

# Non-linearity in Parametric Bandits: the Logistic Bandit case

joint work with M. Abeille<sup>1</sup>, C. Calauzènes<sup>1</sup> and O. Fercoq<sup>2</sup>

<sup>1</sup> Criteo AI Lab

<sup>2</sup> LTCI TelecomParis

# Presentation Outline

- **Goal.** Study **non-linearity** in sequential decision making problem.
  - ▶ Logistic Bandit: theoretical qualities.
    - ↪ Simple extension of the Linear Bandit.
    - ↪ Isolates the effect of non-linearity.
  - ▶ Logistic Bandit: practical relevance.
    - ↪ Model real-life problems with **binary** feedback.
    - ↪ news recommendation, clinical trials, ..
- **Logistic Bandit: high-level contributions.**
  - ▶ Improved algorithms with enhanced performances.
  - ▶ New theoretical insights: non-linearity makes the problem **easier**.

# Warm-Up: Linear Bandits

- Repeated game with **structured** feedback.



- **Motivation.** Generalizes the classical Multi-Arm Bandit setting
  - ▶ encode similarities between actions.
  - ▶ handle infinite number of actions.
  - ▶ handle contextual information  $x_t$ :  $\mathbf{f}(\phi(\mathbf{a}_t, x_t))^T \boldsymbol{\theta}_*$ .

# Warm-up: Linear Bandit (ctn'd)

- **Goal.** Minimize cumulative **regret**; with  $\mathbf{a}_\star = \operatorname{argmax}_{\mathbf{a} \in \mathcal{A}} \mathbf{a}^\top \boldsymbol{\theta}_\star$ :

$$\operatorname{Regret}_{\boldsymbol{\theta}_\star}(T) := T \mathbf{a}_\star^\top \boldsymbol{\theta}_\star - \sum_{t=1}^T \mathbf{a}_t^\top \boldsymbol{\theta}_\star .$$

↪ objective:  $\operatorname{Regret}_{\boldsymbol{\theta}_\star}(T) \leq T^\alpha$  for  $\alpha < 1$ .

↪ balance exploitation and exploration.

# Warm-up: Linear Bandit (ctn'd)

- **Goal.** Minimize cumulative **regret**; with  $\mathbf{a}_* = \operatorname{argmax}_{\mathbf{a} \in \mathcal{A}} \mathbf{a}^\top \boldsymbol{\theta}_*$ :

$$\operatorname{Regret}_{\boldsymbol{\theta}_*}(T) := T \mathbf{a}_*^\top \boldsymbol{\theta}_* - \sum_{t=1}^T \mathbf{a}_t^\top \boldsymbol{\theta}_* .$$

↪ objective:  $\operatorname{Regret}_{\boldsymbol{\theta}_*}(T) \leq T^\alpha$  for  $\alpha < 1$ .

↪ balance exploitation and exploration.

- **Solved:** **minimax-optimal** and efficient algorithms.

$$\operatorname{Regret}(T) = \tilde{O}(d\sqrt{T}) ,$$

where  $\tilde{O}$  hides only logarithmic dependencies.

# Warm-up: Linear Bandit (ctn'd)

- Exploration/exploitation trade-off vs. **optimism** in face of uncertainty.
  - ▶ **Learning** is performed via ordinary least-squares:

$$\hat{\theta}_t := V_t^{-1} \left( \sum_{s=1}^{t-1} r_s a_s \right) \quad \text{where} \quad V_t^{-1} = \sum_{s=1}^t a_s a_s^\top + \lambda I_d .$$

# Warm-up: Linear Bandit (ctn'd)

- Exploration/exploitation trade-off vs. **optimism** in face of uncertainty.
  - ▶ **Learning** is performed via ordinary least-squares:

$$\hat{\theta}_t := V_t^{-1} \left( \sum_{s=1}^{t-1} r_s a_s \right) \quad \text{where} \quad V_t^{-1} = \sum_{s=1}^t a_s a_s^T + \lambda I_d .$$

