Beyond UCB: statistical complexity and optimal algorithm for non-linear ridge bandits

Yanjun Han (MIT IDSS)

Joint work with:

Jiantao Jiao Nived Rajaraman Kannan Ramchandran Berkeley EECS Berkeley EECS Berkeley EECS

AMCS Colloquium, UPenn October 21, 2022

General setting of stochastic bandit

Input parameters:

- $\bullet\,$ parameter set $\Theta\,$
- \bullet action space ${\cal A}$
- reward function class $\mathcal{F} = (f_{ heta})_{ heta \in \Theta}$
- time horizon T

General setting of stochastic bandit

Input parameters:

- $\bullet\,$ parameter set $\Theta\,$
- \bullet action space ${\cal A}$
- reward function class $\mathcal{F} = (f_{ heta})_{ heta \in \Theta}$
- time horizon T

Stochastic bandit environment:

- nature chooses $\theta^{\star} \in \Theta$, fixed across time and unknown to the learner
- at time $t = 1, \cdots, T$, learner chooses action $a_t \in A$ and observes a random reward r_t with $\mathbb{E}[r_t \mid a_t = a] = f_{\theta^*}(a)$
- learner aims to minimize the worst-case (pseudo) regret

$$\mathsf{MinmaxReg}(\Theta, \mathcal{A}, \mathcal{F}, T) = \inf_{a^T} \sup_{\theta^{\star} \in \Theta} \mathbb{E}_{\theta^{\star}} \left[T \cdot \max_{a \in \mathcal{A}} f_{\theta^{\star}}(a) - \sum_{t=1}^T f_{\theta^{\star}}(a_t) \right].$$

General setting of stochastic bandit

Input parameters:

- $\bullet\,$ parameter set $\Theta\,$
- \bullet action space ${\cal A}$
- reward function class $\mathcal{F} = (f_{ heta})_{ heta \in \Theta}$
- time horizon T

Stochastic bandit environment:

- nature chooses $\theta^{\star} \in \Theta$, fixed across time and unknown to the learner
- at time $t = 1, \cdots, T$, learner chooses action $a_t \in A$ and observes a random reward r_t with $\mathbb{E}[r_t \mid a_t = a] = f_{\theta^*}(a)$
- learner aims to minimize the worst-case (pseudo) regret

$$\mathsf{MinmaxReg}(\Theta, \mathcal{A}, \mathcal{F}, T) = \inf_{a^T} \sup_{\theta^{\star} \in \Theta} \mathbb{E}_{\theta^{\star}} \left[T \cdot \max_{a \in \mathcal{A}} f_{\theta^{\star}}(a) - \sum_{t=1}^T f_{\theta^{\star}}(a_t) \right].$$

Linear bandit

 $f_{ heta}(a) = \langle heta, \phi(a)
angle$ with a known feature map $\phi: \mathcal{A} o \mathbb{R}^d$

$$f_{ heta}(a) = \langle heta, a
angle^3: \qquad heta \in \mathbb{S}^{d-1}, \quad a \in \mathbb{B}^d.$$



$$f_{ heta}(a) = \langle heta, a
angle^3: \qquad heta \in \mathbb{S}^{d-1}, \quad a \in \mathbb{B}^d.$$



$$f_{ heta}(a) = \langle heta, a
angle^3: \qquad heta \in \mathbb{S}^{d-1}, \quad a \in \mathbb{B}^d.$$



$$f_{ heta}(a) = \langle heta, a
angle^3: \qquad heta \in \mathbb{S}^{d-1}, \quad a \in \mathbb{B}^d.$$



$$f_{ heta}(a) = \langle heta, a
angle^3: \qquad heta \in \mathbb{S}^{d-1}, \quad a \in \mathbb{B}^d.$$



$$f_{ heta}(a) = \langle heta, a
angle^3: \qquad heta \in \mathbb{S}^{d-1}, \quad a \in \mathbb{B}^d.$$











Curious phenomena in non-linear bandits:

- phase transition in the regret
- initialization phase: regret grows linearly and results in a fixed cost
 - $\rightarrow\,$ find a good "initial action" to start learning
- learning phase: regret grows sublinearly and looks like a linear bandit
 - $\rightarrow\,$ bandit learning starts from the good initial action

Curious phenomena in non-linear bandits:

- phase transition in the regret
- initialization phase: regret grows linearly and results in a fixed cost
 - $\rightarrow\,$ find a good "initial action" to start learning
- learning phase: regret grows sublinearly and looks like a linear bandit
 - $\rightarrow\,$ bandit learning starts from the good initial action

Aim of this talk:



Curious phenomena in non-linear bandits:

- phase transition in the regret
- initialization phase: regret grows linearly and results in a fixed cost
 - $\rightarrow\,$ find a good "initial action" to start learning
- learning phase: regret grows sublinearly and looks like a linear bandit
 - $\rightarrow\,$ bandit learning starts from the good initial action

Aim of this talk:

Questions

• what is the optimal fixed cost in the initialization phase?

Curious phenomena in non-linear bandits:

- phase transition in the regret
- initialization phase: regret grows linearly and results in a fixed cost
 - $\rightarrow\,$ find a good "initial action" to start learning
- learning phase: regret grows sublinearly and looks like a linear bandit
 - $\rightarrow\,$ bandit learning starts from the good initial action

Aim of this talk:

Questions

- what is the optimal fixed cost in the initialization phase?
- what algorithms should we use in different phases?

Curious phenomena in non-linear bandits:

- phase transition in the regret
- initialization phase: regret grows linearly and results in a fixed cost
 - $\rightarrow\,$ find a good "initial action" to start learning
- learning phase: regret grows sublinearly and looks like a linear bandit
 - $\rightarrow\,$ bandit learning starts from the good initial action

Aim of this talk:

Questions

- what is the optimal fixed cost in the initialization phase?
- what algorithms should we use in different phases?
- how to explore when learner has not started learning?

Plan of this talk

- setting and main results
- proof of upper bound
- proof of lower bound
- discussions and extensions

- parameter space $\Theta = \mathbb{S}^{d-1} = \{ \theta \in \mathbb{R}^d : \|\theta\|_2 = 1 \}$
- action space $\mathcal{A} = \mathbb{B}^d = \{ a \in \mathbb{R}^d : \|a\|_2 \leq 1 \}$
- reward function $f_{\theta}(a) = f(\langle \theta, a \rangle)$ with a known link function f

- parameter space $\Theta = \mathbb{S}^{d-1} = \{ \theta \in \mathbb{R}^d : \|\theta\|_2 = 1 \}$
- action space $\mathcal{A} = \mathbb{B}^d = \{ a \in \mathbb{R}^d : \|a\|_2 \leq 1 \}$
- reward function $f_{\theta}(a) = f(\langle \theta, a \rangle)$ with a known link function f

Assumptions

- parameter space $\Theta = \mathbb{S}^{d-1} = \{ \theta \in \mathbb{R}^d : \|\theta\|_2 = 1 \}$
- action space $\mathcal{A} = \mathbb{B}^d = \{ a \in \mathbb{R}^d : \|a\|_2 \leq 1 \}$
- reward function $f_{\theta}(a) = f(\langle \theta, a \rangle)$ with a known link function f

Assumptions

- monotonicity: $f : [-1,1] \rightarrow [-1,1]$ is increasing (or f(-x) = f(x) and f is increasing on [0,1]) with $f(0) = 0, f(1) \approx 1$
 - \rightarrow best action is $a = \theta^*$

- parameter space $\Theta = \mathbb{S}^{d-1} = \{ \theta \in \mathbb{R}^d : \|\theta\|_2 = 1 \}$
- action space $\mathcal{A} = \mathbb{B}^d = \{ a \in \mathbb{R}^d : \|a\|_2 \leq 1 \}$
- reward function $f_{ heta}(a) = f(\langle heta, a \rangle)$ with a known link function f

Assumptions

- monotonicity: $f : [-1,1] \rightarrow [-1,1]$ is increasing (or f(-x) = f(x) and f is increasing on [0,1]) with $f(0) = 0, f(1) \asymp 1$
 - \rightarrow best action is $a = \theta^*$
- local linearity near 1: $\max_{x \in [0.1,1]} f'(x) / \min_{x \in [0.1,1]} f'(x) \le c < \infty$
 - \rightarrow essentially linear reward when $\langle \theta^{\star}, a \rangle$ becomes large

Literature review

Literature review

Ridge bandit $f_{\theta}(a) = f(\langle \theta, a \rangle)$:

- linear bandit f(x) = x: optimal regret $\widetilde{\Theta}(d\sqrt{T})$ [Dani et al. 2008, Chu et al. 2011, Abbasi-Yadkori et al. 2011]
- generalized linear bandit with $c_1 \le |f'(x)| \le c_2$: same as linear bandit [Filippi et al. 2010, Russo and Van Roy 2014]
- concave bandit (f is concave): same as linear bandit [Lattimore, 2021]
- bandit phase retrieval $(f(x) = x^2)$: same as linear bandit [Lattimore and Hao, 2021]
- polynomial bandit $(f(x) = x^p, p \ge 2)$: optimal regret achieved by noisy gradient method [Huang et al. 2021]

Literature review

Ridge bandit $f_{\theta}(a) = f(\langle \theta, a \rangle)$:

- linear bandit f(x) = x: optimal regret $\widetilde{\Theta}(d\sqrt{T})$ [Dani et al. 2008, Chu et al. 2011, Abbasi-Yadkori et al. 2011]
- generalized linear bandit with $c_1 \le |f'(x)| \le c_2$: same as linear bandit [Filippi et al. 2010, Russo and Van Roy 2014]
- concave bandit (f is concave): same as linear bandit [Lattimore, 2021]
- bandit phase retrieval $(f(x) = x^2)$: same as linear bandit [Lattimore and Hao, 2021]
- polynomial bandit $(f(x) = x^p, p \ge 2)$: optimal regret achieved by noisy gradient method [Huang et al. 2021]

