

# Variance-Sensitive Confidence Regions for Parametric Bandits

**Louis Faury**

TélécomParis and Criteo

supervised by **Olivier Fercoq** and co-supervised by **Marc Abeille**.

*Paris, October 11, 2021.*

**CRITEO**  
**AI Lab**



# Outline

- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).

# Outline

- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).
  - ▶ Sequential decision making problems

# Outline

- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).
  - ▶ Sequential decision making problems
  - ▶ Isolate the interaction non-linearity  $\leftrightarrow$  exploration/exploitation trade-off.

# Outline

- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).
  - ▶ Sequential decision making problems
  - ▶ Isolate the interaction non-linearity  $\leftrightarrow$  exploration/exploitation trade-off.
- **Previous studies.** Non-linearity is harmful.

# Outline

- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).
  - ▶ Sequential decision making problems
  - ▶ Isolate the interaction non-linearity  $\leftrightarrow$  exploration/exploitation trade-off.
- **Previous studies.** Non-linearity is harmful.
  - ▶ The more non-linear, the **harder** the problem.

# Outline

- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).
  - ▶ Sequential decision making problems
  - ▶ Isolate the interaction non-linearity  $\leftrightarrow$  exploration/exploitation trade-off.
- **Previous studies.** Non-linearity is harmful.
  - ▶ The more non-linear, the **harder** the problem.
  - ▶ Status-quo since  $\approx 10$  years.

# Outline

- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).
  - ▶ Sequential decision making problems
  - ▶ Isolate the interaction non-linearity  $\leftrightarrow$  exploration/exploitation trade-off.
- **Previous studies.** Non-linearity is harmful.
  - ▶ The more non-linear, the **harder** the problem.
  - ▶ Status-quo since  $\approx$  10 years.
- **Today.** A different story:



# Outline

- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).
  - ▶ Sequential decision making problems
  - ▶ Isolate the interaction non-linearity  $\leftrightarrow$  exploration/exploitation trade-off.
- **Previous studies.** Non-linearity is harmful.
  - ▶ The more non-linear, the **harder** the problem.
  - ▶ Status-quo since  $\approx 10$  years.
- **Today.** A different story:
  - ▶ Improved algorithms and refined analysis.

# Outline

- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).
  - ▶ Sequential decision making problems
  - ▶ Isolate the interaction non-linearity  $\leftrightarrow$  exploration/exploitation trade-off.
- **Previous studies.** Non-linearity is harmful.
  - ▶ The more non-linear, the **harder** the problem.
  - ▶ Status-quo since  $\approx 10$  years.
- **Today.** A different story:
  - ▶ Improved algorithms and refined analysis.
  - ▶ Effects of non-linearity are **short-term** (and not always detrimental).

# Outline

- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).
  - ▶ Sequential decision making problems
  - ▶ Isolate the interaction non-linearity  $\leftrightarrow$  exploration/exploitation trade-off.
- **Previous studies.** Non-linearity is harmful.
  - ▶ The more non-linear, the **harder** the problem.
  - ▶ Status-quo since  $\approx 10$  years.
- **Today.** A different story:
  - ▶ Improved algorithms and refined analysis.
  - ▶ Effects of non-linearity are **short-term** (and not always detrimental).
  - ▶ Even better! in some cases:

# Outline

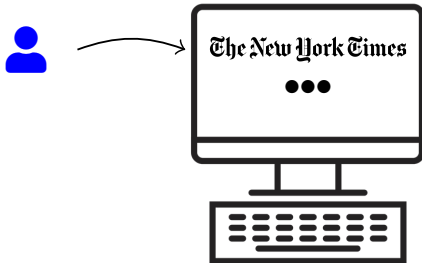
- **Motivation.** Study **non-linearity** in Generalized Linear Bandits (GLBs).
  - ▶ Sequential decision making problems
  - ▶ Isolate the interaction non-linearity  $\leftrightarrow$  exploration/exploitation trade-off.
- **Previous studies.** Non-linearity is harmful.
  - ▶ The more non-linear, the **harder** the problem.
  - ▶ Status-quo since  $\approx 10$  years.
- **Today.** A different story:
  - ▶ Improved algorithms and refined analysis.
  - ▶ Effects of non-linearity are **short-term** (and not always detrimental).
  - ▶ Even better! in some cases:

The more non-linear, the **easier** the problem.

# Motivation and Setting

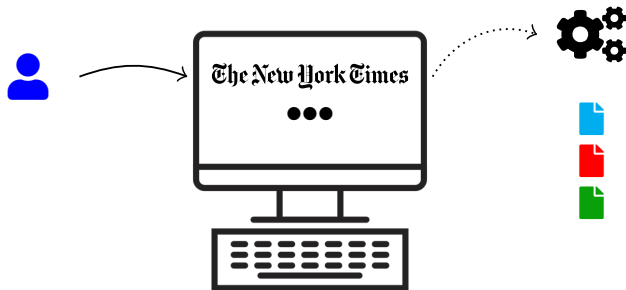
# A motivating example (1/3)

- **News recommendation.** Among others (e.g clinical trials, ..).



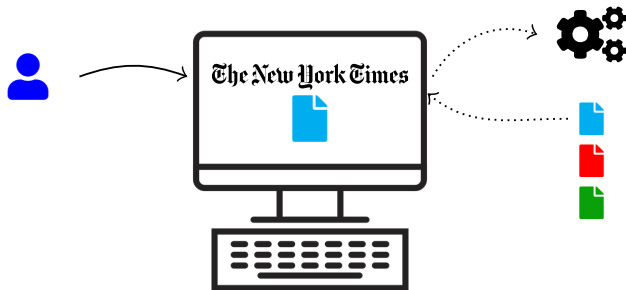
# A motivating example (1/3)

- News recommendation. Among others (e.g clinical trials, ..).



# A motivating example (1/3)

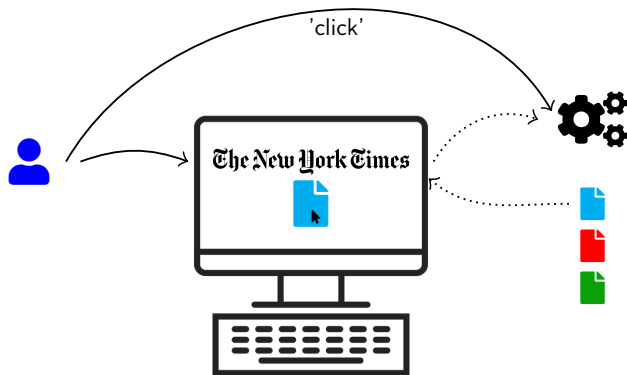
- News recommendation. Among others (e.g clinical trials, ..).





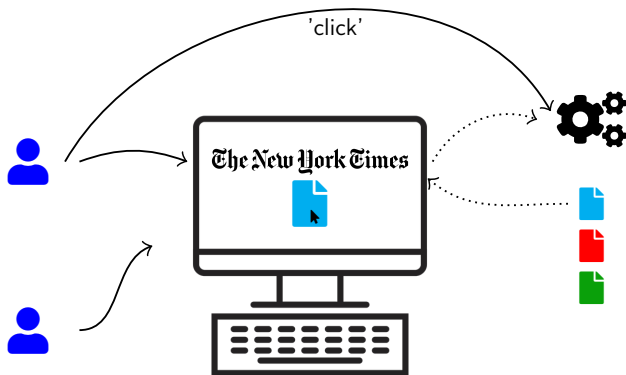
# A motivating example (1/3)

- News recommendation. Among others (e.g clinical trials, ..).



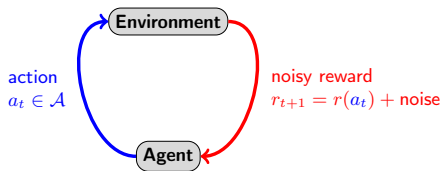
# A motivating example (1/3)

- News recommendation. Among others (e.g clinical trials, ..).



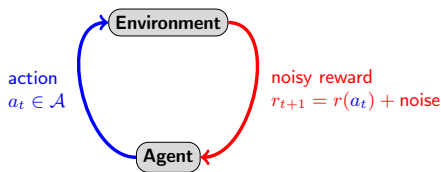
# A motivating example (2/3)

- **Formalization:** stochastic bandit framework.



# A motivating example (2/3)

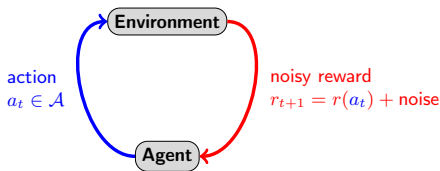
- **Formalization:** stochastic bandit framework.



- **Goal:** minimize  $\text{Regret}(T) = T \max_a r(a) - \sum_{t=1}^T r(a_t)$ .

# A motivating example (2/3)

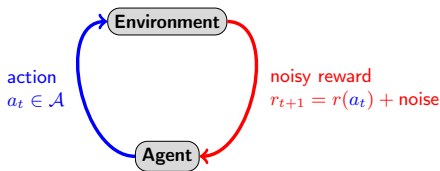
- **Formalization:** stochastic bandit framework.



- **Goal:** minimize  $\text{Regret}(T) = T \max_a r(a) - \sum_{t=1}^T r(a_t)$ .
- **Challenge:** observe (noisy) reward only for the action we play.

# A motivating example (2/3)

- **Formalization:** stochastic bandit framework.



- **Goal:** minimize  $\text{Regret}(T) = T \max_a r(a) - \sum_{t=1}^T r(a_t)$ .
- **Challenge:** observe (noisy) reward only for the action we play.
  - ▶ exploration/exploitation dilemma.

# A motivating example (3/3)

Modelling the reward function.

# A motivating example (3/3)

Modelling the reward function.

- **Challenges.**

(1) Large action space  $\mathcal{A}$  but interrelated payoffs.



# A motivating example (3/3)

Modelling the reward function.

- **Challenges.**

- (1) Large action space  $\mathcal{A}$  but interrelated payoffs.
- (2) Learn from continuous / binary / ordinal / categorical feedback.

# A motivating example (3/3)

Modelling the reward function.

- **Challenges.**

- (1) Large action space  $\mathcal{A}$  but interrelated payoffs.
- (2) Learn from continuous / binary / ordinal / categorical feedback.
- (3) **Theoretical guarantees:**  $\text{Regret}(T) = o(T)$ .

# A motivating example (3/3)

Modelling the reward function.

- **Challenges.**

- (1) Large action space  $\mathcal{A}$  but interrelated payoffs.
- (2) Learn from continuous / binary / ordinal / categorical feedback.
- (3) **Theoretical guarantees:**  $\text{Regret}(T) = o(T)$ .

- Solution for (1): establish **structure** through parametric model.

# A motivating example (3/3)

Modelling the reward function.