- ▶ **Planning** by resorting to confidence sets:

$$\theta_{\star} \in C_t(\delta) = \left\{ \theta, \left\| \theta - \hat{\theta}_t \right\|_{V_t}^2 \leq d \log(t/\delta) \right\} \quad \text{with proba. at least } 1 - \delta$$

and enforcing optimism:

$$\text{play } a_{t+1} \in \operatorname{argmax}_{a \in \mathcal{A}} \max_{\theta \in C_t(\delta)} a^T \theta .$$

the hard part

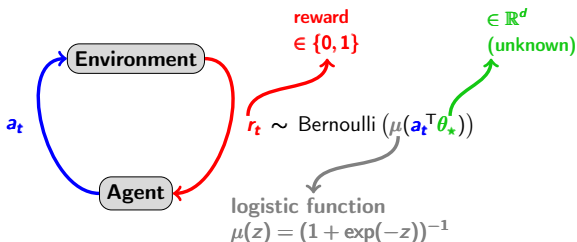
# The Logistic Bandit

- **Motivations.** The Linear Bandit setting has (many) limitations;
  - ▶ **Theoretical:** towards rich reward models.
    - ↪ The real world is fundamentally non-linear.
    - ↪ Does the same principle work?
    - ↪ Will it be optimal?
  - ▶ **Practical:**
    - ↪ The Linear Bandit covers only continuous rewards.
    - ↪ What about binary rewards (click, sale, success)?



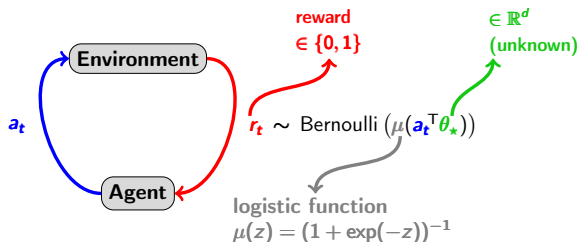
# The Logistic Bandit

- Repeated game with structured **binary** feedback.



# The Logistic Bandit

- Repeated game with structured **binary** feedback.

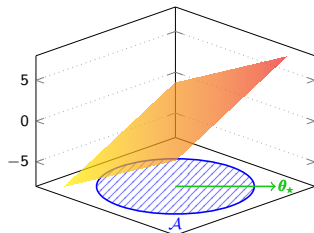


- Regret.** The agent tries to minimize its cumulative pseudo-regret:

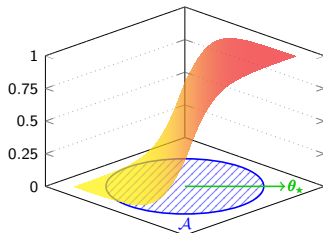
$$\text{Regret}_{\theta_\star}(T) := T\mu(\mathbf{a}_\star^\top \theta_\star) - \sum_{t=1}^T \mu(\mathbf{a}_t^\top \theta_\star).$$

# The Learning Problem (ctn'd)

- **Reward model.** Minimalist non-linear extension from the linear bandit.



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

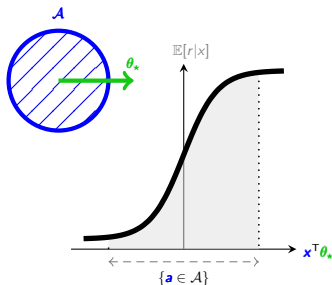


$$\mathbb{E}[r_t | a_t] = (1 + \exp(-a_t^\top \theta_*))^{-1}$$

- **Exploration-exploitation.** Same recipe:
  - ▶ Learning: maximum likelihood (**logistic regression**).
  - ▶ Planning: Optimism through confidence sets.
- **Additional challenge.** Non-linearity: information vs. regret.

# Quantifying non-linearity

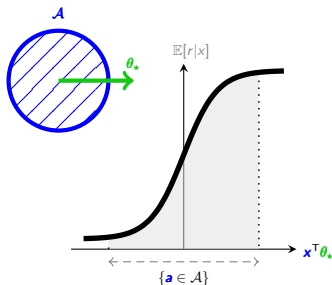
- Level of non-linearity = conditioning.
  - ▶ How **flat** are the tails.



- Important quantities. The level of non-linearity is problem-dependent.