General complexity measures for bandits:

- decision-estimation coefficient (DEC) [Foster et al. 2021, 2022]
- information ratio [Lattimore, 2022]
- often do not lead to tight regret dependence on d

Main Results

Theorem (main upper bound, informal)

Under monotonicity and local linearity of f:

$$\mathsf{MinmaxReg}(T, d, f) \lesssim \min \left\{ d^2 \cdot \int_{1/\sqrt{d}}^{1/2} \frac{\mathsf{d}(x^2)}{\max_{1/\sqrt{d} \leq y \leq x} f'(y)^2} + d\sqrt{T}, T \right\}.$$

Theorem (main upper bound, informal)

Under monotonicity and local linearity of f:

$$\mathsf{MinmaxReg}(T, d, f) \lesssim \min \left\{ d^2 \cdot \int_{1/\sqrt{d}}^{1/2} \frac{\mathsf{d}(x^2)}{\max_{1/\sqrt{d} \leq y \leq x} f'(y)^2} + d\sqrt{T}, T \right\}.$$

• a useful corollary:

$$\mathsf{MinmaxReg}(T, d, f) \lesssim \min\left\{d^2 \cdot \int_{1/\sqrt{d}}^{1/2} \frac{\mathsf{d}(x^2)}{f'(x)^2} + d\sqrt{T}, T\right\}$$

Theorem (main upper bound, informal)

Under monotonicity and local linearity of f:

$$\mathsf{MinmaxReg}(T, d, f) \lesssim \min \left\{ d^2 \cdot \int_{1/\sqrt{d}}^{1/2} \frac{\mathsf{d}(x^2)}{\max_{1/\sqrt{d} \leq y \leq x} f'(y)^2} + d\sqrt{T}, T \right\}.$$

• a useful corollary:

$$\mathsf{MinmaxReg}(T, d, f) \lesssim \min\left\{d^2 \cdot \int_{1/\sqrt{d}}^{1/2} \frac{\mathsf{d}(x^2)}{f'(x)^2} + d\sqrt{T}, T\right\}$$

• the formal version:

$$\mathsf{MinmaxReg}(T,d,f) \lesssim \sum_{m=1}^{d/4} \frac{1}{\max_{0 \leq y \leq \sqrt{m/d}} \min_{z \in [y,2y]} (f(z+1/\sqrt{d})-f(z))^2} + d\sqrt{T}.$$

Theorem (main upper bound, informal)

Under monotonicity and local linearity of f:

$$\mathsf{MinmaxReg}(T, d, f) \lesssim \min\left\{ d^2 \cdot \int_{1/\sqrt{d}}^{1/2} \frac{\mathsf{d}(x^2)}{\max_{1/\sqrt{d} \leq y \leq x} f'(y)^2} + d\sqrt{T}, T \right\}.$$

Theorem (main lower bound)

Under monotonicity and local linearity of f:

$$\mathsf{MinmaxReg}(T, d, f) \gtrsim \min \left\{ d \cdot \int_{1/\sqrt{d}}^{1/2} \frac{\mathsf{d}(x^2)}{f(x)^2} + d\sqrt{T}, T \right\}.$$

Theorem (main upper bound, informal)

Under monotonicity and local linearity of f:

$$\mathsf{MinmaxReg}(T, d, f) \lesssim \min\left\{ d^2 \cdot \int_{1/\sqrt{d}}^{1/2} \frac{\mathsf{d}(x^2)}{\max_{1/\sqrt{d} \leq y \leq x} f'(y)^2} + d\sqrt{T}, T \right\}.$$

Theorem (main lower bound)

Under monotonicity and local linearity of f:

$$\mathsf{MinmaxReg}(T, d, f) \gtrsim \min \left\{ d \cdot \int_{1/\sqrt{d}}^{1/2} \frac{\mathsf{d}(x^2)}{f(x)^2} + d\sqrt{T}, T \right\}.$$

- both results within poly-logarithmic factors
- pointwise upper and lower bounds
- fixed cost depends on the entire function f

$$x_t = \langle \theta^\star, a_t \rangle$$

Theorem (learning trajectory)

$$x_{t} = \langle \theta^{\star}, a_{t} \rangle$$

$$1/\sqrt{d}$$

Theorem (learning trajectory)



Theorem (learning trajectory)

• there is an algorithm attaining the UB learning curve



Theorem (learning trajectory)

• there is an algorithm attaining the UB learning curve
Main results: learning trajectory in the initialization phase



Theorem (learning trajectory)

- there is an algorithm attaining the UB learning curve
- for any algorithm, its learning trajectory lies below the LB learning curve with probability at least $1 T\delta$ under $\theta^* \sim \text{Unif}(\mathbb{S}^{d-1})$

Main results: learning trajectory in the initialization phase



Theorem (learning trajectory)

- there is an algorithm attaining the UB learning curve
- for any algorithm, its learning trajectory lies below the LB learning curve with probability at least $1 T\delta$ under $\theta^* \sim \text{Unif}(\mathbb{S}^{d-1})$
- UCB algorithm makes no progress whenever $t < d/f(1/\sqrt{d})^2!$

Examples

• polynomial bandit $f(x) = x^p$:

$$\mathsf{MinmaxReg} \asymp \begin{cases} \min\{d\sqrt{T}, T\} & \text{if } 0 2. \end{cases}$$

 $\rightarrow\,$ both Eluder-UCB and information-directed sampling give an additional $O(d^{p+1})$ term when p>1

Examples

• polynomial bandit $f(x) = x^p$:

$$\mathsf{MinmaxReg} \asymp egin{cases} \min\{d\sqrt{T}, T\} & ext{if } 0 2. \end{cases}$$

- $\rightarrow\,$ both Eluder-UCB and information-directed sampling give an additional $O(d^{p+1})$ term when p>1
- ReLU bandit $f(x) = (x 0.1)_+$: $T = e^{\Omega(d)}$ is necessary for sublinear regret

Examples

• polynomial bandit $f(x) = x^p$:

$$\mathsf{MinmaxReg} \asymp egin{cases} \min\{d\sqrt{T}, T\} & ext{if } 0 2. \end{cases}$$

- $\rightarrow\,$ both Eluder-UCB and information-directed sampling give an additional $O(d^{p+1})$ term when p>1
- ReLU bandit $f(x) = (x 0.1)_+$: $T = e^{\Omega(d)}$ is necessary for sublinear regret
- importance of f at every point:



Upper Bounds

Key feature in the learning phase

The learner has found a good "initial action" a_0 such that $\langle a_0, \theta^* \rangle \geq \text{const.}$

Key feature in the learning phase

The learner has found a good "initial action" a_0 such that $\langle a_0, \theta^* \rangle \geq \text{const.}$

A simple explore-then-commit (ETC) algorithm:

Key feature in the learning phase

The learner has found a good "initial action" a_0 such that $\langle a_0, \theta^* \rangle \geq \text{const.}$

A simple explore-then-commit (ETC) algorithm:

• for the first m rounds, uniformly explore the following 2d directions:

$$\left\{\lambda a_0 \pm \sqrt{1-\lambda^2} e_1, \cdots, \lambda a_0 \pm \sqrt{1-\lambda^2} e_d\right\}, \quad \lambda = \lambda (ext{const});$$

Key feature in the learning phase

The learner has found a good "initial action" a_0 such that $\langle a_0, \theta^{\star} \rangle \geq \text{const.}$

A simple explore-then-commit (ETC) algorithm:

• for the first m rounds, uniformly explore the following 2d directions:

$$\left\{\lambda a_0 \pm \sqrt{1-\lambda^2}e_1, \cdots, \lambda a_0 \pm \sqrt{1-\lambda^2}e_d\right\}, \quad \lambda = \lambda (ext{const});$$

• find the least squares estimate of θ^* :

$$\widehat{\theta} = \arg\min_{\theta:\langle \theta, a_0 \rangle \ge \text{const}} \frac{1}{2} \sum_{t=1}^m (r_m - f(\langle \theta, a_t \rangle))^2;$$

Key feature in the learning phase

The learner has found a good "initial action" a_0 such that $\langle a_0, \theta^* \rangle \geq \text{const.}$

A simple explore-then-commit (ETC) algorithm:

• for the first m rounds, uniformly explore the following 2d directions:

$$\left\{\lambda a_0 \pm \sqrt{1-\lambda^2} e_1, \cdots, \lambda a_0 \pm \sqrt{1-\lambda^2} e_d\right\}, \quad \lambda = \lambda ({\sf const});$$

• find the least squares estimate of θ^* :

$$\widehat{\theta} = \arg\min_{\theta: \langle \theta, a_0 \rangle \ge \text{const}} \frac{1}{2} \sum_{t=1}^m (r_m - f(\langle \theta, a_t \rangle))^2;$$

• for the remaining rounds, greedily play $a_t = \hat{\theta}$.

