

# Variance-Aware Sparse Linear Bandits

(Published as a conference paper at ICLR 2023)

Yan Dai<sup>1</sup>



Ruosong Wang<sup>2</sup>



Simon S. Du<sup>2</sup>



<sup>1</sup>IIS, Tsinghua University

<sup>2</sup>Paul G. Allen School, University of Washington

# Table of Contents

## 1 Introduction

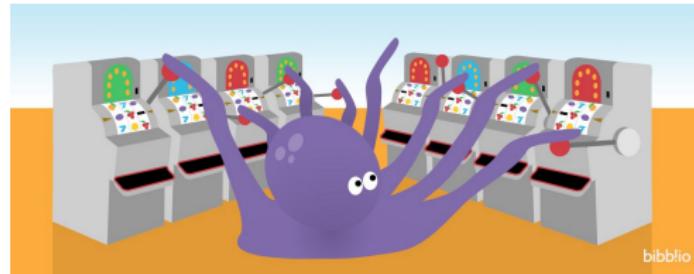
- Preliminaries
- Related Work

## 2 Algorithm

- Classical Design
- Our Design

# Linear Bandit

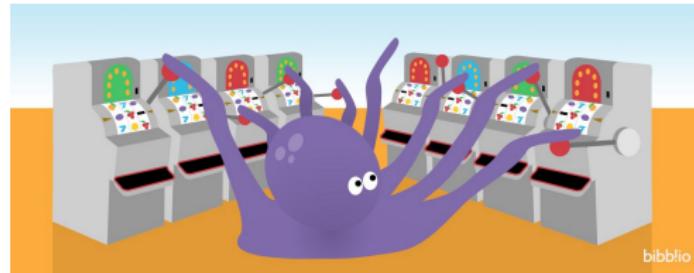
- A  $T$ -round game between an **agent** and the **environment**.<sup>1</sup>



<sup>1</sup>Figure from *Reinforcement Learning – Multi-Arm Bandit Implementation*, Jeremy Zhang.

# Linear Bandit

- A  $T$ -round game between an **agent** and the **environment**.<sup>1</sup>



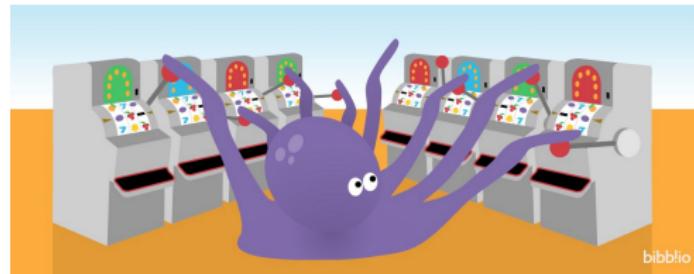
- For each round  $t = 1, 2, \dots, T$ , the agent plays an **action**  $a_t$  from the unit sphere  $\mathbb{S}^{d-1}$  (our assumption).

---

<sup>1</sup>Figure from *Reinforcement Learning – Multi-Arm Bandit Implementation*, Jeremy Zhang.

# Linear Bandit

- A  $T$ -round game between an **agent** and the **environment**.<sup>1</sup>

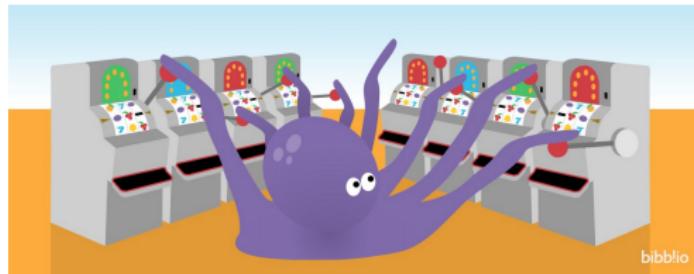


- For each round  $t = 1, 2, \dots, T$ , the agent plays an **action**  $a_t$  from the unit sphere  $\mathbb{S}^{d-1}$  (our assumption).
- For this round, she gains **reward**  $r(a_t) = \langle a_t, \theta^* \rangle$  where  $\theta^* \in \mathbb{R}^d$  is a *fixed but unknown* parameter.

<sup>1</sup>Figure from *Reinforcement Learning – Multi-Arm Bandit Implementation*, Jeremy Zhang.

# Linear Bandit

- A  $T$ -round game between an **agent** and the **environment**.<sup>1</sup>



- For each round  $t = 1, 2, \dots, T$ , the agent plays an **action**  $a_t$  from the unit sphere  $\mathbb{S}^{d-1}$  (our assumption).
- For this round, she gains **reward**  $r(a_t) = \langle a_t, \theta^* \rangle$  where  $\theta^* \in \mathbb{R}^d$  is a *fixed but unknown* parameter.
- She cannot directly access  $r(a_t)$ , but only observes noisy feedback  $r(a_t) + \eta_t$  where  $\eta_t$  is a zero-mean *random noise*. We assume  $\eta_t \sim \mathcal{N}(0, \sigma_t^2)$  where  $\sigma_t$ 's are *fixed but unknown*.

