

JOINTLY EFFICIENT AND OPTIMAL ALGORITHMS FOR LOGISTIC BANDITS

LOUIS FAURY¹, MARC ABEILLE¹, KWANG-SUNG JUN², CLÉMENT CALAUZÈNES¹

¹Criteo AI Lab, ²University of Arizona

MOTIVATION

Logistic Bandits. *Structured* and *binary* bandit games;

- neat study of non-linearity in parametric bandits,
- highly relevant in practice (sequential decision-making under binary feedback).

Context.

- fruitful stream of research on the impact of *non-linearity* in Logistic Bandits,
- led to the development of *statistically* efficient (minimax-optimal) algorithms,
- however (prohibitively) computationally *inefficient*.

↪ although well understood from a learning-theoretic standpoint, we are still missing fast algorithms for Logistic Bandits.

Can we achieve computational efficiency without sacrificing statistical tightness?

LOGISTIC BANDITS

Reward model

- $\mathcal{A} \subset \mathbb{R}^d$ is the arm set,
- $r_{t+1} = r(a_t) \in \{0, 1\}$ is the associated reward,
- $\theta_* \in \mathbb{R}^d$ *unknown* parameter.

[Binary reward]

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

[Non-linear link function]

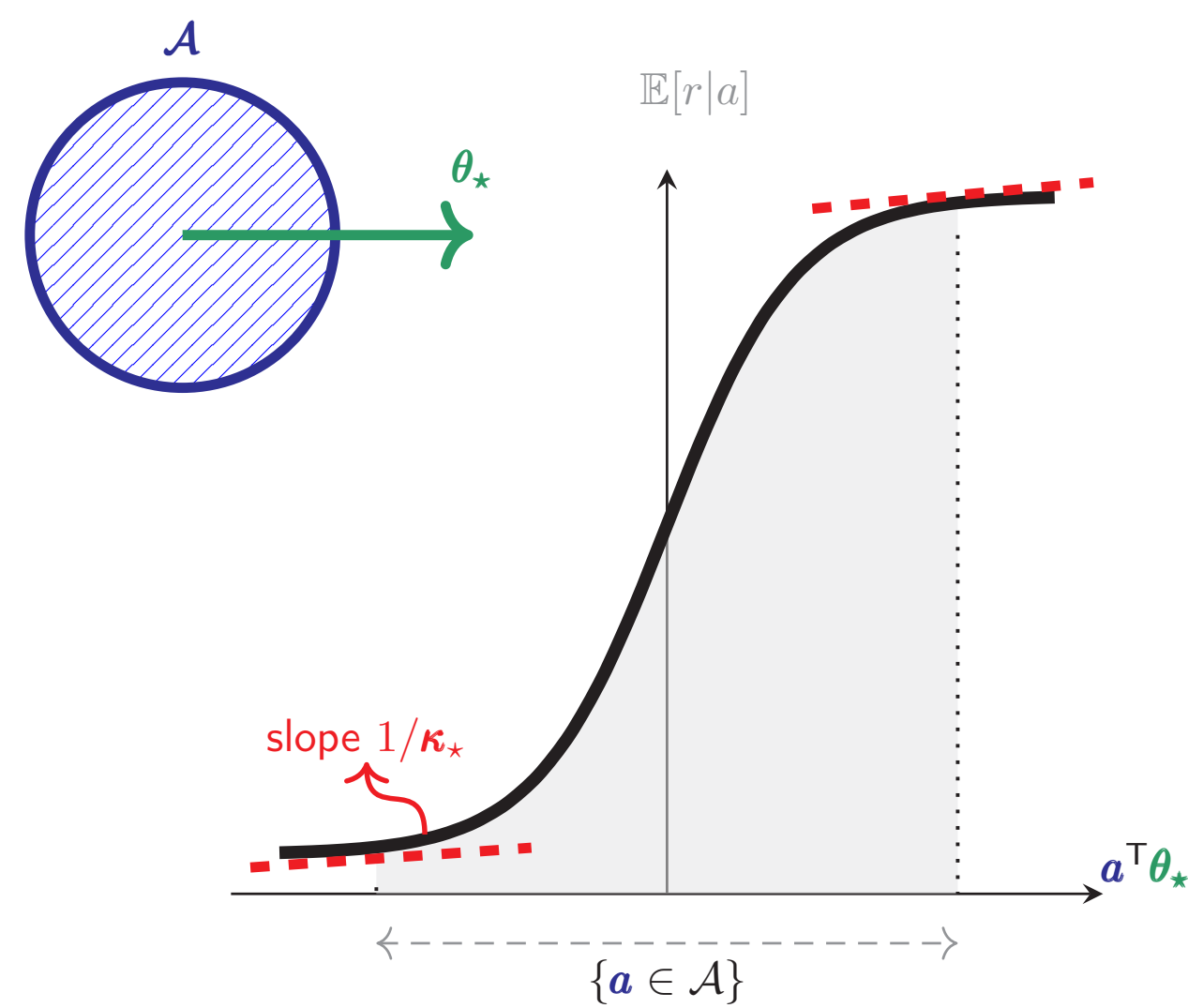
$$\mu(z) = (1 + \exp(-z))^{-1}$$

Learning problem

At each step $t \leq T$:

- choose arm $a_t \in \mathcal{A}$,
- receive reward r_{t+1} .

Objective is to minimize regret: $\text{Regret}(T) = \sum_{t=1}^T [\max_{a \in \mathcal{A}} \mu(a^\top \theta_*) - \mu(a_t^\top \theta_*)]$.



Level of non-linearity

Measured by a *problem-dependent* constant:

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

- (inverse of) minimal variance,