• standard least squares analysis gives w.h.p.

$$\sum_{t=1}^{m} (f(\langle \widehat{\theta}, a_t \rangle) - f(\langle \theta^{\star}, a_t \rangle))^2 = \widetilde{O}(d);$$

• standard least squares analysis gives w.h.p.

$$\sum_{t=1}^{m} (f(\langle \widehat{\theta}, a_t \rangle) - f(\langle \theta^{\star}, a_t \rangle))^2 = \widetilde{O}(d);$$

• local linearity of f near 1 implies that

$$\|\widehat{ heta} - heta^{\star}\|_2^2 = \widetilde{O}\left(rac{d^2}{m \cdot f'(1)^2}
ight);$$

• standard least squares analysis gives w.h.p.

$$\sum_{t=1}^{m} (f(\langle \widehat{\theta}, a_t \rangle) - f(\langle \theta^{\star}, a_t \rangle))^2 = \widetilde{O}(d);$$

• local linearity of f near 1 implies that

$$\|\widehat{ heta} - heta^{\star}\|_2^2 = \widetilde{O}\left(rac{d^2}{m \cdot f'(1)^2}
ight);$$

• instantaneous regret when greedily plays $\hat{\theta}$:

$$f(1) - f(\langle heta^{\star}, \widehat{ heta}
angle) \lesssim f'(1)(1 - \langle heta^{\star}, \widehat{ heta}
angle) \lesssim rac{d^2}{m \cdot f'(1)};$$

• standard least squares analysis gives w.h.p.

$$\sum_{t=1}^{m} (f(\langle \widehat{\theta}, a_t \rangle) - f(\langle \theta^{\star}, a_t \rangle))^2 = \widetilde{O}(d);$$

• local linearity of f near 1 implies that

$$\|\widehat{\theta} - \theta^{\star}\|_{2}^{2} = \widetilde{O}\left(rac{d^{2}}{m \cdot f'(1)^{2}}
ight);$$

• instantaneous regret when greedily plays $\hat{\theta}$:

$$f(1) - f(\langle heta^{\star}, \widehat{ heta}
angle) \lesssim f'(1)(1 - \langle heta^{\star}, \widehat{ heta}
angle) \lesssim rac{d^2}{m \cdot f'(1)};$$

• total regret in the learning phase:

$$m \cdot f'(1) + (T-m) \cdot rac{d^2}{m \cdot f'(1)} \stackrel{m symp d \sqrt{T}/f'(1)}{symp} d\sqrt{T}.$$

Target in the initialization phase

Target in the initialization phase



Target in the initialization phase



Target in the initialization phase



Target in the initialization phase



Target in the initialization phase



Certify that $\langle \theta^{\star}, a \rangle \in [r - \delta, r + \delta]$ can be done with $\widetilde{O}(1/[\delta f'(r)]^2)$ samples

Recursive step

Recursive step

Given an action a_{pre} with $\langle \theta^*, a_{\text{pre}} \rangle \in [x_{\text{pre}}, 2x_{\text{pre}}]$, where x_{pre} is known, how to find a_{now} and certify that $\langle \theta^*, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} > x_{\text{pre}}$?

• idea: find $a \perp a_{\rm pre}$ with $\langle \theta^{\star}, a \rangle \asymp 1/\sqrt{d}$ and play $a_{\rm now} = \lambda a_{\rm pre} + \sqrt{1-\lambda^2}a$

Recursive step

- idea: find $a \perp a_{\rm pre}$ with $\langle \theta^{\star}, a \rangle \asymp 1/\sqrt{d}$ and play $a_{\rm now} = \lambda a_{\rm pre} + \sqrt{1 \lambda^2} a$
- for proper λ , if $\langle \theta^{\star}, a \rangle \in [1/\sqrt{d}, 2/\sqrt{d}]$, then $\langle \theta^{\star}, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} = \sqrt{x_{\text{pre}}^2 + 1/d}$

Recursive step

- idea: find $a \perp a_{pre}$ with $\langle \theta^{\star}, a \rangle \asymp 1/\sqrt{d}$ and play $a_{now} = \lambda a_{pre} + \sqrt{1 \lambda^2} a$
- for proper λ , if $\langle \theta^{\star}, a \rangle \in [1/\sqrt{d}, 2/\sqrt{d}]$, then $\langle \theta^{\star}, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} = \sqrt{x_{\text{pre}}^2 + 1/d}$
- exploration: easy, as $\mathbb{P}(\langle heta^{\star}, a
 angle \in [1/\sqrt{d}, 2/\sqrt{d}]) = \Omega(1)$ for uniform a

Recursive step

- idea: find $a \perp a_{\rm pre}$ with $\langle \theta^{\star}, a \rangle \asymp 1/\sqrt{d}$ and play $a_{\rm now} = \lambda a_{\rm pre} + \sqrt{1 \lambda^2} a$
- for proper λ , if $\langle \theta^{\star}, a \rangle \in [1/\sqrt{d}, 2/\sqrt{d}]$, then $\langle \theta^{\star}, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} = \sqrt{x_{\text{pre}}^2 + 1/d}$
- exploration: easy, as $\mathbb{P}(\langle heta^{\star}, a
 angle \in [1/\sqrt{d}, 2/\sqrt{d}]) = \Omega(1)$ for uniform a
- certification: should make use of apre!

Recursive step

Given an action a_{pre} with $\langle \theta^*, a_{\text{pre}} \rangle \in [x_{\text{pre}}, 2x_{\text{pre}}]$, where x_{pre} is known, how to find a_{now} and certify that $\langle \theta^*, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} > x_{\text{pre}}$?

- idea: find $a \perp a_{\rm pre}$ with $\langle \theta^{\star}, a \rangle \asymp 1/\sqrt{d}$ and play $a_{\rm now} = \lambda a_{\rm pre} + \sqrt{1 \lambda^2} a$
- for proper λ , if $\langle \theta^{\star}, a \rangle \in [1/\sqrt{d}, 2/\sqrt{d}]$, then $\langle \theta^{\star}, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} = \sqrt{x_{\text{pre}}^2 + 1/d}$
- exploration: easy, as $\mathbb{P}(\langle heta^{\star}, a
 angle \in [1/\sqrt{d}, 2/\sqrt{d}]) = \Omega(1)$ for uniform a

• certification: should make use of apre!

$$\frac{\langle \theta^{\star}, \mathbf{a} \rangle}{\sqrt{2}} = \left\langle \theta^{\star}, \frac{\mathbf{a} + \mathbf{a}_{\mathsf{pre}}}{\sqrt{2}} \right\rangle - \left\langle \theta^{\star}, \frac{\mathbf{a}_{\mathsf{pre}}}{\sqrt{2}} \right\rangle$$

Recursive step

Given an action a_{pre} with $\langle \theta^*, a_{\text{pre}} \rangle \in [x_{\text{pre}}, 2x_{\text{pre}}]$, where x_{pre} is known, how to find a_{now} and certify that $\langle \theta^*, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} > x_{\text{pre}}$?

- idea: find $a \perp a_{\rm pre}$ with $\langle \theta^{\star}, a \rangle \asymp 1/\sqrt{d}$ and play $a_{\rm now} = \lambda a_{\rm pre} + \sqrt{1 \lambda^2} a$
- for proper λ , if $\langle \theta^{\star}, a \rangle \in [1/\sqrt{d}, 2/\sqrt{d}]$, then $\langle \theta^{\star}, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} = \sqrt{x_{\text{pre}}^2 + 1/d}$
- exploration: easy, as $\mathbb{P}(\langle heta^{\star}, a
 angle \in [1/\sqrt{d}, 2/\sqrt{d}]) = \Omega(1)$ for uniform a

certification: should make use of apre!

$$\frac{\langle \theta^{\star}, \mathbf{a} \rangle}{\sqrt{2}} = \left\langle \theta^{\star}, \frac{\mathbf{a} + \mathbf{a}_{\mathsf{pre}}}{\sqrt{2}} \right\rangle - \left\langle \theta^{\star}, \frac{\mathbf{a}_{\mathsf{pre}}}{\sqrt{2}} \right\rangle$$

ightarrow each terms uses $\widetilde{O}(d/f'(x_{
m pre})^2)$ samples for certification

Recursive step

Given an action a_{pre} with $\langle \theta^*, a_{\text{pre}} \rangle \in [x_{\text{pre}}, 2x_{\text{pre}}]$, where x_{pre} is known, how to find a_{now} and certify that $\langle \theta^*, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} > x_{\text{pre}}$?

- idea: find $a \perp a_{\rm pre}$ with $\langle \theta^{\star}, a \rangle \asymp 1/\sqrt{d}$ and play $a_{\rm now} = \lambda a_{\rm pre} + \sqrt{1 \lambda^2} a$
- for proper λ , if $\langle \theta^{\star}, a \rangle \in [1/\sqrt{d}, 2/\sqrt{d}]$, then $\langle \theta^{\star}, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} = \sqrt{x_{\text{pre}}^2 + 1/d}$
- exploration: easy, as $\mathbb{P}(\langle heta^{\star}, a
 angle \in [1/\sqrt{d}, 2/\sqrt{d}]) = \Omega(1)$ for uniform a

certification: should make use of apre!

$$\frac{\langle \theta^{\star}, \mathbf{a} \rangle}{\sqrt{2}} = \left\langle \theta^{\star}, \frac{\mathbf{a} + \mathbf{a}_{\mathsf{pre}}}{\sqrt{2}} \right\rangle - \left\langle \theta^{\star}, \frac{\mathbf{a}_{\mathsf{pre}}}{\sqrt{2}} \right\rangle$$

ightarrow each terms uses $\widetilde{O}(d/f'(x_{
m pre})^2)$ samples for certification

 \rightarrow total sample complexity is roughly $d^2 \int_{1/\sqrt{d}}^{1/2} \frac{\mathrm{d}x^2}{f'(x)^2}$

Recursive step

Given an action a_{pre} with $\langle \theta^{\star}, a_{\text{pre}} \rangle \in [x_{\text{pre}}, 2x_{\text{pre}}]$, where x_{pre} is known, how to find a_{now} and certify that $\langle \theta^{\star}, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} > x_{\text{pre}}$?