<sup>1</sup>Figure from *Reinforcement Learning – Multi-Arm Bandit Implementation*, Jeremy Zhang.

# Agent's Goal?

**Maximize** the (expected) total reward

$$\mathbb{E} \left[ \sum_{t=1}^T r(a_t) \right] = \mathbb{E} \left[ \sum_{t=1}^T \langle a_t, \theta^* \rangle \right],$$

# Agent's Goal?

**Maximize** the (expected) total reward

$$\mathbb{E} \left[ \sum_{t=1}^T r(a_t) \right] = \mathbb{E} \left[ \sum_{t=1}^T \langle a_t, \theta^* \rangle \right],$$

or equivalently, minimize the **regret**

$$\begin{aligned} \mathcal{R}_T &\triangleq \max_{a \in \mathbb{S}^{d-1}} \mathbb{E} \left[ \sum_{t=1}^T \langle a - a_t, \theta^* \rangle \right]. \\ &= \mathbb{E} \left[ \sum_{t=1}^T \langle \theta^* - a_t, \theta^* \rangle \right]. \end{aligned}$$

# Sparse Linear Bandit

- We know  $\theta^*$  only has a few non-zero coordinates. In other words,  $\|\theta^*\|_0 \triangleq s \ll d$ . However,  $s$  is *unknown* to the agent.

# Sparse Linear Bandit

- We know  $\theta^*$  only has a few non-zero coordinates. In other words,  $\|\theta^*\|_0 \triangleq s \ll d$ . However,  $s$  is *unknown* to the agent.
- **Known Results on “Worst-Case” ( $\sigma_t \equiv 1$ ) SLB:**
  - Upper Bound:  $\tilde{\mathcal{O}}(\sqrt{sdT})$  [Abbasi-Yadkori et al., 2012].
  - Lower Bound:  $\Omega(\sqrt{dT})$  [Antos and Szepesvári, 2009] even when sparsity factor  $s = 1$  and the action set is  $\mathbb{S}^{d-1}$ .

# Sparse Linear Bandit

- We know  $\theta^*$  only has a few non-zero coordinates. In other words,  $\|\theta^*\|_0 \triangleq s \ll d$ . However,  $s$  is *unknown* to the agent.
- **Known Results on “Worst-Case” ( $\sigma_t \equiv 1$ ) SLB:**
  - Upper Bound:  $\tilde{\mathcal{O}}(\sqrt{sdT})$  [Abbasi-Yadkori et al., 2012].
  - Lower Bound:  $\Omega(\sqrt{dT})$  [Antos and Szepesvári, 2009] even when sparsity factor  $s = 1$  and the action set is  $\mathbb{S}^{d-1}$ .
  - **Conclusion:**  $\tilde{\mathcal{O}}(\sqrt{sdT})$  is optimal for “worst-case” SLB!

# Variance-Aware Sparse Linear Bandit?

- **Worst Case** ( $\sigma_t \equiv 1$ ):  $\tilde{\mathcal{O}}(\sqrt{sdT})$  is shown to be optimal.

# Variance-Aware Sparse Linear Bandit?

- **Worst Case** ( $\sigma_t \equiv 1$ ):  $\tilde{\mathcal{O}}(\sqrt{sdT})$  is shown to be optimal.
- **Deterministic case** ( $\sigma_t \equiv 0$ ): Divide-and-Conquer gets  $\tilde{\mathcal{O}}(s)$ .

# Variance-Aware Sparse Linear Bandit?

- **Worst Case** ( $\sigma_t \equiv 1$ ):  $\tilde{\mathcal{O}}(\sqrt{sdT})$  is shown to be optimal.
- **Deterministic case** ( $\sigma_t \equiv 0$ ): Divide-and-Conquer gets  $\tilde{\mathcal{O}}(s)$ .
- **In Between?**

# Variance-Aware Sparse Linear Bandit?

- **Worst Case** ( $\sigma_t \equiv 1$ ):  $\tilde{\mathcal{O}}(\sqrt{sDT})$  is shown to be optimal.
- **Deterministic case** ( $\sigma_t \equiv 0$ ): Divide-and-Conquer gets  $\tilde{\mathcal{O}}(s)$ .
- **In Between? Objective of this paper!**

Design an algorithm whose regret is **variance-aware**:

$$\mathcal{R}_T = \tilde{\mathcal{O}} \left( \text{poly}(s) \sqrt{d \sum_{t=1}^T \sigma_t^2} + \text{poly}(s) \right),$$

where  $\sigma_t^2 = \text{Var}(\eta_t) \in [0, 1]$  is the noise variance ( $\sigma_t$ 's are all *unknown*) and  $s = \|\theta^*\|_0$  is the sparsity ( $s$  is also *unknown*).