- **Challenges.**

- (1) Large action space  $\mathcal{A}$  but interrelated payoffs.
- (2) Learn from continuous / binary / ordinal / categorical feedback.
- (3) **Theoretical guarantees:**  $\text{Regret}(T) = o(T)$ .

- Solution for (1): establish **structure** through parametric model.
  - ▶ Embed  $\mathcal{A}$  in  $\mathbb{R}^d$  with  $d \ll |\mathcal{A}|$  (feature map)

# A motivating example (3/3)

Modelling the reward function.

- **Challenges.**

- (1) Large action space  $\mathcal{A}$  but interrelated payoffs.
- (2) Learn from continuous / binary / ordinal / categorical feedback.
- (3) **Theoretical guarantees:**  $\text{Regret}(T) = o(T)$ .

- Solution for (1): establish **structure** through parametric model.

- ▶ Embed  $\mathcal{A}$  in  $\mathbb{R}^d$  with  $d \ll |\mathcal{A}|$  (feature map)
- ▶ Reward function belongs to known parametric family:

$$\mathbb{E}[r_{t+1}|a_t] = f_{\theta_*}(a_t) \quad \text{where } f_{\theta_*} \in \left\{ f_{\theta} : \mathbb{R}^d \mapsto \mathbb{R}, \theta \in \Theta \right\},$$

where  $\theta_*$  is shared but **unknown**.

# Generalized Linear Bandits *[Filippi et al. 2010]*

- Expected reward follows a Generalized Linear model:

$$\mathbb{E}[r_{t+1} | a_t] = \mu(a_t^\top \theta_*)$$

where  $\mu$  is continuously differentiable, strictly increasing.

# Generalized Linear Bandits [Filippi et al. 2010]

- Expected reward follows a Generalized Linear model:

$$\mathbb{E}[r_{t+1}|a_t] = \mu(a_t^\top \theta_*)$$

where  $\mu$  is continuously differentiable, strictly increasing.

- **Reward distribution.** Exponential family with underlying **linear** structure:

$$d\mathbb{P}(r|a) \propto \exp(r a^\top \theta_* - b(a^\top \theta_*)) d\nu(r)$$

covers Gaussian, Bernoulli, Poisson, .. distributions.

# Generalized Linear Bandits [Filippi et al. 2010]

- Expected reward follows a Generalized Linear model:

$$\mathbb{E}[r_{t+1}|a_t] = \mu(a_t^\top \theta_*)$$

where  $\mu$  is continuously differentiable, strictly increasing.

- **Reward distribution.** Exponential family with underlying **linear** structure:

$$d\mathbb{P}(r|a) \propto \exp(r a^\top \theta_* - b(a^\top \theta_*)) d\nu(r)$$

covers Gaussian, Bernoulli, Poisson, .. distributions.  $\rightarrow$  challenge (2).



# Generalized Linear Bandits [Filippi et al. 2010]

- Expected reward follows a Generalized Linear model:

$$\mathbb{E}[r_{t+1}|a_t] = \mu(a_t^\top \theta_*)$$

where  $\mu$  is continuously differentiable, strictly increasing.

- **Reward distribution.** Exponential family with underlying **linear** structure:

$$d\mathbb{P}(r|a) \propto \exp(r a^\top \theta_* - b(a^\top \theta_*)) d\nu(r)$$

covers Gaussian, Bernoulli, Poisson, .. distributions.  $\rightarrow$  challenge (2).

- **Learnability.** Maximum-likelihood principle

$$\hat{\theta}_t := \operatorname{argmin}_{\theta} \sum_{s=1}^{t-1} -\log d\mathbb{P}(r_{s+1}|a_s) + \lambda \|\theta\|^2 / 2,$$

# Generalized Linear Bandits [Filippi et al. 2010]

- **Expected reward** follows a Generalized Linear model:

$$\mathbb{E}[r_{t+1}|a_t] = \mu(a_t^\top \theta_*)$$

where  $\mu$  is continuously differentiable, strictly increasing.

- **Reward distribution.** Exponential family with underlying **linear** structure:

$$d\mathbb{P}(r|a) \propto \exp(r a^\top \theta_* - b(a^\top \theta_*)) d\nu(r)$$

covers Gaussian, Bernoulli, Poisson, .. distributions.  $\rightarrow$  challenge (2).

- **Learnability.** Maximum-likelihood principle  $\rightarrow$  challenge (3)

$$\hat{\theta}_t := \operatorname{argmin}_{\theta} \sum_{s=1}^{t-1} -\log d\mathbb{P}(r_{s+1}|a_s) + \lambda \|\theta\|^2 / 2,$$

# An illustration: the Logistic Bandit

- **Logistic Bandit.** Structured **binary** feedback:

$$r_{t+1} \sim \text{Bernoulli}(\mu(\mathbf{a}_t^\top \boldsymbol{\theta}_*))$$

where  $\mu(z) = (1 + \exp(-z))^{-1}$  is the logistic function.

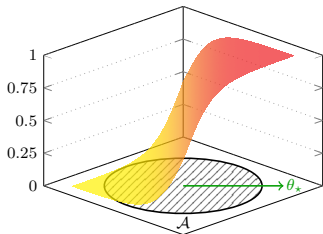
# An illustration: the Logistic Bandit

- **Logistic Bandit.** Structured **binary** feedback:

$$r_{t+1} \sim \text{Bernoulli}(\mu(a_t^\top \theta_*))$$

where  $\mu(z) = (1 + \exp(-z))^{-1}$  is the logistic function.

- Two-dimensional illustration:



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

# Generalized Linear Bandits: beyond linearity

- **Linear Bandit** (LB). Special case with  $\mu = \text{Id}$ :

$$\mathbb{E}[r_{t+1}|a_t] = a_t^\top \theta_*$$

# Generalized Linear Bandits: beyond linearity

- **Linear Bandit** (LB). Special case with  $\mu = \text{Id}$ :

$$\mathbb{E}[r_{t+1}|a_t] = a_t^\top \theta_* .$$

- ▶ Well-understood: [Auer. 2002, Dani et al. 2008, Abbasi-Yadkori et al. 2011, ..].
- ▶ Minimax-optimal and efficient algorithms:

$$\text{Regret}(T) = \tilde{O}(d\sqrt{T}) .$$

# Generalized Linear Bandits: beyond linearity

- **Linear Bandit (LB)**. Special case with  $\mu = \text{Id}$ :

$$\mathbb{E}[r_{t+1}|a_t] = a_t^\top \theta_*$$

- ▶ Well-understood: [Auer. 2002, Dani et al. 2008, Abbasi-Yadkori et al. 2011, ..].
- ▶ Minimax-optimal and efficient algorithms:

$$\text{Regret}(T) = \tilde{O}(d\sqrt{T}) .$$

- **GLBs**. Beyond linearity;

$$\mathbb{E}[r_{t+1}|a_t] = \mu(a_t^\top \theta_*) .$$

# Generalized Linear Bandits: beyond linearity

- **Linear Bandit (LB)**. Special case with  $\mu = \text{Id}$ :

$$\mathbb{E}[r_{t+1}|a_t] = a_t^\top \theta_*$$

- ▶ Well-understood: [Auer. 2002, Dani et al. 2008, Abbasi-Yadkori et al. 2011, ..].
- ▶ Minimax-optimal and efficient algorithms:

$$\text{Regret}(T) = \tilde{O}(d\sqrt{T}) .$$

- **GLBs**. Beyond linearity;

$$\mathbb{E}[r_{t+1}|a_t] = \mu(a_t^\top \theta_*) .$$

- ▶ minimalistic **non-linear** extension of LB.



# Generalized Linear Bandits: beyond linearity

- **Linear Bandit (LB)**. Special case with  $\mu = \text{Id}$ :

$$\mathbb{E}[r_{t+1}|a_t] = a_t^\top \theta_*$$

- ▶ Well-understood: [Auer. 2002, Dani et al. 2008, Abbasi-Yadkori et al. 2011, ..].
- ▶ Minimax-optimal and efficient algorithms:

$$\text{Regret}(T) = \tilde{O}(d\sqrt{T}) .$$

- **GLBs**. Beyond linearity;

$$\mathbb{E}[r_{t+1}|a_t] = \mu(a_t^\top \theta_*) .$$

- ▶ minimalistic **non-linear** extension of LB.
- ▶ first step towards richer reward signal.

# GLBs: quantifying non-linearity

- Level of non-linearity = conditioning of the reward signal =  $\kappa_{\mu}(\theta_*, \mathcal{A})$ .

# GLBs: quantifying non-linearity

- Level of non-linearity = conditioning of the reward signal =  $\kappa_\mu(\theta_*, \mathcal{A})$ .

$$\kappa_\mu(\theta_*, \mathcal{A}) := \frac{\max_{a \in \mathcal{A}} \dot{\mu}(a^\top \theta_*)}{\min_{a \in \mathcal{A}} \dot{\mu}(a^\top \theta_*)} =: \frac{\mathbf{L}_\mu}{\mathbf{l}_\mu}$$

# GLBs: quantifying non-linearity

- Level of non-linearity = conditioning of the reward signal =  $\kappa_\mu(\theta_*, \mathcal{A})$ .

$$\kappa_\mu(\theta_*, \mathcal{A}) := \frac{\max_{a \in \mathcal{A}} \dot{\mu}(a^\top \theta_*)}{\min_{a \in \mathcal{A}} \dot{\mu}(a^\top \theta_*)} =: \frac{\mathbf{L}_\mu}{\mathbf{l}_\mu}$$

- ▶ the more non-linear the reward signal, the larger  $\kappa_\mu$ .

# GLBs: quantifying non-linearity

- Level of non-linearity = conditioning of the reward signal =  $\kappa_\mu(\theta_*, \mathcal{A})$ .

$$\kappa_\mu(\theta_*, \mathcal{A}) := \frac{\max_{a \in \mathcal{A}} \dot{\mu}(a^\top \theta_*)}{\min_{a \in \mathcal{A}} \dot{\mu}(a^\top \theta_*)} =: \frac{\mathbf{L}_\mu}{\mathbf{l}_\mu}$$

- ▶ the more non-linear the reward signal, the larger  $\kappa_\mu$ .
- ▶ “distance” from the linear model ( $\kappa_\mu = 1$  for Linear Bandit).

# GLBs: quantifying non-linearity

- Level of non-linearity = conditioning of the reward signal =  $\kappa_\mu(\theta_*, \mathcal{A})$ .

$$\kappa_\mu(\theta_*, \mathcal{A}) := \frac{\max_{a \in \mathcal{A}} \dot{\mu}(a^\top \theta_*)}{\min_{a \in \mathcal{A}} \dot{\mu}(a^\top \theta_*)} =: \frac{\mathbf{L}_\mu}{\mathbf{l}_\mu}$$

- ▶ the more non-linear the reward signal, the larger  $\kappa_\mu$ .
- ▶ “distance” from the linear model ( $\kappa_\mu = 1$  for Linear Bandit).
- ▶ numerically very large ( $\kappa_\mu \propto \exp(\|\theta_*\|) \approx 10^3!$ )

Previous work, limitations,  
contributions.