- idea: find $a \perp a_{\rm pre}$ with $\langle \theta^{\star}, a \rangle \asymp 1/\sqrt{d}$ and play $a_{\rm now} = \lambda a_{\rm pre} + \sqrt{1 \lambda^2} a$
- for proper λ , if $\langle \theta^{\star}, a \rangle \in [1/\sqrt{d}, 2/\sqrt{d}]$, then $\langle \theta^{\star}, a_{\text{now}} \rangle \in [x_{\text{now}}, 2x_{\text{now}}]$ with $x_{\text{now}} = \sqrt{x_{\text{pre}}^2 + 1/d}$
- exploration: easy, as $\mathbb{P}(\langle heta^\star, a
 angle \in [1/\sqrt{d}, 2/\sqrt{d}]) = \Omega(1)$ for uniform a

• certification: should make use of apre!

$$\frac{\langle \theta^{\star}, \mathbf{a} \rangle}{\sqrt{2}} = \left\langle \theta^{\star}, \frac{\mathbf{a} + \mu \mathbf{a}_{\mathsf{pre}}}{\sqrt{2}} \right\rangle - \left\langle \theta^{\star}, \frac{\mu \mathbf{a}_{\mathsf{pre}}}{\sqrt{2}} \right\rangle$$

 \rightarrow each terms uses $\widetilde{O}(d/\max_{y \leq x} f'(y)^2)$ samples for certification;

 \rightarrow total sample complexity is roughly $d^2 \int_{1/\sqrt{d}}^{1/2} \frac{dx^2}{\max_{y < x} f'(y)^2}$.

Target of certification

Given actions a and a + b with $\langle \theta^{\star}, a \rangle \in [x, 2x]$, find a test which

- outputs "failure" w.h.p. if $\langle \theta^{\star}, b \rangle \notin [z, 2z];$
- outputs "success" w.h.p. if $\langle \theta^{\star}, b \rangle \in [1.2z, 1.8z]$.

Target of certification

Given actions a and a + b with $\langle \theta^{\star}, a \rangle \in [x, 2x]$, find a test which

- outputs "failure" w.h.p. if $\langle \theta^{\star}, b \rangle \notin [z, 2z]$;
- outputs "success" w.h.p. if $\langle \theta^{\star}, b \rangle \in [1.2z, 1.8z]$.

• pull both actions $\widetilde{O}(1/\delta^2)$ times to obtain

$$|\widehat{f_1} - f(\langle heta^\star, a \rangle)| \leq \delta, \qquad |\widehat{f_2} - f(\langle heta^\star, a + b
angle)| \leq \delta;$$

Target of certification

Given actions a and a + b with $\langle \theta^*, a \rangle \in [x, 2x]$, find a test which

- outputs "failure" w.h.p. if $\langle \theta^{\star}, b \rangle \notin [z, 2z];$
- outputs "success" w.h.p. if $\langle \theta^{\star}, b \rangle \in [1.2z, 1.8z]$.

• pull both actions $\widetilde{O}(1/\delta^2)$ times to obtain

$$|\widehat{f_1} - f(\langle heta^\star, a \rangle)| \leq \delta, \qquad |\widehat{f_2} - f(\langle heta^\star, a + b
angle)| \leq \delta;$$

• test returns "success" iff $\exists u \in [x, 2x], v \in [1.2z, 1.8z]$ such that $|\widehat{f_1} - f(u)| \leq \delta$ and $|\widehat{f_2} - f(u+v)| \leq \delta$:

Target of certification

Given actions a and a + b with $\langle \theta^*, a \rangle \in [x, 2x]$, find a test which

- outputs "failure" w.h.p. if $\langle \theta^{\star}, b \rangle \notin [z, 2z]$;
- outputs "success" w.h.p. if $\langle \theta^{\star}, b \rangle \in [1.2z, 1.8z]$.

• pull both actions $\widetilde{O}(1/\delta^2)$ times to obtain

$$|\widehat{f_1} - f(\langle heta^\star, a \rangle)| \leq \delta, \qquad |\widehat{f_2} - f(\langle heta^\star, a + b
angle)| \leq \delta;$$

• test returns "success" iff $\exists u \in [x, 2x], v \in [1.2z, 1.8z]$ such that $|\widehat{f_1} - f(u)| \leq \delta$ and $|\widehat{f_2} - f(u+v)| \leq \delta$: \rightarrow if $\langle \theta^*, b \rangle \in [1.2z, 1.8z]$, then $(u, v) = (\langle \theta^*, a \rangle, \langle \theta^*, b \rangle)$ passes the test;

Target of certification

Given actions a and a + b with $\langle \theta^*, a \rangle \in [x, 2x]$, find a test which

- outputs "failure" w.h.p. if $\langle \theta^{\star}, b \rangle \notin [z, 2z]$;
- outputs "success" w.h.p. if $\langle \theta^{\star}, b \rangle \in [1.2z, 1.8z]$.

• pull both actions $\widetilde{O}(1/\delta^2)$ times to obtain

$$|\widehat{f_1} - f(\langle heta^\star, a \rangle)| \leq \delta, \qquad |\widehat{f_2} - f(\langle heta^\star, a + b
angle)| \leq \delta;$$

• test returns "success" iff $\exists u \in [x, 2x], v \in [1.2z, 1.8z]$ such that $|\widehat{f_1} - f(u)| \leq \delta$ and $|\widehat{f_2} - f(u+v)| \leq \delta$: \rightarrow if $\langle \theta^*, b \rangle \in [1.2z, 1.8z]$, then $(u, v) = (\langle \theta^*, a \rangle, \langle \theta^*, b \rangle)$ passes the test; \rightarrow if $\langle \theta^*, b \rangle \notin [z, 2z]$, then existence of (u, v) implies

$$|u - \langle \theta^{\star}, a \rangle| \ge 0.2z, \quad \text{ or } \quad |u + v - \langle \theta^{\star}, a + b \rangle| \ge 0.2z;$$

Target of certification

Given actions a and a + b with $\langle \theta^*, a \rangle \in [x, 2x]$, find a test which

- outputs "failure" w.h.p. if $\langle \theta^{\star}, b \rangle \notin [z, 2z]$;
- outputs "success" w.h.p. if $\langle \theta^{\star}, b \rangle \in [1.2z, 1.8z]$.

• pull both actions $\widetilde{O}(1/\delta^2)$ times to obtain

$$|\widehat{f_1} - f(\langle heta^\star, a \rangle)| \leq \delta, \qquad |\widehat{f_2} - f(\langle heta^\star, a + b
angle)| \leq \delta;$$

• test returns "success" iff $\exists u \in [x, 2x], v \in [1.2z, 1.8z]$ such that $|\widehat{f_1} - f(u)| \leq \delta$ and $|\widehat{f_2} - f(u+v)| \leq \delta$: \rightarrow if $\langle \theta^*, b \rangle \in [1.2z, 1.8z]$, then $(u, v) = (\langle \theta^*, a \rangle, \langle \theta^*, b \rangle)$ passes the test; \rightarrow if $\langle \theta^*, b \rangle \notin [z, 2z]$, then existence of (u, v) implies

$$|u - \langle \theta^{\star}, a \rangle| \geq 0.2z, \quad \text{ or } \quad |u + v - \langle \theta^{\star}, a + b \rangle| \geq 0.2z;$$

 \rightarrow test works if $\delta < \min_{y \in [x, 2x+2z]} [f(y+0.2z) - f(y)]/2.$
Lower Bounds

Theorem (formal lower bound)

Let $\delta>0$ be any parameter, and c>0 be a large absolute constant. Define a sequence $\{\varepsilon_t\}_{t\geq 1}$ with

$$arepsilon_1 = \sqrt{rac{c\log(1/\delta)}{d}}, \quad arepsilon_{t+1}^2 = arepsilon_t^2 + rac{c}{d}f(arepsilon_t)^2, \quad t \geq 1$$

Theorem (formal lower bound)

Let $\delta>0$ be any parameter, and c>0 be a large absolute constant. Define a sequence $\{\varepsilon_t\}_{t\geq 1}$ with

$$arepsilon_1 = \sqrt{rac{c\log(1/\delta)}{d}}, \quad arepsilon_{t+1}^2 = arepsilon_t^2 + rac{c}{d}f(arepsilon_t)^2, \quad t \geq 1$$

Then if $\theta^{\star} \sim \mathsf{Unif}(\mathbb{S}^{d-1})$, any learner $\{a_t\}_{t\geq 1}$ satisfies that

$$\mathbb{P}\left(\bigcap_{1\leq t\leq T}\left\{\langle \theta^{\star}, a_t\rangle \leq \varepsilon_t\right\}\right) \geq 1 - T\delta$$

Theorem (formal lower bound)

Let $\delta>0$ be any parameter, and c>0 be a large absolute constant. Define a sequence $\{\varepsilon_t\}_{t\geq 1}$ with

$$arepsilon_1 = \sqrt{rac{c\log(1/\delta)}{d}}, \quad arepsilon_{t+1}^2 = arepsilon_t^2 + rac{c}{d}f(arepsilon_t)^2, \quad t \geq 1$$

Then if $\theta^{\star} \sim \mathsf{Unif}(\mathbb{S}^{d-1})$, any learner $\{a_t\}_{t\geq 1}$ satisfies that

$$\mathbb{P}\left(igcap_{1\leq t\leq \mathcal{T}}\left\{\langle heta^{\star}, oldsymbol{a}_t
angle\leqarepsilon_t
ight\}
ight)\geq 1-\mathcal{T}\delta.$$

• the continuous-time version of $\{\varepsilon_t\}$ gives the differential equation

Theorem (formal lower bound)

Let $\delta>0$ be any parameter, and c>0 be a large absolute constant. Define a sequence $\{\varepsilon_t\}_{t\geq 1}$ with

$$arepsilon_1 = \sqrt{rac{c\log(1/\delta)}{d}}, \quad arepsilon_{t+1}^2 = arepsilon_t^2 + rac{c}{d}f(arepsilon_t)^2, \quad t \geq 1$$