- "distance to linearity" over the decision set,
- typically large as $\kappa_* \approx \exp(\|\theta_*\|)$.

MINIMAX-OPTIMALITY

Minimax-Optimal algorithms w.r.t T , d and κ_* :

[Abeille et al., AISTATS'21]

$$\text{Regret}(T) \lesssim d\sqrt{T/\kappa_*} + \kappa_*$$

- ↪ exponential improvement over previous approaches.
- ↪ the more non-linear the problem, the easier!
- ↪ matching lower-bound.

Optimistic algorithm based on confidence region $\mathcal{C}_t(\delta)$; play $a_t = \arg \max_{\mathcal{A}} \max_{\theta \in \mathcal{C}_t(\delta)} a^\top \theta$.

Main teachings from those optimal algorithms:

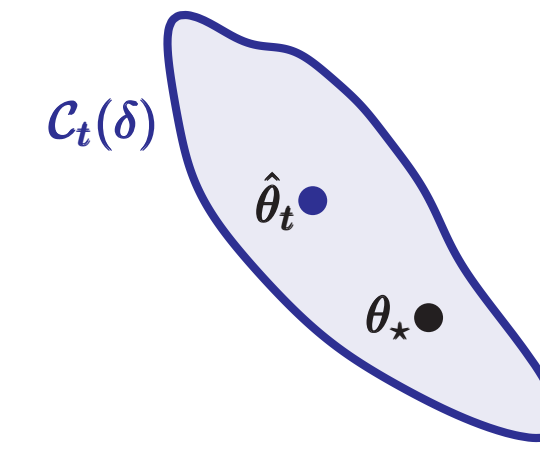
- the learning procedure must leverage the *local* behavior of non-linearity,
- the learner might suffer an inevitably large regret during an initial *transitory* regime,
- in the *permanent* regime, the local behavior around θ_* dominates.

COMPUTATIONAL EFFICIENCY

Learning. [Abeille et al., 2021] use *variance-sensitive* confidence sets:

$$\mathcal{C}_t(\delta) := \left\{ \theta, \|\theta - \hat{\theta}_t\|_{\mathbf{H}_t(\theta)}^2 \lesssim d \log(t) \right\}$$

$$\begin{cases} \hat{\theta}_t := \arg \min_{\theta} \sum_{s=1}^t \ell_s(\theta) \text{ the MLE estimator,} \\ \mathbf{H}_t(\theta) := \sum_{s=1}^t \dot{\mu}(a_s^\top \theta) a_s a_s^\top + \lambda \mathbf{I}_d. \end{cases}$$



Expensive to maintain;

- testing membership $\theta \in \mathcal{C}_t(\delta)$ costs $\Omega(t)$ operations!
- batch computation of the MLE,
- non-ellipsoidal confidence sets impacts the complexity of the *planning* mechanism.