# Related Work

## ① “Worst-Case” ( $\sigma_t \equiv 1$ ) Sparse Linear Bandit:

- Upper Bound:  $\tilde{\mathcal{O}}(\sqrt{sdT})$  [Abbasi-Yadkori et al., 2012].
- Lower Bound:  $\Omega(\sqrt{dT})$  [Antos and Szepesvári, 2009].

# Related Work

## ① “Worst-Case” ( $\sigma_t \equiv 1$ ) Sparse Linear Bandit:

- Upper Bound:  $\tilde{\mathcal{O}}(\sqrt{sdT})$  [Abbasi-Yadkori et al., 2012].
- Lower Bound:  $\Omega(\sqrt{dT})$  [Antos and Szepesvári, 2009].

## ② “Worst-Case” ( $\sigma_t \equiv 1$ ) Linear Bandits:

- Upper Bound:  $\tilde{\mathcal{O}}(d\sqrt{T})$  [Dani et al., 2008].
- Lower Bound:  $\Omega(d\sqrt{T})$  [Dani et al., 2008].

# Related Work

## ① “Worst-Case” ( $\sigma_t \equiv 1$ ) Sparse Linear Bandit:

- Upper Bound:  $\tilde{\mathcal{O}}(\sqrt{sdT})$  [Abbasi-Yadkori et al., 2012].
- Lower Bound:  $\Omega(\sqrt{dT})$  [Antos and Szepesvári, 2009].

## ② “Worst-Case” ( $\sigma_t \equiv 1$ ) Linear Bandits:

- Upper Bound:  $\tilde{\mathcal{O}}(d\sqrt{T})$  [Dani et al., 2008].
- Lower Bound:  $\Omega(d\sqrt{T})$  [Dani et al., 2008].

## ③ “Variance-Aware” Linear Bandits:

- $\tilde{\mathcal{O}}(d^{1.5} \sqrt{\sum \sigma_t^2} + d^2)$  [Kim et al., 2022].
- $\tilde{\mathcal{O}}(d \sqrt{\sum \sigma_t^2} + \sqrt{dT})$  [Zhou et al., 2021].

# Related Work

## ① “Worst-Case” ( $\sigma_t \equiv 1$ ) Sparse Linear Bandit:

- Upper Bound:  $\tilde{\mathcal{O}}(\sqrt{sdT})$  [Abbasi-Yadkori et al., 2012].
- Lower Bound:  $\Omega(\sqrt{dT})$  [Antos and Szepesvári, 2009].

## ② “Worst-Case” ( $\sigma_t \equiv 1$ ) Linear Bandits:

- Upper Bound:  $\tilde{\mathcal{O}}(d\sqrt{T})$  [Dani et al., 2008].
- Lower Bound:  $\Omega(d\sqrt{T})$  [Dani et al., 2008].

## ③ “Variance-Aware” Linear Bandits:

- $\tilde{\mathcal{O}}(d^{1.5} \sqrt{\sum \sigma_t^2} + d^2)$  [Kim et al., 2022].
- $\tilde{\mathcal{O}}(d \sqrt{\sum \sigma_t^2} + \sqrt{dT})$  [Zhou et al., 2021].

**This paper:** convert any LB algorithm  $\mathcal{A}$  to SLB algorithm  $\mathcal{B}$ , s.t.:

if  $\mathcal{A}$  ensures  $\mathcal{R}_T^{\text{LB}} = \tilde{\mathcal{O}}\left(f(d)\sqrt{\sum \sigma_t^2} + g(d)\right)$  for some  $f, g$ ,

then  $\mathcal{B}$  ensures  $\mathcal{R}_T^{\text{SLB}} = \tilde{\mathcal{O}}\left((sf(s) + s\sqrt{d})\sqrt{\sum \sigma_t^2} + sg(s)\right)$ .

# Classical “Explore-then-Commit” Idea for Worst-Case SLBs

- ① **Explore:** Find coordinates with “large enough” magnitudes.
- ② **Commit:** Play “wisely” on these coordinates (ignore others).

# Classical “Explore-then-Commit” Idea for Worst-Case SLBs

- ① **Explore:** Find coordinates with “large enough” magnitudes.
- ② **Commit:** Play “wisely” on these coordinates (ignore others).

**Example [Carpentier and Munos, 2012]:**

- ① **Explore:** Identify all  $i$  with  $|\theta_i^*| = \Omega((Ts/d)^{-1/4})$  (call this threshold  $\Delta$ ).