Then if $\theta^{\star} \sim \text{Unif}(\mathbb{S}^{d-1})$, any learner $\{a_t\}_{t \geq 1}$ satisfies that

$$\mathbb{P}\left(igcap_{1\leq t\leq \mathcal{T}}\left\{\langle heta^{\star}, oldsymbol{a}_t
angle\leqarepsilon_t
ight\}
ight)\geq 1-\mathcal{T}\delta.$$

- the continuous-time version of $\{\varepsilon_t\}$ gives the differential equation
- hard(?) to prove via usual arguments of hypothesis testing

Let $I_t = I(\theta^*; \mathcal{H}_t)$ be the mutual information between the true parameter θ^* and the history \mathcal{H}_t up to time t, then

$$egin{aligned} & I_{t+1} - I_t = I(heta^\star; r_{t+1} \mid a_{t+1}, \mathcal{H}_t) \ & \leq \mathbb{E}\left[rac{1}{2}\log\left(1 + \mathbb{E}[f(\langle heta^\star, a_{t+1}
angle))^2]
ight)
ight] \ & \leq rac{1}{2}\mathbb{E}[f(\langle heta^\star, a_{t+1}
angle)^2]. \end{aligned}$$

Let $I_t = I(\theta^*; \mathcal{H}_t)$ be the mutual information between the true parameter θ^* and the history \mathcal{H}_t up to time t, then

$$egin{aligned} &I_{t+1}-I_t=I(heta^\star;r_{t+1}\mid a_{t+1},\mathcal{H}_t)\ &\leq \mathbb{E}\left[rac{1}{2}\log\left(1+\mathbb{E}[f(\langle heta^\star,a_{t+1}
angle))^2]
ight)
ight]\ &\leq rac{1}{2}\mathbb{E}[f(\langle heta^\star,a_{t+1}
angle)^2]. \end{aligned}$$

To argue that $\langle heta^{\star}, a_{t+1}
angle$ should not be large, note that

$$I(\theta^*; a_{t+1}) \leq I(\theta^*; \mathcal{H}_t) = I_t.$$

Let $I_t = I(\theta^*; \mathcal{H}_t)$ be the mutual information between the true parameter θ^* and the history \mathcal{H}_t up to time t, then

$$egin{aligned} &I_{t+1}-I_t=I(heta^\star;r_{t+1}\mid a_{t+1},\mathcal{H}_t)\ &\leq \mathbb{E}\left[rac{1}{2}\log\left(1+\mathbb{E}[f(\langle heta^\star,a_{t+1}
angle))^2]
ight)
ight]\ &\leq rac{1}{2}\mathbb{E}[f(\langle heta^\star,a_{t+1}
angle)^2]. \end{aligned}$$

To argue that $\langle heta^{\star}, a_{t+1}
angle$ should not be large, note that

$$I(\theta^{\star}; a_{t+1}) \leq I(\theta^{\star}; \mathcal{H}_t) = I_t.$$

Key insight

$$I(\theta^{\star}; a) \leq I \Longrightarrow \langle \theta^{\star}, a \rangle \lesssim \sqrt{I/d}$$
 with high probability.

Let $I_t = I(\theta^*; \mathcal{H}_t)$ be the mutual information between the true parameter θ^* and the history \mathcal{H}_t up to time t, then

$$egin{aligned} &I_{t+1}-I_t=I(heta^\star;r_{t+1}\mid a_{t+1},\mathcal{H}_t)\ &\leq \mathbb{E}\left[rac{1}{2}\log\left(1+\mathbb{E}[f(\langle heta^\star,a_{t+1}
angle))^2]
ight)
ight]\ &\leq rac{1}{2}\mathbb{E}[f(\langle heta^\star,a_{t+1}
angle)^2]. \end{aligned}$$

To argue that $\langle heta^{\star}, a_{t+1}
angle$ should not be large, note that

$$I(\theta^{\star}; a_{t+1}) \leq I(\theta^{\star}; \mathcal{H}_t) = I_t.$$

Key insight

$$I(heta^{\star}; a) \leq I \Longrightarrow \langle heta^{\star}, a \rangle \lesssim \sqrt{I/d}$$
 with high probability.

Applying the insight gives the desired recursion

$$d(\varepsilon_{t+1}^2 - \varepsilon_t^2) \lesssim f(\varepsilon_t)^2$$

More on the above insights

• reasoning behind the insight:

$$\mathsf{a} \mid \theta^{\star} \sim \mathsf{Unif}(\{\mathsf{a} \in \mathbb{S}^{d-1} : \langle \mathsf{a}, \theta^{\star} \rangle \geq \varepsilon\}) \Longrightarrow \mathsf{I}(\mathsf{a}; \theta^{\star}) \asymp \mathsf{d}\varepsilon^{2}$$

• reasoning behind the insight:

$$\mathsf{a} \mid \theta^{\star} \sim \mathsf{Unif}(\{\mathsf{a} \in \mathbb{S}^{d-1} : \langle \mathsf{a}, \theta^{\star} \rangle \geq \varepsilon\}) \Longrightarrow \mathsf{I}(\mathsf{a}; \theta^{\star}) \asymp \mathsf{d}\varepsilon^2$$

• however, it does not hold with high probability: Fano's inequality only gives

$$\mathbb{P}(\langle heta^{\star}, extbf{a}
angle \leq arepsilon) \geq 1 - rac{I(heta^{\star}; extbf{a}) + \log 2}{\Theta(darepsilon^2)},$$

which is tight for the worst-case distribution of (θ^{\star}, a)

• reasoning behind the insight:

$$\mathsf{a} \mid \theta^{\star} \sim \mathsf{Unif}(\{\mathsf{a} \in \mathbb{S}^{d-1} : \langle \mathsf{a}, \theta^{\star} \rangle \geq \varepsilon\}) \Longrightarrow \mathsf{I}(\mathsf{a}; \theta^{\star}) \asymp \mathsf{d}\varepsilon^{2}$$

• however, it does not hold with high probability: Fano's inequality only gives

$$\mathbb{P}(\langle heta^{\star}, extbf{a}
angle \leq arepsilon) \geq 1 - rac{I(heta^{\star}; extbf{a}) + \log 2}{\Theta(darepsilon^2)},$$

which is tight for the worst-case distribution of (θ^*, a) • our solution: use χ^2 -informativity instead

• χ^2 -informativity between X and Y:

$$I_{\chi^2}(X; Y) = \inf_{Q_Y} \chi^2(P_{XY} || P_X \times Q_Y).$$

• χ^2 -informativity between X and Y:

$$I_{\chi^2}(X;Y) = \inf_{Q_Y} \chi^2(P_{XY} || P_X \times Q_Y).$$

• error probability lower bound using χ^2 -informativity:

$$\mathbb{P}(\langle heta^\star, extbf{a}
angle \leq arepsilon) \geq 1 - e^{-\Theta(darepsilon^2)} \cdot \sqrt{I_{\chi^2}(heta^\star; extbf{a}) + 1}$$

• χ^2 -informativity between X and Y:

$$I_{\chi^2}(X;Y) = \inf_{Q_Y} \chi^2(P_{XY} || P_X \times Q_Y).$$

• error probability lower bound using χ^2 -informativity:

$$\mathbb{P}(\langle heta^\star, extbf{a}
angle \leq arepsilon) \geq 1 - e^{-\Theta(darepsilon^2)} \cdot \sqrt{I_{\chi^2}(heta^\star; extbf{a}) + 1}.$$

• suffices to upper bound $I_{\chi^2}(heta^\star;a_{t+1}) \leq I_{\chi^2}(heta^\star;\mathcal{H}_t)$ for each t

• χ^2 -informativity between X and Y:

$$I_{\chi^2}(X;Y) = \inf_{Q_Y} \chi^2(P_{XY} || P_X \times Q_Y).$$

• error probability lower bound using χ^2 -informativity:

$$\mathbb{P}(\langle heta^\star, extbf{a}
angle \leq arepsilon) \geq 1 - e^{-\Theta(darepsilon^2)} \cdot \sqrt{I_{\chi^2}(heta^\star; extbf{a}) + 1}.$$

• suffices to upper bound $I_{\chi^2}(\theta^*; a_{t+1}) \leq I_{\chi^2}(\theta^*; \mathcal{H}_t)$ for each t

 ${\, \bullet \,}$ issue: $\chi^2 {\rm -informativity}$ does not satisfy the chain rule or subadditivity

• let $\mathcal{E}_t = \cap_{s \leq t} \{ \langle \theta^\star, a_s \rangle \leq \varepsilon_s \}$ be the error event

- let $\mathcal{E}_t = \cap_{s \leq t} \{ \langle \theta^\star, a_s \rangle \leq \varepsilon_s \}$ be the error event
- upper bound of conditioned χ^2 -informativity:

$$I_{\chi^2}(\theta^\star; \mathcal{H}_t \mid \mathcal{E}_t) + 1$$

• let $\mathcal{E}_t = \cap_{s \leq t} \{ \langle \theta^\star, a_s \rangle \leq \varepsilon_s \}$ be the error event

$$I_{\chi^{2}}(\theta^{\star};\mathcal{H}_{t} \mid \mathcal{E}_{t}) + 1 \leq \min_{\mathbb{Q}_{t-1}} \int \underbrace{\frac{\left[\frac{\mathbb{I}(\mathcal{E}_{t})}{\mathbb{P}(\mathcal{E}_{t})}\pi(\theta^{\star})\prod_{s\leq t}\varphi(r_{s}-\langle\theta^{\star},a_{s}\rangle)\right]^{2}}{\pi(\theta^{\star})\mathbb{Q}_{t-1}(r^{t-1})\cdot\varphi(r_{t})}}_{\pi(\theta^{\star})\mathbb{Q}_{t}(\mathcal{H}_{t})} d\theta^{\star} dr^{t}$$

