

INSTANCE-WISE MINIMAX-OPTIMAL ALGORITHMS FOR LOGISTIC BANDITS

MARC ABEILLE¹, LOUIS FAURY^{1,2}, CLÉMENT CALAUZÈNES¹

¹Criteo AI Lab, ²LTCI Télécom Paris

MOTIVATION

Toward non-linear reward model.

- Parametric bandit results mostly concern the linear setting,
- Non-linearity often arises in real-world application,
- Impact of non-linearity on the exploration-exploitation tradeoff is poorly understood.

The logistic bandit setting.

- Non-linear reward signal,
- Compact and minimal setting,
- Widely used for practical applications.

We characterize the impact of non-linearity for Logistic Bandit:

- First problem-dependent lower-bound,
- Minimax-optimal algorithm.

THE LOGISTIC BANDIT PROBLEM

The reward model.

- $\mathcal{X} \subset \mathbb{R}^d$ is the arm set,
- $r(x) \in \{0, 1\}$ is the reward associated with arm $x \in \mathcal{X}$,
- $\theta_* \in \mathbb{R}^d$ *unknown* parameter.

[Binary reward]

$$r(x) \sim \text{Bernoulli}(\mu(x^\top \theta_*))$$

[Non-linear link function]

$$\mu(z) = (1 + \exp(-z))^{-1}$$

The learning problem.

At each step $t \leq T$:

- Choose a arm $x_t \in \mathcal{X}$,
- Receive $r(x_t)$,

Objective: Minimize Regret

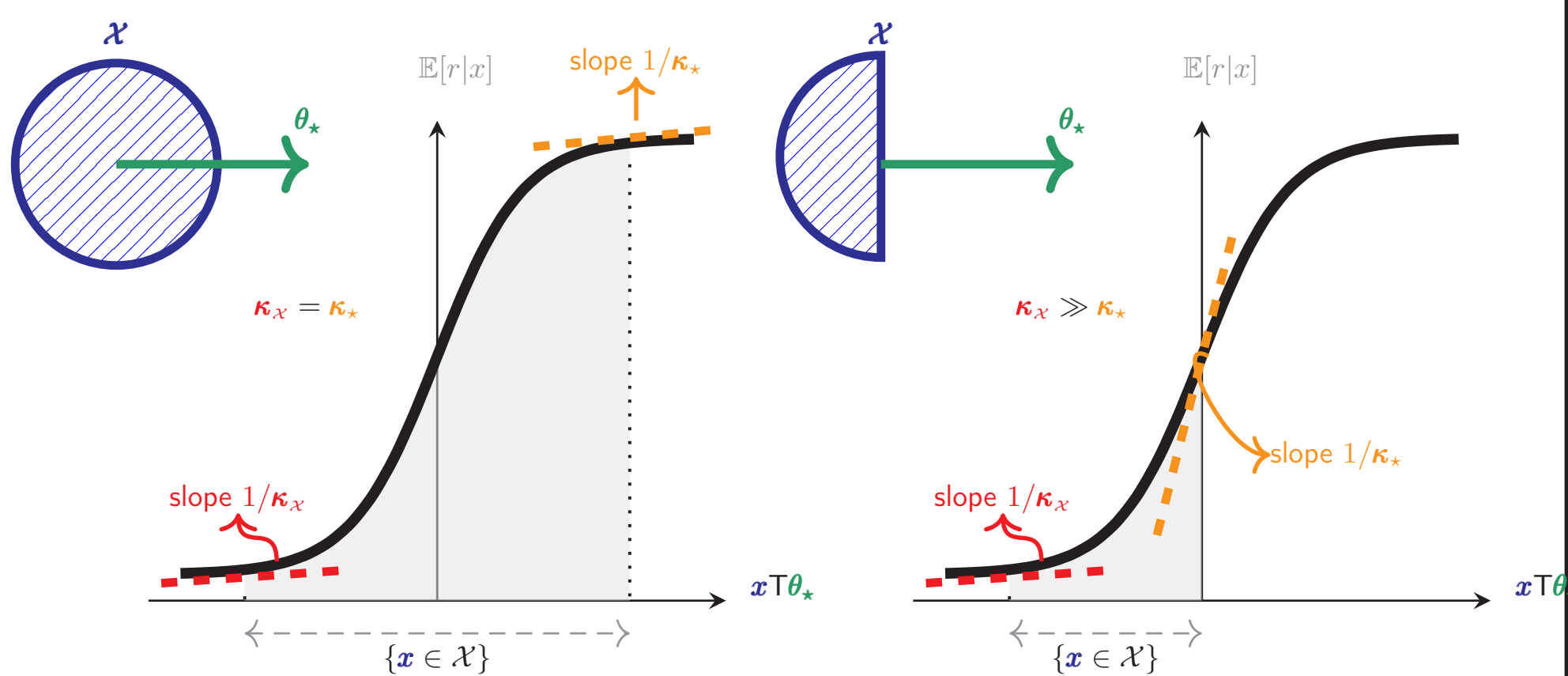
$$R_{\theta_*}(T) = \sum_{t=1}^T \left[\max_{x \in \mathcal{X}} \mu(x^\top \theta_*) - \mu(x_t^\top \theta_*) \right].$$