# Classical “Explore-then-Commit” Idea for Worst-Case SLBs

- ① **Explore:** Find coordinates with “large enough” magnitudes.
- ② **Commit:** Play “wisely” on these coordinates (ignore others).

**Example [Carpentier and Munos, 2012]:**

- ① **Explore:** Identify all  $i$  with  $|\theta_i^*| = \Omega((Ts/d)^{-1/4})$  (call this threshold  $\Delta$ ). Takes  $N = \tilde{\mathcal{O}}(\Delta^{-2}d) = \tilde{\mathcal{O}}(\sqrt{sdT})$  rounds to make the confidence radius  $\sqrt{d/n}$  smaller than  $\Delta/2$ .

# Classical “Explore-then-Commit” Idea for Worst-Case SLBs

- ① **Explore:** Find coordinates with “large enough” magnitudes.
- ② **Commit:** Play “wisely” on these coordinates (ignore others).

**Example [Carpentier and Munos, 2012]:**

- ① **Explore:** Identify all  $i$  with  $|\theta_i^*| = \Omega((Ts/d)^{-1/4})$  (call this threshold  $\Delta$ ). Takes  $N = \tilde{\mathcal{O}}(\Delta^{-2}d) = \tilde{\mathcal{O}}(\sqrt{sdT})$  rounds to make the confidence radius  $\sqrt{d/n}$  smaller than  $\Delta/2$ .
- ② **Commit:** For the remaining  $T - N$  rounds, execute a linear bandit algorithm on these coordinates (i.e., only consider an  $\mathcal{O}(s)$ -dimensional subspace) and play 0 on the other ones.

# Classical “Explore-then-Commit” Idea for Worst-Case SLBs

- ① **Explore:** Find coordinates with “large enough” magnitudes.
- ② **Commit:** Play “wisely” on these coordinates (ignore others).

**Example [Carpentier and Munos, 2012]:**

- ① **Explore:** Identify all  $i$  with  $|\theta_i^*| = \Omega((Ts/d)^{-1/4})$  (call this threshold  $\Delta$ ). Takes  $N = \tilde{\mathcal{O}}(\Delta^{-2}d) = \tilde{\mathcal{O}}(\sqrt{sdT})$  rounds to make the confidence radius  $\sqrt{d/n}$  smaller than  $\Delta/2$ .
- ② **Commit:** For the remaining  $T - N$  rounds, execute a linear bandit algorithm on these coordinates (i.e., only consider an  $\mathcal{O}(s)$ -dimensional subspace) and play 0 on the other ones.

**Regret Analysis:** The regret  $\mathcal{R}_T = \tilde{\mathcal{O}}(\sqrt{sdT})$ , as:

- **Exploration** causes no more than  $N = \tilde{\mathcal{O}}(\sqrt{sdT})$  regret.
- **Linear bandit** on  $s$  coordinates has regret  $\tilde{\mathcal{O}}(s\sqrt{T})$ .
- Each **un-explored coordinate**  $i$  (which is “small”) incurs regret  $\leq T\Delta^2 = \sqrt{dT/s}$ . There are no more than  $s$  such  $i$ ’s.

# Generalize to Variance-Aware SLB?

**Question 1: How to get  $\sqrt{\sum \sigma_t^2}$ -style regret in “commit”?**

# Generalize to Variance-Aware SLB?

**Question 1: How to get  $\sqrt{\sum \sigma_t^2}$ -style regret in “commit”?**

- **Answer:** Use variance-aware LB algorithms instead.

# Generalize to Variance-Aware SLB?

**Question 1: How to get  $\sqrt{\sum \sigma_t^2}$ -style regret in “commit”?**

- **Answer:** Use variance-aware LB algorithms instead.

**Question 2: How to get  $\sqrt{\sum \sigma_t^2}$ -style regret in “explore”?**

# Generalize to Variance-Aware SLB?

**Question 1: How to get  $\sqrt{\sum \sigma_t^2}$ -style regret in “commit”?**

- **Answer:** Use variance-aware LB algorithms instead.

**Question 2: How to get  $\sqrt{\sum \sigma_t^2}$ -style regret in “explore”?**

- ① **Worst-Case:** Exploration threshold  $\Delta \sim T^{-1/4}$ .
- ② **Deterministic-Case:** Exploration threshold  $\Delta \sim 0$ .

# Generalize to Variance-Aware SLB?

**Question 1: How to get  $\sqrt{\sum \sigma_t^2}$ -style regret in “commit”?**

- **Answer:** Use variance-aware LB algorithms instead.

**Question 2: How to get  $\sqrt{\sum \sigma_t^2}$ -style regret in “explore”?**

- ① **Worst-Case:** Exploration threshold  $\Delta \sim T^{-1/4}$ .
  - ② **Deterministic-Case:** Exploration threshold  $\Delta \sim 0$ .
- **Answer:** Decide the “threshold”  $\Delta$  **adaptively!**

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
-

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 3: How to do exploration?

