes '09)     *av reg.* $\frac{d}{\sqrt{T}}$

# Outline

# Confidence

In establishing the upper bounds there are two main propositions from which the upper bounds follow. The first is in showing that the confidence region is valid.

## Proposition

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

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

# Sum of Squares Regret Bound

Assuming the confidence event holds, the following controls on the growth of the regret.

## Proposition

*(Sum of Squares Regret Bound) Define:*

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

*instantaneous regret*

*Suppose $\|x\| \leq B$ for $x \in D$. Suppose $\beta_t$ is set as in Theorem 1. Suppose $\mu^\star \in \text{BALL}_t$ for all t, then*

$O(\cdot \log T)$

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

$$\beta_T \approx O\left(d \log T\right)$$

**Proof:**[Proof of Theorem 1] With the two previous Propositions, along with the Cauchy-Schwarz inequality, we have, with probability at least $1 - \delta$,

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

The remainder of the proof follows from using our chosen value of $\beta_T$ and algebraic manipulations. ∎

# Outline

# "Width" of Confidence Ball

## Lemma

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

$$|(\mu - \widehat{\mu}_t)^\top x| \leq \sqrt{\beta_t x^\top \Sigma_t^{-1} x}$$
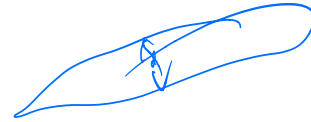
# "Width" of Confidence Ball

## Lemma

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

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

**Proof:** By Cauchy-Schwarz, we have:

$$|(\mu - \widehat{\mu}_t)^\top x| = |(\mu - \widehat{\mu}_t)^\top \Sigma_t^{1/2} \Sigma_t^{-1/2} x| = |(\Sigma_t^{1/2}(\mu - \widehat{\mu}_t))^\top \Sigma_t^{-1/2} x|$$

$$\le \|\Sigma_t^{1/2}(\mu - \widehat{\mu}_t)\| \|\Sigma_t^{-1/2} x\| = \|\Sigma_t^{1/2}(\mu - \widehat{\mu}_t)\| \sqrt{x^\top \Sigma_t^{-1} x} \le \sqrt{\beta_t x^\top \Sigma_t^{-1} x}$$

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

# Instantaneous Regret Lemma

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.

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 $\mu^\star \in \mathrm{BALL}_t$, then

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

$\beta_t$ is increasing

# Instantaneous Regret Lemma

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

*(handwritten: $\Sigma_{t+1} \leftarrow \Sigma_t + x_t x_t^\top$)*

which is the "normalized width" at time $t$ in the direction of our decision.

## Lemma

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

*(handwritten: $\mathbb{E}[r_t | x_t] \in [-1, 1]$)*

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

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

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

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

$$\text{regret}_t = (\mu^\star)^\top x^* - (\mu^\star)^\top x_t \leq (\widetilde{\mu} - \mu^\star)^\top x_t$$
$$= (\widetilde{\mu} - \widehat{\mu}_t)^\top x_t + (\widehat{\mu}_t - \mu^\star)^\top x_t \leq 2 \sqrt{\beta_t} w_t$$

# Geometric Argument: Part 1

The next two lemmas give us 'geometric' potential function argument, where can bound the sum of widths independently of the choices made by the algorithm.

# Geometric Argument: Part 1

The next two lemmas give us 'geometric' potential function argument, where can bound the sum of widths independently of the choices made by the algorithm.

## Lemma

*We have:*

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

# Geometric Argument: Part 1

The next two lemmas give us 'geometric' potential function argument, where can bound the sum of widths independently of the choices made by the algorithm.

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

## Lemma

*We have:*

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

*volume update equality*

**Proof:** By the definition of $\Sigma_{t+1}$, we have

$$\det \Sigma_{t+1} = \det(\Sigma_t + x_t x_t^\top) = \det(\Sigma_t^{1/2}(I + \Sigma_t^{-1/2} x_t x_t^\top \Sigma_t^{-1/2})\Sigma_t^{1/2})$$

$$= \det(\Sigma_t) \det(I + \Sigma_t^{-1/2} x_t (\Sigma_t^{-1/2} x_t)^\top) = \det(\Sigma_t) \det(I + v_t v_t^\top),$$

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

## Lemma

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

$$\Sigma_0 \doteq \lambda I$$

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

# Geometric Argument: Part 2

## Lemma

*For any sequence $x_0, \ldots x_{T-1}$ 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 $\sigma_1, \ldots \sigma_d$, and note:

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

Using the AM-GM inequality,

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

$$= d\log\left(\prod_{i=1}^{d}\left(1 + \sigma_i/\lambda\right)\right)^{1/d} \leq d\log\left(\frac{1}{d}\sum_{i=1}^{d}(1 + \sigma_i/\lambda)\right) \leq d\log\left(1 + \frac{TB^2}{d\lambda}\right)$$

