

# Linear Bandits

Sham M. Kakade and Wen Sun

## 1 Recap

## 2 Linear Bandits

- Setting
- LinUCB
- An Optimal Regret Bound

## 3 Analysis

- Regret Analysis
- Confidence Analysis

# Intro to MAB

## Setting:

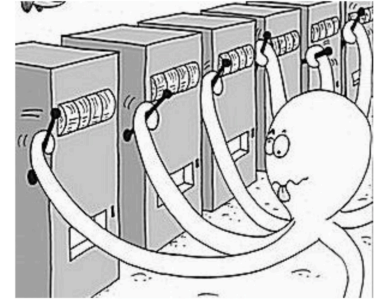
We have  $K$  many arms:  $a_1, \dots, a_K$

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:

$$\text{Regret}_T = T\mu^* - \sum_{t=0}^{T-1} \mu_{I_t} \quad \mu^* = \max_{i \in [K]} \mu_i$$

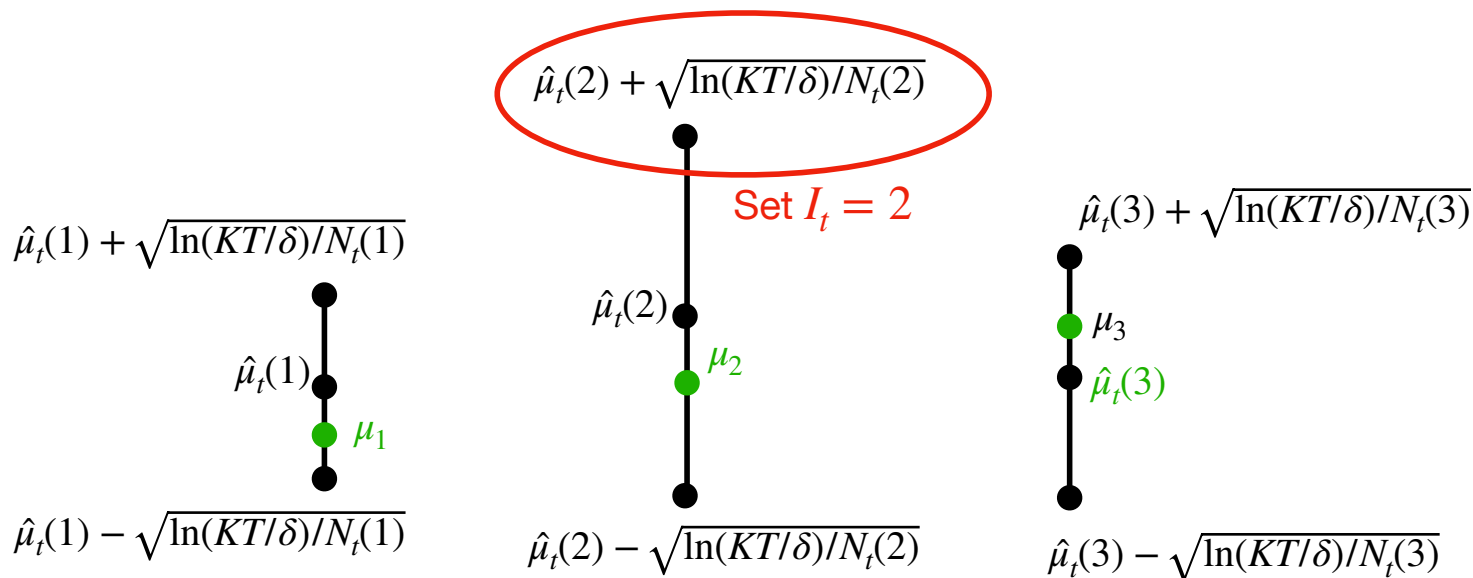
Total expected reward if we pulled best arm over T rounds

Total expected reward of the arms we pulled over T rounds

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

# UCB: Optimism in the face of Uncertainty

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



# UCB Regret:

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

$$\text{Regret}_T = \tilde{O}(\sqrt{KT}) \approx \frac{\sqrt{K}}{\sqrt{T}}$$

- (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

## 1 Recap

## 2 Linear Bandits

- Setting
- LinUCB
- An Optimal Regret Bound

## 3 Analysis

- Regret Analysis
- Confidence Analysis



# Handling Large Actions Spaces

- 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

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

# 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

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

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

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

$\mu^* \leftarrow$  edge cost vector.  
is i.i.d noise.