# Quantifying non-linearity

- Level of non-linearity = conditioning.
  - ▶ How flat are the tails.

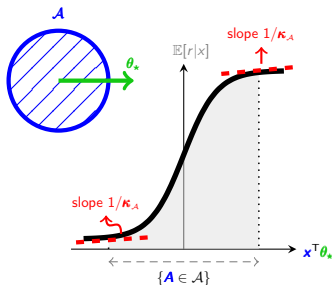


- Important quantities. The level of non-linearity is problem-dependent.
  - ▶ Historically characterized by a constant  $\kappa_{\mathcal{A}}$ :

$$\kappa_{\mathcal{A}} := \frac{1}{\min_{\mathbf{a} \in \mathcal{A}} \dot{\mu}(\mathbf{a}^T \boldsymbol{\theta}_*)} .$$

# Quantifying non-linearity

- Level of non-linearity = conditioning.
  - ▶ How **flat** are the tails.



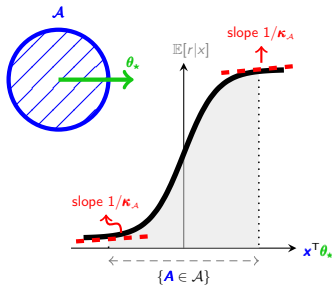
- **Important quantities.** The level of non-linearity is problem-dependent.
  - ▶ Historically characterized by a constant  $\kappa_{\mathcal{A}}$ :

$$\kappa_{\mathcal{A}} := \frac{1}{\min_{\mathbf{a} \in \mathcal{A}} \dot{\mu}(\mathbf{a}^T \boldsymbol{\theta}_*)}.$$

- ▶ The more non-linear the reward, the **bigger**.

# Quantifying non-linearity

- Level of non-linearity = conditioning.
  - ▶ How **flat** are the tails.



- **Important quantities.** The level of non-linearity is problem-dependent.
  - ▶ Historically characterized by a constant  $\kappa_{\mathcal{A}}$ :

$$\kappa_{\mathcal{A}} := \frac{1}{\min_{\mathbf{a} \in \mathcal{A}} \dot{\mu}(\mathbf{a}^T \boldsymbol{\theta}_*)} .$$

- ▶ The more non-linear the reward, the **bigger**.
- ▶ Typically  $\kappa_{\mathcal{A}} \geq \exp(\|\boldsymbol{\theta}_*\|)$  ! In practical case;  $\kappa_{\mathcal{A}} \sim 10^3$ .

# Previous approaches

- A lot of existing work on the logistic bandit: [Filippi et al. 2010; Li et al 2017; Kveton et al. 2019; Dong et al. 2019];
- All rely on a global **linearization** approach.



# Previous approaches

- A lot of existing work on the logistic bandit: [Filippi et al. 2010; Li et al 2017; Kveton et al. 2019; Dong et al. 2019];
- All rely on a global **linearization** approach.
- Information vs. regret : worst of both world!
  - ▶ Confidence set (at algorithmic design time):

$$\theta_{\star} \in \mathcal{C}_t(\delta) = \left\{ \theta, \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{V}_t}^2 \leq \kappa_{\mathcal{A}} d \log(t/\delta) \right\} \quad \text{with proba. at least } 1 - \delta$$

# Previous approaches

- A lot of existing work on the logistic bandit: [Filippi et al. 2010; Li et al 2017; Kveton et al. 2019; Dong et al. 2019];
- All rely on a global **linearization** approach.
- Information vs. regret : worst of both world!
  - ▶ Confidence set (at algorithmic design time):

$$\theta_{\star} \in \mathcal{C}_t(\delta) = \left\{ \theta, \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{V}_t}^2 \leq \kappa_{\mathcal{A}} d \log(t/\delta) \right\} \quad \text{with proba. at least } 1 - \delta$$

- ▶ Prediction error (at analysis time):

$$\mu(\mathbf{a}^T \theta_{\star}) - \mu(\mathbf{a}^T \theta) \leq \mathbf{a}^T (\theta_{\star} - \theta) / 4 .$$

# Previous approaches (ctn'd)

- Global linearization  $\Rightarrow$  disappointing results!
  - ▶ Poor regret guarantees.

$$\text{Regret}_{\theta_*}(T) = \tilde{O}\left(\kappa_{\mathcal{A}} d \sqrt{T}\right).$$

- ▶ Because the algorithms are **over-explorative**.
- ▶ Disappointing story about the effects of non-linearity.
- ↪ The more non-linear the problem, the larger the regret!

