

# 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 Telecom Paris

# Presentation Outline

- Goal.
  - ▶ Study non-linearity in sequential decision making.
  - ▶ A simple problem: the Logistic Bandit.
    - ↪ Compact non-linear extension to the Linear Bandit.
    - ↪ Very relevant in practical problems with **binary** feedback.

# Presentation Outline

- **Goal.**

- ▶ Study non-linearity in sequential decision making.
- ▶ A simple problem: the Logistic Bandit.
  - ↪ Compact non-linear extension to the Linear Bandit.
  - ↪ Very relevant in practical problems with **binary** feedback.

- **Logistic Bandit: high-level contributions.**

- ▶ [Filippi et al. 2010, Faury et al. 2020]: non-linearity is harmful. Actually:

Non-linearity can make the problem **easier**.

# Presentation Outline

- **Goal.**

- ▶ Study non-linearity in sequential decision making.
- ▶ A simple problem: the Logistic Bandit.
  - ↪ Compact non-linear extension to the Linear Bandit.
  - ↪ Very relevant in practical problems with **binary** feedback.

- **Logistic Bandit: high-level contributions.**

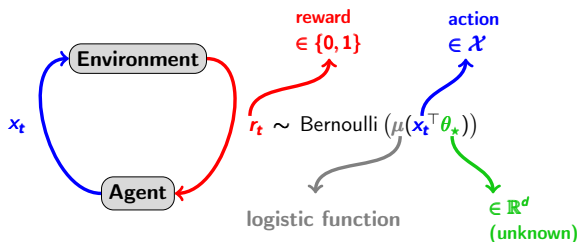
- ▶ [Filippi et al. 2010, Faury et al. 2020]: non-linearity is harmful. Actually:

Non-linearity can make the problem **easier**.

- ▶ Identify two distinct regimes:
  - ↪ Short-term  $\leftrightarrow$  early exploration phase: **neutral** (most often).
  - ↪ Long-term  $\leftrightarrow$  exploration-exploitation phase: **beneficial**.

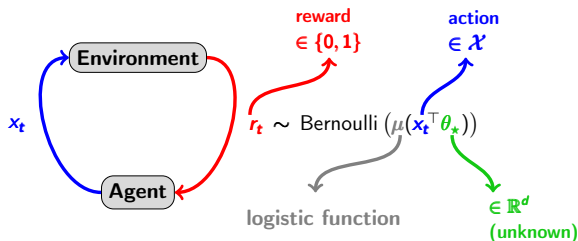
# The Learning Problem

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



# The Learning Problem

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

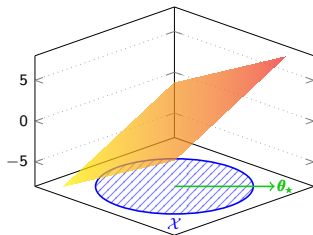


- Regret.** The agent tries to minimize its cumulative pseudo-regret:

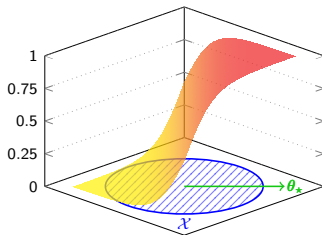
$$\text{Regret}_{\theta_*}(T) := T \max_{x \in \mathcal{X}} \mu(x^\top \theta_*) - \sum_{t=1}^T \mu(x_t^\top \theta_*).$$

# The Learning Problem (ctn'd)

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



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

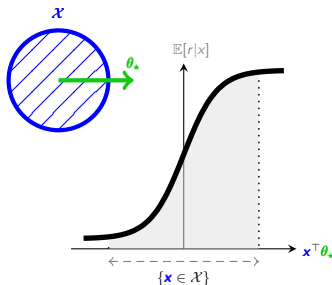


$$\mathbb{E}[r_t | x_t] = (1 + \exp(-x_t^\top \theta_*))^{-1}$$

- **Exploration-exploitation.** Same recipe:
  - ▶ Learning: maximum likelihood.
  - ▶ 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.

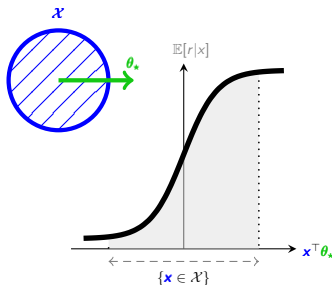


- **Important quantities.** The level of non-linearity is problem-dependent.