• let $\mathcal{E}_t=\cap_{s\leq t}\{\langle heta^\star,a_s
angle\leq arepsilon_s\}$ be the error event

$$I_{\chi^{2}}(\theta^{*};\mathcal{H}_{t} \mid \mathcal{E}_{t}) + 1 \leq \min_{\mathbb{Q}_{t-1}} \int \underbrace{\frac{\left[\frac{\mathbb{I}(\mathcal{E}_{t})}{\mathbb{P}(\mathcal{E}_{t})}\pi(\theta^{*})\prod_{s\leq t}\varphi(r_{s}-\langle\theta^{*},a_{s}\rangle)\right]^{2}}{\pi(\theta^{*})\mathbb{Q}_{t-1}(r^{t-1})\cdot\varphi(r_{t})}}_{\pi(\theta^{*})\mathbb{Q}_{t}(\mathcal{H}_{t})} d\theta^{*}dr^{t}$$
$$= \min_{\mathbb{Q}_{t-1}} \int \frac{\left[\frac{\mathbb{I}(\mathcal{E}_{t})}{\mathbb{P}(\mathcal{E}_{t})}\pi(\theta^{*})\prod_{s\leq t-1}\varphi(r_{s}-\langle\theta^{*},a_{s}\rangle)\right]^{2}}{\pi(\theta^{*})\mathbb{Q}_{t-1}(r^{t-1})} \cdot \exp(\langle\theta^{*},a_{t}\rangle^{2})d\theta^{*}dr^{t-1}$$

• let $\mathcal{E}_t = \cap_{s \leq t} \{ \langle \theta^\star, a_s \rangle \leq \varepsilon_s \}$ be the error event

$$I_{\chi^{2}}(\theta^{\star};\mathcal{H}_{t} \mid \mathcal{E}_{t}) + 1 \leq \min_{\mathbb{Q}_{t-1}} \int \underbrace{\frac{\left[\frac{\mathbb{I}(\mathcal{E}_{t})}{\mathbb{P}(\mathcal{E}_{t})}\pi(\theta^{\star})\prod_{s\leq t}\varphi(r_{s}-\langle\theta^{\star},a_{s}\rangle)\right]^{2}}{\pi(\theta^{\star})\mathbb{Q}_{t}-1(r^{t-1})\cdot\varphi(r_{t})}}_{\pi(\theta^{\star})\mathbb{Q}_{t}(\mathcal{H}_{t})} d\theta^{\star}dr^{t}$$

$$= \min_{\mathbb{Q}_{t-1}} \int \frac{\left[\frac{\mathbb{I}(\mathcal{E}_{t})}{\mathbb{P}(\mathcal{E}_{t})}\pi(\theta^{\star})\prod_{s\leq t-1}\varphi(r_{s}-\langle\theta^{\star},a_{s}\rangle)\right]^{2}}{\pi(\theta^{\star})\mathbb{Q}_{t-1}(r^{t-1})} \cdot \exp(\langle\theta^{\star},a_{t}\rangle^{2})d\theta^{\star}dr^{t-1}$$

$$\leq \exp(\varepsilon_{t}^{2})\cdot\min_{\mathbb{Q}_{t-1}} \int \frac{\left[\frac{\mathbb{I}(\mathcal{E}_{t})}{\mathbb{P}(\mathcal{E}_{t})}\pi(\theta^{\star})\prod_{s\leq t-1}\varphi(r_{s}-\langle\theta^{\star},a_{s}\rangle)\right]^{2}}{\pi(\theta^{\star})\mathbb{Q}_{t-1}(r^{t-1})} dr^{t-1}$$

• let $\mathcal{E}_t = \cap_{s \leq t} \{ \langle \theta^\star, a_s \rangle \leq \varepsilon_s \}$ be the error event

$$\begin{split} & \frac{\mathbb{P}(\theta^{\star},\mathcal{H}_{t}|\mathcal{E}_{t})^{2}}{I_{\chi^{2}}(\theta^{\star};\mathcal{H}_{t}\mid\mathcal{E}_{t})+1\leq\min_{\mathbb{Q}_{t-1}}\int\frac{\left[\frac{\mathbb{I}(\mathcal{E}_{t})}{\mathbb{P}(\mathcal{E}_{t})}\pi(\theta^{\star})\prod_{s\leq t}\varphi(r_{s}-\langle\theta^{\star},a_{s}\rangle)\right]^{2}}{\pi(\theta^{\star})\mathbb{Q}_{t-1}(r^{t-1})\cdot\varphi(r_{t})}d\theta^{\star}dr^{t}\\ &=\min_{\mathbb{Q}_{t-1}}\int\frac{\left[\frac{\mathbb{I}(\mathcal{E}_{t})}{\mathbb{P}(\mathcal{E}_{t})}\pi(\theta^{\star})\prod_{s\leq t-1}\varphi(r_{s}-\langle\theta^{\star},a_{s}\rangle)\right]^{2}}{\pi(\theta^{\star})\mathbb{Q}_{t-1}(r^{t-1})}\cdot\exp(\langle\theta^{\star},a_{t}\rangle^{2})d\theta^{\star}dr^{t-1}\\ &\leq\exp(\varepsilon_{t}^{2})\cdot\min_{\mathbb{Q}_{t-1}}\int\frac{\left[\frac{\mathbb{I}(\mathcal{E}_{t})}{\mathbb{P}(\mathcal{E}_{t})}\pi(\theta^{\star})\prod_{s\leq t-1}\varphi(r_{s}-\langle\theta^{\star},a_{s}\rangle)\right]^{2}}{\pi(\theta^{\star})\mathbb{Q}_{t-1}(r^{t-1})}dr^{t-1}\\ &\leq\frac{\exp(\varepsilon_{t}^{2})}{\mathbb{P}(\mathcal{E}_{t}\mid\mathcal{E}_{t-1})^{2}}\cdot\min_{\mathbb{Q}_{t-1}}\int\frac{\left[\frac{\mathbb{I}(\mathcal{E}_{t-1})}{\mathbb{P}(\mathcal{E}_{t-1})}\pi(\theta^{\star})\prod_{s\leq t-1}\varphi(r_{s}-\langle\theta^{\star},a_{s}\rangle)\right]^{2}}{\pi(\theta^{\star})\mathbb{Q}_{t-1}(r^{t-1})}dr^{t-1} \end{split}$$

• let $\mathcal{E}_t=\cap_{s\leq t}\{\langle heta^\star,a_s
angle\leq arepsilon_s\}$ be the error event

$$I_{\chi^2}(\theta^\star;\mathcal{H}_t\mid \mathcal{E}_t)+1 \leq \frac{\exp(\varepsilon_t^2)}{\mathbb{P}(\mathcal{E}_t\mid \mathcal{E}_{t-1})^2}\left(I_{\chi^2}(\theta^\star;\mathcal{H}_{t-1}\mid \mathcal{E}_{t-1})+1\right).$$

• let $\mathcal{E}_t=\cap_{s\leq t}\{\langle heta^\star,a_s
angle\leq arepsilon_s\}$ be the error event

• upper bound of conditioned $\chi^2\text{-informativity:}$

$$I_{\chi^2}(\theta^\star;\mathcal{H}_t\mid \mathcal{E}_t)+1 \leq \frac{\exp(\varepsilon_t^2)}{\mathbb{P}(\mathcal{E}_t\mid \mathcal{E}_{t-1})^2}\left(I_{\chi^2}(\theta^\star;\mathcal{H}_{t-1}\mid \mathcal{E}_{t-1})+1\right).$$

• continuing this process gives

$$I_{\chi^2}(\theta^{\star}; \mathcal{H}_t \mid \mathcal{E}_t) + 1 \leq \frac{\exp(\sum_{s \leq t} \varepsilon_s^2)}{\mathbb{P}(\mathcal{E}_t)^2}.$$

• let $\mathcal{E}_t=\cap_{s\leq t}\{\langle\theta^\star,a_s\rangle\leq\varepsilon_s\}$ be the error event

• upper bound of conditioned χ^2 -informativity:

$$I_{\chi^2}(\theta^\star;\mathcal{H}_t\mid \mathcal{E}_t)+1 \leq \frac{\exp(\varepsilon_t^2)}{\mathbb{P}(\mathcal{E}_t\mid \mathcal{E}_{t-1})^2}\left(I_{\chi^2}(\theta^\star;\mathcal{H}_{t-1}\mid \mathcal{E}_{t-1})+1\right).$$

• continuing this process gives

$$I_{\chi^2}(\theta^{\star}; \mathcal{H}_t \mid \mathcal{E}_t) + 1 \leq \frac{\exp(\sum_{s \leq t} \varepsilon_s^2)}{\mathbb{P}(\mathcal{E}_t)^2}.$$

.