# GLBs: previous work

- First studied in the seminal work of *[Filippi et al. 2010]*.
  - ▶ many extensions: *[Li et al. 2017, Jun et al. 2017, Kveton et al. 2019, ...]*



# GLBs: previous work

- First studied in the seminal work of [Filippi et al. 2010].
  - ▶ many extensions: [Li et al. 2017, Jun et al. 2017, Kveton et al. 2019, ..]
- **Regret upper-bound.** With high probability:

$$\text{Regret}(T) = \tilde{O}(d\sqrt{T}) .$$

# GLBs: previous work

- First studied in the seminal work of [Filippi et al. 2010].
  - ▶ many extensions: [Li et al. 2017, Jun et al. 2017, Kveton et al. 2019, ..]
- **Regret upper-bound.** With high probability:

$$\text{Regret}(T) = \tilde{O}(d\sqrt{T}) .$$

- ✔ Extend LB tools to generic GLBs.

# GLBs: previous work

- First studied in the seminal work of [Filippi et al. 2010].
  - ▶ many extensions: [Li et al. 2017, Jun et al. 2017, Kveton et al. 2019, ..]
- **Regret upper-bound.** With high probability:

$$\text{Regret}(T) = \tilde{O}(\kappa_{\mu} d \sqrt{T}) .$$

- ✓ Extend LB tools to generic GLBs.

# GLBs: previous work

- First studied in the seminal work of [Filippi et al. 2010].
  - ▶ many extensions: [Li et al. 2017, Jun et al. 2017, Kveton et al. 2019, ..]
- **Regret upper-bound.** With high probability:

$$\text{Regret}(T) = \tilde{O}(\kappa_{\mu} d \sqrt{T}) .$$

- ✓ Extend LB tools to generic GLBs.
- ✗ Large regret upper bound.

# GLBs: previous work

- First studied in the seminal work of [Filippi et al. 2010].
  - ▶ many extensions: [Li et al. 2017, Jun et al. 2017, Kveton et al. 2019, ..]
- **Regret upper-bound.** With high probability:

$$\text{Regret}(T) = \tilde{O}(\kappa_{\mu} d \sqrt{T}) .$$

- ✓ Extend LB tools to generic GLBs.
- ✗ Large regret upper bound.
- ✗ Over-exploratory algorithms, poor empirical performance.

# GLBs: previous work

- First studied in the seminal work of [Filippi et al. 2010].
  - ▶ many extensions: [Li et al. 2017, Jun et al. 2017, Kveton et al. 2019, ..]

- **Regret upper-bound.** With high probability:

$$\text{Regret}(T) = \tilde{O}(\kappa_{\mu} d \sqrt{T}) .$$

- ✓ Extend LB tools to generic GLBs.
  - ✗ Large regret upper bound.
  - ✗ Over-exploratory algorithms, poor empirical performance.
- **Learning-theoretic:** non-linearity is **detrimental!**

# GLBs: previous work

- First studied in the seminal work of [Filippi et al. 2010].
  - ▶ many extensions: [Li et al. 2017, Jun et al. 2017, Kveton et al. 2019, ..]

- **Regret upper-bound.** With high probability:

$$\text{Regret}(T) = \tilde{O}(\kappa_{\mu} d \sqrt{T}) .$$

- ✓ Extend LB tools to generic GLBs.
  - ✗ Large regret upper bound.
  - ✗ Over-exploratory algorithms, poor empirical performance.
- **Learning-theoretic:** non-linearity is **detrimental!**
  - ▶ the more non-linear the problem, the worse the performance.

# GLBs: previous approach (1/2)

- **Algorithmic design.** Two main ingredients; at each round  $t$ :



# GLBs: previous approach (1/2)

- **Algorithmic design.** Two main ingredients; at each round  $t$ :

1. Confidence set  $\mathcal{E}_t(\delta)$  for  $\theta_*$ ;

$$\mathbb{P}\left(\forall t \geq 1, \theta_* \in \mathcal{E}_t(\delta)\right) \geq 1 - \delta.$$

# GLBs: previous approach (1/2)

- **Algorithmic design.** Two main ingredients; at each round  $t$ :

1. Confidence set  $\mathcal{E}_t(\delta)$  for  $\theta_*$ ;

$$\mathbb{P}\left(\forall t \geq 1, \theta_* \in \mathcal{E}_t(\delta)\right) \geq 1 - \delta.$$

2. Optimism in face of uncertainty:

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

# GLBs: previous approach (1/2)

- **Algorithmic design.** Two main ingredients; at each round  $t$ :

1. Confidence set  $\mathcal{E}_t(\delta)$  for  $\theta_*$ ;

$$\mathbb{P}(\forall t \geq 1, \theta_* \in \mathcal{E}_t(\delta)) \geq 1 - \delta.$$

2. Optimism in face of uncertainty:

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

- **Confidence set** from previous works; with  $\mathbf{V}_t = \sum_{s=1}^{t-1} a_s a_s^\top + \lambda \mathbf{I}_d$ :

$$\mathcal{E}_t(\delta) = \left\{ \theta, \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{V}_t} \leq \sqrt{d \log(t/\delta)} / \ell_\mu \right\}$$

# GLBs: previous approach (1/2)

- **Algorithmic design.** Two main ingredients; at each round  $t$ :

1. Confidence set  $\mathcal{E}_t(\delta)$  for  $\theta_*$ ;

$$\mathbb{P}\left(\forall t \geq 1, \theta_* \in \mathcal{E}_t(\delta)\right) \geq 1 - \delta.$$

2. Optimism in face of uncertainty:

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

- **Confidence set** from previous works; with  $\mathbf{V}_t = \sum_{s=1}^{t-1} a_s a_s^\top + \lambda \mathbf{I}_d$ :

$$\mathcal{E}_t(\delta) = \left\{ \theta, \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{V}_t} \leq \sqrt{d \log(t/\delta)} / \ell_\mu \right\}$$

- ▶ Radius  $\propto \ell_\mu^{-1} \Rightarrow$  large confidence region  $\Rightarrow$  aggressive exploration.

# GLBs: previous approach (1/2)

- **Algorithmic design.** Two main ingredients; at each round  $t$ :

1. Confidence set  $\mathcal{E}_t(\delta)$  for  $\theta_*$ ;

$$\mathbb{P}(\forall t \geq 1, \theta_* \in \mathcal{E}_t(\delta)) \geq 1 - \delta.$$

2. Optimism in face of uncertainty:

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

- **Confidence set** from previous works; with  $\mathbf{V}_t = \sum_{s=1}^{t-1} a_s a_s^\top + \lambda \mathbf{I}_d$ :

$$\mathcal{E}_t(\delta) = \left\{ \theta, \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{V}_t} \leq \sqrt{d \log(t/\delta)} / \ell_\mu \right\}$$

- ▶ Radius  $\propto \ell_\mu^{-1} \Rightarrow$  large confidence region  $\Rightarrow$  aggressive exploration.
- ▶ Learn as slow as in the flattest region, in every direction.

## GLBs: previous approach (2/2)

- **Analysis.** Upper linear bound:

$$\text{Regret}(T) = \sum_{t=1}^T \mu(a_{\star}^{\top} \theta_{\star}) - \mu(a_t^{\top} \theta_{\star})$$

## GLBs: previous approach (2/2)

- **Analysis.** Upper linear bound:

$$\text{Regret}(T) \leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_*) \quad (\theta_t \text{ opt. param})$$

## GLBs: previous approach (2/2)

- **Analysis.** Upper linear bound:

$$\begin{aligned} \text{Regret}(T) &\leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_\star) && (\theta_t \text{ opt. param}) \\ &\leq \mathbf{L}_\mu \sum_{t=1}^T a_t^\top (\theta_t - \theta_\star) \end{aligned}$$



# GLBs: previous approach (2/2)

- **Analysis.** Upper linear bound:

$$\begin{aligned} \text{Regret}(T) &\leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_\star) && (\theta_t \text{ opt. param}) \\ &\leq \mathbf{L}_\mu \sum_{t=1}^T a_t^\top (\theta_t - \theta_\star) \\ &\leq \mathbf{L}_\mu \sum_{t=1}^T \|a_t\|_{\mathbf{V}_t^{-1}} \|\theta_t - \theta_\star\|_{\mathbf{V}_t} && \text{(C.S ineq.)} \end{aligned}$$

# GLBs: previous approach (2/2)

- **Analysis.** Upper linear bound:

$$\begin{aligned}\text{Regret}(T) &\leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_\star) && (\theta_t \text{ opt. param}) \\ &\leq \mathbf{L}_\mu \sum_{t=1}^T a_t^\top (\theta_t - \theta_\star) \\ &\leq \mathbf{L}_\mu \sum_{t=1}^T \|a_t\|_{\mathbf{V}_t^{-1}} \|\theta_t - \theta_\star\|_{\mathbf{V}_t} && \text{(C.S ineq.)} \\ &\leq \mathbf{L}_\mu / \check{\ell}_\mu \sqrt{d \log(T/\delta)} \sum_{t=1}^T \|a_t\|_{\mathbf{V}_t^{-1}}\end{aligned}$$

# GLBs: previous approach (2/2)

- **Analysis.** Upper linear bound:

$$\text{Regret}(T) \leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_\star) \quad (\theta_t \text{ opt. param})$$

$$\leq \mathbf{L}_\mu \sum_{t=1}^T a_t^\top (\theta_t - \theta_\star)$$

$$\leq \mathbf{L}_\mu \sum_{t=1}^T \|a_t\|_{\mathbf{V}_t^{-1}} \|\theta_t - \theta_\star\|_{\mathbf{V}_t} \quad (\text{C.S ineq.})$$

$$\leq \mathbf{L}_\mu / \ell_\mu \sqrt{d \log(T/\delta)} \sum_{t=1}^T \|a_t\|_{\mathbf{V}_t^{-1}}$$

$$= \kappa_\mu d \sqrt{T} \log(T/\delta)$$

## GLBs: previous approach (2/2)

- **Analysis.** Upper linear bound:

$$\text{Regret}(T) \leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_\star) \quad (\theta_t \text{ opt. param})$$

$$\leq \mathbf{L}_\mu \sum_{t=1}^T a_t^\top (\theta_t - \theta_\star)$$

$$\leq \mathbf{L}_\mu \sum_{t=1}^T \|a_t\|_{\mathbf{V}_t^{-1}} \|\theta_t - \theta_\star\|_{\mathbf{V}_t} \quad (\text{C.S ineq.})$$

$$\leq \mathbf{L}_\mu / \ell_\mu \sqrt{d \log(T/\delta)} \sum_{t=1}^T \|a_t\|_{\mathbf{V}_t^{-1}}$$

$$= \kappa_\mu d \sqrt{T} \log(T/\delta)$$

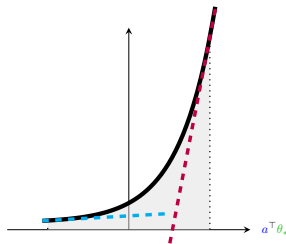
- We pay errors as in the **sharpest** linear case.