# Quantifying non-linearity

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

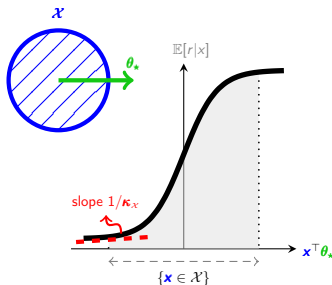


- **Important quantities.** The level of non-linearity is problem-dependent.
  - ▶ Historically characterized by a constant  $\kappa_{\mathcal{X}}$ :

$$\kappa_{\mathcal{X}} := \frac{1}{\min_{\mathbf{x} \in \mathcal{X}} \dot{\mu}(\mathbf{x}^\top \boldsymbol{\theta}_*)}.$$

# Quantifying non-linearity

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

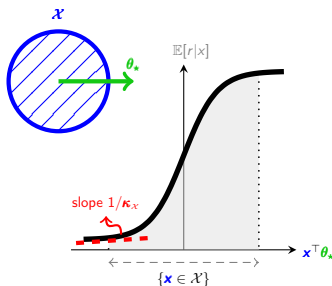


- **Important quantities.** The level of non-linearity is problem-dependent.
  - ▶ Historically characterized by a constant  $\kappa_X$ :

$$\kappa_X := \frac{1}{\min_{x \in X} \dot{\mu}(x^T \theta_*)}.$$

# 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  $\kappa_{\mathcal{X}}$ :

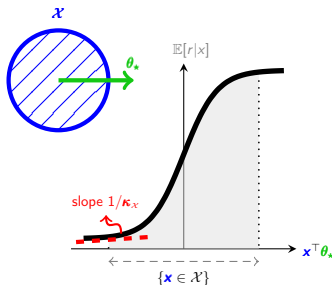
$$\kappa_{\mathcal{X}} := \frac{1}{\min_{\mathbf{x} \in \mathcal{X}} \dot{\mu}(\mathbf{x}^\top \boldsymbol{\theta}_*)}$$

the more non-linear  
the bigger

$\propto \exp(\|\boldsymbol{\theta}_*\|)$

# 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  $\kappa_{\mathcal{X}}$ :

$$\kappa_{\mathcal{X}} := \frac{1}{\min_{\mathbf{x} \in \mathcal{X}} \dot{\mu}(\mathbf{x}^\top \boldsymbol{\theta}_*)}.$$

the more non-linear  
the bigger

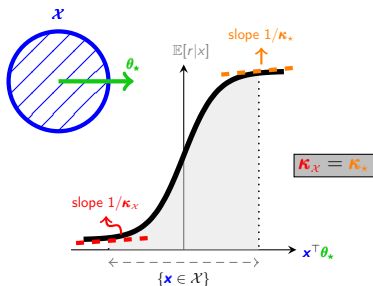
$\propto \exp(\|\boldsymbol{\theta}_*\|)$

- ▶ Inverse slope at the optimum; letting  $\mathbf{x}_* = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \mathbf{x}^\top \boldsymbol{\theta}_*$ :

$$\kappa_* := \frac{1}{\dot{\mu}(\mathbf{x}_*^\top \boldsymbol{\theta}_*)}.$$

# 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  $\kappa_X$ :

$$\kappa_X := \frac{1}{\min_{x \in X} \dot{\mu}(x^\top \theta_*)}.$$

the more non-linear  
the bigger

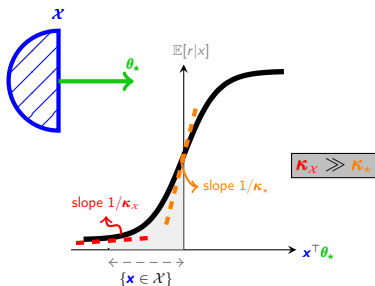
$\propto \exp(\|\theta_*\|)$

- ▶ Inverse slope at the optimum; letting  $x_* = \operatorname{argmax}_{x \in X} x^\top \theta_*$ :

$$\kappa_* := \frac{1}{\dot{\mu}(x_*^\top \theta_*)}.$$

# 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  $\kappa_{\mathcal{X}}$ :

$$\kappa_{\mathcal{X}} := \frac{1}{\min_{\mathbf{x} \in \mathcal{X}} \dot{\mu}(\mathbf{x}^\top \boldsymbol{\theta}_*)}.$$

the more non-linear  
the bigger

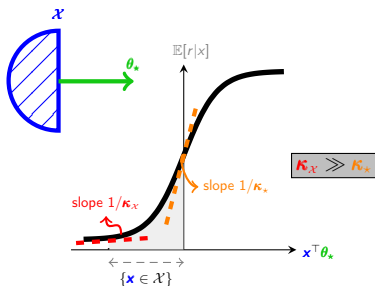
$\propto \exp(\|\boldsymbol{\theta}_*\|)$

