

Toward Better Use of Data in Contextual and Linear Bandits

Nima Hamidi and Mohsen Bayati

Stanford University

October 2, 2020

References: [arXiv 2002.05152](#) & [arXiv 2006.06790](#)

Overview

1 Motivation

2 Confidence-based Policies

3 Sieved-Greedy

How to Test New Medical Interventions?

- A hospital wants to reduce post-discharge complications:
 - Use one of two newly designed telehealth (A or B)
- Should select one of A or B per patient
- A/B test or Randomized Control Trial (RCT) have high opportunity cost
 - In healthcare, experimentation is costly or unethical¹



A



B

¹Sibbald, Bonnie. 1998. Understanding controlled trials: Why are randomized controlled trials important?, British Medical Journal (Clinical Research Ed.) 316(201).

Beyond Healthcare

“Today, Microsoft and several other leading companies, including Amazon, Booking.com, Facebook, and Google, each conduct more than 10,000 online controlled experiments annually, with many tests engaging millions of users.”

Kohavi and Thompke, Harvard Business Review, 2017

Beyond Healthcare

“Today, Microsoft and several other leading companies, including Amazon, Booking.com, Facebook, and Google, each conduct more than 10,000 online controlled experiments annually, with many tests engaging millions of users.”

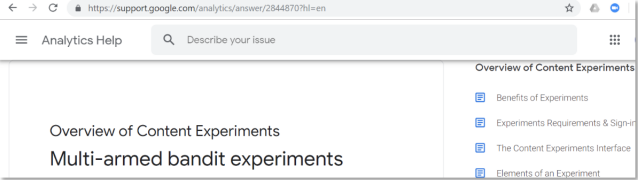
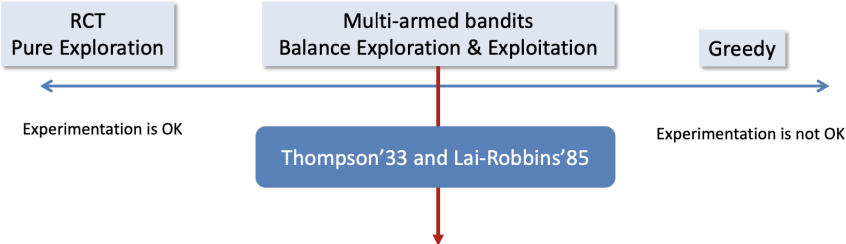
Kohavi and Thompke, Harvard Business Review, 2017

Also,



The screenshot shows a web browser displaying a news article from The Guardian. The page header includes navigation links for 'News', 'Opinion', 'Sport', 'Culture', 'Lifestyle', and 'More'. The main headline reads 'Facebook emotion study breached ethical guidelines, researchers say'. Below the headline, a sub-headline states: 'Lack of 'informed consent' means that Facebook experiment on nearly 700,000 news feeds broke rules on tests on human subjects, say scientists'. A poll question is visible: 'Poll: Facebook's secret mood experiment: have you lost trust in the social network?'. The author's name, Charles Arthur, is listed with a profile picture and social media links. A small image of a Facebook sign-in screen is partially visible at the bottom of the article. On the right side, there is a dark blue advertisement for 'MODE' with the text 'Start with SQL. Go anywhere.' and a 'Try for Free' button.

Multi-armed Bandit Experiments



Example (Google Analytics)²

- A/B testing
 - Website configurations A and B with conversion rates 4% and 5% respectively
- Using **Thompson Sampling**, instead of A/B testing, can run experiment with 78.5% less data → 97.5 conversions saved (on avg.)

²Source: Google Analytics Support Page

Stochastic Linear Bandit Problem

- Let $\Theta^* \in \mathbb{R}^d$ be fixed (and unknown).
- At time t , the **action set** $\mathcal{A}_t \subseteq \mathbb{R}^d$ is revealed to a **policy** π .
- The policy chooses $\tilde{A}_t \in \mathcal{A}_t$.
- It observes a **reward** $r_t = \langle \Theta^*, \tilde{A}_t \rangle + \varepsilon_t$.
- Conditional on the history, ε_t has zero mean.

Evaluation Metric

- The objective is to **improve using past experiences**.
- The **cumulative regret** is defined as

$$\text{Regret}(T, \Theta^*, \pi) := \mathbb{E} \left[\sum_{i=1}^T \sup_{A \in \mathcal{A}_t} \langle \Theta^*, A \rangle - \langle \Theta^*, \tilde{A}_t \rangle \mid \Theta^* \right].$$

Evaluation Metric

- The objective is to **improve using past experiences**.
- The **cumulative regret** is defined as

$$\text{Regret}(T, \Theta^*, \pi) := \mathbb{E} \left[\sum_{i=1}^T \sup_{A \in \mathcal{A}_t} \langle \Theta^*, A \rangle - \langle \Theta^*, \tilde{A}_t \rangle \mid \Theta^* \right].$$

- In the Bayesian setting, the **Bayesian regret** is given by

$$\text{BayesRegret}(T, \pi) := \mathbb{E}_{\Theta^* \sim \mathcal{P}} [\text{Regret}(T, \Theta^*, \pi)].$$

Special Cases

- Standard multi-armed bandit problem

Special Cases

- Standard multi-armed bandit problem
- k -armed contextual bandit problem

Special Cases

- Standard multi-armed bandit problem
- k -armed contextual bandit problem
- Dynamic-pricing with demand covariates

$$\text{Expected Demand} = \alpha + \beta p + \langle \Gamma, X \rangle$$

$$\text{Expected Revenue} = \alpha p + \beta p^2 + \langle \Gamma, X \rangle p$$

Special Cases

- Standard multi-armed bandit problem
- k -armed contextual bandit problem
- Dynamic-pricing with demand covariates

$$\text{Expected Demand} = \alpha + \beta p + \langle \Gamma, X \rangle$$

$$\text{Expected Revenue} = \alpha p + \beta p^2 + \langle \Gamma, X \rangle p$$

can be mapped to a linear bandit by setting

$$\mathcal{A} = \{(p, p^2, pX) \mid p \in [p_{\min}, p_{\max}]\} \text{ and } \Theta^* = \begin{bmatrix} \alpha \\ \beta \\ \Gamma \end{bmatrix}$$

Related Literature

- **UCB/OFUL:** Auer, Cesa-Bianchi, and Fischer 2002; Dani, Hayes, and Kakade 2008; Rusmevichientong and Tsitsiklis 2010; Abbasi-Yadkori, Pál, and Szepesvári 2011
- **Thompson sampling:** Agrawal and Goyal 2013; Russo and Van Roy 2014, 2016; Abeille and Lazaric 2017
- **ϵ -Greedy and variants:** Langford and Zhang 2008; Goldenshluger and Zeevi 2013