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 3: How to do exploration?

- If explore on all coordinates, what's the point of halving?

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 3: How to do exploration?

- If explore on all coordinates, what's the point of halving?
- If ignore identified coordinates, they contribute large regret?

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 3: How to do exploration?

- If explore on all coordinates, what's the point of halving?
- If ignore identified coordinates, they contribute large regret?
- **Solution:** For those with magnitude  $> 2\Delta$  (i.e., identified before), always put their estimation  $\hat{\theta}_i$  on those coordinates.  
Use **remaining mass**  $1 - \sum \hat{\theta}_i^2$  to explore remaining ones.

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

**Question 4: When to stop exploration?**

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 4: When to stop exploration?

- Calculate confidence radius using Chernoff / Bernstein?

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 4: When to stop exploration?

- Calculate confidence radius using Chernoff / Bernstein?
- The confidence radius  $\frac{1}{n} \sqrt{d \sum_{k=1}^n \sigma_k^2}$  contains unknown  $\sigma_k$ 's!

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 4: When to stop exploration?

- Calculate confidence radius using Chernoff / Bernstein?
- The confidence radius  $\frac{1}{n} \sqrt{d \sum_{k=1}^n \sigma_k^2}$  contains unknown  $\sigma_k$ 's!
- Can we use some “empirical” observations to replace  $\sigma_k^2$ ?

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 4: When to stop exploration?

**Lemma:** For common-mean, independent & symmetric  $\{X_i\}_{i=1}^n$ ,

$$|\bar{X} - \mu| \leq \frac{1}{n} \sqrt{2 \sum_{i=1}^n (X_i - \bar{X})^2 \ln \frac{4}{\delta}} \quad \text{w.p. } 1 - \delta,$$

where  $n < \infty$  is stopping time,  $\mu = \mathbb{E}[X_i]$ , and  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ .

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

**Question 5: When to stop commit?**

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 5: When to stop commit?

- Recall that we need  $\hat{\theta}_i$  for all identified  $i$  in “explore”?

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 5: When to stop commit?

- Recall that we need  $\hat{\theta}_i$  for all identified  $i$  in “explore”?
- Recall LB algorithm can give  $\mathcal{O}(\sqrt{T})$  regret: the longer it executes, the closer the action  $a_t$  and the parameter  $\theta^*$  is.

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:     **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:     **Commit:** Deploy VA LB on all identified coordinates.
  - 4:     **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 5: When to stop commit?

- Recall that we need  $\hat{\theta}_i$  for all identified  $i$  in “explore”?
- Recall LB algorithm can give  $\mathcal{O}(\sqrt{T})$  regret: the longer it executes, the closer the action  $a_t$  and the parameter  $\theta^*$  is.
- Can we use this property for an accurate enough estimation?

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:   **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:   **Commit:** Deploy VA LB on all identified coordinates.
  - 4:   **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 5: When to stop commit?

“Regret-to-Sample-Complexity”: if its per-round regret  $< \Delta^2$ , i.e.,

$$\mathcal{R}_n^{\text{LB}} = \sum_{k=1}^n \langle \theta^* - a_k, \theta^* \rangle \leq n\Delta^2 \implies \hat{\theta} \triangleq \frac{1}{n} \sum_{k=1}^n a_k \text{ satisfies } \langle \theta^* - \hat{\theta}, \theta^* \rangle \leq \Delta^2.$$

# Our Idea: “Adaptive” Exploration Threshold

---

## Algorithm “Explore-then-Commit” with Adaptive Threshold

---

- 1: **for**  $\Delta = \frac{1}{2}, \dots$  **do**
  - 2:   **Explore:** Identify all coordinates with magnitude  $[\Delta, 2\Delta]$ .
  - 3:   **Commit:** Deploy VA LB on all identified coordinates.
  - 4:   **Continue:** Halve  $\Delta$  and repeat.
- 

### Question 5: When to stop commit?

“Regret-to-Sample-Complexity”: if its per-round regret  $< \Delta^2$ , i.e.,

$$\mathcal{R}_n^{\text{LB}} = \sum_{k=1}^n \langle \theta^* - a_k, \theta^* \rangle \leq n\Delta^2 \implies \hat{\theta} \triangleq \frac{1}{n} \sum_{k=1}^n a_k \text{ satisfies } \langle \theta^* - \hat{\theta}, \theta^* \rangle \leq \Delta^2.$$

So waiting until  $\mathcal{R}_n^{\text{LB}} \leq n\Delta^2$  can give “good” estimation  $\hat{\theta}$ !