# GLBs: previous approach (2/2)

- **Analysis.** Upper linear bound:

$$\begin{aligned}\text{Regret}(T) &\leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_\star) && (\theta_t \text{ opt. param}) \\ &\leq \mathbf{L}_\mu \sum_{t=1}^T a_t^\top (\theta_t - \theta_\star) \\ &\leq \mathbf{L}_\mu \sum_{t=1}^T \|a_t\|_{\mathbf{V}_t^{-1}} \|\theta_t - \theta_\star\|_{\mathbf{V}_t} && \text{(C.S ineq.)} \\ &\leq \mathbf{L}_\mu / \underline{\ell}_\mu \sqrt{d \log(T/\delta)} \sum_{t=1}^T \|a_t\|_{\mathbf{V}_t^{-1}} \\ &= \kappa_\mu d \sqrt{T} \log(T/\delta)\end{aligned}$$

- We pay errors as in the **sharpest** linear case.
- Worst-case **errors** / worst-case **learning**.



# GLBs: our approach

**Local** treatment of non-linearity for **improved** regret bounds.

# GLBs: our approach

**Local** treatment of non-linearity for **improved** regret bounds.

- **New confidence set.**
  - ▶ sensitive to effective reward sensitivity ( $\neq$  worst-case)
  - ▶ provably tighter.

# GLBs: our approach

**Local** treatment of non-linearity for **improved** regret bounds.

- **New confidence set.**
  - ▶ sensitive to effective reward sensitivity ( $\neq$  worst-case)
  - ▶ provably tighter.
- **Locality-sensitive analysis** under generalized self-concordance [*Bach. 2010*]:

$$|\ddot{\mu}| \leq c\dot{\mu}$$

- ▶ allows exact Taylor control with local quantities.



# GLBs: our approach

**Local** treatment of non-linearity for **improved** regret bounds.

- **New confidence set.**
  - ▶ sensitive to effective reward sensitivity ( $\neq$  worst-case)
  - ▶ provably tighter.
- **Locality-sensitive analysis** under generalized self-concordance [*Bach. 2010*]:

$$|\ddot{\mu}| \leq c\dot{\mu}$$

- ▶ allows exact Taylor control with local quantities.
- ▶ not restrictive: Logistic and Poisson Bandits ( $c=1$ ).

# Variance-Sensitive Confidence Sets for GLBs

# Improved confidence set: asymptotic intuition

- **Objective.** Dependence to effective reward sensitivity:
  - ▶ measured through the **variance** of the reward signal:

$$\text{Var}(r_{t+1}|a_t) = \dot{\mu}(a_t^\top \theta_\star) .$$

# Improved confidence set: asymptotic intuition

- **Objective.** Dependence to effective reward sensitivity:
  - ▶ measured through the **variance** of the reward signal:

$$\text{Var}(r_{t+1}|a_t) = \dot{\mu}(a_t^\top \theta_*) .$$

- ▶ Variance-sensitive concentration tools.

# Improved confidence set: asymptotic intuition

- **Objective.** Dependence to effective reward sensitivity:
  - ▶ measured through the **variance** of the reward signal:

$$\text{Var}(r_{t+1}|a_t) = \dot{\mu}(a_t^\top \theta_\star).$$

- ▶ Variance-sensitive concentration tools.
- **Asymptotic intuition.** Let  $\mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^\top \theta) a_s a_s^\top$ .

# Improved confidence set: asymptotic intuition

- **Objective.** Dependence to effective reward sensitivity:
  - ▶ measured through the **variance** of the reward signal:

$$\text{Var}(r_{t+1}|a_t) = \dot{\mu}(a_t^\top \theta_\star).$$

- ▶ Variance-sensitive concentration tools.
- **Asymptotic intuition.** Let  $\mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^\top \theta) a_s a_s^\top$ .

$$\lim_{t \rightarrow \infty} \mathbb{P} \left( \|\hat{\theta}_t - \theta_\star\|_{\mathbf{H}_t(\theta_\star)}^2 \leq d \log(1/\delta) \right) \geq 1 - \delta.$$

under **random** design.

# Improved confidence set: asymptotic intuition

- **Objective.** Dependence to effective reward sensitivity:
  - ▶ measured through the **variance** of the reward signal:

$$\text{Var}(r_{t+1}|a_t) = \dot{\mu}(a_t^\top \theta_\star) .$$

- ▶ Variance-sensitive concentration tools.
- **Asymptotic intuition.** Let  $\mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^\top \theta) a_s a_s^\top$ .

$$\lim_{t \rightarrow \infty} \mathbb{P} \left( \|\hat{\theta}_t - \theta_\star\|_{\mathbf{H}_t(\theta_\star)}^2 \leq d \log(1/\delta) \right) \geq 1 - \delta .$$

under **random** design.

- **Challenge.** Generalization for:
  - ▶ finite-time (non-asymptotic).
  - ▶ **adaptive** design ( $\{a_1, \dots, a_s\}_s$  are not independent).

# A novel concentration inequality (1/2)

Theorem (F., Abeille, Calauzènes and Fercoq, 2020.)

Let  $\{\mathcal{F}_t\}_{t \in \mathbb{N}}$  be a filtration and:

- $\{a_t\}_{t \in \mathbb{N}}$  a  $\mathcal{F}_t$ -measurable stochastic process.



# A novel concentration inequality (1/2)

Theorem (F., Abeille, Calauzènes and Fercoq, 2020.)

Let  $\{\mathcal{F}_t\}_{t \in \mathbb{N}}$  be a filtration and:

- $\{a_t\}_{t \in \mathbb{N}}$  a  $\mathcal{F}_t$ -measurable stochastic process.
- $\{\eta_{t+1}\}_{t \in \mathbb{N}}$  a  $\mathcal{F}_{t+1}$ -measurable martingale difference sequence s.t:
  - ▶  $|\eta_{t+1}| \leq \sigma$  almost surely and  $v_t^2 := \text{Var}(\eta_{t+1} | \mathcal{F}_t)$ .

# A novel concentration inequality (1/2)

Theorem (F., Abeille, Calauzènes and Fercoq, 2020.)

Let  $\{\mathcal{F}_t\}_{t \in \mathbb{N}}$  be a filtration and:

- $\{a_t\}_{t \in \mathbb{N}}$  a  $\mathcal{F}_t$ -measurable stochastic process.
- $\{\eta_{t+1}\}_{t \in \mathbb{N}}$  a  $\mathcal{F}_{t+1}$ -measurable martingale difference sequence s.t:
  - ▶  $|\eta_{t+1}| \leq \sigma$  almost surely and  $\mathbf{v}_t^2 := \mathbb{V}\text{ar}(\eta_{t+1} | \mathcal{F}_t)$ .
- Let  $\lambda > 0$  and define for  $t \in \mathbb{N}$ :

$$S_{t+1} := \sum_{s=1}^t \eta_{s+1} a_s \quad \text{and} \quad H_t := \sum_{s=1}^t \mathbf{v}_s^2 a_s a_s^\top + \lambda \mathbf{I}_d$$

# A novel concentration inequality (1/2)

Theorem (F., Abeille, Calauzènes and Fercoq, 2020.)

For  $\delta \in (0, 1]$  the event:

$$\forall t \geq 1, \|S_{t+1}\|_{\mathbf{H}_t^{-1}} \leq \frac{\sqrt{\lambda}}{2\sigma} + \frac{2\sigma}{\sqrt{\lambda}} d \log \left( \frac{4(1 + \sigma^2 t / (d\lambda))}{\delta} \right),$$

holds with probability at least  $1 - \delta$ .

Let  $\{\mathcal{F}_t\}_{t \in \mathbb{N}}$  be a filtration and:

- $\{a_t\}_{t \in \mathbb{N}}$  a  $\mathcal{F}_t$ -measurable stochastic process.
- $\{\eta_{t+1}\}_{t \in \mathbb{N}}$  a  $\mathcal{F}_{t+1}$ -measurable martingale difference sequence s.t:
  - ▶  $|\eta_{t+1}| \leq \sigma$  almost surely and  $\mathbf{v}_t^2 := \mathbb{V}\text{ar}(\eta_{t+1} | \mathcal{F}_t)$ .
- Let  $\lambda > 0$  and define for  $t \in \mathbb{N}$ :

$$S_{t+1} := \sum_{s=1}^t \eta_{s+1} a_s \text{ and } \mathbf{H}_t := \sum_{s=1}^t \mathbf{v}_s^2 a_s a_s^\top + \lambda \mathbf{I}_d$$

# A novel concentration inequality (1/2)

Theorem (F., Abeille, Calauzènes and Fercoq, 2020.)

For  $\delta \in (0, 1]$  the event:

$$\forall t \geq 1, \|S_{t+1}\|_{\mathbf{H}_t^{-1}} \leq \mathcal{O}\left(\sqrt{d \log(t/\delta)}\right),$$

holds with probability at least  $1 - \delta$  (setting  $\lambda_t = d \log(t/\delta)$ ).

Let  $\{\mathcal{F}_t\}_{t \in \mathbb{N}}$  be a filtration and:

- $\{a_t\}_{t \in \mathbb{N}}$  a  $\mathcal{F}_t$ -measurable stochastic process.
- $\{\eta_{t+1}\}_{t \in \mathbb{N}}$  a  $\mathcal{F}_{t+1}$ -measurable martingale difference sequence s.t:
  - ▶  $|\eta_{t+1}| \leq \sigma$  almost surely and  $\mathbf{v}_t^2 := \mathbb{V}\text{ar}(\eta_{t+1} | \mathcal{F}_t)$ .
- Let  $\lambda > 0$  and define for  $t \in \mathbb{N}$ :

$$S_{t+1} := \sum_{s=1}^t \eta_{s+1} a_s \quad \text{and} \quad \mathbf{H}_t := \sum_{s=1}^t \mathbf{v}_s^2 a_s a_s^\top + \lambda \mathbf{I}_d$$

## A novel concentration inequality (2/2)

Theorem (F., Abeille, Calauzènes and Fercoq, 2020.)

For  $\delta \in (0, 1]$  the event:

$$\forall t \geq 1, \|S_{t+1}\|_{\mathbf{H}_t^{-1}} \leq \mathcal{O}\left(\sqrt{d \log(t/\delta)}\right),$$

holds with probability at least  $1 - \delta$  (setting  $\lambda_t = d \log(t/\delta)$ ).

- Sketch of proof.

## A novel concentration inequality (2/2)

Theorem (F., Abeille, Calauzènes and Fercoq, 2020.)