Related Literature

- **UCB/OFUL:** Auer, Cesa-Bianchi, and Fischer 2002; Dani, Hayes, and Kakade 2008; Rusmevichientong and Tsitsiklis 2010; Abbasi-Yadkori, Pál, and Szepesvári 2011
- **Thompson sampling:** Agrawal and Goyal 2013; Russo and Van Roy 2014, 2016; Abeille and Lazaric 2017
- **ϵ -Greedy and variants:** Langford and Zhang 2008; Goldenshluger and Zeevi 2013

Learning and earning in operations: Carvalho and Puterman05, Araman and Caldentey09, Besbes and Zeevi0911, Harrison et al.12, den Boer and Zwart14-16, Keskin and Zeevi14-16, Gur et al.'14, Johnson et al.15, Chen et al.15, Cohen et al 16, Bayati and Bastani'15, Kallus and Udell16, Javanmard and Nazerzadeh16, Javanmard17, Elmachtoub et. al.'17, Ban and Keskin17, Cheung et al '18, Bastani et al.'19, **and many more!**

Algorithms

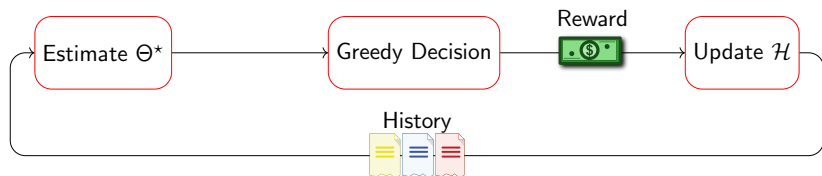
Greedy

At time $t = 1, 2, \dots, T$:

- Using the set of observations

$$\mathcal{H}_{t-1} := \{(\tilde{A}_1, r_1), \dots, (\tilde{A}_{t-1}, r_{t-1})\},$$

- Construct an **estimate** $\hat{\Theta}_{t-1}$ for Θ^* ,
- Choose the action $A \in \mathcal{A}_t$ with **largest** $\langle A, \hat{\Theta}_{t-1} \rangle$.



Greedy

The **ridge estimator** is used to obtain $\hat{\Theta}_t$ (for a fixed λ):

$$\mathbf{V}_t := \lambda \mathbf{I} + \sum_{i=1}^t \tilde{\mathbf{A}}_i \tilde{\mathbf{A}}_i^\top \in \mathbb{R}^{d \times d}, \quad (1)$$

and

$$\hat{\Theta}_t := \mathbf{V}_t^{-1} \left(\sum_{i=1}^t \tilde{\mathbf{A}}_i r_i \right) \in \mathbb{R}^d. \quad (2)$$

Algorithm 1 Greedy algorithm

- 1: **for** $t = 1$ to T **do**
 - 2: Pull $\tilde{A}_t := \arg \max_{A \in \mathcal{A}_t} \langle A, \hat{\Theta}_{t-1} \rangle$
 - 3: Observe the reward r_t
 - 4: Compute $\mathbf{V}_t = \lambda \mathbf{I} + \sum_{i=1}^t \tilde{A}_i \tilde{A}_i^\top$
 - 5: Compute $\hat{\Theta}_t = \mathbf{V}_t^{-1} \left(\sum_{i=1}^t \tilde{A}_i r_i \right)$
 - 6: **end for**
-

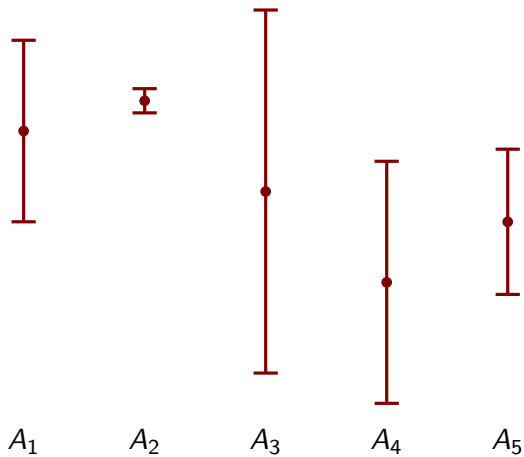
Algorithm 2 Greedy algorithm

- 1: **for** $t = 1$ to T **do**
 - 2: Pull $\tilde{A}_t := \arg \max_{A \in \mathcal{A}_t} \langle A, \hat{\Theta}_{t-1} \rangle$
 - 3: Observe the reward r_t
 - 4: Compute $\mathbf{V}_t = \lambda \mathbf{I} + \sum_{i=1}^t \tilde{A}_i \tilde{A}_i^\top$
 - 5: Compute $\hat{\Theta}_t = \mathbf{V}_t^{-1} \left(\sum_{i=1}^t \tilde{A}_i r_i \right)$
 - 6: **end for**
-

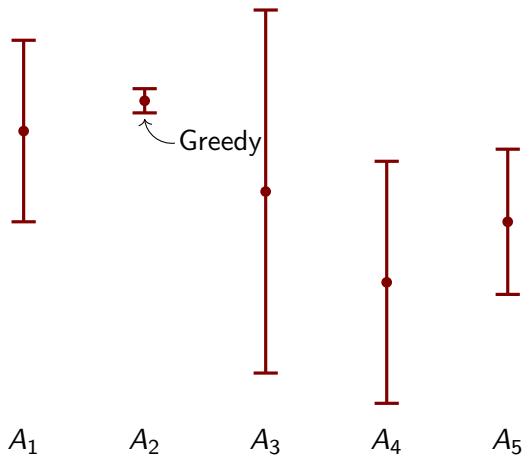
Greedy makes wrong decisions due to **over-** or **under-estimating** the true rewards.

- The over-estimation is **automatically** corrected.
- The under-estimation can cause **linear regret**.

Greedy



Greedy



Optimism in Face of Uncertainty (OFU) Algorithm

- Key idea: **be optimistic** when estimating the reward of actions.