• recursion of error probability:

$$\mathbb{P}(\mathcal{E}_{t+1}) = \mathbb{P}(\mathcal{E}_t) \cdot \mathbb{P}(\langle \theta^{\star}, \boldsymbol{a}_{t+1} \rangle \leq \varepsilon_{t+1} \mid \mathcal{E}_t)$$

• let $\mathcal{E}_t=\cap_{s\leq t}\{\langle heta^\star,a_s
angle\leq arepsilon_s\}$ be the error event

• upper bound of conditioned χ^2 -informativity:

$$I_{\chi^2}(\theta^\star;\mathcal{H}_t\mid \mathcal{E}_t)+1 \leq \frac{\exp(\varepsilon_t^2)}{\mathbb{P}(\mathcal{E}_t\mid \mathcal{E}_{t-1})^2}\left(I_{\chi^2}(\theta^\star;\mathcal{H}_{t-1}\mid \mathcal{E}_{t-1})+1\right).$$

• continuing this process gives

$$I_{\chi^2}(\theta^{\star}; \mathcal{H}_t \mid \mathcal{E}_t) + 1 \leq \frac{\exp(\sum_{s \leq t} \varepsilon_s^2)}{\mathbb{P}(\mathcal{E}_t)^2}.$$

• recursion of error probability:

$$\mathbb{P}(\mathcal{E}_{t+1}) = \mathbb{P}(\mathcal{E}_t) \cdot \mathbb{P}(\langle \theta^{\star}, a_{t+1} \rangle \leq \varepsilon_{t+1} \mid \mathcal{E}_t)$$

 $\geq \mathbb{P}(\mathcal{E}_t) \left(1 - e^{-\Theta(d\varepsilon_{t+1}^2)} \sqrt{I_{\chi^2}(\theta^{\star}; \mathcal{H}_t \mid \mathcal{E}_t) + 1} \right)$

٠

• let $\mathcal{E}_t=\cap_{s\leq t}\{\langle heta^\star,a_s
angle\leq arepsilon_s\}$ be the error event

• upper bound of conditioned χ^2 -informativity:

$$I_{\chi^2}(\theta^\star;\mathcal{H}_t\mid \mathcal{E}_t)+1 \leq \frac{\exp(\varepsilon_t^2)}{\mathbb{P}(\mathcal{E}_t\mid \mathcal{E}_{t-1})^2}\left(I_{\chi^2}(\theta^\star;\mathcal{H}_{t-1}\mid \mathcal{E}_{t-1})+1\right).$$

• continuing this process gives

$$I_{\chi^2}(\theta^{\star}; \mathcal{H}_t \mid \mathcal{E}_t) + 1 \leq \frac{\exp(\sum_{s \leq t} \varepsilon_s^2)}{\mathbb{P}(\mathcal{E}_t)^2}.$$

• recursion of error probability:

$$egin{aligned} \mathbb{P}(\mathcal{E}_{t+1}) &= \mathbb{P}(\mathcal{E}_t) \cdot \mathbb{P}(\langle heta^\star, oldsymbol{a}_{t+1}
angle \leq arepsilon_{t+1} \mid \mathcal{E}_t) \ &\geq \mathbb{P}(\mathcal{E}_t) \left(1 - e^{-\Theta(darepsilon_{t+1}^2)} \sqrt{I_{\chi^2}(heta^\star; \mathcal{H}_t \mid \mathcal{E}_t) + 1}
ight) \ &\geq \mathbb{P}(\mathcal{E}_t) - \underbrace{e^{-\Theta(darepsilon_{t+1}^2) + rac{1}{2}\sum_{s \leq t}arepsilon_s^2}_{=\delta}}_{=\delta}. \end{aligned}$$

Discussions and Further Questions

• Eluder-UCB algorithm [Russo and Van Roy 2014]:

• Eluder-UCB algorithm [Russo and Van Roy 2014]: at each time t,

 \rightarrow form the least squares estimator $\hat{\theta}_t = \arg \min_{\theta} \sum_{s < t} (r_s - f(\langle \theta, a_s \rangle))^2$;

• Eluder-UCB algorithm [Russo and Van Roy 2014]: at each time t,

 \rightarrow form the least squares estimator $\hat{\theta}_t = \arg \min_{\theta} \sum_{s < t} (r_s - f(\langle \theta, a_s \rangle))^2$;

 $\rightarrow\,$ construct the confidence set of $\theta^{\star}:$

$${\mathcal C}_t = \left\{ heta \in \mathbb{S}^{d-1} : \sum_{s < t} (f(\langle heta, {\sf a}_s \rangle) - f(\langle \widehat{ heta}_t, {\sf a}_s \rangle))^2 = \widetilde{O}(d)
ight\};$$

- Eluder-UCB algorithm [Russo and Van Roy 2014]: at each time t,
 - \rightarrow form the least squares estimator $\hat{\theta}_t = \arg \min_{\theta} \sum_{s < t} (r_s f(\langle \theta, a_s \rangle))^2$;
 - $\rightarrow\,$ construct the confidence set of $\theta^{\star}:$

$$\mathcal{C}_t = \left\{ heta \in \mathbb{S}^{d-1} : \sum_{s < t} (f(\langle heta, \mathbf{a}_s \rangle) - f(\langle \widehat{ heta}_t, \mathbf{a}_s \rangle))^2 = \widetilde{O}(d)
ight\};$$

 \rightarrow play action $a_t = \arg \max_a \max_{\theta \in C_t} f(\langle \theta, a \rangle).$

- Eluder-UCB algorithm [Russo and Van Roy 2014]: at each time t,
 - \rightarrow form the least squares estimator $\widehat{\theta}_t = \arg \min_{\theta} \sum_{s < t} (r_s f(\langle \theta, a_s \rangle))^2$;
 - $\rightarrow\,$ construct the confidence set of $\theta^{\star}:$

$$\mathcal{C}_t = \left\{ \theta \in \mathbb{S}^{d-1} : \sum_{s < t} (f(\langle \theta, \mathbf{a}_s \rangle) - f(\langle \widehat{\theta}_t, \mathbf{a}_s \rangle))^2 = \widetilde{O}(d) \right\};$$

 \rightarrow play action $a_t = \arg \max_a \max_{\theta \in C_t} f(\langle \theta, a \rangle).$

Theorem (lower bound for Eluder-UCB)

For every f, there exist a bandit instance such that for (a certain tie-breaking rule of) Eluder-UCB, achieving a sublinear regret requires

$$T \gtrsim \max_{K} \min\left\{K, \frac{d}{f(\sqrt{(\log K)/d})^2}\right\}$$

• Eluder-UCB algorithm [Russo and Van Roy 2014]: at each time t,

 \rightarrow form the least squares estimator $\hat{\theta}_t = \arg \min_{\theta} \sum_{s < t} (r_s - f(\langle \theta, a_s \rangle))^2$;

 $\rightarrow\,$ construct the confidence set of $\theta^{\star}:$

$$\mathcal{C}_t = \left\{ \theta \in \mathbb{S}^{d-1} : \sum_{s < t} (f(\langle \theta, \mathbf{a}_s \rangle) - f(\langle \widehat{\theta}_t, \mathbf{a}_s \rangle))^2 = \widetilde{O}(d) \right\};$$

 \rightarrow play action $a_t = \arg \max_a \max_{\theta \in C_t} f(\langle \theta, a \rangle).$

Theorem (lower bound for Eluder-UCB)

For every f, there exist a bandit instance such that for (a certain tie-breaking rule of) Eluder-UCB, achieving a sublinear regret requires

$$T \gtrsim \max_{K} \min\left\{K, \frac{d}{f(\sqrt{(\log K)/d})^2}\right\}$$

• for $f(x) = x^3$, Eluder-UCB requires $T \gtrsim d^4$, but optimal is $T \gtrsim d^3$
Online regression oracle model [Foster et al. 2020]: for any adversarial sequence {a_t}, oracle outputs {θ
_t} such that

$$\sum_{t=1}^{T} (f(\langle \theta^{\star}, a_t \rangle) - f(\langle \widehat{\theta}_t, a_t \rangle))^2 \leq \operatorname{Reg}_{\mathsf{Sq}}(T)$$

Online regression oracle model [Foster et al. 2020]: for any adversarial sequence {a_t}, oracle outputs {θ
_t} such that

$$\sum_{t=1}^{T} (f(\langle \theta^{\star}, a_t \rangle) - f(\langle \widehat{\theta}_t, a_t \rangle))^2 \leq \mathsf{Reg}_{\mathsf{Sq}}(T)$$

 \rightarrow learner only observes $\widehat{ heta}_t$ in the oracle model, but not r_t ;

Online regression oracle model [Foster et al. 2020]: for any adversarial sequence {a_t}, oracle outputs {θ
_t} such that

$$\sum_{t=1}^{T} (f(\langle \theta^{\star}, a_t \rangle) - f(\langle \widehat{\theta}_t, a_t \rangle))^2 \leq \mathsf{Reg}_{\mathsf{Sq}}(T)$$

- ightarrow learner only observes $\widehat{ heta}_t$ in the oracle model, but not $r_t;$
- \rightarrow a natural choice of $\operatorname{Reg}_{\operatorname{Sq}}(T)$ is $\widetilde{O}(d)$.

Online regression oracle model [Foster et al. 2020]: for any adversarial sequence {a_t}, oracle outputs {θ
_t} such that

$$\sum_{t=1}^{T} (f(\langle \theta^{\star}, \boldsymbol{a}_t \rangle) - f(\langle \widehat{\theta}_t, \boldsymbol{a}_t \rangle))^2 \leq \mathsf{Reg}_{\mathsf{Sq}}(T)$$

→ learner only observes $\hat{\theta}_t$ in the oracle model, but not r_t ; → a natural choice of $\operatorname{Reg}_{S_q}(T)$ is $\widetilde{O}(d)$.

Theorem (lower bound for the oracle model)

For every f, there exist a bandit instance under the oracle model such that for every algorithm, achieving a sublinear regret requires

$$T \gtrsim \max_{K} \min\left\{K, \frac{d}{f(\sqrt{(\log K)/d})^2}\right\}$$

Online regression oracle model [Foster et al. 2020]: for any adversarial sequence {a_t}, oracle outputs {θ_t} such that

$$\sum_{t=1}^{T} (f(\langle \theta^{\star}, \boldsymbol{a}_t \rangle) - f(\langle \widehat{\theta}_t, \boldsymbol{a}_t \rangle))^2 \leq \mathsf{Reg}_{\mathsf{Sq}}(T)$$

→ learner only observes $\hat{\theta}_t$ in the oracle model, but not r_t ; → a natural choice of $\operatorname{Reg}_{S_q}(T)$ is $\widetilde{O}(d)$.