# Previous approaches (ctn'd)

- Global linearization  $\Rightarrow$  disappointing results!
  - ▶ Poor regret guarantees.

$$\text{Regret}_{\theta_*}(T) = \tilde{O}\left(\kappa_{\mathcal{A}} d \sqrt{T}\right).$$

- ▶ Because the algorithms are **over-explorative**.
  - ▶ Disappointing story about the effects of non-linearity.  
 $\rightsquigarrow$  The more non-linear the problem, the larger the regret!
- Our goal is to improve this with:
  - ▶ Enhanced confidence sets for  $\theta_*$ .
  - ▶ Improvement treatment of the **local** behavior of the reward signal.

# Improved confidence set

- Let  $H_t(\theta) = \sum_{s=1}^{t-1} \dot{\mu}(a_s^\top \theta) a_s a_s^\top + \lambda I_d$ ; then

$$\theta_\star \in \mathcal{E}_t(\delta) := \left\{ \theta, \left\| \theta - \hat{\theta}_t \right\|_{H_t(\theta)}^2 \leq d \log(t/\delta) \right\} \text{ with proba } \geq 1 - \delta.$$

- Based on a **new concentration inequality** for self-normalized process.
- Smaller than  $\mathcal{C}_t(\delta)$  by **at least**  $\sqrt{\kappa_{\mathcal{A}}}$ .
- Undergoes convex relaxation for tractability.

# Improved confidence set

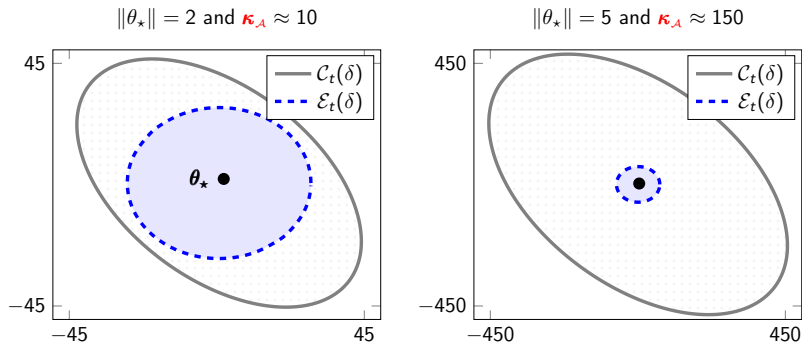


Figure: Visualization of two-dimensional Logistic bandit confidence sets.

- Smaller confidence set  $\Rightarrow$  less explorative algorithm, better performance.

# Improved algorithm and analysis

- We use the same recipe for enforcing **optimism**:

$$\text{play } a_t = \operatorname{argmax}_{a \in \mathcal{A}} \max_{\theta \in \mathcal{E}_t(\delta)} \mu(a^\top \theta).$$

# Improved algorithm and analysis

- We use the same recipe for enforcing **optimism**:

$$\text{play } a_t = \operatorname{argmax}_{a \in \mathcal{A}} \max_{\theta \in \mathcal{E}_t(\delta)} \mu(a^\top \theta).$$

- We introduce a new analysis for the Logistic Bandit:



# Improved algorithm and analysis

- We use the same recipe for enforcing **optimism**:

$$\text{play } a_t = \operatorname{argmax}_{a \in \mathcal{A}} \max_{\theta \in \mathcal{E}_t(\delta)} \mu(a^\top \theta).$$

- We introduce a new analysis for the Logistic Bandit:
  - ▶ Leverage the **self-concordance** property of the logistic function.
  - ▶ Allows for a **local** treatment of the non-linearity.
  - ▶ Strikes the right balance between information and regret.