Optimism in Face of Uncertainty (OFU) Algorithm

- Key idea: **be optimistic** when estimating the reward of actions.
- For $\rho > 0$, define the **confidence set** $\mathcal{C}_{t-1}(\rho)$ to be

$$\mathcal{C}_{t-1}(\rho) := \{\Theta \mid \|\Theta - \hat{\Theta}_{t-1}\|_{\mathbf{V}_{t-1}} \leq \rho\},$$

where

$$\|\mathbf{X}\|_{\mathbf{V}_{t-1}}^2 = \mathbf{X}^\top \mathbf{V}_{t-1} \mathbf{X} \in \mathbb{R}^+.$$

Optimism in Face of Uncertainty (OFU) Algorithm

- Key idea: **be optimistic** when estimating the reward of actions.
- For $\rho > 0$, define the **confidence set** $\mathcal{C}_{t-1}(\rho)$ to be

$$\mathcal{C}_{t-1}(\rho) := \{\Theta \mid \|\Theta - \hat{\Theta}_{t-1}\|_{\mathbf{V}_{t-1}} \leq \rho\},$$

where

$$\|X\|_{\mathbf{V}_{t-1}}^2 = X^\top \mathbf{V}_{t-1} X \in \mathbb{R}^+.$$

Theorem (Informal, Abbasi-Yadkori, Pál, and Szepesvári 2011)

Letting $\rho := \tilde{O}(\sqrt{d})$, we have $\Theta^* \in \mathcal{C}_{t-1}(\rho)$ with high probability.

Optimism in Face of Uncertainty (OFU) Algorithm

Algorithm 3 OFUL algorithm

- 1: **for** $t = 1$ to T **do**
 - 2: Pull $\tilde{A}_t := \arg \max_{A \in \mathcal{A}_t} \sup_{\Theta \in \mathcal{C}_{t-1}(\rho)} \langle A, \Theta \rangle$
 - 3: Observe the reward r_t
 - 4: Compute $\mathbf{V}_t = \lambda \mathbf{I} + \sum_{i=1}^t \tilde{A}_i \tilde{A}_i^\top$
 - 5: Compute $\hat{\Theta}_t = \mathbf{V}_t^{-1} \left(\sum_{i=1}^t \tilde{A}_i r_i \right)$
 - 6: **end for**
-

Optimism in Face of Uncertainty (OFU) Algorithm

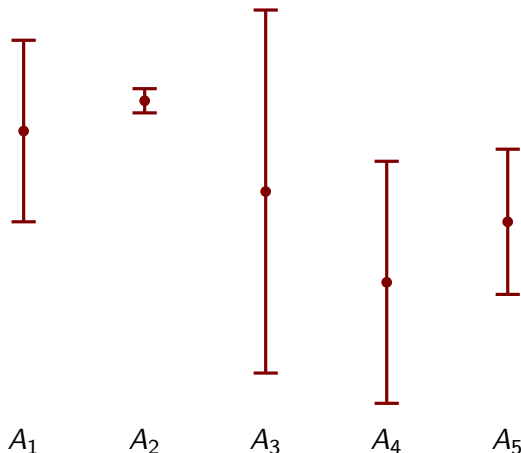
Algorithm 3 OFUL algorithm

- 1: **for** $t = 1$ to T **do**
 - 2: Pull $\tilde{A}_t := \arg \max_{A \in \mathcal{A}_t} \sup_{\Theta \in \mathcal{C}_{t-1}(\rho)} \langle A, \Theta \rangle$
 - 3: Observe the reward r_t
 - 4: Compute $\mathbf{V}_t = \lambda \mathbf{I} + \sum_{i=1}^t \tilde{A}_i \tilde{A}_i^\top$
 - 5: Compute $\hat{\Theta}_t = \mathbf{V}_t^{-1} \left(\sum_{i=1}^t \tilde{A}_i r_i \right)$
 - 6: **end for**
-

It can be shown that

$$\sup_{\Theta \in \mathcal{C}_{t-1}(\rho)} \langle A, \Theta \rangle = \langle A, \hat{\Theta}_{t-1} \rangle + \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}.$$

Optimism in Face of Uncertainty (OFU) Algorithm



Linear Thompson Sampling (LinTS) Algorithm

- Key idea: use **randomization** to address under-estimation.

Linear Thompson Sampling (LinTS) Algorithm

- Key idea: use **randomization** to address under-estimation.
- LinTS samples from the **posterior** distribution of Θ^* .

Algorithm 4 LinTS algorithm

- 1: **for** $t = 1$ to T **do**
 - 2: Sample $\tilde{\Theta}_{t-1} \sim \mathbb{P}(\Theta^* \mid \mathcal{H}_{t-1})$
 - 3: Pull $A_t := \arg \max_{A \in \mathcal{A}_t} \langle A, \tilde{\Theta}_{t-1} \rangle$
 - 4: Observe the reward r_t
 - 5: Update $\mathcal{H}_t \leftarrow \mathcal{H}_{t-1} \cup \{(A_t, r_t)\}$
 - 6: **end for**
-

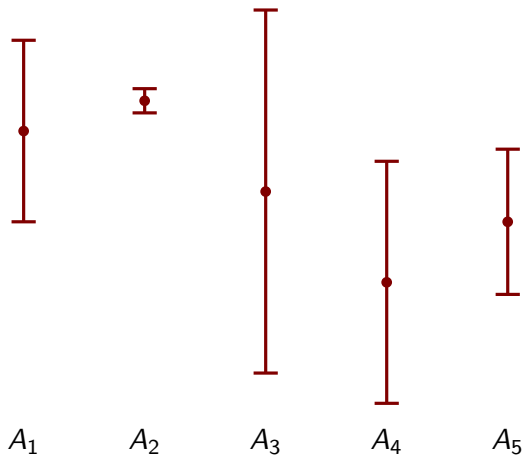
Linear Thompson Sampling (LinTS) Algorithm

- Under **normality**, LinTS becomes:

Algorithm 5 LinTS algorithm under normality

- 1: **for** $t = 1$ to T **do**
 - 2: Sample $\tilde{\Theta}_{t-1} \sim \mathcal{N}(\hat{\Theta}_{t-1}, \mathbf{V}_{t-1}^{-1})$
 - 3: Pull $A_t := \arg \max_{A \in \mathcal{A}_t} \langle A, \tilde{\Theta}_{t-1} \rangle$
 - 4: Observe the reward r_t
 - 5: Compute $\mathbf{V}_t = \lambda \mathbf{I} + \sum_{i=1}^t \tilde{A}_i \tilde{A}_i^\top$
 - 6: Compute $\hat{\Theta}_t = \mathbf{V}_t^{-1} \left(\sum_{i=1}^t \tilde{A}_i r_i \right)$
 - 7: **end for**
-

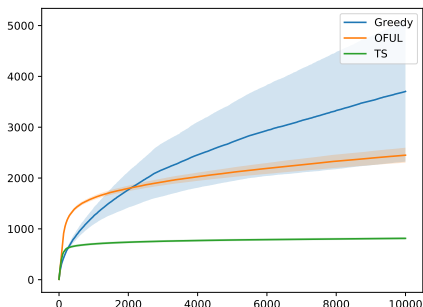
Linear Thompson Sampling (LinTS) Algorithm



Why Is LinTS Popular?

