

# INSTANCE-WISE MINIMAX-OPTIMAL ALGORITHMS FOR LOGISTIC BANDITS

MARC ABEILLE<sup>1</sup>, LOUIS FAURY<sup>1,2</sup>, CLÉMENT CALAUZÈNES<sup>1</sup>

<sup>1</sup>Criteo AI Lab, <sup>2</sup>LTCI 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

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

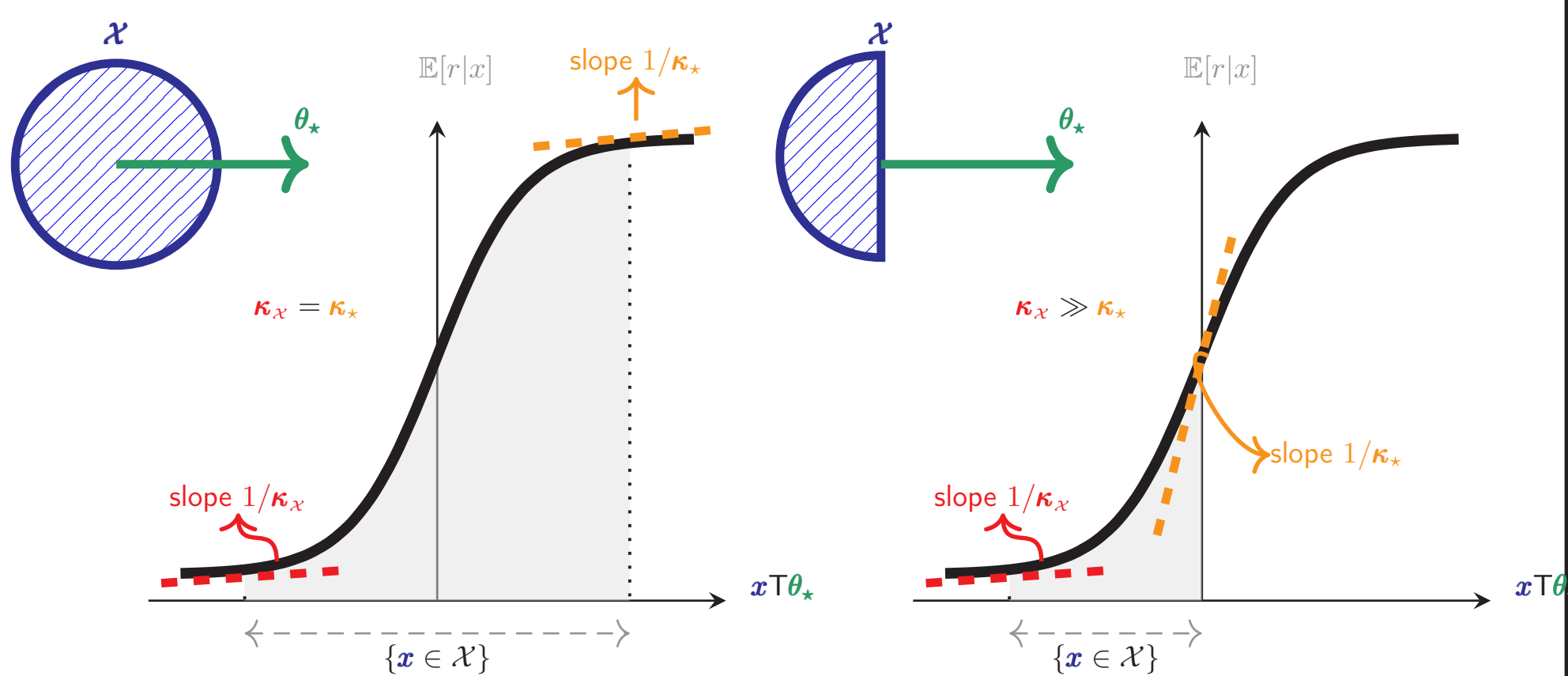
### Quantifying non-linearity

We consider two important *problem-dependent* constants:

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

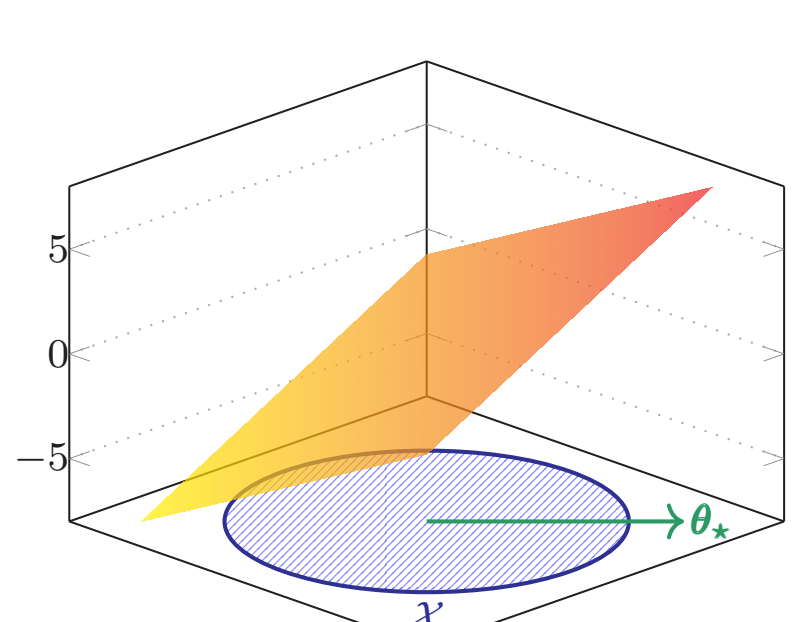
$$\kappa_{\mathcal{X}} := 1/\min_{x \in \mathcal{X}} \mu(x^\top \theta_*)$$