# Final Algorithm

---

**Algorithm** Final Algorithm Using VA LB Algorithm  $\mathcal{A}$ 

---

1: **for**  $\Delta = \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots$  (i.e., halve until  $T$  rounds) **do**

# Final Algorithm

---

**Algorithm** Final Algorithm Using VA LB Algorithm  $\mathcal{A}$ 

---

- 1: **for**  $\Delta = \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots$  (i.e., halve until  $T$  rounds) **do**
- 2:     For each round, put  $\hat{\theta}_i$  on  $i$  for all identified  $i$ , and use remaining mass to explore like [Carpentier and Munos, 2012].

# Final Algorithm

---

**Algorithm** Final Algorithm Using VA LB Algorithm  $\mathcal{A}$ 

---

- 1: **for**  $\Delta = \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots$  (i.e., halve until  $T$  rounds) **do**
- 2:     For each round, put  $\hat{\theta}_i$  on  $i$  for all identified  $i$ , and use remaining mass to explore like [Carpentier and Munos, 2012].
- 3:     Terminate until the “explore” rounds  $n_{\Delta}^b$  satisfies

$$2\sqrt{2\sum_{k=1}^{n_{\Delta}^b}(r_{k,i} - \bar{r}_i)^2 \ln \frac{4}{\delta}} < n_{\Delta}^b \cdot \frac{\Delta}{4}, \quad \forall i \text{ unidentified},$$

where  $r_{k,i}$  is the  $k$ -th estimate of  $\theta_i^*$  and  $\bar{r}_i$  is the average of all  $r_{k,i}$ 's. Then mark all  $i$  with  $|\bar{r}_i| > \Delta$  as “identified”.

# Final Algorithm

---

**Algorithm** Final Algorithm Using VA LB Algorithm  $\mathcal{A}$ 

---

- 1: **for**  $\Delta = \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots$  (i.e., halve until  $T$  rounds) **do**
- 2:     For each round, put  $\hat{\theta}_i$  on  $i$  for all identified  $i$ , and use remaining mass to explore like [Carpentier and Munos, 2012].
- 3:     Terminate until the “explore” rounds  $n_{\Delta}^b$  satisfies

$$2\sqrt{2\sum_{k=1}^{n_{\Delta}^b}(r_{k,i} - \bar{r}_i)^2 \ln \frac{4}{\delta}} < n_{\Delta}^b \cdot \frac{\Delta}{4}, \quad \forall i \text{ unidentified},$$

where  $r_{k,i}$  is the  $k$ -th estimate of  $\theta_i^*$  and  $\bar{r}_i$  is the average of all  $r_{k,i}$ 's. Then mark all  $i$  with  $|\bar{r}_i| > \Delta$  as “identified”.

- 4:     Deploy  $\mathcal{A}$  on all identified coordinates until the “commit” rounds  $n_{\Delta}^a$  satisfies  $\mathcal{R}_{n_{\Delta}^a}^{\text{LB}} < n_{\Delta}^a \cdot \Delta^2$ . Calculate  $\hat{\theta}_i$  for all identified coordinates from “Regret-to-Sample-Complexity” conversion.
-

# Analysis Sketch

**Recap:** For each  $\Delta$ ,  $n_{\Delta}^b$  and  $n_{\Delta}^a$  are defined as (ignore constants)

$$n_{\Delta}^b \approx \Delta^{-1} \sqrt{\sum_{k=1}^{n_{\Delta}^b} (r_{k,i} - \bar{r}_i)^2}, \quad n_{\Delta}^a \approx \Delta^{-2} \mathcal{R}_{n_{\Delta}^a}^{\text{LB}}.$$

# Analysis Sketch

**Recap:** For each  $\Delta$ ,  $n_{\Delta}^b$  and  $n_{\Delta}^a$  are defined as (ignore constants)

$$n_{\Delta}^b \approx \Delta^{-1} \sqrt{\sum_{k=1}^{n_{\Delta}^b} (r_{k,i} - \bar{r}_i)^2}, \quad n_{\Delta}^a \approx \Delta^{-2} \mathcal{R}_{n_{\Delta}^a}^{\text{LB}}.$$

## ① “Explore” Regret:

- ① Identified ones contribute regret  $n_{\Delta}^b \langle \theta^* - \hat{\theta}, \theta^* \rangle \leq n_{\Delta}^b \cdot \Delta^2$ .