- **Empirical superiority:**

- $d = 120$, $\Theta^* \sim \mathcal{N}(0, \mathbf{I}_d)$,
- $k = 10$, $X \sim \mathcal{N}(0, \mathbf{I}_{12})$,
- Each A_t contains X as a block³.



³This is the 10-armed contextual bandit with 12 dimensional covariates.

Comparison of Regret Bounds

Theorem (Abbasi-Yadkori, Pál, and Szepesvári 2011)

Under some conditions, the regret of OFUL is bounded by

$$\text{Regret}(T, \Theta^*, \pi^{OFUL}) \leq \tilde{O}(d\sqrt{T}).$$

Comparison of Regret Bounds

Theorem (Abbasi-Yadkori, Pál, and Szepesvári 2011)

Under some conditions, the regret of OFUL is bounded by

$$\text{Regret}(T, \Theta^*, \pi^{\text{OFUL}}) \leq \tilde{O}(d\sqrt{T}).$$

Theorem (Russo and Van Roy 2014)

Under minor assumptions, the Bayesian regret of LinTS is bounded by

$$\text{BayesRegret}(T, \pi^{\text{LinTS}}) \leq \tilde{O}(d\sqrt{T}).$$

Comparison of Regret Bounds

Theorem (Abbasi-Yadkori, Pál, and Szepesvári 2011)

Under some conditions, the regret of OFUL is bounded by

$$\text{Regret}(T, \Theta^*, \pi^{\text{OFUL}}) \leq \tilde{O}(d\sqrt{T}).$$

Theorem (Russo and Van Roy 2014)

Under minor assumptions, the Bayesian regret of LinTS is bounded by

$$\text{BayesRegret}(T, \pi^{\text{LinTS}}) \leq \tilde{O}(d\sqrt{T}).$$

Theorem (Dani, Hayes, and Kakade 2008)

There is a Bayesian linear bandit problem that satisfies

$$\inf_{\pi} \text{BayesRegret}(T, \pi) \geq \Omega(d\sqrt{T}).$$

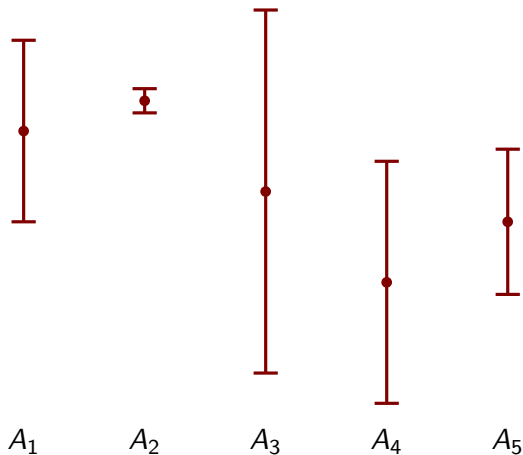
A Worst-Case Regret Bound for LinTS

- Question: can one prove a similar worst-case regret bound for LinTS?
- The only known results require **inflating** the posterior variance.

Algorithm 6 LinTS algorithm under normality

- 1: **for** $t = 1$ to T **do**
 - 2: Sample $\tilde{\Theta}_{t-1} \sim \mathcal{N}(\hat{\Theta}_{t-1}, \beta^2 \mathbf{V}_{t-1}^{-1})$
 - 3: Pull $A_t := \arg \max_{A \in \mathcal{A}_t} \langle A, \tilde{\Theta}_{t-1} \rangle$
 - 4: Update \mathbf{V}_t and $\hat{\Theta}_t$
 - 5: **end for**
-

Inflated Linear Thompson Sampling (LinTS) Algorithm



A Worst-Case Regret Bound for LinTS

Theorem (Agrawal and Goyal 2013; Abeille and Lazaric 2017)

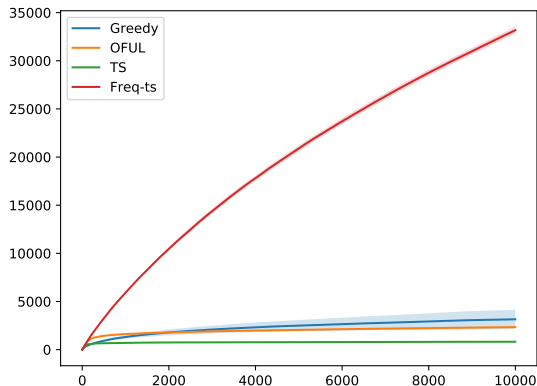
If $\beta \propto \sqrt{d}$, then

$$\text{Regret}(T, \Theta^*, \pi^{\text{LinTS}}) \leq \tilde{O}(d\sqrt{dT}).$$

This result is far from optimal by a \sqrt{d} factor.

Empirical Performance of Inflated LinTS

- Unfortunately, the inflated variant of LinTS performs poorly...



Bayesian Analyses are Brittle

We prove that the inflation is **necessary** for LinTS to work.

Theorem

There exists a linear bandit problem such that for $T \leq \exp(\Omega(d))$, we have

$$\text{BayesRegret}(T, \pi^{\text{LinTS}}) = \Omega(T).$$

Bayesian Analyses are Brittle

We prove that the inflation is **necessary** for LinTS to work.

Theorem

There exists a linear bandit problem such that for $T \leq \exp(\Omega(d))$, we have

$$\text{BayesRegret}(T, \pi^{\text{LinTS}}) = \Omega(T).$$

The counter-example satisfies the following properties:

- $\Theta^* \sim \mathcal{N}(0, \mathbf{I}_d)$,
- LinTS uses the right prior,
- LinTS assumes noises are standard normal,
- $r_t = \langle \Theta^*, A_t \rangle$. (i.e., **noiseless** data!)

Some Remarks on LinTS

- Under some assumptions on the action set we can prove that LinTS with only a logarithmic inflation works.
- Bastani, Simchi-Levi, and Zhu (2019) showed, in the dynamic pricing case, when only the prior mean is unknown but prior variance is known, regret loss can be only up to a constant.
- Jin, Xu, Shi, Xiao, and Gu (2020) showed a variant of TS achieves optimal regret bound in the standard k -armed bandit setting.

Improving OFUL

Weaker Optimism and Improving OFUL

- Unlike its name, OFUL is very **pessimistic**.
- It assumes that $\langle A, \hat{\Theta}_t \rangle$ is **as small as it can be** for all arms.
- In practice, this may not be the case.
- One may be able to **use data more efficiently**.
- This is challenging since even choosing **the second most optimistic action** can lead to **linear regret**.