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

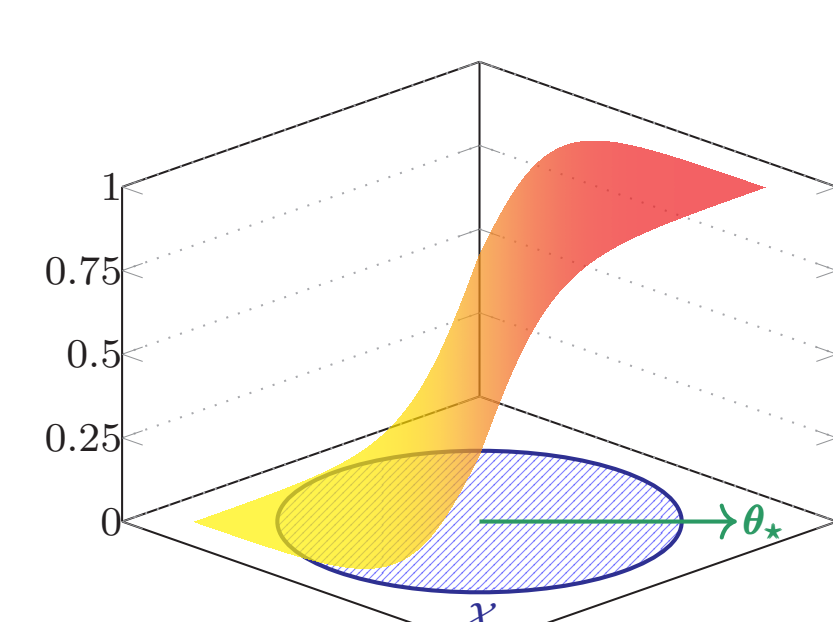


## NON-LINEARITY: BLESSING OR CURSE ?

### From LB to LogB



$$\mathbb{E}[r_t | x_t] = x_t^\top \theta_*$$



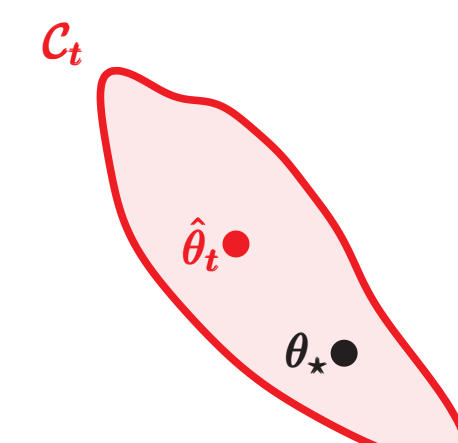
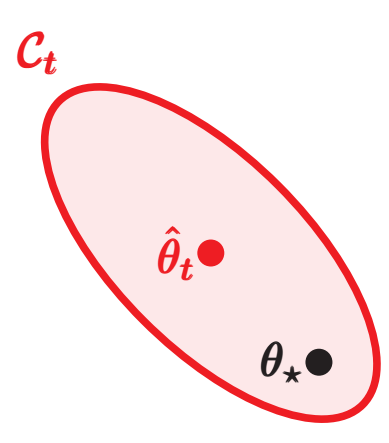
$$\mathbb{E}[r_t | x_t] = (1 + \exp(-x_t^\top \theta_*))^{-1}$$

### 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 than others → 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*

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

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:

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

Illustration: if  $\mathcal{X} = \{\|x\| \leq 1\}$  then  $\kappa_* = \kappa_{\mathcal{X}} \approx \exp(\|\theta_*\|)$ :

$$R_{\theta_*}(T) \lesssim d \sqrt{T/\kappa_*}, \lesssim d \exp(-\|\theta_*\|/2) \sqrt{T}$$

- ↗ the more non-linear the model, the smaller the regret!
- ↗ exponential improvement over existing bounds.

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

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

where  $\epsilon$  is small enough that  $\forall \theta \in \{\|\theta - \theta_*\| \leq \epsilon\}$  we have  $\kappa_*(\theta) = \Theta(\kappa_*)$ .

- ↗ the upper-bound is *optimal* for large  $T$ .
- ↗ the lower-bound holds for all instances  $\theta_*$ .

## IDEAS BEHIND THE LOWER BOUND

### Objective and approach

- We shoot for a *problem-dependent* lower-bound,
- usual approaches consider worst-case over *all possible instances*,
- inspired by [Simchowitz et al., ICML'20] → *local* lower-bound,
- worst-case over nearby alternatives around a given *problem instance*.

### High-level idea

- We consider a given instance parametrized by  $\theta_*$ ,
- let  $\pi$  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  $\kappa_*$  → smaller impact when deviating from  $x_*(\theta_*)$ !**

### Low exploration ↔ large set of plausible alternative

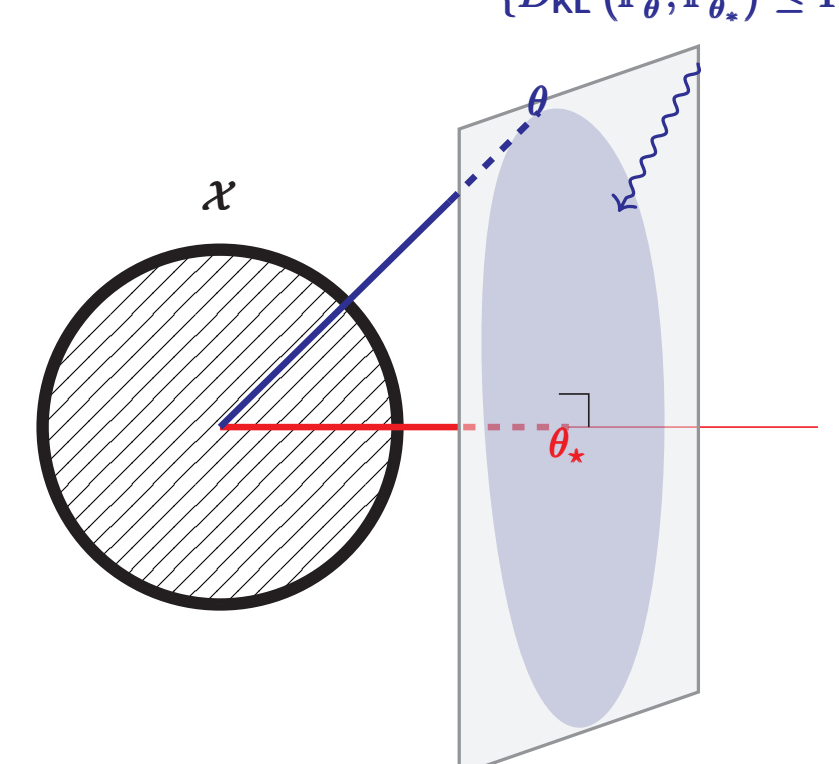
- We quantify the *similarity* between instances  $\theta, \theta_*$  under policy  $\pi$  by the *discrepancy*