- ▶ Inverse slope at the optimum; letting  $\mathbf{x}_* = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \mathbf{x}^\top \boldsymbol{\theta}_*$ :

$$\kappa_* := \frac{1}{\dot{\mu}(\mathbf{x}_*^\top \boldsymbol{\theta}_*)}.$$

# 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  $\kappa_{\mathcal{X}}$ :

$$\kappa_{\mathcal{X}} := \frac{1}{\min_{\mathbf{x} \in \mathcal{X}} \dot{\mu}(\mathbf{x}^\top \boldsymbol{\theta}_*)}$$

the more non-linear  
the bigger

$$\propto \exp(\|\boldsymbol{\theta}_*\|)$$

- ▶ Inverse slope at the optimum; letting  $\mathbf{x}_* = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \mathbf{x}^\top \boldsymbol{\theta}_*$ :

$$\kappa_* := \frac{1}{\dot{\mu}(\mathbf{x}_*^\top \boldsymbol{\theta}_*)}$$

$$\in [4, \kappa_{\mathcal{X}}]$$

# Non-linearity vs. regret: previous work

Approach	Regret
[Filippi et al. 2010] Linearization (global)	$\tilde{O}(\kappa_{\mathcal{X}} d\sqrt{T})$
[Faury et al. 2020] Self-concordance (local)	$\tilde{O}(d\sqrt{T} + \kappa_{\mathcal{X}})$
<b>This work</b> Refined local approach	$\tilde{O}(d\sqrt{T/\kappa_{\star}} (+ \kappa_{\mathcal{X}}))$



# Non-linearity vs. regret: previous work

Approach	Regret
[Filippi et al. 2010] Linearization (global)	$\tilde{O}(\kappa_{\mathcal{X}} d\sqrt{T})$
[Faury et al. 2020] Self-concordance (local)	$\tilde{O}(d\sqrt{T} + \kappa_{\mathcal{X}})$
<b>This work</b> Refined local approach	$\tilde{O}(d\sqrt{T/\kappa_*} + \kappa_{\mathcal{X}})$

- **Exponential improvement.** If  $\mathcal{X} = \{\|x\| \leq 1\}$  then  $\kappa_{\mathcal{X}} = \kappa_* \geq e^{\|\theta_*\|}$  then regret:

$$\tilde{O}(e^{\|\theta_*\|} d\sqrt{T}) \rightarrow \tilde{O}(d\sqrt{T} + e^{\|\theta_*\|}) \rightarrow \tilde{O}(e^{-\|\theta_*\|/2} d\sqrt{T})$$

# Regret Upper-Bound

- Effects of non-linearity: transitory and permanent regime.

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

# Regret Upper-Bound

- **Effects of non-linearity:** transitory and permanent regime.

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

- **Permanent regime.** For  $t \gg 1$ , only the local slope around  $x_*$  matters.

▶ Conceptually:

- Sub-linear regret  $\rightsquigarrow$  play mostly  $x_t \approx x_*$  for large  $t$ .
- Linear bandit with slope  $\dot{\mu}(x_*^\top \theta_*) = \frac{1}{\kappa_*}$  (potentially  $\ll 1$ ).

# Regret Upper-Bound

- **Effects of non-linearity:** transitory and permanent regime.

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

- **Permanent regime.** For  $t \gg 1$ , only the local slope around  $x_*$  matters.

- ▶ Conceptually:

- Sub-linear regret  $\rightsquigarrow$  play mostly  $x_t \approx x_*$  for large  $t$ .
- Linear bandit with slope  $\dot{\mu}(x_*^\top \theta_*) = \frac{1}{\kappa_*}$  (potentially  $\ll 1$ ).

- ▶ The smaller this local slope, the easier the problem:

$$R^{\text{perm}}(T) = \tilde{O}\left(d\sqrt{T/\kappa_*}\right).$$

- Formal proof: self-concordance.

# Regret Upper-Bound

- **Effects of non-linearity:** transitory and permanent regime.

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

- **Permanent regime.** For  $t \gg 1$ , only the local slope around  $x_*$  matters.