Optimism Baseline

Let $L_t(A)$ be the **lower confidence-bound** for action A defined as

$$L_t(A) := \langle A, \hat{\Theta}_{t-1} \rangle - \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}.$$

Optimism Baseline

Let $L_t(A)$ be the **lower confidence-bound** for action A defined as

$$L_t(A) := \langle A, \hat{\Theta}_{t-1} \rangle - \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}.$$

A simple observation:

Lemma

The following inequality holds with high probability:

$$\sup_{A \in \mathcal{A}_t} \langle A, \Theta^* \rangle \geq \sup_{A \in \mathcal{A}_t} L_t(A).$$

Then, define the **baseline** as $B_t := \sup_{A \in \mathcal{A}_t} L_t(A)$.

Sieved-Greedy

Let $L_t(A)$ and $U_t(A)$ be the **confidence-bounds** for action A given by

$$U_t(A) := \langle A, \hat{\Theta}_{t-1} \rangle + \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}$$

$$L_t(A) := \langle A, \hat{\Theta}_{t-1} \rangle - \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}.$$

Sieved-Greedy

Let $L_t(A)$ and $U_t(A)$ be the **confidence-bounds** for action A given by

$$U_t(A) := \langle A, \hat{\Theta}_{t-1} \rangle + \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}$$

$$L_t(A) := \langle A, \hat{\Theta}_{t-1} \rangle - \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}.$$

For a fixed **sieving-rate** $\alpha \in (0, 1]$, define

$$\mathcal{A}'_t := \left\{ A \in \mathcal{A}_t : U_t(A) - B_t \geq \alpha \left(\sup_{A' \in \mathcal{A}_t} U_t(A') - B_t \right) \right\}.$$

Sieved-Greedy

Let $L_t(A)$ and $U_t(A)$ be the **confidence-bounds** for action A given by

$$U_t(A) := \langle A, \hat{\Theta}_{t-1} \rangle + \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}$$
$$L_t(A) := \langle A, \hat{\Theta}_{t-1} \rangle - \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}.$$

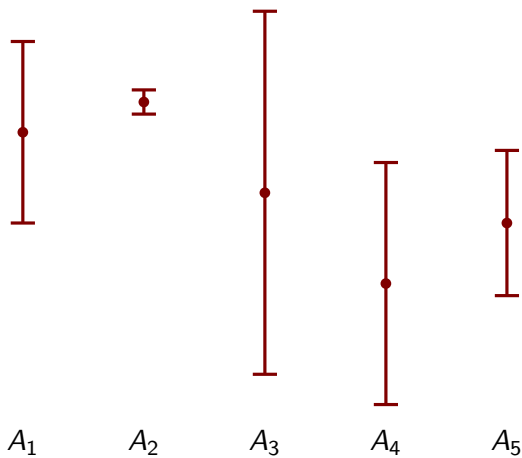
For a fixed **sieving-rate** $\alpha \in (0, 1]$, define

$$\mathcal{A}'_t := \left\{ A \in \mathcal{A}_t : U_t(A) - B_t \geq \alpha \left(\sup_{A' \in \mathcal{A}_t} U_t(A') - B_t \right) \right\}.$$

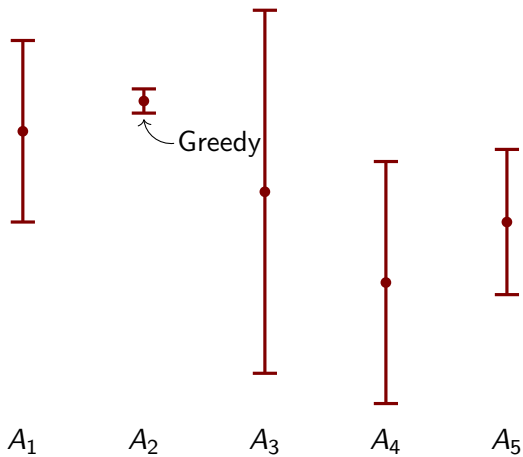
Sieved-greedy (SG) then chooses

$$\tilde{A}_t \in \arg \max_{A \in \mathcal{A}'_t} \langle A, \hat{\Theta}_{t-1} \rangle.$$