Quantifying non-linearity We consider two important *problem-dependent* constants:

$$\kappa_* := 1/\dot{\mu}(\max_{x \in \mathcal{X}} x^\top \theta_*)$$

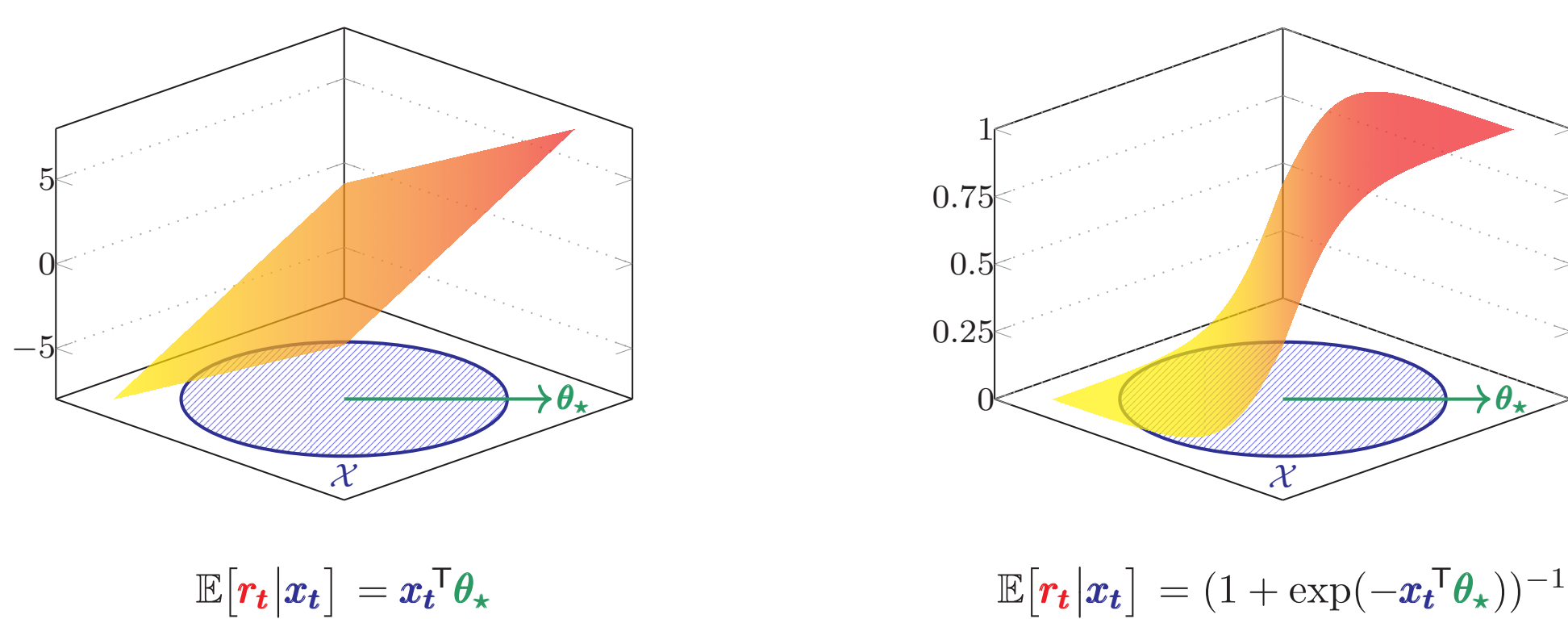
$$\kappa_{\mathcal{X}} := 1/\min_{x \in \mathcal{X}} \dot{\mu}(x^\top \theta_*)$$