Goal. New confidence sets for Logistic Bandits such that:

- *sufficient statistics* require $\mathcal{O}(1)$ operations,
- without sacrificing tightness.

ECOLOG PROCEDURE

Iterative Procedure. Given some convex set Θ , compute estimators $\{\theta_t\}_t$ of θ_* through:

ECOLOG

1. $\theta_t \leftarrow \arg \min_{\theta \in \Theta} \|\theta - \theta_{t-1}\|_{\mathbf{W}_{t-1}}^2 + \eta \ell_t(\theta)$
2. $\mathbf{W}_t \leftarrow \mathbf{W}_{t-1} + \dot{\mu}(a_t^\top \theta_t) a_t a_t^\top$

where $\eta = \text{diam}(\Theta)$ is the *learning-rate* and $\ell_t(\cdot)$ is the immediate log-loss.

Main Idea. Inspired from Online Convex Optimization literature [Jézéquel et al., 2020];

- based on *local* quadratic lower-bounds for the logistic-loss:

$$\ell_t(\theta_*) \geq \ell_t(\theta) + \nabla \ell_t(\theta)^\top (\theta_* - \theta) + \dot{\mu}(a_t^\top \theta) (a_t^\top (\theta_* - \theta))^2.$$

- decompose the total log-loss $\mathcal{L}_t(\theta) := \sum_{s=1}^t \ell_s(\theta)$ as:

$$\mathcal{L}_t(\theta) = \tilde{\mathcal{L}}_{t-1}(\theta) + l_t(\theta) = \eta^{-1} (\theta - \theta_{t-1})^\top \mathbf{W}_{t-1} \mathbf{T} (\theta - \theta_{t-1}) + l_t(\theta).$$

- minimize this *strongly convex* proxy to obtain θ_t at $\tilde{\mathcal{O}}(d^2)$ cost.

Confidence Regions. Emulates the confidence set $\mathcal{C}_t(\delta)$ through \mathbf{W}_t ;

Confidence Set

Let $\delta \in (0, 1]$. If $\theta_* \in \Theta$ then:

$$\mathcal{E}_t(\delta) = \left\{ \|\theta - \theta_t\|_{\mathbf{W}_t}^2 \lesssim \exp(\eta) d \log(t/\delta) \right\}$$

is a confidence region for θ_* ; $\mathbb{P}(\theta_* \in \mathcal{E}_t(\delta)) \geq 1 - \delta$.

↪ *local* behavior captured through Θ , an *admissible* parameter-set.

Forced-Exploration phase allows to identify a *small* admissible set Θ ;

- such that $\eta = \text{diam}(\Theta) \approx 1$,
- for at most κ_* rounds (e.g. second-order regret term),
- hardcodes the transitory phase of [Abeille et al., 2021].

ALGORITHM AND REGRET BOUND

Algorithm. Combines forced-exploration and ECOLog;

OFU-ECOLOG

- Perform $\tau = \kappa_*$ rounds of forced exploration to obtain Θ such that: $\theta_* \in \Theta$ with high proba. and $\text{diam}(\Theta) \leq 1$.
- For $t \geq \tau$:
 1. play the optimistic arm $a_t = \arg \max_{a \in \mathcal{A}} \max_{\theta \in \mathcal{E}_t(\delta)} a^\top \theta$,
 2. receive r_{t+1} , compute $(\theta_{t+1}, \mathbf{W}_{t+1})$ by running ECOLog.

Computational Cost. At each round at most $\mathcal{O}(d^2 |\mathcal{A}| \log(t))$ operations, since:

- ECOLog can be solved to arbitrary precision ε at cost $d \log(1/\varepsilon)$,
- the confidence region is *ellipsoidal*; closed-form for the optimistic arm (exploration bonus).

Regret Bounds. Minimax-optimal rates;

Regret

The regret of OFU-ECOLOG satisfies with high probability:

$$\text{Regret}(T) \lesssim d\sqrt{T/\kappa_*} + \kappa_*.$$

Adaptive Version. ada-OFU-ECOLOG dilutes the warm-up throughout the learning.