For  $\delta \in (0, 1]$  the event:

$$\forall t \geq 1, \|S_{t+1}\|_{\mathbf{H}_t^{-1}} \leq \mathcal{O}\left(\sqrt{d \log(t/\delta)}\right),$$

holds with probability at least  $1 - \delta$  (setting  $\lambda_t = d \log(t/\delta)$ ).

- **Sketch of proof.**

- ▶ Pseudo-maximization (methods of mixture) [de la Peña. 2007].

## A novel concentration inequality (2/2)

Theorem (F., Abeille, Calauzènes and Fercoq, 2020.)

For  $\delta \in (0, 1]$  the event:

$$\forall t \geq 1, \|S_{t+1}\|_{\mathbf{H}_t^{-1}} \leq \mathcal{O}\left(\sqrt{d \log(t/\delta)}\right),$$

holds with probability at least  $1 - \delta$  (setting  $\lambda_t = d \log(t/\delta)$ ).

- **Sketch of proof.**

- ▶ Pseudo-maximization (methods of mixture) [de la Peña. 2007].
- ▶ Similar to [Abbasi-Yadkori et al. 2011], different base super-martingale:

$$M_t(\xi) = \xi^\top S_{t+1} - \|\xi\|_{\mathbf{H}_t}^2 \quad \text{for } \|\xi\| \leq 1.$$

- ▶ Bernstein vs. Hoeffding conditions.

# Application to GLBs (1/2)

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda t}$ :

$$\forall t \geq 1, \quad \left\| \theta_\star - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_\star)} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)}$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^\top \theta_\star) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^\top \theta) a_s a_s^\top + \lambda_t \mathbf{I}_d$$



# Application to GLBs (1/2)

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda t}$ :

$$\forall t \geq 1, \quad \left\| \theta_\star - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_\star)} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)} = \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right)$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^\top \theta_\star) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^\top \theta) a_s a_s^\top + \lambda_t \mathbf{I}_d$$

# Application to GLBs (1/2)

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_\star - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_\star)} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)} = \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right)$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^\top \theta_\star) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^\top \theta) a_s a_s^\top + \lambda_t \mathbf{I}_d$$

- New confidence set:

Proposition (F., Abeille, Calauzènes and Fercoq, 2020.)

For  $\delta \in (0, 1]$  let:

$$\mathcal{C}_t(\delta) := \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\} .$$

Then  $\mathbb{P}(\forall t \geq 1, \theta_\star \in \mathcal{C}_t(\delta)) \geq 1 - \delta$ .

## Application to GLBs (2/2)

$$\mathcal{C}_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\}, \quad (\text{ours})$$

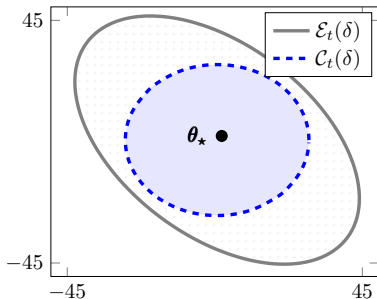
$$\mathcal{E}_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{V}_t} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} / \ell_\mu \right) \right\}. \quad [\text{Filippi et al.}]$$

# Application to GLBs (2/2)

$$C_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\}, \quad (\text{ours})$$

$$\mathcal{E}_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{V}_t} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} / \ell_\mu \right) \right\}. \quad [\text{Filippi et al.}]$$

- Illustration for Logistic Bandit:



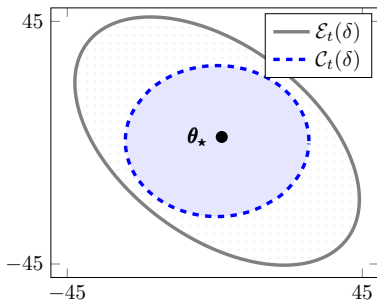
$$\|\theta_*\| = 2 \text{ and } 1/\ell_\mu \approx 10$$

# Application to GLBs (2/2)

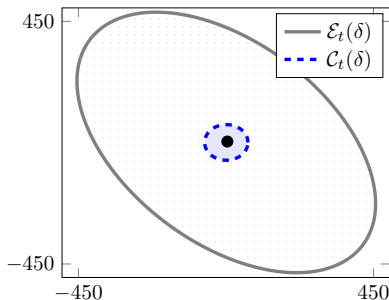
$$C_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\}, \quad (\text{ours})$$

$$\mathcal{E}_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{V}_t} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} / \ell_\mu \right) \right\}. \quad [\text{Filippi et al.}]$$

- Illustration for Logistic Bandit:



$$\|\theta_*\| = 2 \text{ and } 1/\ell_\mu \approx 10$$

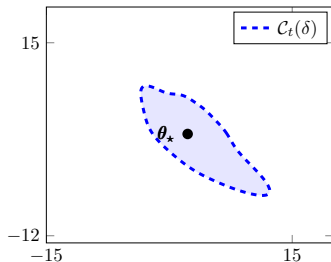


$$\|\theta_*\| = 5 \text{ and } 1/\ell_\mu \approx 150$$

## Extension: convex relaxation

$$\mathcal{C}_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\}$$

- Non-convex, burdensome to manipulate.

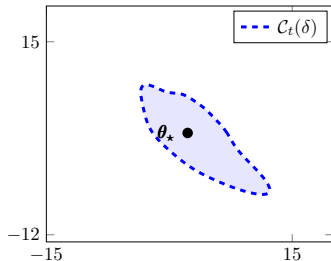


# Extension: convex relaxation

$$\mathcal{C}_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\}$$

- Non-convex, burdensome to manipulate.
- **Convex** relaxation based on log-loss  $\mathcal{L}_t$ :

$$\mathcal{C}_t^c(\delta) = \left\{ \mathcal{L}_t(\theta) - \mathcal{L}_t(\hat{\theta}_t) \leq \mathcal{O} (d \log(t/\delta)) \right\} .$$

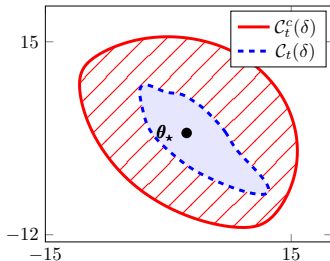


# Extension: convex relaxation

$$\mathcal{C}_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\}$$

- Non-convex, burdensome to manipulate.
- **Convex** relaxation based on log-loss  $\mathcal{L}_t$ :

$$\mathcal{C}_t^c(\delta) = \left\{ \mathcal{L}_t(\theta) - \mathcal{L}_t(\hat{\theta}_t) \leq \mathcal{O} (d \log(t/\delta)) \right\} .$$



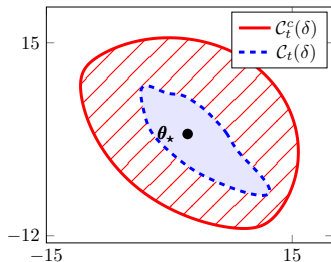


# Extension: convex relaxation

$$\mathcal{C}_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\}$$

- Non-convex, burdensome to manipulate.
- **Convex** relaxation based on log-loss  $\mathcal{L}_t$ :

$$\mathcal{C}_t^c(\delta) = \left\{ \mathcal{L}_t(\theta) - \mathcal{L}_t(\hat{\theta}_t) \leq \mathcal{O} \left( d \log(t/\delta) \right) \right\} .$$



Proposition (Abeille, F. and Calauzènes, 2021)

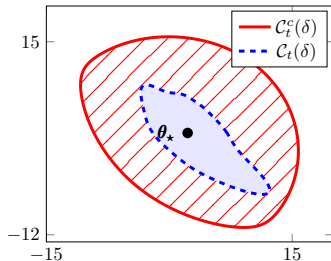
*The following holds:*

# Extension: convex relaxation

$$\mathcal{C}_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\}$$

- Non-convex, burdensome to manipulate.
- **Convex** relaxation based on log-loss  $\mathcal{L}_t$ :

$$\mathcal{C}_t^c(\delta) = \left\{ \mathcal{L}_t(\theta) - \mathcal{L}_t(\hat{\theta}_t) \leq \mathcal{O} (d \log(t/\delta)) \right\} .$$



## Proposition (Abeille, F. and Calauzènes, 2021)

*The following holds:*

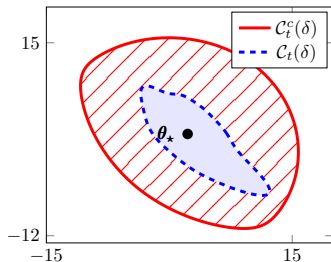
- For all  $t \geq 1$ ,  $\mathcal{C}_t(\delta) \subseteq \mathcal{C}_t^c(\delta)$  i.e  $\mathcal{C}_t^c(\delta)$  is a confidence set for  $\theta_*$ .

# Extension: convex relaxation

$$\mathcal{C}_t(\delta) = \left\{ \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\}$$

- Non-convex, burdensome to manipulate.
- **Convex** relaxation based on log-loss  $\mathcal{L}_t$ :

$$\mathcal{C}_t^c(\delta) = \left\{ \mathcal{L}_t(\theta) - \mathcal{L}_t(\hat{\theta}_t) \leq \mathcal{O} (d \log(t/\delta)) \right\} .$$



## Proposition (Abeille, F. and Calauzènes, 2021)

*The following holds:*

- For all  $t \geq 1$ ,  $\mathcal{C}_t(\delta) \subseteq \mathcal{C}_t^c(\delta)$  i.e  $\mathcal{C}_t^c(\delta)$  is a confidence set for  $\theta_*$ .
- With proba. at least  $1 - \delta$ :

$$\forall \theta \in \mathcal{C}_t^c(\delta), \forall t \geq 1, \left\| \theta - \theta_* \right\|_{\mathbf{H}_t(\theta_*)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right)$$

Algorithm and regret bounds

# OFU-GLB

- **Algorithm.** New ingredients, same recipe.

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

# OFU-GLB

- **Algorithm.** New ingredients, same recipe.

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

- Pseudo-code.

---

## Algorithm OFU-GLB

---

**input:** Arm set  $\mathcal{A}$ , regularizations  $\{\lambda_t\}_t$ , failure level  $\delta$ , norm upper-bound  $S$ .

Set  $\mathbf{H}_1 \leftarrow \lambda_1 \mathbf{I}_d$ ,  $\hat{\theta}_1 \leftarrow 0_d$ .

**for**  $t \in [1, T]$  **do**

Solve  $a_t \in \operatorname{argmax}_{\mathcal{A}} \max_{\theta \in \mathcal{C}_t(\delta)} a^\top \theta$ .

▷ *planning*

Play the arm  $a_t$  and observe reward  $r_{t+1}$ .

Update the estimator  $\hat{\theta}_{t+1}$  and the confidence interval  $\mathcal{C}_t(\delta)$ .