**Proof:**[Proof of Proposition 3] Assume $\mu^\star \in \text{BALL}_t$ for all $t$. We have:

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

$$\leq 4\beta_T \sum_{t=0}^{T-1} \ln(1 + w_\tau^2) \leq 4\beta_T \log\left(\det \Sigma_{T-1}/\det \Sigma_0\right)$$

$$= 4\beta_T d \log\left(1 + \frac{TB^2}{d\lambda}\right)$$

where the first inequality follow from by Lemma 5; the second from that $\beta_t$ is an increasing function of $t$; the third uses that for $0 \leq y \leq 1$, $\ln(1 + y) \geq y/2$; the final two inequalities follow by Lemmas 6 and 7. $\blacksquare$

# Outline

$$Pr\left(\forall t, \; \mu^* \in Ball_t\right)$$

$$\leq S$$

**Proof:** Since $r_\tau = x_\tau \cdot \mu^\star + \eta_\tau$, we have:

$$\widehat{\mu}_t - \mu^\star = \Sigma_t^{-1} \sum_{\tau=0}^{t-1} r_\tau x_\tau - \mu^\star = \Sigma_t^{-1} \sum_{\tau=0}^{t-1} x_\tau (x_\tau \cdot \mu^\star + \eta_\tau) - \mu^\star$$

$$= \Sigma_t^{-1} \left( \sum_{\tau=0}^{t-1} x_\tau (x_\tau)^\top \right) \mu^\star - \mu^\star + \Sigma_t^{-1} \sum_{\tau=0}^{t-1} \eta_\tau x_\tau$$

$$= \lambda \Sigma_t^{-1} \mu^\star + \Sigma_t^{-1} \sum_{\tau=0}^{t-1} \eta_\tau x_\tau$$

$+ \lambda I - \lambda I$

$\Sigma_t = \Sigma x_\tau x_\tau^\top + \lambda I$

# Confidence [Proof of Proposition 2]

**Proof:** Since $r_\tau = x_\tau \cdot \mu^\star + \eta_\tau$, we have:

$$\widehat{\mu}_t - \mu^\star = \Sigma_t^{-1} \sum_{\tau=0}^{t-1} r_\tau x_\tau - \mu^\star = \Sigma_t^{-1} \sum_{\tau=0}^{t-1} x_\tau (x_\tau \cdot \mu^\star + \eta_\tau) - \mu^\star$$

$$= \Sigma_t^{-1} \left( \sum_{\tau=0}^{t-1} x_\tau (x_\tau)^\top \right) \mu^\star - \mu^\star + \Sigma_t^{-1} \sum_{\tau=0}^{t-1} \eta_\tau x_\tau$$

$$= \lambda \Sigma_t^{-1} \mu^\star + \Sigma_t^{-1} \sum_{\tau=0}^{t-1} \eta_\tau x_\tau$$

By the triangle inequality,

$$(\widehat{\mu}_t - \mu^\star)^\top \Sigma_t (\widehat{\mu}_t - \mu^\star) \leq \left\| \lambda (\Sigma_t)^{-1/2} \mu^\star \right\| + \left\| \Sigma_t^{-1} \sum_{\tau=0}^{t-1} \eta_\tau x_\tau \right\|$$

$$\leq \sqrt{\lambda} \| \mu^\star \| \qquad + \qquad ??.$$

How can we bound "??" To be continued...

## 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 $\varepsilon_i$ is bounded by $\sigma$. Let $\{X_i\}_{i=1}^{\infty}$ be a stochastic process. Define $V_t = V_0 + \sum_{i=1}^{t} X_i X_i^{\top}$. With probability at least $1 - \delta$, we have for all $t \geq 1$:*

$$\left\| \sum_{i=1}^{t} X_i \varepsilon_i \right\|_{V_t^{-1}}^2 \leq 2\sigma^2 \log \left( \frac{\det(V_t)^{1/2} \det(V_0)^{-1/2}}{\delta} \right).$$

**Proof:**

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

We seek to lower bound $\Pr(\forall t, \mu^\star \in \mathrm{BALL}_t)$. Assign failure probability $\delta_t = (3/\pi^2)/t^2$ for the $t$-th event, which gives us:

$$1 - \Pr(\forall t, \mu^\star \in \mathrm{BALL}_t) = \Pr(\exists t, \mu^\star \notin \mathrm{BALL}_t) \leq \sum_{t=1}^{\infty} \Pr(\mu^\star \notin \mathrm{BALL}_t)$$

$$< \sum_{t=1}^{\infty} (1/t^2)(3/\pi^2) = 1/2.$$

This along with Lemma 7 completes the proof. ∎