# Analysis Sketch

**Recap:** For each  $\Delta$ ,  $n_{\Delta}^b$  and  $n_{\Delta}^a$  are defined as (ignore constants)

$$n_{\Delta}^b \approx \Delta^{-1} \sqrt{\sum_{k=1}^{n_{\Delta}^b} (r_{k,i} - \bar{r}_i)^2}, \quad n_{\Delta}^a \approx \Delta^{-2} \mathcal{R}_{n_{\Delta}^a}^{\text{LB}}.$$

## ① “Explore” Regret:

- ① Identified ones contribute regret  $n_{\Delta}^b \langle \theta^* - \hat{\theta}, \theta^* \rangle \leq n_{\Delta}^b \cdot \Delta^2$ .
- ② Unidentified ones contribute regret  $n_{\Delta}^b \sum_i (\theta_i^*)^2 \leq n_{\Delta}^b \cdot s \Delta^2$ .

# Analysis Sketch

**Recap:** For each  $\Delta$ ,  $n_{\Delta}^b$  and  $n_{\Delta}^a$  are defined as (ignore constants)

$$n_{\Delta}^b \approx \Delta^{-1} \sqrt{\sum_{k=1}^{n_{\Delta}^b} (r_{k,i} - \bar{r}_i)^2}, \quad n_{\Delta}^a \approx \Delta^{-2} \mathcal{R}_{n_{\Delta}^a}^{\text{LB}}.$$

## ① “Explore” Regret:

- ① Identified ones contribute regret  $n_{\Delta}^b \langle \theta^* - \hat{\theta}, \theta^* \rangle \leq n_{\Delta}^b \cdot \Delta^2$ .
- ② Unidentified ones contribute regret  $n_{\Delta}^b \sum_i (\theta_i^*)^2 \leq n_{\Delta}^b \cdot s \Delta^2$ .

## ② “Commit” Regret:

- ① Identified ones contribute regret  $\mathcal{R}_{n_{\Delta}^a}^{\text{LB}} < n_{\Delta}^a \cdot \Delta^2$ .

# Analysis Sketch

**Recap:** For each  $\Delta$ ,  $n_{\Delta}^b$  and  $n_{\Delta}^a$  are defined as (ignore constants)

$$n_{\Delta}^b \approx \Delta^{-1} \sqrt{\sum_{k=1}^{n_{\Delta}^b} (r_{k,i} - \bar{r}_i)^2}, \quad n_{\Delta}^a \approx \Delta^{-2} \mathcal{R}_{n_{\Delta}^a}^{\text{LB}}.$$

## ① “Explore” Regret:

- ① Identified ones contribute regret  $n_{\Delta}^b \langle \theta^* - \hat{\theta}, \theta^* \rangle \leq n_{\Delta}^b \cdot \Delta^2$ .
- ② Unidentified ones contribute regret  $n_{\Delta}^b \sum_i (\theta_i^*)^2 \leq n_{\Delta}^b \cdot s \Delta^2$ .

## ② “Commit” Regret:

- ① Identified ones contribute regret  $\mathcal{R}_{n_{\Delta}^a}^{\text{LB}} < n_{\Delta}^a \cdot \Delta^2$ .
- ② Unidentified ones contribute regret  $n_{\Delta}^a \sum_i (\theta_i^*)^2 \leq n_{\Delta}^a \cdot s \Delta^2$ .

# Analysis Sketch

**Recap:** For each  $\Delta$ ,  $n_\Delta^b$  and  $n_\Delta^a$  are defined as (ignore constants)

$$n_\Delta^b \approx \Delta^{-1} \sqrt{\sum_{k=1}^{n_\Delta^b} (r_{k,i} - \bar{r}_i)^2}, \quad n_\Delta^a \approx \Delta^{-2} \mathcal{R}_{n_\Delta^a}^{\text{LB}}.$$

## ① “Explore” Regret:

- ① Identified ones contribute regret  $n_\Delta^b \langle \theta^* - \hat{\theta}, \theta^* \rangle \leq n_\Delta^b \cdot \Delta^2$ .
- ② Unidentified ones contribute regret  $n_\Delta^b \sum_i (\theta_i^*)^2 \leq n_\Delta^b \cdot s\Delta^2$ .

## ② “Commit” Regret:

- ① Identified ones contribute regret  $\mathcal{R}_{n_\Delta^a}^{\text{LB}} < n_\Delta^a \cdot \Delta^2$ .
- ② Unidentified ones contribute regret  $n_\Delta^a \sum_i (\theta_i^*)^2 \leq n_\Delta^a \cdot s\Delta^2$ .