- κ_* : "distance to linearity" around the optimal action,
- $\kappa_{\mathcal{X}}$: worst-case "distance to linearity" over the decision set.



NON-LINEARITY: BLESSING OR CURSE ?

From LB to LogB

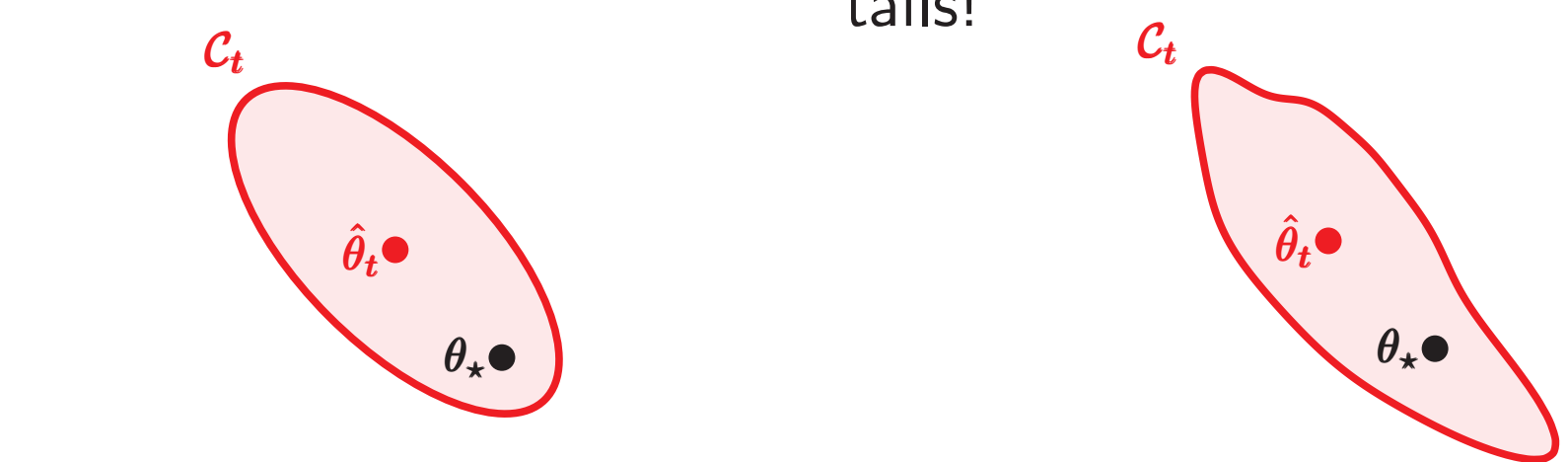


Impact on the learning.

Different richness of information associated with sampling an arm:

LB Same everywhere,

LogB High in the center, low in the tails!



✓ Despite non-linearity → available conf. set. C_t for **LogB**

[Faury et al., Improved Optimistic Algorithms for Logistic Bandits, ICML'20]

✗ Some regions are *harder* to learn that other → the conf. set. C_t is *not* an ellipsoid!

Impact on the predicted performance

✓ **LogB** deviation in parameters → little to no deviation in performance *in the tails*

$$\|\theta - \theta_*\| = \delta \Rightarrow \mu(x^\top \theta) \simeq \mu(x^\top \theta_*)$$

Open question: does *easy* prediction cancel out *hard* learning?