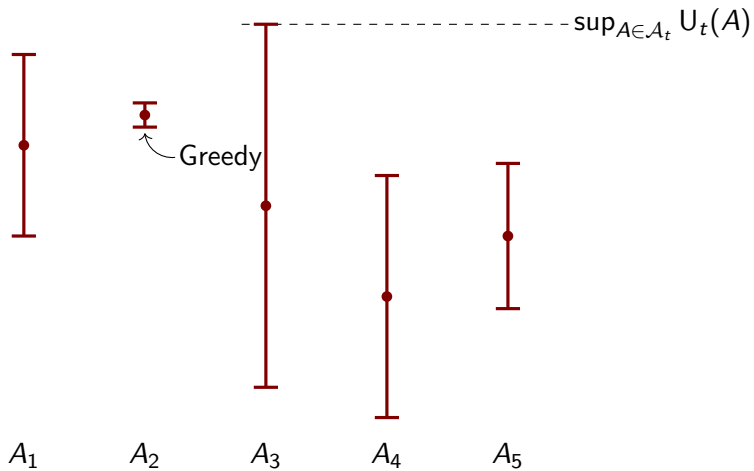
Sieved-Greedy



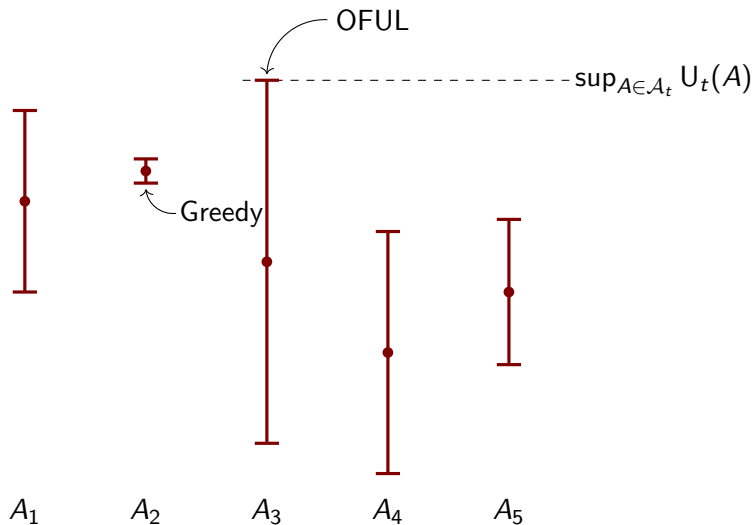
Sieved-Greedy



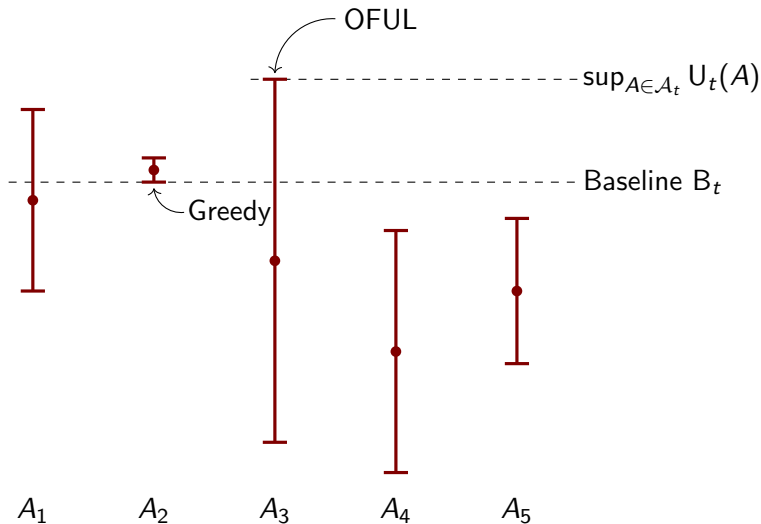
Sieved-Greedy



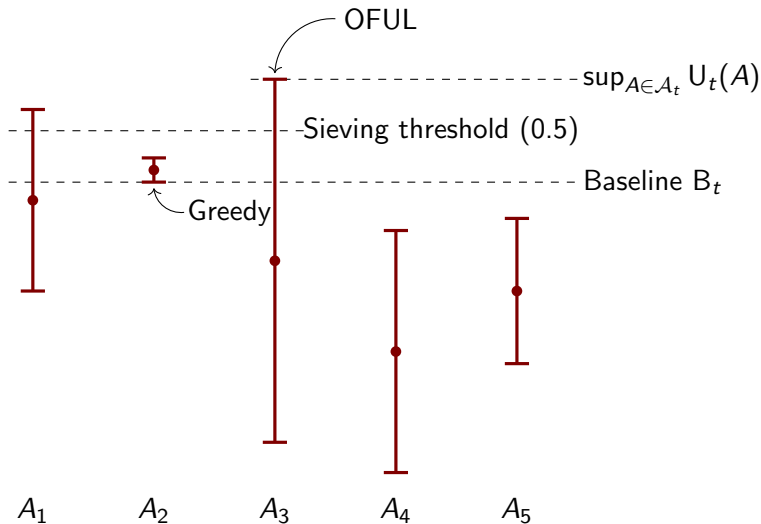
Sieved-Greedy



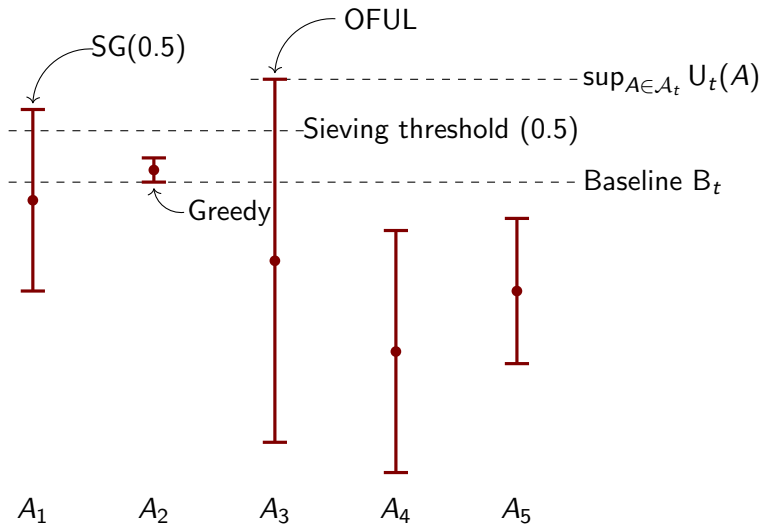
Sieved-Greedy



Sieved-Greedy



Sieved-Greedy



Sieved-Greedy

Two instances of Sieved-greedy:

- For $\alpha = 1$, Sieved-greedy is equivalent to OFUL.
- For $\alpha = 0$, Sieved-greedy is the same as Greedy.

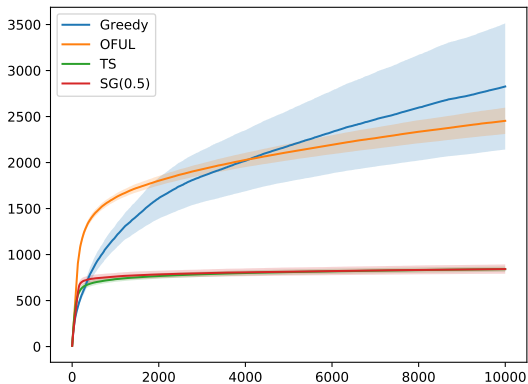
Sieved-Greedy

Theorem

For $\alpha > 0$, the regret of Sieved-greedy is bounded by

$$\text{Regret}(T, \Theta^*, \pi^{\text{SG}(\alpha)}) \leq \tilde{O}\left(\frac{d\sqrt{T}}{\alpha}\right).$$

Simulations



A General Regret Bound

- By a **worth function**, we mean a function \tilde{M}_t that maps each $A \in \mathcal{A}_t$ to \mathbb{R} such that

$$|\tilde{M}_t(A) - \langle A, \hat{\Theta}_{t-1} \rangle| \leq \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}$$

with probability at least $1 - \frac{1}{T^2}$.

A General Regret Bound

- By a **worth function**, we mean a function \tilde{M}_t that maps each $A \in \mathcal{A}_t$ to \mathbb{R} such that

$$|\tilde{M}_t(A) - \langle A, \hat{\Theta}_{t-1} \rangle| \leq \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}$$

with probability at least $1 - \frac{1}{T^2}$.