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

- *large*  $D_{\text{KL}}(\mathbb{P}_{\theta}^{\pi}, \mathbb{P}_{\theta_*}^{\pi})$  → *easy* to distinguish  $\theta$  and  $\theta_*$  under  $\pi$ ,
- *small*  $D_{\text{KL}}(\mathbb{P}_{\theta}^{\pi}, \mathbb{P}_{\theta_*}^{\pi})$  → *hard* to distinguish  $\theta$  and  $\theta_*$  under  $\pi$ .

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

- *large*  $\kappa_*$  degrades the richness of acquired information,
- $D_{\text{KL}}(\mathbb{P}_{\theta}^{\pi}, \mathbb{P}_{\theta_*}^{\pi})$  decreases with  $\kappa_*$ .



### Tension and trade-off

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

### Trade-off

- Let  $\pi$  perform well for  $\theta_*$ ,
- consider an alternative instance  $\theta$  such that  $\|\theta - \theta_*\|^2 \approx \sqrt{\kappa_*/T}$ ,
- the regret of  $\pi$  for the instance  $\theta$  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

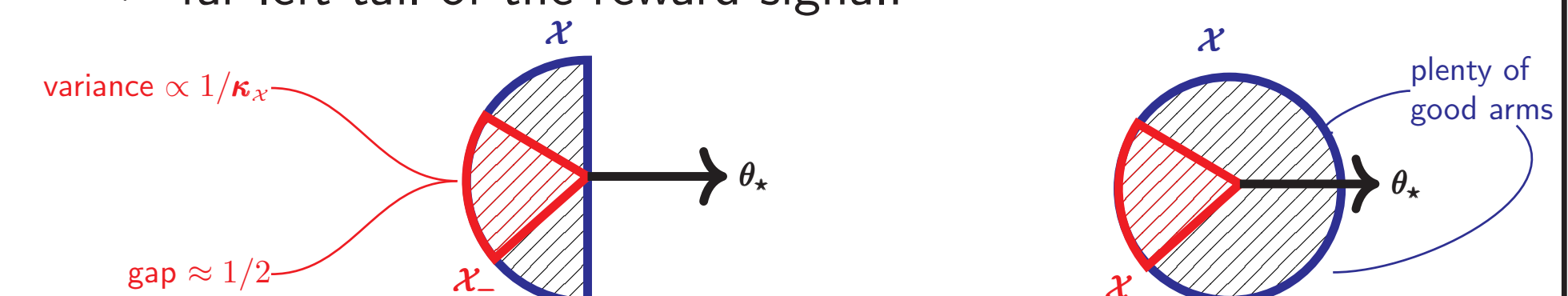
- Sublinear regret ⇒ play mostly around the best arm  $x_*$ .
- ↗ Almost a linear bandit with slope  $1/\kappa_*$ .
- A finer analysis is coherent with this conceptual argument:

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

- 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 arms is small:

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

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

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

- ↗ independent of  $\kappa_{\mathcal{X}}$  for reasonable configurations!

## ALGORITHM AND EXPERIMENTS

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

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

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

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.