RELATED WORK AND CONTRIBUTIONS

Related work.

[Filippi et al., NIPS'10]

$$R_{\theta_*}(T) \lesssim \kappa_{\mathcal{X}} d \sqrt{T}$$

[Faury et al., ICML'20]

$$R_{\theta_*}(T) \lesssim d \sqrt{T} + \kappa_{\mathcal{X}}$$

[Dong et al., COLT'19]

In the worst case, $R_{\theta_*}(T)$ must increase with $\kappa_{\mathcal{X}}$

Contributions.

Theorem 1. (Regret Upper Bound) The regret of OFU-Log satisfies with high-probability:

$$R(T) \lesssim d \sqrt{\frac{T}{\kappa_*}} + (\kappa_{\mathcal{X}}).$$

Theorem 2. (Local Lower Bound) Let $\mathcal{X} = \mathcal{S}_d(0, 1)$, for any θ_* and T large enough, it exists $\epsilon > 0$ small enough s.t.

$$\min_{\pi} \max_{\|\theta - \theta_*\| \leq \epsilon} \mathbb{E}[R_{\theta}^{\pi}(T)] = \Omega\left(d \sqrt{\frac{T}{\kappa_*}}\right).$$

OPTIMISTIC ALGORITHM OFULog

for $t = \{0, \dots, T\}$ **do**

Set $\lambda_t \leftarrow d \log(t)$.

(Learning) Solve $\hat{\theta}_t = \arg \min_{\theta} \mathcal{L}_t(\theta)$.

(Planning) Solve $(x_t, \theta_t) \in \arg \max_{\mathcal{X}, \mathcal{C}_t(\delta)} \mu(x^\top \theta)$.

Play x_t and observe reward r_{t+1} .

end for

where $\mathcal{L}_t(\theta)$ and $\mathcal{C}_t(\delta)$ are the log-likelihood function and confidence set associated with the learning problem.

IDEAS BEHIND THE LOWER BOUND

Objective and approach

- we shoot for a *problem-dependent* lower-bound
- standard approach consider worst-case over *all possible instance*
- inspired by [Simchowitz et al., ICML'20] → *local* lower-bound
- consider worst-case over all nearby alternative around a given *problem instance*.

ideas

- we consider a given instance parametrized by θ_* ,
- let π denote a policy that outputs a sequence of arms, and $R_{\theta_*}^{\pi}(T)$ the induced expected regret.

Small regret ↔ low exploration

$$R_{\theta_*}^{\pi}(T) \propto 1/\kappa_* \sum_{t=1}^T \|x_t - x_*(\theta_*)\|^2, \quad x_*(\theta_*) = \arg \max_{x \in \mathcal{X}} \mu(x^\top \theta_*)$$

- $R_{\theta_*}^{\pi}(T)$ small ↔ $x_t \simeq x_*(\theta_*)$,
- directions orthogonal to $x_*(\theta_*)$ are poorly explored!
- Larger* κ_* → *smaller* impact when deviating from $x_*(\theta_*)$!

Low exploration ↔ large set of plausible alternative

- We quantify the *similarity* between instances θ, θ_* under policy π by the *discrepancy*