- Next, define **Randomized OFUL (ROFUL)** to be:

Algorithm 7 ROFUL algorithm

- 1: **for** $t = 1$ to T **do**
 - 2: Pull $\tilde{A}_t := \arg \max_{A \in \mathcal{A}_t} \tilde{M}_t(A)$
 - 3: Observe the reward r_t
 - 4: Compute $\mathbf{V}_t = \lambda \mathbf{I} + \sum_{i=1}^t \tilde{A}_i \tilde{A}_i^\top$
 - 5: Compute $\hat{\Theta}_t = \mathbf{V}_t^{-1} \left(\sum_{i=1}^t \tilde{A}_i r_i \right)$
 - 6: **end for**
-

A General Regret Bound

Examples of worth functions:

- Greedy: $\tilde{M}_t(A) = \langle A, \hat{\Theta}_{t-1} \rangle$
- OFUL: $\tilde{M}_t(A) = \langle A, \hat{\Theta}_{t-1} \rangle + \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}$
- LinTS: $\tilde{M}_t(A) = \langle A, \tilde{\Theta}_{t-1} \rangle$

A General Regret Bound

Examples of worth functions:

- Greedy: $\tilde{M}_t(A) = \langle A, \hat{\Theta}_{t-1} \rangle$
- OFUL: $\tilde{M}_t(A) = \langle A, \hat{\Theta}_{t-1} \rangle + \rho \|A\|_{\mathbf{V}_{t-1}^{-1}}$
- LinTS: $\tilde{M}_t(A) = \langle A, \tilde{\Theta}_{t-1} \rangle$
- Sieved-greedy:

$$\tilde{M}_t(A) = \begin{cases} U_t(A) & \text{if } A = \tilde{A}_t \\ L_t(A) & \text{otherwise} \end{cases}$$

Weaker Optimism and Regret Bound

Definition (Optimism in expectation)

We say a worth function \tilde{M}_t is **optimistic in expectation (OIE)** if

$$\mathbb{E} \left[\left(\sup_{A \in \mathcal{A}_t} \tilde{M}_t(A) - B_t \right)^2 \right] \geq p \mathbb{E} \left[\left(\sup_{A \in \mathcal{A}_t} \langle A, \Theta^* \rangle - B_t \right)^2 \right].$$

Weaker Optimism and Regret Bound

Definition (Optimism in expectation)

We say a worth function \tilde{M}_t is **optimistic in expectation (OIE)** if

$$\mathbb{E} \left[\left(\sup_{A \in \mathcal{A}_t} \tilde{M}_t(A) - B_t \right)^2 \right] \geq p \mathbb{E} \left[\left(\sup_{A \in \mathcal{A}_t} \langle A, \Theta^* \rangle - B_t \right)^2 \right].$$

Theorem

For a sequence of OIE worth functions $(\tilde{M}_t)_{t=1}^T$, we have

$$\text{Regret}(T, \pi^{\text{ROFUL}}) \leq \tilde{O} \left(\rho \sqrt{\frac{dT}{p}} \right).$$

Conclusion

- Proved that LinTS without inflation can incur linear regret.
- Provided a general regret bound for confidence-based policies.
- Introduced a new algorithm that is less pessimistic than OFUL while enjoying similar regret bounds.

References I



Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. “Finite-time analysis of the multiarmed bandit problem”. In: *Machine learning* 47.2-3 (2002), pp. 235–256.



Varsha Dani, Thomas P. Hayes, and Sham M. Kakade. “Stochastic Linear Optimization under Bandit Feedback”. In: *COLT*. 2008.



John Langford and Tong Zhang. “The Epoch-Greedy Algorithm for Multi-armed Bandits with Side Information”. In: *Advances in Neural Information Processing Systems 20*. Ed. by J. C. Platt et al. Curran Associates, Inc., 2008, pp. 817–824.



Paat Rusmevichientong and John N Tsitsiklis. “Linearly parameterized bandits”. In: *Mathematics of Operations Research* 35.2 (2010), pp. 395–411.

References II

-  Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. “Improved algorithms for linear stochastic bandits”. In: *Advances in Neural Information Processing Systems*. 2011, pp. 2312–2320.
-  Shipra Agrawal and Navin Goyal. “Thompson Sampling for Contextual Bandits with Linear Payoffs.”. In: *ICML (3)*. 2013, pp. 127–135.
-  Alexander Goldenshluger and Assaf Zeevi. “A linear response bandit problem”. In: *Stochastic Systems* 3.1 (2013), pp. 230–261.
-  Daniel Russo and Benjamin Van Roy. “Learning to Optimize via Posterior Sampling”. In: *Mathematics of Operations Research* 39.4 (2014), pp. 1221–1243. DOI: 10.1287/moor.2014.0650.
-  Daniel Russo and Benjamin Van Roy. “An information-theoretic analysis of thompson sampling”. In: *The Journal of Machine Learning Research* 17.1 (2016), pp. 2442–2471.

References III



Marc Abeille, Alessandro Lazaric, et al. “Linear Thompson sampling revisited”. In: *Electronic Journal of Statistics* 11.2 (2017), pp. 5165–5197.

Thank you!

Any questions?

A General Regret Bound

Definition (Strong optimism)

We say a worth function \tilde{M}_t is **optimistic** if

$$\sup_{A \in \mathcal{A}_t} \tilde{M}_t(A) \geq \sup_{A \in \mathcal{A}_t} \langle A, \Theta^* \rangle \quad (3)$$

with probability at least p .

A General Regret Bound

Definition (Strong optimism)

We say a worth function \tilde{M}_t is **optimistic** if

$$\sup_{A \in \mathcal{A}_t} \tilde{M}_t(A) \geq \sup_{A \in \mathcal{A}_t} \langle A, \Theta^* \rangle \quad (3)$$

with probability at least p .

Theorem

Let $(\tilde{M}_t)_{t=1}^T$ be a sequence of optimistic worth functions. Then, the regret of ROFUL with this worth function is bounded by

$$\text{Regret}(T, \pi^{\text{ROFUL}}) \leq \tilde{O} \left(\rho \sqrt{\frac{dT}{p}} \right).$$

A Sufficient Condition for Optimism

- Recall that the worth function for LinTS is given by

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t \rangle.$$

A Sufficient Condition for Optimism

- Recall that the worth function for LinTS is given by

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t \rangle.$$

- We can decompose it as

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t - \hat{\Theta}_{t-1} \rangle + \langle A, \hat{\Theta}_{t-1} - \Theta^* \rangle + \langle A, \Theta^* \rangle.$$

A Sufficient Condition for Optimism

- Recall that the worth function for LinTS is given by

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t \rangle.$$

- We can decompose it as

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t - \hat{\Theta}_{t-1} \rangle + \langle A, \hat{\Theta}_{t-1} - \Theta^* \rangle + \langle A, \Theta^* \rangle.$$

- Hence, we have

$$\sup_{A \in \mathcal{A}_t} \tilde{M}_t(A) - \sup_{A \in \mathcal{A}_t} \langle A, \Theta^* \rangle \geq \tilde{M}_t(A_t^*) - \langle A_t^*, \Theta^* \rangle$$