Theorem (lower bound for the oracle model)

For every f, there exist a bandit instance under the oracle model such that for every algorithm, achieving a sublinear regret requires

$$T \gtrsim \max_{K} \min \left\{ K, \frac{d}{f(\sqrt{(\log K)/d})^2} \right\}$$

• Key modeling difference: in oracle model, choosing repeated action may not reduce the estimation error

• in linear bandit, fewer actions lead to smaller regret

- in linear bandit, fewer actions lead to smaller regret
 - \rightarrow minimax regret decreases from $\Theta(d\sqrt{T})$ to $\Theta(\sqrt{dT\log K})$ with K actions

- in linear bandit, fewer actions lead to smaller regret
 - \rightarrow minimax regret decreases from $\Theta(d\sqrt{T})$ to $\Theta(\sqrt{dT \log K})$ with K actions
 - $\rightarrow\,$ intuition: UCB needs to construct fewer confidence intervals

- in linear bandit, fewer actions lead to smaller regret
 - \rightarrow minimax regret decreases from $\Theta(d\sqrt{T})$ to $\Theta(\sqrt{dT \log K})$ with K actions
 - $\rightarrow\,$ intuition: UCB needs to construct fewer confidence intervals
- does similar phenomenon hold for non-linear bandits?

- in linear bandit, fewer actions lead to smaller regret
 - \rightarrow minimax regret decreases from $\Theta(d\sqrt{T})$ to $\Theta(\sqrt{dT \log K})$ with K actions
 - $\rightarrow\,$ intuition: UCB needs to construct fewer confidence intervals
- does similar phenomenon hold for non-linear bandits?

Theorem (lower bound for finite actions)

For every link function f and K = poly(d), there exist an K-armed ridge bandit instance such that achieving a sublinear regret requires

$$T\gtrsim \min\left\{K,rac{1}{f(1/\sqrt{d})^2}
ight\}.$$

- in linear bandit, fewer actions lead to smaller regret
 - \rightarrow minimax regret decreases from $\Theta(d\sqrt{T})$ to $\Theta(\sqrt{dT \log K})$ with K actions
 - $\rightarrow\,$ intuition: UCB needs to construct fewer confidence intervals
- does similar phenomenon hold for non-linear bandits?

Theorem (lower bound for finite actions)

For every link function f and K = poly(d), there exist an K-armed ridge bandit instance such that achieving a sublinear regret requires

$$T\gtrsim \min\left\{K,rac{1}{f(1/\sqrt{d})^2}
ight\}.$$

• implication: for $f(x) = x^3$, the fixed cost for the finite-action problem is already d^3 , same as the infinite-action problem

- in linear bandit, fewer actions lead to smaller regret
 - \rightarrow minimax regret decreases from $\Theta(d\sqrt{T})$ to $\Theta(\sqrt{dT \log K})$ with K actions
 - $\rightarrow\,$ intuition: UCB needs to construct fewer confidence intervals
- does similar phenomenon hold for non-linear bandits?

Theorem (lower bound for finite actions)

For every link function f and K = poly(d), there exist an K-armed ridge bandit instance such that achieving a sublinear regret requires

$$T\gtrsim \min\left\{K,rac{1}{f(1/\sqrt{d})^2}
ight\}.$$

- implication: for $f(x) = x^3$, the fixed cost for the finite-action problem is already d^3 , same as the infinite-action problem
- reason: the learner cannot explore every direction in the initialization phase

Unit sphere vs unit ball

What happens if we assume that $\theta^* \in \mathbb{B}^d$ instead of $\theta^* \in \mathbb{S}^{d-1}$?

Unit sphere vs unit ball

What happens if we assume that $\theta^{\star} \in \mathbb{B}^d$ instead of $\theta^{\star} \in \mathbb{S}^{d-1}$?

Theorem (modified upper bound)

Under monotonicity and local linearity of f:

$$\operatorname{MinmaxReg}(T, d, f) \lesssim \max_{r \in [0,1]} \min \left\{ d^2 \frac{f(r)}{r^4} \int_{r/\sqrt{d}}^{r/2} \frac{d(x^2)}{\max_{r/\sqrt{d} \leq y \leq x} f'(y)^2} + d\sqrt{T}, Tf(r) \right\}$$

Theorem (modified lower bound)

Under monotonicity and local linearity of f:

$$\mathsf{MinmaxReg}(T, d, f) \gtrsim \max_{r \in [0,1]} \min \left\{ d \frac{f(r)}{r^2} \int_{r/\sqrt{d}}^{r/2} \frac{\mathsf{d}(x^2)}{f(x)^2} + d\sqrt{T}, Tf(r) \right\}.$$

Unit sphere vs unit ball

What happens if we assume that $\theta^* \in \mathbb{B}^d$ instead of $\theta^* \in \mathbb{S}^{d-1}$?

Theorem (modified upper bound)

Under monotonicity and local linearity of f:

$$\operatorname{MinmaxReg}(T, d, f) \lesssim \max_{r \in [0,1]} \min \left\{ d^2 \frac{f(r)}{r^4} \int_{r/\sqrt{d}}^{r/2} \frac{d(x^2)}{\max_{r/\sqrt{d} \le y \le x} f'(y)^2} + d\sqrt{T}, Tf(r) \right\}$$

Theorem (modified lower bound)

Under monotonicity and local linearity of f:

$$\operatorname{MinmaxReg}(T, d, f) \gtrsim \max_{r \in [0,1]} \min \left\{ d \frac{f(r)}{r^2} \int_{r/\sqrt{d}}^{r/2} \frac{d(x^2)}{f(x)^2} + d\sqrt{T}, Tf(r) \right\}.$$

minimax regret often exhibits only one elbow instead of two

Further questions

$$I_t - I_{t-1} \leq \mathsf{Var}(f(\langle \theta^{\star}, a_t \rangle) \mid a_t, \mathcal{H}_{t-1}) \stackrel{?}{\lesssim} \max_{y \leq \varepsilon_t} \frac{f'(y)^2}{d}$$

$$I_t - I_{t-1} \leq \mathsf{Var}(f(\langle \theta^{\star}, a_t \rangle) \mid a_t, \mathcal{H}_{t-1}) \stackrel{?}{\lesssim} \max_{y \leq \varepsilon_t} \frac{f'(y)^2}{d}$$

• analyze the learning trajectory of information-directed sampling $a_t = \arg \max_a I(\theta^*; r_t \mid \mathcal{H}_{t-1}, a_t = a)$

$$I_t - I_{t-1} \leq \mathsf{Var}(f(\langle \theta^{\star}, a_t \rangle) \mid a_t, \mathcal{H}_{t-1}) \stackrel{?}{\lesssim} \max_{y \leq \varepsilon_t} \frac{f'(y)^2}{d}$$

- analyze the learning trajectory of information-directed sampling $a_t = \arg \max_a I(\theta^*; r_t \mid \mathcal{H}_{t-1}, a_t = a)$
- more general reward functions such as $f_{\theta}(a) = \sum_{i=1}^{m} f_i(\langle \theta_i, a \rangle)$

$$I_t - I_{t-1} \leq \mathsf{Var}(f(\langle \theta^{\star}, a_t \rangle) \mid a_t, \mathcal{H}_{t-1}) \stackrel{?}{\lesssim} \max_{y \leq \varepsilon_t} \frac{f'(y)^2}{d}$$

- analyze the learning trajectory of information-directed sampling $a_t = \arg \max_a I(\theta^*; r_t \mid \mathcal{H}_{t-1}, a_t = a)$
- more general reward functions such as $f_{\theta}(a) = \sum_{i=1}^{m} f_i(\langle \theta_i, a \rangle)$
- more systematic methods for exploration in the initialization phase

$$I_t - I_{t-1} \leq \mathsf{Var}(f(\langle \theta^{\star}, a_t \rangle) \mid a_t, \mathcal{H}_{t-1}) \stackrel{?}{\lesssim} \max_{y \leq \varepsilon_t} \frac{f'(y)^2}{d}$$

- analyze the learning trajectory of information-directed sampling $a_t = \arg \max_a I(\theta^*; r_t \mid \mathcal{H}_{t-1}, a_t = a)$
- more general reward functions such as $f_{\theta}(a) = \sum_{i=1}^{m} f_i(\langle \theta_i, a \rangle)$
- more systematic methods for exploration in the initialization phase
- more complicated settings such as contextual bandits and RL

Take-home message:

• there could be a phase transition in the regret of non-linear bandits

- there could be a phase transition in the regret of non-linear bandits
- in the initialization phase, the learner needs algorithms beyond UCB to explore a good initial action, which incurs a fixed cost

- there could be a phase transition in the regret of non-linear bandits
- in the initialization phase, the learner needs algorithms beyond UCB to explore a good initial action, which incurs a fixed cost
- in the learning phase, the learner can employ UCB-type algorithm around the good initial action

- there could be a phase transition in the regret of non-linear bandits
- in the initialization phase, the learner needs algorithms beyond UCB to explore a good initial action, which incurs a fixed cost
- in the learning phase, the learner can employ UCB-type algorithm around the good initial action
- learning trajectory of the initialization phase could be characterized by proper differential equations

- there could be a phase transition in the regret of non-linear bandits
- in the initialization phase, the learner needs algorithms beyond UCB to explore a good initial action, which incurs a fixed cost
- in the learning phase, the learner can employ UCB-type algorithm around the good initial action
- learning trajectory of the initialization phase could be characterized by proper differential equations
- traditional learning algorithms may fail to obtain the optimal initialization cost

Take-home message:

- there could be a phase transition in the regret of non-linear bandits
- in the initialization phase, the learner needs algorithms beyond UCB to explore a good initial action, which incurs a fixed cost
- in the learning phase, the learner can employ UCB-type algorithm around the good initial action
- learning trajectory of the initialization phase could be characterized by proper differential equations
- traditional learning algorithms may fail to obtain the optimal initialization cost

Thank You!