- ▶ Conceptually:

- Sub-linear regret  $\rightsquigarrow$  play mostly  $x_t \approx x_*$  for large  $t$ .
- Linear bandit with slope  $\dot{\mu}(x_*^\top \theta_*) = \frac{1}{\kappa_*}$  (potentially  $\ll 1$ ).

- ▶ The smaller this local slope, the easier the problem:

$$R^{\text{perm}}(T) = \tilde{O}\left(d\sqrt{T/\kappa_*}\right).$$

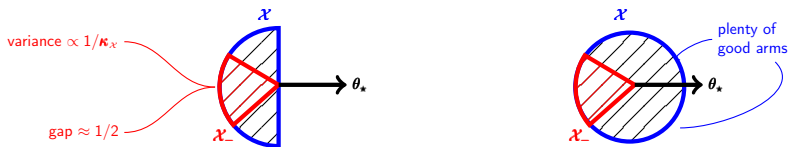
- Formal proof: self-concordance.

- ▶ Question: how long to reach it?


$$\approx \exp(\|\theta_*\|)!$$

# Regret Upper Bounds (ctn'd)

- **Transitory Regret.** Also linked to the problem's geometry..
  - ▶ Proportion of **detrimental** arms: little information and large sub-optimality.

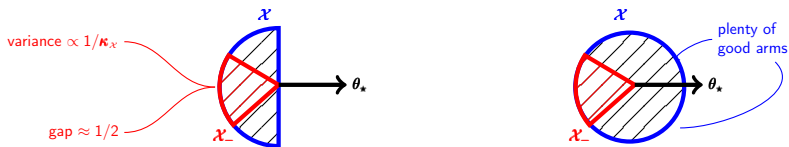


- ▶ Transitory regret = how long are we stuck playing **detrimental** arms?

$$R^{\text{trans}}(T) \propto \sum_{t=1}^T \mathbb{1}(x_t \in \mathcal{X}_-)$$

# Regret Upper Bounds (ctn'd)

- **Transitory Regret.** Also linked to the problem's geometry..
  - ▶ Proportion of **detrimental** arms: little information and large sub-optimality.



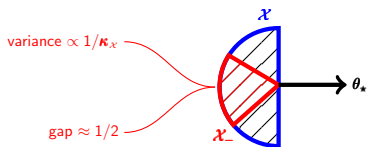
$$R^{\text{trans}}(T) = \tilde{O}(\kappa_X)$$

- ▶ Transitory regret = how long are we stuck playing **detrimental** arms?

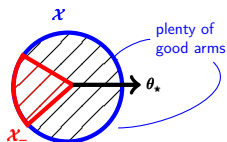
$$R^{\text{trans}}(T) \propto \sum_{t=1}^T \mathbb{1}(x_t \in \mathcal{X}_-)$$

# Regret Upper Bounds (ctn'd)

- **Transitory Regret.** Also linked to the problem's geometry..
  - ▶ Proportion of **detrimental** arms: little information and large sub-optimality.



$$R^{\text{trans}}(T) = \tilde{O}(\kappa_{\mathcal{X}})$$



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

- ▶ Transitory regret = how long are we stuck playing **detrimental** arms?

$$R^{\text{trans}}(T) \propto \sum_{t=1}^T \mathbb{1}(x_t \in \mathcal{X}_-)$$



# Regret Upper Bounds (ctn'd)

- Wrapping up.

## Theorem (Regret upper-bound)

*With high probability:*

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

- Refined problem-dependent bounds:

# Regret Upper Bounds (ctn'd)

- Wrapping up.

## Theorem (Regret upper-bound)

With high probability:

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

- Refined problem-dependent bounds:
  - ▶ **Worst configuration.**

$$\text{Regret}_{\theta_*}(T) = \tilde{O}(d\sqrt{T} + \kappa_{\mathcal{X}})$$

# Regret Upper Bounds (ctn'd)

- Wrapping up.

## Theorem (Regret upper-bound)

With high probability:

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

- Refined problem-dependent bounds:
  - ▶ Worst configuration.

$$\text{Regret}_{\theta_*}(T) = \tilde{O}(d\sqrt{T} + \kappa_X)$$

- ▶ Best configuration.

$$\text{Regret}_{\theta_*}(T) = \tilde{O}(d\sqrt{T/\kappa_X})$$

↪ Is this optimal?

# Problem-dependent lower-bound

- **Challenge.** Study optimality w.r.t problem-dependent constants  $\kappa_{\mathcal{X}}$ .
  - ▶ Lower-bound for a *continuum* of problems, each with different  $\kappa_{\mathcal{X}}$ .
  - ▶ Traditional lower-bound technique fails.