▷ *learning*

**end for**

---

# OFU-GLB: analysis

- Sketch of proof.

$$\text{Regret}(T) = \sum_{t=1}^T \mu(a_{\star}^{\top} \theta_{\star}) - \mu(a_t^{\top} \theta_{\star})$$

# OFU-GLB: analysis

- Sketch of proof.

$$\text{Regret}(T) \leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_*) \quad (\theta_t \text{ optim.})$$



# OFU-GLB: analysis

- Sketch of proof.

$$\text{Regret}(T) \leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_*) \quad (\theta_t \text{ optim.})$$

$$\leq \sum_{t=1}^T \dot{\mu}(a_t^\top \theta_*) a_t^\top (\theta_t - \theta_*) + \mathbf{L}_\mu \sum_{t=1}^T (a_t^\top (\theta_t - \theta_*))^2 \quad (\text{Taylor})$$

# OFU-GLB: analysis

- Sketch of proof.

$$\begin{aligned} \text{Regret}(T) &\leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_*) && (\theta_t \text{ optim.}) \\ &\leq \underbrace{\sum_{t=1}^T \dot{\mu}(a_t^\top \theta_*) a_t^\top (\theta_t - \theta_*)}_{\textcircled{1}} + \mathbf{L}_\mu \underbrace{\sum_{t=1}^T (a_t^\top (\theta_t - \theta_*))^2}_{\textcircled{2}} && (\text{Taylor}) \end{aligned}$$

# OFU-GLB: analysis

- Sketch of proof.

$$\begin{aligned} \text{Regret}(T) &\leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_*) && (\theta_t \text{ optim.}) \\ &\leq \underbrace{\sum_{t=1}^T \dot{\mu}(a_t^\top \theta_*) a_t^\top (\theta_t - \theta_*)}_{\textcircled{1}} + \kappa_\mu d^2 \log(T/\delta)^2 \end{aligned}$$

# OFU-GLB: analysis

- Sketch of proof.

$$\begin{aligned} \text{Regret}(T) &\leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_*) && (\theta_t \text{ optim.}) \\ &\leq \underbrace{\sum_{t=1}^T \dot{\mu}(a_t^\top \theta_*) a_t^\top (\theta_t - \theta_*)}_{\textcircled{1}} + \kappa_\mu d^2 \log(T/\delta)^2 \end{aligned}$$

- Bounding  $\textcircled{1}$ :

# OFU-GLB: analysis

- Sketch of proof.

$$\begin{aligned} \text{Regret}(T) &\leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_*) && (\theta_t \text{ optim.}) \\ &\leq \underbrace{\sum_{t=1}^T \dot{\mu}(a_t^\top \theta_*) a_t^\top (\theta_t - \theta_*)}_{\textcircled{1}} + \kappa_\mu d^2 \log(T/\delta)^2 \end{aligned}$$

- Bounding  $\textcircled{1}$ :

$$\textcircled{1} \leq d \log(T/\delta) \sqrt{\sum_{t=1}^T \dot{\mu}(a_t^\top \theta_*)} \quad (\text{Ell. Pot.})$$

# OFU-GLB: analysis

- Sketch of proof.

$$\begin{aligned} \text{Regret}(T) &\leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_*) && (\theta_t \text{ optim.}) \\ &\leq \underbrace{\sum_{t=1}^T \dot{\mu}(a_t^\top \theta_*) a_t^\top (\theta_t - \theta_*)}_{\textcircled{1}} + \kappa_\mu d^2 \log(T/\delta)^2 \end{aligned}$$

- Bounding  $\textcircled{1}$ :

$$\begin{aligned} \textcircled{1} &\leq d \log(T/\delta) \sqrt{\sum_{t=1}^T \dot{\mu}(a_t^\top \theta_*)} && (\text{Ell. Pot.}) \\ &\leq d \log(T/\delta) \sqrt{T \dot{\mu}(a_*^\top \theta_*) + \text{Regret}(T)} && (\text{s.c}) \end{aligned}$$

where  $a_* = \operatorname{argmax}_{\mathcal{A}} a^\top \theta_*$  (best action in hindsight)

# OFU-GLB: analysis

- Sketch of proof.

$$\begin{aligned} \text{Regret}(T) &\leq \sum_{t=1}^T \mu(a_t^\top \theta_t) - \mu(a_t^\top \theta_*) && (\theta_t \text{ optim.}) \\ &\leq \underbrace{\sum_{t=1}^T \dot{\mu}(a_t^\top \theta_*) a_t^\top (\theta_t - \theta_*)}_{\textcircled{1}} + \kappa_\mu d^2 \log(T/\delta)^2 \\ &\leq d \log(T/\delta) \sqrt{T \dot{\mu}(a_*^\top \theta_*) + \text{Regret}(T)} + \kappa_\mu d^2 \log(T/\delta)^2 \end{aligned}$$

- Bounding  $\textcircled{1}$ :

$$\begin{aligned} \textcircled{1} &\leq d \log(T/\delta) \sqrt{\sum_{t=1}^T \dot{\mu}(a_t^\top \theta_*)} && (\text{Ell. Pot.}) \\ &\leq d \log(T/\delta) \sqrt{T \dot{\mu}(a_*^\top \theta_*) + \text{Regret}(T)} && (\text{s.c}) \end{aligned}$$

where  $a_* = \operatorname{argmax}_{\mathcal{A}} a^\top \theta_*$  (best action in hindsight)

# OFU-GLB: new regret upper-bound

Theorem (Extends [F., Abeille, Calauzènes, Fercoq (2020)])

*For all self-concordant GLBs, OFU-GLB satisfies:*

$$\text{Regret}(T) = \tilde{O} \left( d \sqrt{\dot{\mu}(a_{\star}^{\top} \theta_{\star}) T} + \kappa_{\mu} d^2 \right),$$

*with probability at least  $1 - \delta$ .*



# OFU-GLB: new regret upper-bound

Theorem (Extends [F., Abeille, Calauzènes, Fercoq (2020)])

For all self-concordant GLBs, OFU-GLB satisfies:

$$\text{Regret}(T) = \tilde{O} \left( d \sqrt{\dot{\mu}(a_{\star}^{\top} \theta_{\star}) T} + \kappa_{\mu} d^2 \right),$$

with probability at least  $1 - \delta$ .

- Non-linearity deferred to **second-order term**.

$$\text{for } T \gg \kappa_{\mu}^2, \quad \text{Regret}(T) = \tilde{O} \left( d \sqrt{\dot{\mu}(a_{\star}^{\top} \theta_{\star}) T} \right).$$

# OFU-GLB: new regret upper-bound

Theorem (Extends [F., Abeille, Calauzènes, Fercoq (2020)])

For all self-concordant GLBs, OFU-GLB satisfies:

$$\text{Regret}(T) = \tilde{O} \left( d \sqrt{\dot{\mu}(a_{\star}^{\top} \theta_{\star}) T} + \kappa_{\mu} d^2 \right),$$

with probability at least  $1 - \delta$ .

- **Non-linearity** deferred to **second-order term**.

$$\text{for } T \gg \kappa_{\mu}^2, \quad \text{Regret}(T) = \tilde{O} \left( d \sqrt{\dot{\mu}(a_{\star}^{\top} \theta_{\star}) T} \right).$$

- **Exponential improvement** over previous work: e.g Logistic Bandit:

$$\text{(before)} \quad \text{Regret}(T) \lesssim \kappa_{\mu} d \sqrt{T},$$

$$\text{(now)} \quad \text{Regret}(T) \lesssim \exp(-\|\theta_{\star}\|/2) d \sqrt{T}.$$

# Transitory and permanent regimes

- Making sense of the regret bound:

$$\text{Regret}(T) = \tilde{O} \left( d \sqrt{\dot{\mu}(a_*^\top \theta_*) T} + \kappa_\mu d^2 \right) ,$$

# Transitory and permanent regimes

- Making sense of the regret bound:

$$\begin{aligned}\text{Regret}(T) &= \tilde{O} \left( d\sqrt{\dot{\mu}(a_{\star}^{\top} \theta_{\star})T} + \kappa_{\mu} d^2 \right) , \\ &= R_T^{\text{perm}} + R_T^{\text{trans}} .\end{aligned}$$

- ▶ each term associated to a different **regime** of algorithm behavior.

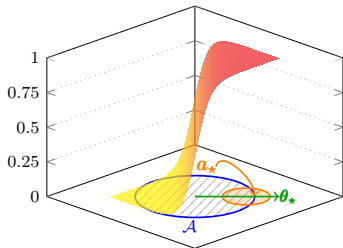
# Transitory and permanent regimes

- Making sense of the regret bound:

$$\begin{aligned}\text{Regret}(T) &= \tilde{O} \left( d\sqrt{\hat{\mu}(a_\star^\top \theta_\star)T} + \kappa_\mu d^2 \right), \\ &= R_T^{\text{perm}} + R_T^{\text{trans}}.\end{aligned}$$

- ▶ each term associated to a different **regime** of algorithm behavior.

- **Permanent** regret  $\longleftrightarrow a_\star \approx$  located.

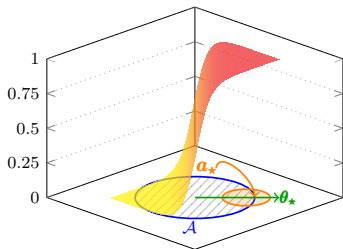


# Transitory and permanent regimes

- Making sense of the regret bound:

$$\begin{aligned}\text{Regret}(T) &= \tilde{O} \left( d \sqrt{\dot{\mu}(a_*^\top \theta_*) T} + \kappa_\mu d^2 \right), \\ &= R_T^{\text{perm}} + R_T^{\text{trans}}.\end{aligned}$$

- ▶ each term associated to a different **regime** of algorithm behavior.
- **Permanent** regret  $\longleftrightarrow a_* \approx$  located.
- ▶ locally *linear* with slope  $\dot{\mu}(a_*^\top \theta_*)$ .

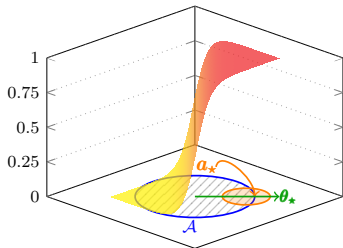


# Transitory and permanent regimes

- Making sense of the regret bound:

$$\begin{aligned}\text{Regret}(T) &= \tilde{O} \left( d \sqrt{\dot{\mu}(a_*^\top \theta_*) T} + \kappa_\mu d^2 \right), \\ &= R_T^{\text{perm}} + R_T^{\text{trans}}.\end{aligned}$$

- ▶ each term associated to a different **regime** of algorithm behavior.
- **Permanent** regret  $\longleftrightarrow a_* \approx$  located.
  - ▶ locally *linear* with slope  $\dot{\mu}(a_*^\top \theta_*)$ .
  - ▶ e.g **flat**  $\Rightarrow$  small regret.