$$D_{\text{KL}}(\mathbb{P}_{\theta}^{\pi}, \mathbb{P}_{\theta_*}^{\pi})$$

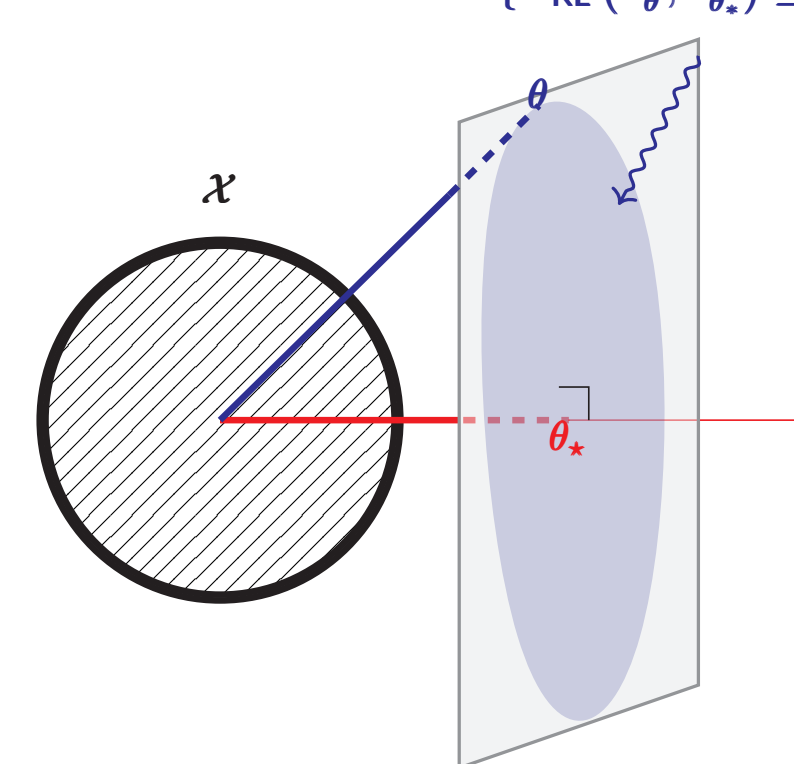
large $D_{\text{KL}}(\mathbb{P}_{\theta}^{\pi}, \mathbb{P}_{\theta_*}^{\pi}) \rightarrow$ *easy* to distinguish θ and θ_* under π ,

small $D_{\text{KL}}(\mathbb{P}_{\theta}^{\pi}, \mathbb{P}_{\theta_*}^{\pi}) \rightarrow$ *hard* to distinguish θ and θ_* under π ,

$$D_{\text{KL}}(\mathbb{P}_{\theta}^{\pi}, \mathbb{P}_{\theta_*}^{\pi}) \propto \sqrt{\frac{T}{\kappa_*}} \|\theta - \theta_*\|^2$$

- large* κ_* degrades the richness of acquired information,

→ $D_{\text{KL}}(\mathbb{P}_{\theta}^{\pi}, \mathbb{P}_{\theta_*}^{\pi})$ decreases with κ_* ,



Tension and trade-off

- Policy π cannot perform well on two *distinct* instances,
- but may not yield *similar* information.

Trade-off

- Let π perform well for θ_* ,
- consider an alternative instance θ such that $\|\theta - \theta_*\|^2 \approx \sqrt{\frac{\kappa_*}{T}}$,
- the regret of π for the instance θ must be large:

$$R_{\theta}^{\pi}(T) \approx 1/\kappa_* \sum_{t=1}^T \|x_t - x_*(\theta)\|^2 \approx 1/\kappa_* \sum_{t=1}^T \|x_*(\theta_*) - x_*(\theta)\|^2 \approx T \|\theta_* - \theta\|^2 / \kappa_* \approx \sqrt{T/\kappa_*}.$$

IDEAS BEHIND THE UPPER BOUND

Permanent and transitory regimes

- Regret decomposition:

$$R_{\theta_*}(T) = \underbrace{R^{\text{perm}}(T)}_{\tilde{O}(\sqrt{T})} + \underbrace{R^{\text{trans}}(T)}_{\tilde{O}(1)}$$

Permanent regime: intuition.