# Regret guarantees

- Enhanced regret guarantee; denote  $\mathbf{a}_*$  the **best** action:

$$\text{Regret}_{\theta_*}(T) = \tilde{O} \left( d \sqrt{\dot{\mu}(\mathbf{a}_*^\top \theta_*) T} \right) .$$

# Regret guarantees

- Enhanced regret guarantee; denote  $\mathbf{a}_*$  the **best** action:

$$\text{Regret}_{\theta_*}(T) = \tilde{O}\left(d\sqrt{\dot{\mu}(\mathbf{a}_*^T\theta_*)T}\right).$$

- Illustration for the unit ball arm-set:  $\dot{\mu}(\mathbf{a}_*^T\theta_*) = 1/\kappa_{\mathcal{A}} \approx \exp(-\|\theta_*\|)$ .
  - ▶ The regret is:

$$\text{Regret}_{\theta_*}(T) = \tilde{O}\left(d\sqrt{T/\kappa_{\mathcal{A}}}\right).$$

- ▶ Improvement by  $\kappa_{\mathcal{A}}^{3/2} \approx \exp(3\|\theta_*\|/2)$ !

# Regret guarantees

- Enhanced regret guarantee; denote  $\mathbf{a}_*$  the **best** action:

$$\text{Regret}_{\theta_*}(T) = \tilde{O}\left(d\sqrt{\dot{\mu}(\mathbf{a}_*^T\theta_*)T}\right).$$

- Illustration for the unit ball arm-set:  $\dot{\mu}(\mathbf{a}_*^T\theta_*) = 1/\kappa_{\mathcal{A}} \approx \exp(-\|\theta_*\|)$ .
  - The regret is:

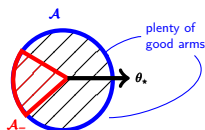
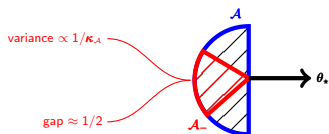
$$\text{Regret}_{\theta_*}(T) = \tilde{O}\left(d\sqrt{T/\kappa_{\mathcal{A}}}\right).$$

- Improvement by  $\kappa_{\mathcal{A}}^{3/2} \approx \exp(3\|\theta_*\|/2)$ !
- This rate is **minimax-optimal** w.r.t  $d$ ,  $T$  and  $\kappa_{\mathcal{A}}$ .

# Effects of non-linearity

- Non-linearity seems to be beneficial!
  - ▶ The larger  $\kappa_{\mathcal{A}}$ , the smaller the regret!
- Not entirely true; non-linearity can impact a **transitory phase**.
  - ▶ Second-order term of the regret.
  - ▶ Happens before highly rewarding areas of  $\mathcal{A}$  are identified.

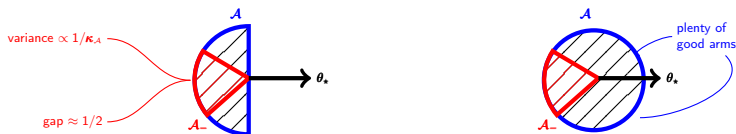
$$\text{Regret}_{\theta_*}(T) = \tilde{O} \left( d \sqrt{\dot{\mu}(a_*^\top \theta_*)} T + R^{\text{transitory}}(T) \right).$$



# Effects of non-linearity

- Non-linearity seems to be beneficial!
  - ▶ The larger  $\kappa_{\mathcal{A}}$ , the smaller the regret!
- Not entirely true; non-linearity can impact a **transitory phase**.
  - ▶ Second-order term of the regret.
  - ▶ Happens before highly rewarding areas of  $\mathcal{A}$  are identified.

$$\text{Regret}_{\theta_*}(T) = \tilde{O} \left( d \sqrt{\dot{\mu}(a_*^\top \theta_*)} T + R^{\text{transitory}}(T) \right).$$