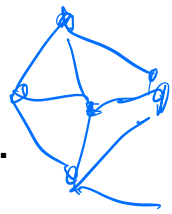
"paths on graph"

model due to Abe & Long '99

SEE cumulative path length.  
edge taken.

B

A



# Our Objective

"Data set" at time  $t$   $\{(\vec{x}_0, r_0), \dots, (\vec{x}_{t-1}, r_{t-1})\}$

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

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

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

$$x^* \in \operatorname{argmax}_{x \in D} \mu^* \cdot x$$

## 1 Recap

## 2 Linear Bandits

- Setting
- **LinUCB**
- An Optimal Regret Bound

## 3 Analysis

- Regret Analysis
- Confidence Analysis

# The “Confidence Ball”

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

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

# The “Confidence Ball”

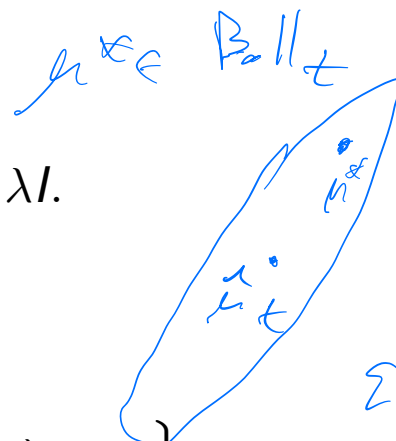
After  $t$  rounds, define our uncertainty region  $\text{BALL}_t$ : with center,  $\hat{\mu}_t$ , and shape,  $\Sigma_t$ , using the  $\lambda$ -regularized least squares solution:

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

hope

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

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



Define the uncertainty region:

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

shape of the ellipse

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

# LinUCB (the algo)

- 1 Input:  $\lambda, \beta_t$
- 2 For  $t = 0, 1, \dots$ 
  - 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:  $R_T \leq O^*(d\sqrt{T})$

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

MAB  $\leftrightarrow$  OJDIT

# LinUCB Regret Bound

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

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

does  $\hat{\mu}_t \rightarrow \mu$  ??  
very slowly

## Theorem

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

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

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

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

where  $c$  is an absolute constant.

due to Dani, Hayes, K. '09

## 1 Recap

## 2 Linear Bandits

- Setting
- LinUCB
- An Optimal Regret Bound

## 3 Analysis

- Regret Analysis
- Confidence Analysis

# 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*

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

$\forall T$   $\mathbb{E}[R_T] \leq 2\sqrt{T}$

# 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:*

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

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

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

*OC*  $(\beta_T \int \log T)$

# 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  $1 - \delta$ ,

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

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

## 1 Recap

## 2 Linear Bandits

- Setting
- LinUCB
- An Optimal Regret Bound

## 3 Analysis

- Regret Analysis
- Confidence Analysis

# “Width” of Confidence Ball

## Lemma

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

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



# “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:

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

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

# Instantaneous Regret 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  $\mu^* \in \text{BALL}_t$ , then

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

Assume all the rewards are bounded by 1

$$\text{regret}_t \leq 2\sqrt{\beta_T} w_t$$

# 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.

## Lemma

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

*True for Lin VCB algo*

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

**Proof:** Let  $\tilde{\mu} \in \text{BALL}_t$  denote the vector which ~~minimizes~~ *maximizes* 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 Argument: Part 1

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

## Lemma

*We have:*

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

$$\Sigma_t \stackrel{?}{=} \lambda I + \sum_{\tau=0}^t x_\tau x_\tau^\top$$

# 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

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

where  $v_t := \Sigma_t^{-1/2} x_t$ . Now observe that  $v_t^\top v_t = \underline{w_t^2}$  and ... det(I + a a^\top) = 1 + ||a||^2 ■

# 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:*

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



# 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,

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

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

$$U \text{Diag}(\sigma_1, \dots, \sigma_d) U^\top$$

# Proving “sum of squares regret” Proposition

$$\text{regret} \leq 2\sqrt{\beta_T} \min(w_{T-1}, 1)$$

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

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

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. ■

## 1 Recap

## 2 Linear Bandits

- Setting
- LinUCB
- An Optimal Regret Bound

## 3 Analysis

- Regret Analysis
- Confidence Analysis

# Confidence [Proof of Proposition 2]

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

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

# Confidence [Proof of Proposition 2]

**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-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:**

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

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

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

This along with Lemma 7 completes the proof. ■