## ③ Conclusion: Total Regret is $\mathcal{O}(\mathbb{E}[\sum_\Delta s\Delta^2(n_\Delta^b + n_\Delta^a)])$ .

# Analysis Sketch (Cont'd)

Omitting expectations,  $\mathcal{R}_T = \mathcal{O}(\sum_{\Delta} s\Delta^2(n_{\Delta}^b + n_{\Delta}^a))$ , i.e.,

$$\mathcal{R}_T = \widetilde{\mathcal{O}}(s) \sum_{\Delta} \Delta^2 \left( \frac{1}{\Delta} \sqrt{\sum_{k=1}^{n_{\Delta}^b} (r_{k,i} - \bar{r}_i)^2} + \Delta^{-2} \mathcal{R}_{n_{\Delta}^a}^{\text{LB}} \right).$$

# Analysis Sketch (Cont'd)

Omitting expectations,  $\mathcal{R}_T = \mathcal{O}(\sum_{\Delta} s\Delta^2(n_{\Delta}^b + n_{\Delta}^a))$ , i.e.,

$$\mathcal{R}_T = \widetilde{\mathcal{O}}(s) \sum_{\Delta} \Delta^2 \left( \frac{1}{\Delta} \sqrt{\sum_{k=1}^{n_{\Delta}^b} (r_{k,i} - \bar{r}_i)^2} + \Delta^{-2} \mathcal{R}_{n_{\Delta}^a}^{\text{LB}} \right).$$

Plugging in  $\mathcal{R}_n^{\text{LB}} = \widetilde{\mathcal{O}}\left(s^{1.5} \sqrt{\sum_{k=1}^{n_{\Delta}^a} \sigma_k^2} + s^2\right)$  [Kim et al., 2022]

and  $\sum_{k=1}^{n_{\Delta}^b} (r_{k,i} - \bar{r}_i)^2 \approx \sum_{k=1}^{n_{\Delta}^b} \mathbb{E}[(r_{k,i} - \bar{r}_i)^2] = \sum_{k=1}^{n_{\Delta}^b} (1 + \frac{d}{\Delta^2} \sigma_k^2)$ ,

# Analysis Sketch (Cont'd)

Omitting expectations,  $\mathcal{R}_T = \mathcal{O}(\sum_{\Delta} s\Delta^2(n_{\Delta}^b + n_{\Delta}^a))$ , i.e.,

$$\mathcal{R}_T = \widetilde{\mathcal{O}}(s) \sum_{\Delta} \Delta^2 \left( \frac{1}{\Delta} \sqrt{\sum_{k=1}^{n_{\Delta}^b} (r_{k,i} - \bar{r}_i)^2} + \Delta^{-2} \mathcal{R}_{n_{\Delta}^a}^{\text{LB}} \right).$$

Plugging in  $\mathcal{R}_n^{\text{LB}} = \widetilde{\mathcal{O}}\left(s^{1.5} \sqrt{\sum_{k=1}^{n_{\Delta}^a} \sigma_k^2} + s^2\right)$  [Kim et al., 2022]

and  $\sum_{k=1}^{n_{\Delta}^b} (r_{k,i} - \bar{r}_i)^2 \approx \sum_{k=1}^{n_{\Delta}^b} \mathbb{E}[(r_{k,i} - \bar{r}_i)^2] = \sum_{k=1}^{n_{\Delta}^b} (1 + \frac{d}{\Delta^2} \sigma_k^2)$ ,

$$\mathcal{R}_T = \widetilde{\mathcal{O}}(s) \sum_{\Delta} \left( \sqrt{\sum_{k=1}^{n_{\Delta}^b} (\Delta^2 + d\sigma_k^2)} + s^{1.5} \sqrt{\sum_{k=1}^{n_{\Delta}^a} \sigma_k^2 + s^2} \right)$$

$$= \widetilde{\mathcal{O}}\left((s^{2.5} + s\sqrt{d}) \sqrt{\sum_{t=1}^T \sigma_t^2} + s^3\right).$$

□

*Thank you for listening!*

Questions are more than welcomed.

# References

-  Abbasi-Yadkori, Y., Pal, D., and Szepesvari, C. (2012).  
Online-to-confidence-set conversions and application to sparse stochastic bandits.  
In *Artificial Intelligence and Statistics*, pages 1–9. PMLR.
-  Antos, A. and Szepesvári, C. (2009).  
Stochastic bandits with large action sets revisited.  
Personal communication.
-  Carpentier, A. and Munos, R. (2012).  
Bandit theory meets compressed sensing for high dimensional stochastic linear bandit.  
In *Artificial Intelligence and Statistics*, pages 190–198. PMLR.
-  Dani, V., Hayes, T. P., and Kakade, S. M. (2008).  
Stochastic linear optimization under bandit feedback.  
In *21st Annual Conference on Learning Theory*, pages 355–366.
-  Kim, Y., Yang, I., and Jun, K.-S. (2022).  
Improved regret analysis for variance-adaptive linear bandits and horizon-free linear mixture mdps.  
In *Advances in Neural Information Processing Systems 35*.
-  Zhou, D., Gu, Q., and Szepesvari, C. (2021).  
Nearly minimax optimal reinforcement learning for linear mixture markov decision processes.  
In *Conference on Learning Theory*, pages 4532–4576. PMLR.