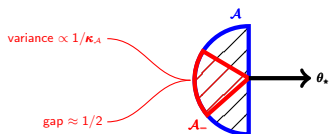


$$R^{\text{transitory}}(T) = \tilde{O}(\kappa_{\mathcal{A}})$$

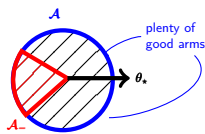
# Effects of non-linearity

- Non-linearity seems to be beneficial!
  - ▶ The larger  $\kappa_{\mathcal{A}}$ , the smaller the regret!
- Not entirely true; non-linearity can impact a **transitory phase**.
  - ▶ Second-order term of the regret.
  - ▶ Happens before highly rewarding areas of  $\mathcal{A}$  are identified.

$$\text{Regret}_{\theta_*}(T) = \tilde{O} \left( d \sqrt{\dot{\mu}(a_*^\top \theta_*)} T + R^{\text{transitory}}(T) \right).$$



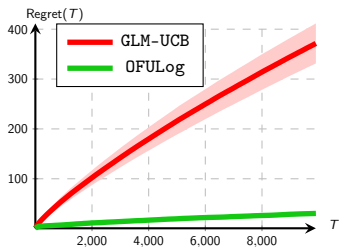
$$R^{\text{transitory}}(T) = \tilde{O}(\kappa_{\mathcal{A}})$$



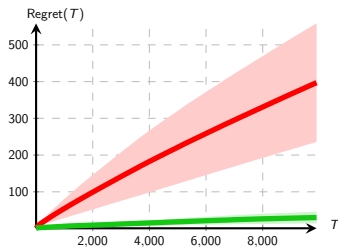
$$R^{\text{transitory}}(T) = \tilde{O}(1)$$

# Empirical performances

- Compared with the GLM-UCB of [Filippi et al. 2010].



(a)  $\kappa_{\mathcal{A}} = 50$



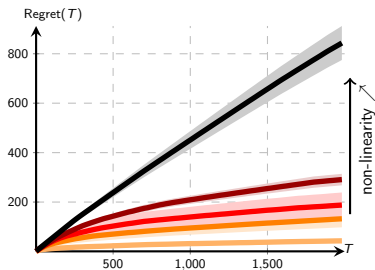
(b)  $\kappa_{\mathcal{A}} = 400$

Empirical comparison of GLM-UCB and OFULog on two LogB toy experiments. The regret curves are averaged over 50 independent runs. Standard-deviation is reported in shaded colors around the averaged cumulative regret. The arm-set  $\mathcal{A}$  is composed of 40 arms drawn uniformly at random in the 2-dimensional ball at the beginning of each run.

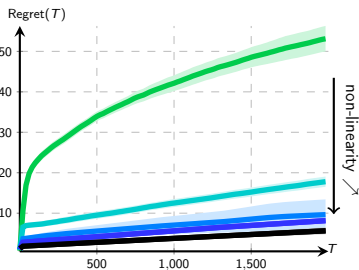


# Empirical performances (ctn'd)

- Check the impact of non-linearity:



(a) GLM-UCB



(b) OFULog

Figure: Comparing the effect of non-linearity on GLM-UCB and OFULog by varying the level of non-linearity in a Logistic Bandit setting.

Thank you!

# Some references

- **Previous work (most relevant)**
  - ▶ Filippi et al. Parametric Bandits: the Generalized Linear case. *NeurIPS*, 2010.
  - ▶ Li et al. Provably Optimal Algorithms for Generalized Linear Bandits. *ICML*, 2017.
  - ▶ Dong et al. On the Performance of Thompson Sampling on Logistic Bandits. *COLT*, 2019.
- **Material for this talk was taken from:**
  - ▶ F., Abeille, Calauzènes and Fercoq. Improved Optimistic Algorithms for Logistic Bandits. *ICML*, 2020.
  - ▶ Abeille, F. and Calauzènes. Instance-Wise Minimax-Optimal Algorithms for Logistic Bandits. *AISTATS*, 2021.
- **Extension** to non-stationary settings.
  - ▶ Russac, F., Cappé, Garivier. Self-Concordant Analysis of Generalized Linear Bandits with Forgetting *AISTATS*, 2021.
  - ▶ F., Russac, Abeille and Calauzènes. Regret Bounds for Generalized Linear Bandits under Parameter Drift. *ALT*, 2021.