# Problem-dependent lower-bound

- **Challenge.** Study optimality w.r.t problem-dependent constants  $\kappa_{\mathcal{X}}$ .
  - ▶ Lower-bound for a *continuum* of problems, each with different  $\kappa_{\mathcal{X}}$ .
  - ▶ Traditional lower-bound technique fails.

## Theorem (A local lower-bound)

Let  $\mathcal{X} = \{\|x\| = 1\}$ , fix  $\theta_* \in \mathbb{R}^d$  and denote  $\kappa = \kappa_*(\theta_*)$ . For any policy

$$\max_{\|\theta' - \theta_*\| \leq \varepsilon} \text{Regret}_{\theta'}(T) = \Omega\left(d\sqrt{T/\kappa}\right) \quad \text{if } T \geq \kappa$$

where  $\varepsilon$  is such that  $\forall \theta' \in \{\|\theta' - \theta\| \leq \varepsilon\}$  we have  $\kappa_*(\theta') = \Theta(\kappa)$ .

# Problem-dependent lower-bound

- **Challenge.** Study optimality w.r.t problem-dependent constants  $\kappa_{\mathcal{X}}$ .
  - ▶ Lower-bound for a *continuum* of problems, each with different  $\kappa_{\mathcal{X}}$ .
  - ▶ Traditional lower-bound technique fails.

## Theorem (A local lower-bound)

Let  $\mathcal{X} = \{\|x\| = 1\}$ , fix  $\theta_{\star} \in \mathbb{R}^d$  and denote  $\kappa = \kappa_{\star}(\theta_{\star})$ . For any policy

$$\max_{\|\theta' - \theta_{\star}\| \leq \varepsilon} \text{Regret}_{\theta'}(T) = \Omega\left(d\sqrt{T/\kappa}\right) \quad \text{if } T \geq \kappa$$

where  $\varepsilon$  is such that  $\forall \theta' \in \{\|\theta' - \theta_{\star}\| \leq \varepsilon\}$  we have  $\kappa_{\star}(\theta') = \Theta(\kappa)$ .

- **Interpretation.** For any problem:
  - ▶ Consider the hardest alternative in nearby instances.
  - ▶ That share the same problem-dependent constant  $\kappa_{\star}$ .
- **Conclusion.** The long-term regret is **tight**.

# Algorithm

- **Algorithm.** OFULog:

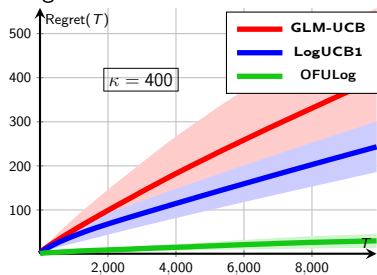
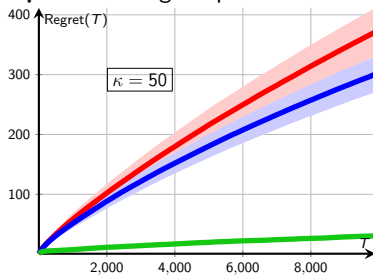
- ▶ Relies on the confidence set  $\mathcal{C}_t(\delta)$  of [Fauray et al. 2020].<sup>1</sup>
- ▶ Parameter-based optimism (vs. bonus-based)

$$x_t = \max_{x \in \mathcal{X}} \max_{\theta \in \mathcal{C}_t(\delta)} x^\top \theta$$

$$(\max_{x \in \mathcal{X}} \mu(x^\top \hat{\theta}_t) + \varepsilon_t(x))$$

- More adaptive to the problem effective's hardness.
- Tractable algorithm (no non-convex optimization routines).

- **In practice.** Large improvement on the regret.



<sup>1</sup>We also introduce a convex relaxation which leads to a fully tractable algorithm

# Bibliography

- Sarah Filippi, Olivier Cappé, Aurélien Garivier, Csaba Szepesvári. *Parametric Bandits: The Generalized Linear Case*, 2010.
- Francis Bach. *Self-Concordant Analysis for Logistic Regression*, 2010.
- Shi Dong, Tengyu Ma, Benjamin Van Roy. *On the Performance of Thompson Sampling on Logistic Bandits*, 2019.
- Louis Fauray, Marc Abeille, Clément Calauzènes, Olivier Fercoq. *Improved Optimistic Algorithms for Logistic Bandits*, 2020.



See you at the Q&A!