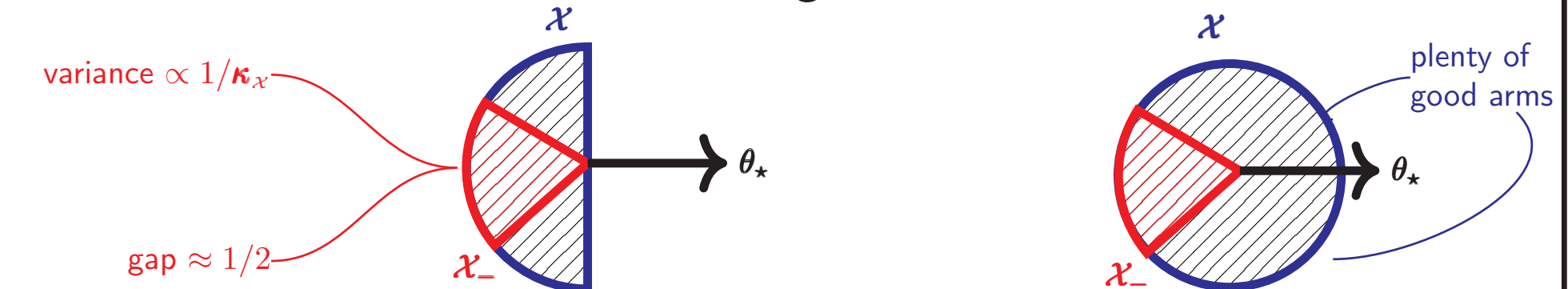
- Sublinear regret ⇒ play mostly around the best arm x_* .
↪ Almost a linear bandit with slope $1/\kappa_*$.
- A finer analysis is coherent with this conceptual argument:

$$R^{\text{perm}}(T) \leq d \sqrt{\sum_{t=1}^T \dot{\mu}(x_t^\top \theta_*)} \approx d \sqrt{T/\kappa_*}.$$

- Formal proof: thanks to self-concordance property.

Transitory regime and detrimental arms.

- Detrimental arm* \mathcal{X}_- : low-information and large gap
↪ far left tail of the reward signal:



- Transitory regime: how long before discarding detrimental arms:

$$R^{\text{trans}}_{\theta_*}(T) \leq \min\left(\kappa_{\mathcal{X}}, \sum_{t=1}^T \mathbb{1}(x_t \in \mathcal{X}_-)\right)$$

- Fast if the proportion of detrimental arm is small:

Proposition 1. (Transitory regret) With h.p :

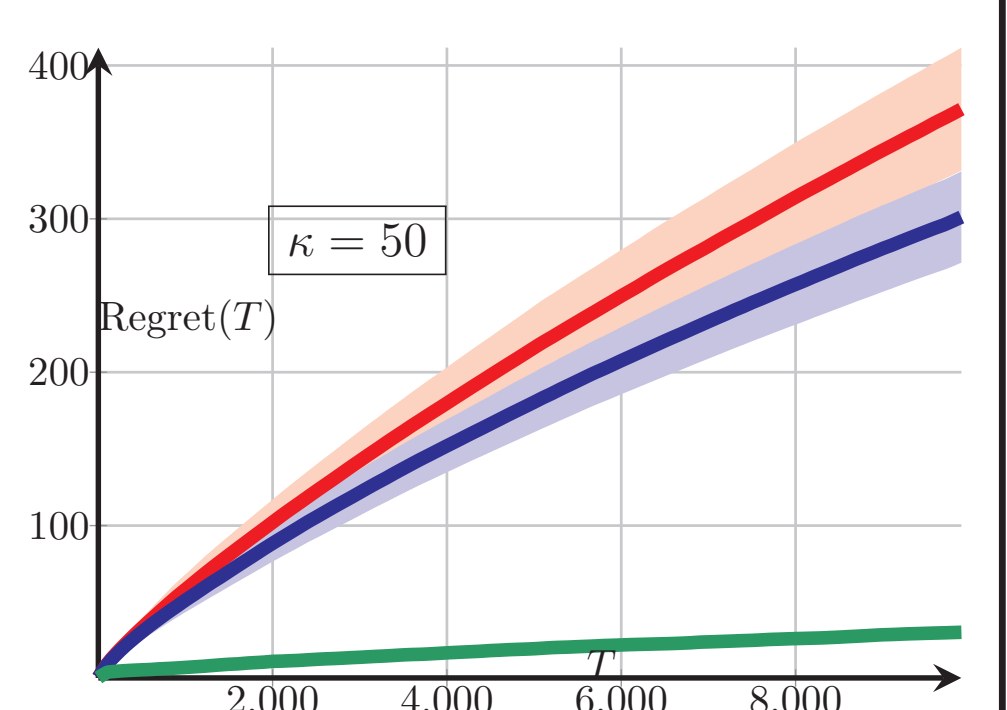
$$R^{\text{trans}}_{\theta_*}(T) \lesssim_T d^2 + dK \quad \text{if } |\mathcal{X}_-| \leq K,$$

$$R^{\text{trans}}_{\theta_*}(T) \lesssim_T d^3 \quad \text{if } \mathcal{X} = \mathcal{B}_d(0, 1).$$