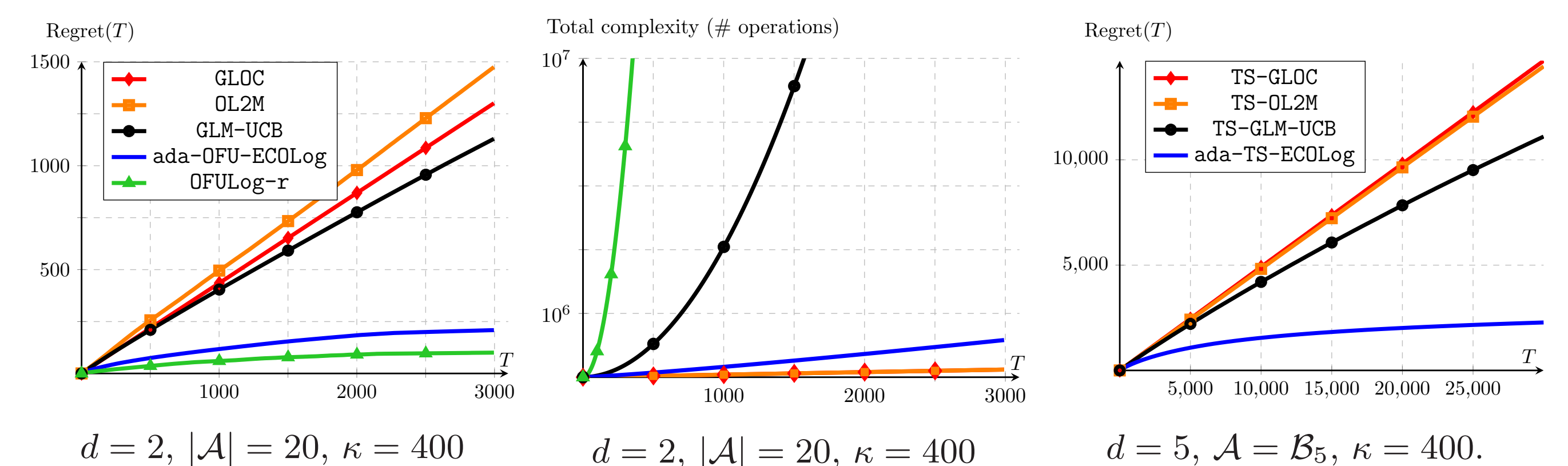
- based on a *data-dependent* width for the confidence regions,
- ✓ preserves statistical efficiency while ultimately removing the need for forced-exploration,
- ✓ coherent with [Abeille et al. 2021]: low-order κ_* dependencies can sometimes be avoided.

CONCLUSION

Joint *statistical* and *computational* efficiency;

Algorithm	Regret Bound	Cost Per-Round	Minimax	Efficient
GLM-UCB [Filippi et al. 2010]	$\tilde{\mathcal{O}}(\kappa_* d \sqrt{T})$	$\mathcal{O}(d^2 \mathcal{A} T)$	✗	✗
GLOC, OL2M [Jun et al. 2017] [Zhang et al. 2016]	$\tilde{\mathcal{O}}(\kappa_* d \sqrt{T})$	$\mathcal{O}(d^2 \mathcal{A})$	✗	✓
OFULog-r [Abeille et al. 2021]	$\tilde{\Theta}(d\sqrt{T/\kappa_*})$	$\mathcal{O}(d^2 \mathcal{A} T)$	✓	✗
(ada-)OFU-ECOLOG (this paper)	$\tilde{\Theta}(d\sqrt{T/\kappa_*})$	$\tilde{\mathcal{O}}(d^2 \mathcal{A})$	✓	✓

Numerical simulations corroborates theoretical results;



REFERENCES

- S. Filippi, O. Cappé, A. Garivier and C. Szepesvári. Parametric Bandits: The Generalized Linear Case. *NeurIPS*, 2010.
- L. Faury, M. Abeille, C. Calauzènes and O. Fercoq. Improved Optimistic Algorithms for Logistic Bandits. *ICML*, 2020.
- M. Abeille, L. Faury and C. Calauzènes. Instance-Wise Minimax-Optimal Algorithms for Logistic Bandits. *AISTATS*, 2021.