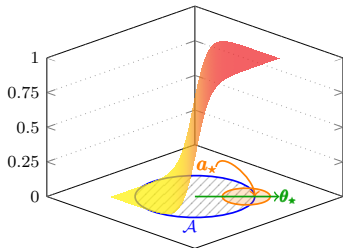
# Transitory and permanent regimes

- Making sense of the regret bound:

$$\begin{aligned}\text{Regret}(T) &= \tilde{O} \left( d\sqrt{\dot{\mu}(a_\star^\top \theta_\star)T} + \kappa_\mu d^2 \right), \\ &= R_T^{\text{perm}} + R_T^{\text{trans}}.\end{aligned}$$

- ▶ each term associated to a different **regime** of algorithm behavior.

- **Permanent** regret  $\longleftrightarrow a_\star \approx$  located.
  - ▶ locally *linear* with slope  $\dot{\mu}(a_\star^\top \theta_\star)$ .
  - ▶ e.g **flat**  $\Rightarrow$  small regret.



- **Transitory** regret: how long to find “good” regions of  $\mathcal{A}$ .
  - ▶ can be hard because of non-linearity since  $R_T^{\text{trans}} \propto \kappa_\mu$ .



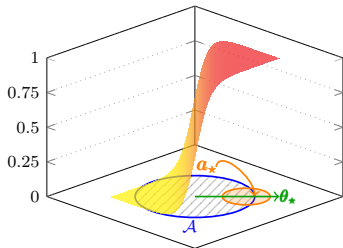
# Transitory and permanent regimes

- Making sense of the regret bound:

$$\begin{aligned}\text{Regret}(T) &= \tilde{O} \left( d\sqrt{\dot{\mu}(a_\star^\top \theta_\star)T} + \kappa_\mu d^2 \right), \\ &= R_T^{\text{perm}} + R_T^{\text{trans}}.\end{aligned}$$

- ▶ each term associated to a different **regime** of algorithm behavior.

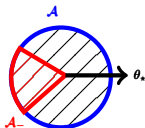
- **Permanent** regret  $\longleftrightarrow a_\star \approx$  located.
  - ▶ locally *linear* with slope  $\dot{\mu}(a_\star^\top \theta_\star)$ .
  - ▶ e.g **flat**  $\Rightarrow$  small regret.



- **Transitory** regret: how long to find “good” regions of  $\mathcal{A}$ .
  - ▶ can be hard because of non-linearity since  $R_T^{\text{trans}} \propto \kappa_\mu$ .
  - ▶ coherent with the Bayesian lower-bound of [Dong et al. 2019].

# Transitory regret on Logistic Bandit

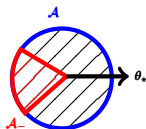
- Detrimental arms: large sub-optimality gap and small information.



# Transitory regret on Logistic Bandit

- Detrimental arms: large sub-optimality gap and small information.
- Transitory regret = how many times  $\mathcal{A}_-$  is played:

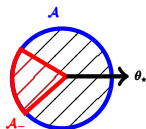
$$R_T^{\text{trans}} \leq \sum_{t=1}^T \mathbb{1} \{a_t \in \mathcal{A}_-\} .$$



# Transitory regret on Logistic Bandit

- Detrimental arms: large sub-optimality gap and small information.
- Transitory regret = how many times  $\mathcal{A}_-$  is played:

$$R_T^{\text{trans}} \leq \sum_{t=1}^T \mathbb{1} \{a_t \in \mathcal{A}_-\} .$$



Proposition (Abeille, F. and Calauzènes (2021))

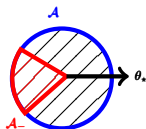
For Logistic Bandit if  $\mathcal{A} = \mathcal{B}_d$  the transitory regret satisfies:

$$R_T^{\text{trans}} = \tilde{O}(d^3)$$

# Transitory regret on Logistic Bandit

- Detrimental arms: large sub-optimality gap and small information.
- Transitory regret = how many times  $\mathcal{A}_-$  is played:

$$R_T^{\text{trans}} \leq \sum_{t=1}^T \mathbb{1} \{a_t \in \mathcal{A}_-\} .$$



## Proposition (Abeille, F. and Calauzènes (2021))

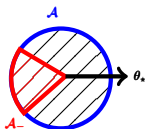
For Logistic Bandit if  $\mathcal{A} = \mathcal{B}_d$  the transitory regret satisfies:

$$R_T^{\text{trans}} = \tilde{O}(d^3) \leftarrow \text{independent of } \kappa_{\mu}!$$

# Transitory regret on Logistic Bandit

- Detrimental arms: large sub-optimality gap and small information.
- Transitory regret = how many times  $\mathcal{A}_-$  is played:

$$R_T^{\text{trans}} \leq \sum_{t=1}^T \mathbb{1} \{a_t \in \mathcal{A}_-\} .$$



## Proposition (Abeille, F. and Calauzènes (2021))

For Logistic Bandit if  $\mathcal{A} = \mathcal{B}_d$  the transitory regret satisfies:

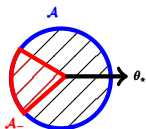
$$R_T^{\text{trans}} = \tilde{O}(d^3) \leftarrow \text{independent of } \kappa_\mu!$$

- ✔ "Good" case where non-linearity has **no effect** on the regret bound.

# Transitory regret on Logistic Bandit

- Detrimental arms: large sub-optimality gap and small information.
- Transitory regret = how many times  $\mathcal{A}_-$  is played:

$$R_T^{\text{trans}} \leq \sum_{t=1}^T \mathbb{1} \{a_t \in \mathcal{A}_-\} .$$



## Proposition (Abeille, F. and Calauzènes (2021))

For Logistic Bandit if  $\mathcal{A} = \mathcal{B}_d$  the transitory regret satisfies:

$$R_T^{\text{trans}} = \tilde{O}(d^3) \leftarrow \text{independent of } \kappa_\mu!$$

- ✓ "Good" case where non-linearity has **no effect** on the regret bound.
- ✓ For **any**  $T$ , we have  $\text{Regret}(T) \lesssim d\sqrt{\dot{\mu}(a_\star^\top \theta_\star)}T + d^3$ .

# Non-linearity: a blessing?

- The Logistic Bandit case on  $\mathcal{A} = \mathcal{B}_d \implies$  no more  $\kappa_\mu$ .

$$\text{Regret}(T) \lesssim d\sqrt{\dot{\mu}(a_\star^\top \theta_\star)T}$$



# Non-linearity: a blessing?

- The Logistic Bandit case on  $\mathcal{A} = \mathcal{B}_d \implies$  no more  $\kappa_\mu$ .

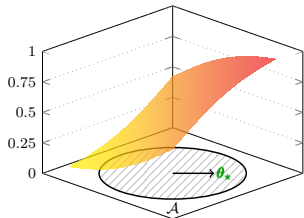
$$\text{Regret}(T) \lesssim d \exp(-\|\theta_\star\|/2) \sqrt{T}$$

# Non-linearity: a blessing?

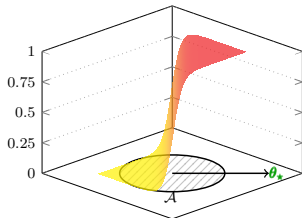
- The Logistic Bandit case on  $\mathcal{A} = \mathcal{B}_d \implies$  no more  $\kappa_\mu$ .

$$\text{Regret}(T) \lesssim d \exp(-\|\theta_\star\|/2) \sqrt{T}$$

$\|\theta_\star\| < 1$



$\|\theta_\star\| \gg 1$

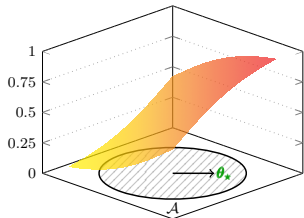


# Non-linearity: a blessing?

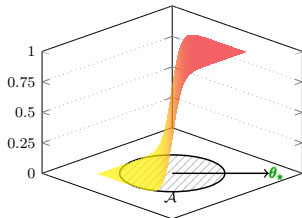
- The Logistic Bandit case on  $\mathcal{A} = \mathcal{B}_d \implies$  no more  $\kappa_\mu$ .

$$\text{Regret}(T) \lesssim d \exp(-\|\theta_*\|/2) \sqrt{T}$$

$\|\theta_*\| < 1$



$\|\theta_*\| \gg 1$



The more non-linear the problem, the smaller the regret.

Numerical simulations

# A tractable algorithm

- The planning step of OFU-GLB is intractable:

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

# A tractable algorithm

- The planning step of OFU-GLB is intractable:

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

- ▶ the constraint  $\theta \in \mathcal{C}_t(\delta)$  is **non-convex**.
- ▶ no principled way to solve (even approximately).

# A tractable algorithm

- The planning step of OFU-GLB is intractable:

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

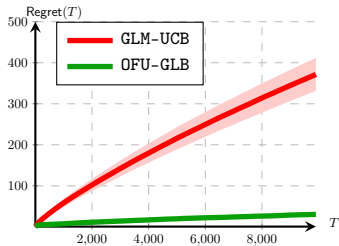
- ▶ the constraint  $\theta \in \mathcal{C}_t(\delta)$  is **non-convex**.
  - ▶ no principled way to solve (even approximately).
- 
- Use the convex relaxation  $\mathcal{C}_t^c(\delta)$ :

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

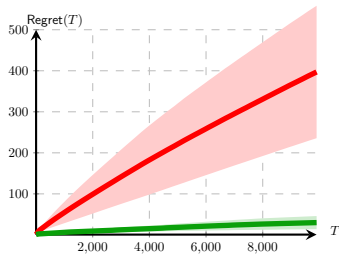
- ▶ tractable when  $|\mathcal{A}| < \infty$  (solve  $|\mathcal{A}|$  convex programs).
- ▶ same theoretical guarantees.

# Empirical performances

- Improved performances compared to GLM-UCB [Filippi et al. 2010]



(a)  $\kappa_\mu = 50$



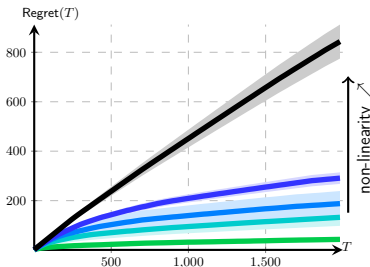
(b)  $\kappa_\mu = 400$

Comparing GLM-UCB and OFU-GLB on toy Logistic Bandit experiments.