A Sufficient Condition for Optimism

- Recall that the worth function for LinTS is given by

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t \rangle.$$

- We can decompose it as

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t - \hat{\Theta}_{t-1} \rangle + \langle A, \hat{\Theta}_{t-1} - \Theta^* \rangle + \langle A, \Theta^* \rangle.$$

- Hence, we have

$$\begin{aligned} \sup_{A \in \mathcal{A}_t} \tilde{M}_t(A) - \sup_{A \in \mathcal{A}_t} \langle A, \Theta^* \rangle &\geq \tilde{M}_t(A_t^*) - \langle A_t^*, \Theta^* \rangle \\ &= \langle A_t^*, \tilde{\Theta}_t - \hat{\Theta}_{t-1} \rangle + \langle A_t^*, \hat{\Theta}_{t-1} - \Theta^* \rangle. \end{aligned}$$

A Sufficient Condition for Optimism

- Recall that the worth function for LinTS is given by

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t \rangle.$$

- We can decompose it as

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t - \hat{\Theta}_{t-1} \rangle + \langle A, \hat{\Theta}_{t-1} - \Theta^* \rangle + \langle A, \Theta^* \rangle.$$

- Hence, we have

$$\begin{aligned} \sup_{A \in \mathcal{A}_t} \tilde{M}_t(A) - \sup_{A \in \mathcal{A}_t} \langle A, \Theta^* \rangle &\geq \tilde{M}_t(A_t^*) - \langle A_t^*, \Theta^* \rangle \\ &= \langle A_t^*, \tilde{\Theta}_t - \hat{\Theta}_{t-1} \rangle + \underbrace{\langle A_t^*, \hat{\Theta}_{t-1} - \Theta^* \rangle}_{\text{Error term}}. \end{aligned}$$

A Sufficient Condition for Optimism

- Recall that the worth function for LinTS is given by

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t \rangle.$$

- We can decompose it as

$$\tilde{M}_t(A) = \langle A, \tilde{\Theta}_t - \hat{\Theta}_{t-1} \rangle + \langle A, \hat{\Theta}_{t-1} - \Theta^* \rangle + \langle A, \Theta^* \rangle.$$

- Hence, we have

$$\begin{aligned} \sup_{A \in \mathcal{A}_t} \tilde{M}_t(A) - \sup_{A \in \mathcal{A}_t} \langle A, \Theta^* \rangle &\geq \tilde{M}_t(A_t^*) - \langle A_t^*, \Theta^* \rangle \\ &= \underbrace{\langle A_t^*, \tilde{\Theta}_t - \hat{\Theta}_{t-1} \rangle}_{\text{Compensation term}} + \underbrace{\langle A_t^*, \hat{\Theta}_{t-1} - \Theta^* \rangle}_{\text{Error term}}. \end{aligned}$$

A Sufficient Condition for Optimism

Define

- Error vector $E := \Theta^* - \hat{\Theta}_{t-1}$
- Compensator vector $C := \tilde{\Theta}_t - \hat{\Theta}_{t-1}$

The optimism assumption holds if, with probability p , the following holds

$$\langle A_t^*, C \rangle \geq \langle A_t^*, E \rangle.$$

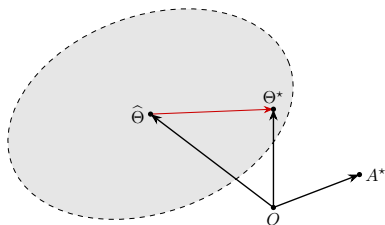
Omniscient Adversary and LinTS

- An **adversary** chooses \mathcal{A}_t at time t .
- The adversary is **omniscient** if he knows $\hat{\Theta}_{t-1}$ and Θ^* .

Omniscient Adversary and LinTS

- An **adversary** chooses \mathcal{A}_t at time t .
- The adversary is **omniscient** if he knows $\hat{\Theta}_{t-1}$ and Θ^* .
- He chooses $A = -c\hat{\Theta}_{t-1} + E$ so that

$$\langle A, \Theta^* \rangle > 0 \quad \text{and} \quad \langle A, \hat{\Theta}_{t-1} \rangle < -\frac{1}{2} \cdot \|A\|_{\mathbf{v}_{t-1}^{-1}} \cdot \underbrace{\|E\|_{\mathbf{v}_{t-1}}}_{\approx \sqrt{d}} \ll 0.$$



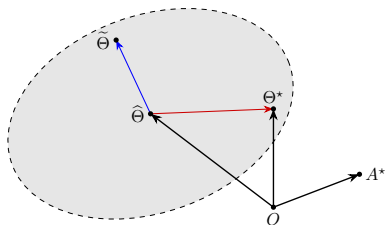
Omniscient Adversary and LinTS

- The adversary sets $\mathcal{A}_t = \{0, A\}$.
- LinTS chooses A if and only if

$$\langle A, \tilde{\Theta}_t \rangle = \langle A, \tilde{\Theta}_t - \hat{\Theta}_{t-1} \rangle + \langle A, \hat{\Theta}_{t-1} \rangle > 0.$$

- This requires

$$\langle A, C \rangle \sim \mathcal{N}(0, \mathbf{V}_{t-1}^{-1}) > \frac{1}{2} \cdot \|A\|_{\mathbf{V}_{t-1}^{-1}} \cdot \underbrace{\|E\|_{\mathbf{V}_{t-1}}}_{\approx \sqrt{d}}.$$



Omniscient Adversary and LinTS

- Next, we have

$$\mathbb{P}(\langle A, \tilde{\Theta}_t \rangle > 0) \leq \exp(-\Omega(d))!$$

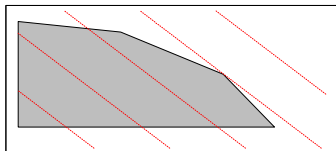
- LinTS chooses the optimal arm A w.p. **exponentially small in $\Omega(d)$** .
- When $\tilde{A}_t = 0$, the reward contains **no new information** about Θ^* .
- The adversary reveals the same action set in the next rounds.
- The regret will grow **linearly**.

Bayesian Analyses are Brittle

- The key point was the **adversary's knowledge of E** .
- This can be relaxed by **slightly modifying** the noise distribution.
- **Reducing the noise variance** reveals information about E .

Why is LinTS Popular?

- **Computation efficiency:** when \mathcal{A}_t is a polytope ...
 - LinTS solves an LP problem,



- OFUL becomes an NP-hard problem!

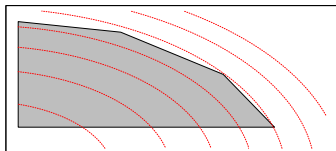


Photo credit: Russo and Van Roy 2014