### Parameter-based optimism

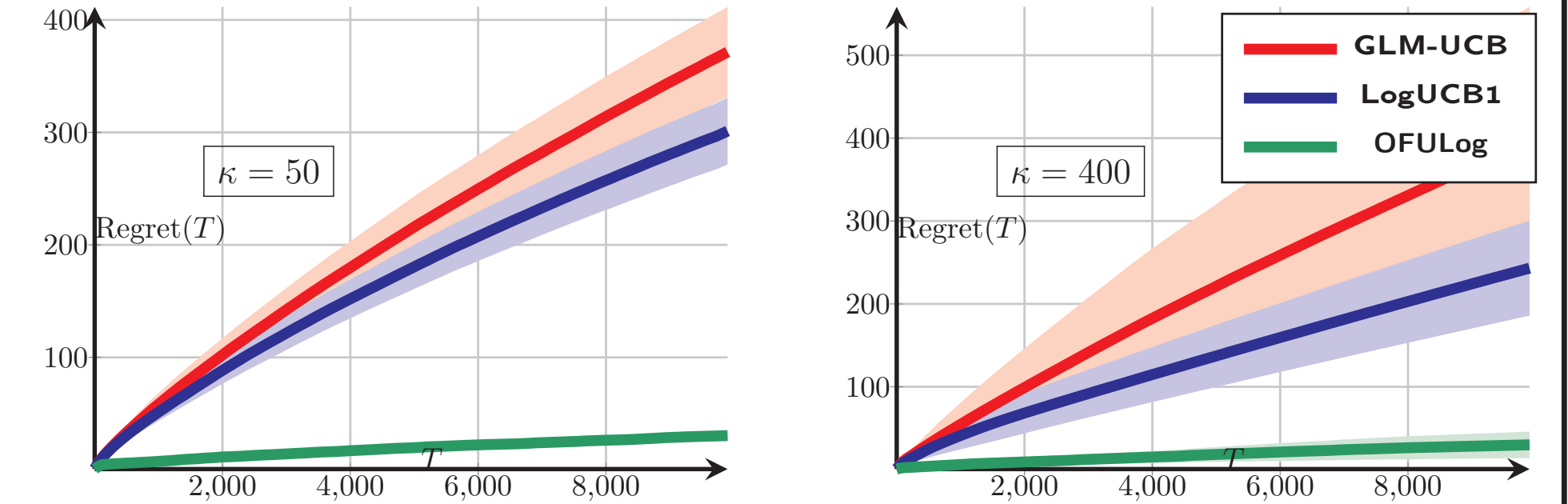
- Enforce optimism through parameter-search (OFUL-like), and not bonus-based approach.
- This yields an *adaptive* algorithm: no tuning needed to adapt to the structure of the decision set.

### Tractable algorithm

- We also introduce a *convex relaxation* of the confidence set  $\mathcal{C}_t(\delta)$  of [Faury et al., ICML'20].
- No non-convex optimization routine (≠ previous work).

### Practical improvements

- Toy experiment: dramatic improvement over GLM-UCB [Filippi et al., NIPS'10] and Log-UCB1 [Faury et al., ICML'20].



## CONCLUSION

- Our conclusion contrasts with previous work:

Logistic Bandit: non-linearity makes the problem **easier!**

- Regret-upper bound with exponential improvement.
- First problem-dependent lower-bound for Logistic Bandit.
- Fully tractable, adaptive algorithm thanks to convex relaxation.

## REFERENCES

- S. Filippi, O. Cappé, A. Garivier and C. Szepesvári. Parametric Bandits: The Generalized Linear Case. *Proceedings of NIPS*, 2010.
- F. Bach. Self-concordant analysis for logistic regression. *Electronic Journal of Statistics*, 2010.
- S. Dong, T. Ma and B. Van Roy. On the Performance of Thompson Sampling on Logistic Bandits. *Proceedings of COLT*, 2019.
- L. Faury, M. Abeille, C. Calauzène and O. Fercoq. Improved Optimistic Algorithms for Logistic Bandits. *Proceedings of ICML*, 2020.
- M. Simchowitz and D. Foster. Naive Exploration is Optimal for Online LQR. *Proceedings of ICML*, 2020.