# Linear Bandits

Sham M. Kakade and Wen Sun

# Outline

# Intro to MAB



**Setting:**
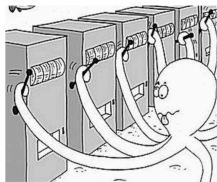
We have K many arms: $a_1, \ldots, a_K$

Each arm has a unknown reward distribution, i.e., $\nu_i \in \Delta([0,1])$, w/ mean $\mu_i = \mathbb{E}_{r \sim \nu_i}[r]$

**Example:** $a_i$ has a Bernoulli distribution $\nu_i$ w/ mean $\mu_i := p$:

Every time we pull arm $a_i$, we observe an i.i.d reward $r = \begin{cases} 1 & \text{w/ prob } p \\ 0 & \text{w/ prob } 1 - p \end{cases}$

# Intro to MAB

**More formally, we have the following learning objective:**

$$\text{Regret}_T = T\mu^\star - \sum_{t=0}^{T-1} \mu_{I_t} \qquad \mu^\star = \max_{i \in [K]} \mu_i$$
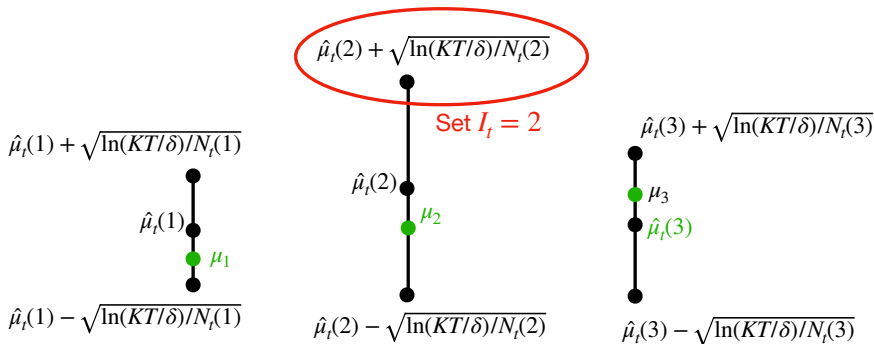
Total expected reward if we pulled best arm over T rounds

Total expected reward of the arms we pulled over T rounds

Goal: no-regret, i.e., $\text{Regret}_T/T \to 0$, as $T \to \infty$

# UCB: Optimism in the face of Uncertainty

Given the confidence interval, we pick arm that has the **highest Upper-Conf-Bound:**



$$\hat{\mu}_t(2) + \sqrt{\ln(KT/\delta)/N_t(2)}$$

Set $I_t = 2$

$\hat{\mu}_t(1) + \sqrt{\ln(KT/\delta)/N_t(1)}$

$\hat{\mu}_t(1)$

$\mu_1$

$\hat{\mu}_t(1) - \sqrt{\ln(KT/\delta)/N_t(1)}$

$\hat{\mu}_t(2)$

$\mu_2$

$\hat{\mu}_t(2) - \sqrt{\ln(KT/\delta)/N_t(2)}$

$\hat{\mu}_t(3) + \sqrt{\ln(KT/\delta)/N_t(3)}$

$\mu_3$

$\hat{\mu}_t(3)$

$\hat{\mu}_t(3) - \sqrt{\ln(KT/\delta)/N_t(3)}$

# UCB Regret:

[Theorem (informal)] With high probability, UCB has the following regret:

$$\text{Regret}_T = \widetilde{O}\left(\sqrt{KT}\right)$$

# Generalization in RL

- (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
  - $\exp(H)$ samples.

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

# Outline

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

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

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

# 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

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

  - 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.

model due to Abe & Long '99

## Our Objective

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

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

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

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

## 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:

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

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

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

## 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:

$$
\begin{aligned}
\widehat{\mu}_t &= \arg\min_{\mu} \sum_{\tau=0}^{t-1} \|\mu \cdot x_\tau - r_\tau\|_2^2 + \lambda\|\mu\|_2^2 \\
&= \Sigma_t^{-1} \sum_{\tau=0}^{t-1} r_\tau x_\tau, \\
\Sigma_t &= \lambda I + \sum_{\tau=0}^{t-1} x_\tau x_\tau^\top, \text{ with } \Sigma_0 = \lambda I.
\end{aligned}
$$

Define the uncertainty region:

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

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

# LinUCB (the algo)

1. Input: $\lambda$, $\beta_t$
2. For $t = 0, 1, \ldots$
   1. Execute

      $$x_t = \operatorname{argmax}_{x \in D} \max_{\mu \in \text{BALL}_t} \mu \cdot x$$

      and observe the reward $r_t$.
   2. Update $\text{BALL}_{t+1}$.

# LinUCB Regret Bound

Sublinear regret: $R_T \leq O^\star(d\sqrt{T})$
poly dependence on $d$ , no dependence on the cardinality $|D|$.

# LinUCB Regret Bound

Sublinear regret: $R_T \leq O^\star(d\sqrt{T})$
poly dependence on $d$, no dependence on the cardinality $|D|$.

## Theorem

*Suppose: $|\mu^\star \cdot x| \leq 1$ and $\|x\| \leq B$ for all $x \in D$; that the noise is $\sigma^2$ sub-Gaussian; and that $\|\mu^\star\| \leq W$. Set $\lambda = \sigma^2/W^2$ and*

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

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

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

*where $c$ is an absolute constant.*

due to Dani, Hayes, K. '09

# 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*

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

# 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:*

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

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

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

## Completing the Proof

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

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

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

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

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

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

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

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

## 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.

## 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.

### Lemma

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

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

## 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.

### Lemma

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

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

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

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

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

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

# 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.

## Lemma

*We have:*

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

**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 ... ∎

# 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:*

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

## Geometric Argument: Part 2

### 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,

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

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

# Proving "sum of squares regret" Proposition

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

$$
\begin{aligned}
\sum_{t=0}^{T-1} \text{regret}_t^2 &\leq \sum_{t=0}^{T-1} 4\beta_t \min(w_t^2, 1) \leq 4\beta_T \sum_{t=0}^{T-1} \min(w_t^2, 1) \\
&\leq 8\beta_T \sum_{t=0}^{T-1} \ln(1 + w_t^2) \leq 8\beta_T \log\left(\det \Sigma_{T-1} / \det \Sigma_0\right) \\
&= 8\beta_T d \log\left(1 + \frac{TB^2}{d\lambda}\right)
\end{aligned}
$$

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. ∎

## Confidence [Proof of Proposition 2]

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

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

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

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

## 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,

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

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

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

# Self-Normalizing Sum

## 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 $\Sigma_t = \Sigma_0 + \sum_{i=1}^{t} X_i X_i^{\top}$. With probability at least $1 - \delta$, we have for all $t \geq 1$:*

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

## Continued... [Proof of Proposition 2]

**Proof:**

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

$$\leq \sqrt{\lambda} \|\mu^\star\| + \sqrt{2\sigma^2 \log\left(\det(\Sigma_t)\det(\Sigma^0)^{-1}/\delta_t\right)}.$$

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. ∎