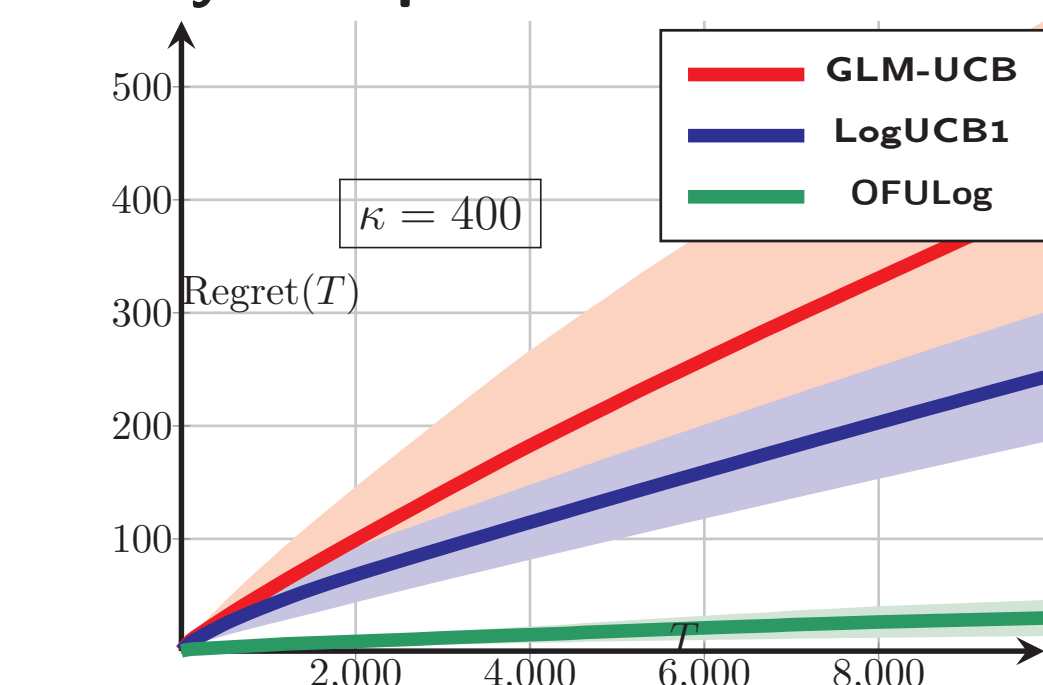
↪ independent of $\kappa_{\mathcal{X}}$ for reasonable configurations.

ALGORITHMIC DESIGN AND EXPERIMENTS

- Convex relaxation.** bla



- Toy experiment.**



CONCLUSION

- Blah
- Blah
- Blah
- Blah

REFERENCES

- Y. Abbasi-Yadkori, Cs. Szepesvári. Regret Bounds for the Adaptive Control of Linear Quadratic Systems. In *Proceedings of COLT*, 2011.
- S. Bittanti and M.C. Campi. Adaptive control of linear time invariant systems: the "bet on the best" principle. *Commun. Inf. Syst. Volume 6*, 2006.
- M.K.S. Faradonbeh, A. Tewari, and G. Michailidis. Finite Time Analysis of Optimal Adaptive Policies for Linear-Quadratic Systems. *arXiv:1711.07230*.
- Y. Ouyang, M. Gagrani, and R. Jain. Learning-based Control of Unknown Linear Systems with Thompson Sampling. *arXiv:1709.04047*.
- I. Osband, and B. Van Roy. Model-based Reinforcement Learning and the Eluder Dimension *NIPS*, 2014.
- M. Abeille, and A. Lazaric. Thompson Sampling for Linear-Quadratic Control Problems. *Proceedings of AISTATS*, 2017.