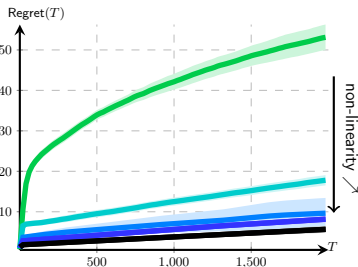


# Empirical performances (ctn'd)

- Check the impact of non-linearity:



(a) GLM-UCB



(b) OFU-GLB

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

# Logistic Bandit

## Regret Lower-Bound

# Optimality (1/3)

- Are these new regret upper-bounds **optimal**?
  - ▶ can we show that for any algorithms, there exist situations where:

$$\text{Regret}(T) \geq d\sqrt{\dot{\mu}(a_{\star}^{\top} \theta_{\star})T} + \kappa_{\mu} d^2 .$$

# Optimality (1/3)

- Are these new regret upper-bounds **optimal**?
  - ▶ can we show that for any algorithms, there exist situations where:

$$\text{Regret}(T) \geq d\sqrt{\dot{\mu}(a_{\star}^{\top} \theta_{\star})T} + \kappa_{\mu} d^2 .$$

- Why is it challenging?
  - ▶ involves **problem-dependent** constants.
  - ▶ describe a **continuum** of hard situations.
  - ▶ existing approaches from LB; typically use  $\|\theta_{\star}\| \propto 1/T$

# Optimality (1/3)

- Are these new regret upper-bounds **optimal**?
  - ▶ can we show that for any algorithms, there exist situations where:

$$\text{Regret}(T) \geq d\sqrt{\dot{\mu}(a_{\star}^{\top} \theta_{\star})T} + \kappa_{\mu} d^2 .$$

- Why is it challenging?
  - ▶ involves **problem-dependent** constants.
  - ▶ describe a **continuum** of hard situations.
  - ▶ existing approaches from LB; typically use  $\|\theta_{\star}\| \propto 1/T$
- Notion of **local** minimax-regret [Simchowitz and Foster. 2021]:

$$\text{MinimaxRegret}_{\theta_{\star}}(T, \varepsilon) := \min_{\pi} \max_{\|\theta - \theta_{\star}\| \leq \varepsilon} \text{Regret}_{\theta}^{\pi}(T)$$

for a given (arbitrary) reference  $\theta_{\star}$ .

## Optimality (2/3)

Theorem (Abeille, F., Calauzènes (2021))

For the Logistic Bandit with  $\mathcal{A} = \mathcal{S}_d$ :  $\forall \theta_*$ ,  $\exists \varepsilon$  small s.t:

$$\text{MinimaxRegret}_{\theta_*}(T, \varepsilon) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T} .$$

whenever  $T \geq d^2 \kappa(\theta_*)$ .

## Optimality (2/3)

Theorem (Abeille, F., Calauzènes (2021))

For the Logistic Bandit with  $\mathcal{A} = \mathcal{S}_d$ :  $\forall \theta_*$ ,  $\exists \varepsilon$  small s.t:

$$\text{MinimaxRegret}_{\theta_*}(T, \varepsilon) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T} .$$

whenever  $T \geq d^2 \kappa(\theta_*)$ .

- **Discussion.**  $\theta_*$  arbitrary reference point,  $\pi$  a given algorithm

## Optimality (2/3)

Theorem (Abeille, F., Calauzènes (2021))

For the Logistic Bandit with  $\mathcal{A} = \mathcal{S}_d$ :  $\forall \theta_*$ ,  $\exists \varepsilon$  small s.t:

$$\text{MinimaxRegret}_{\theta_*}(T, \varepsilon) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T} .$$

whenever  $T \geq d^2 \kappa(\theta_*)$ .

- **Discussion.**  $\theta_*$  arbitrary reference point,  $\pi$  a given algorithm
  - ▶  $\theta$  the “hardest” nearby instance in  $\{\|\theta' - \theta_*\| \leq \varepsilon\}$ .



# Optimality (2/3)

Theorem (Abeille, F., Calauzènes (2021))

For the Logistic Bandit with  $\mathcal{A} = \mathcal{S}_d$ :  $\forall \theta_*$ ,  $\exists \varepsilon$  small s.t:

$$\text{MinimaxRegret}_{\theta_*}(T, \varepsilon) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T}.$$

whenever  $T \geq d^2 \kappa(\theta_*)$ .

- **Discussion.**  $\theta_*$  arbitrary reference point,  $\pi$  a given algorithm
  - ▶  $\theta$  the “hardest” nearby instance in  $\{\|\theta' - \theta_*\| \leq \varepsilon\}$ .
  - ▶ the regret of  $\pi$  against  $\theta$  is:

$$\text{Regret}_{\theta}^{\pi}(T) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T}$$

# Optimality (2/3)

## Theorem (Abeille, F., Calauzènes (2021))

For the Logistic Bandit with  $\mathcal{A} = \mathcal{S}_d$ :  $\forall \theta_*$ ,  $\exists \varepsilon$  small s.t.:

$$\text{MinimaxRegret}_{\theta_*}(T, \varepsilon) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T}.$$

whenever  $T \geq d^2 \kappa(\theta_*)$ .

• **Discussion.**  $\theta_*$  arbitrary reference point,  $\pi$  a given algorithm

- ▶  $\theta$  the “hardest” nearby instance in  $\{\|\theta' - \theta_*\| \leq \varepsilon\}$ .
- ▶ the regret of  $\pi$  against  $\theta$  is:

$$\text{Regret}_{\theta}^{\pi}(T) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T}$$

- ▶ small  $\varepsilon \implies \theta$  and  $\theta_*$  share same problem-dependent constants:

$$\dot{\mu}(a_*^\top \theta_*) \approx \dot{\mu}(a_*(\theta)^\top \theta).$$

# Optimality (2/3)

## Theorem (Abeille, F., Calauzènes (2021))

For the Logistic Bandit with  $\mathcal{A} = \mathcal{S}_d$ :  $\forall \theta_*$ ,  $\exists \varepsilon$  small s.t.:

$$\text{MinimaxRegret}_{\theta_*}(T, \varepsilon) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T}.$$

whenever  $T \geq d^2 \kappa(\theta_*)$ .

• **Discussion.**  $\theta_*$  arbitrary reference point,  $\pi$  a given algorithm

- ▶  $\theta$  the “hardest” nearby instance in  $\{\|\theta' - \theta_*\| \leq \varepsilon\}$ .
- ▶ the regret of  $\pi$  against  $\theta$  is:

$$\text{Regret}_{\theta}^{\pi}(T) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T} \approx d \sqrt{\dot{\mu}(a_*(\theta)^\top \theta) T}$$

- ▶ small  $\varepsilon \implies \theta$  and  $\theta_*$  share same problem-dependent constants:

$$\dot{\mu}(a_*^\top \theta_*) \approx \dot{\mu}(a_*(\theta)^\top \theta).$$

# Optimality (2/3)

## Theorem (Abeille, F., Calauzènes (2021))

For the Logistic Bandit with  $\mathcal{A} = \mathcal{S}_d$ :  $\forall \theta_*$ ,  $\exists \varepsilon$  small s.t.:

$$\text{MinimaxRegret}_{\theta_*}(T, \varepsilon) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T}.$$

whenever  $T \geq d^2 \kappa(\theta_*)$ . The **permanent** regret is **minimax-optimal**

• **Discussion.**  $\theta_*$  arbitrary reference point,  $\pi$  a given algorithm

- ▶  $\theta$  the “hardest” nearby instance in  $\{\|\theta' - \theta_*\| \leq \varepsilon\}$ .
- ▶ the regret of  $\pi$  against  $\theta$  is:

$$\text{Regret}_{\theta}^{\pi}(T) \geq d \sqrt{\dot{\mu}(a_*^\top \theta_*) T} \approx d \sqrt{\dot{\mu}(a_*(\theta)^\top \theta) T}$$

- ▶ small  $\varepsilon \implies \theta$  and  $\theta_*$  share same problem-dependent constants:

$$\dot{\mu}(a_*^\top \theta_*) \approx \dot{\mu}(a_*(\theta)^\top \theta).$$

# Optimality (3/3)

- **Proof sketch.** To find a hard nearby instance  $\theta$  for  $\theta_*$ :

# Optimality (3/3)

- **Proof sketch.** To find a hard nearby instance  $\theta$  for  $\theta_*$ :
  - (1)  $\pi$  must behave similarly on  $\theta_*$  and  $\theta$ .

# Optimality (3/3)

- **Proof sketch.** To find a hard nearby instance  $\theta$  for  $\theta_*$ :
  - (1)  $\pi$  must behave similarly on  $\theta_*$  and  $\theta$ .
  - (2) the best arm differs:  $a_*(\theta) \neq a_*(\theta_*)$ .

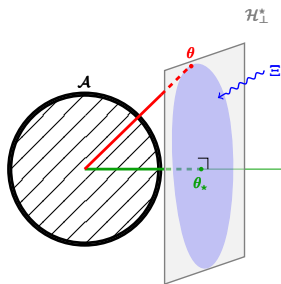
# Optimality (3/3)

• **Proof sketch.** To find a hard nearby instance  $\theta$  for  $\theta_*$ :

- (1)  $\pi$  must behave similarly on  $\theta_*$  and  $\theta$ .
- (2) the best arm differs:  $a_*(\theta) \neq a_*(\theta_*)$ .

► discrepancy measure: for  $\theta' \in \mathcal{H}_\perp^*$ :

$$d(\theta', \theta_*) := \sqrt{T\dot{\mu}(a_*^\top \theta_*)} \|\theta' - \theta_*\|^2$$

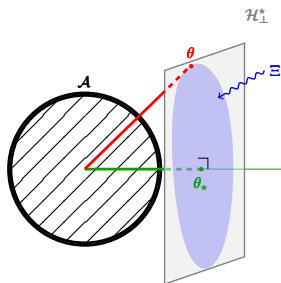




# Optimality (3/3)

• **Proof sketch.** To find a hard nearby instance  $\theta$  for  $\theta_*$ :

- (1)  $\pi$  must behave similarly on  $\theta_*$  and  $\theta$ .
- (2) the best arm differs:  $a_*(\theta) \neq a_*(\theta_*)$ .



► discrepancy measure: for  $\theta' \in \mathcal{H}_\perp^*$ :

$$d(\theta', \theta_*) := \sqrt{T\dot{\mu}(a_*^\top \theta_*)} \|\theta' - \theta_*\|^2$$

► compromise between (1) and (2):

$$\begin{aligned} \theta \in \Xi &:= \{\theta' \in \mathcal{H}_\perp^*, d(\theta', \theta_*) = 1\} . \\ \implies \|\theta - \theta_*\|^2 &= (T\dot{\mu}(a_*^\top \theta_*))^{-1/2} \end{aligned}$$

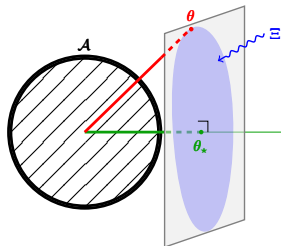
# Optimality (3/3)

- **Proof sketch.** To find a hard nearby instance  $\theta$  for  $\theta_*$ :

- (1)  $\pi$  must behave similarly on  $\theta_*$  and  $\theta$ .
- (2) the best arm differs:  $a_*(\theta) \neq a_*(\theta_*)$ .

- ▶ discrepancy measure: for  $\theta' \in \mathcal{H}_\perp^*$ :

$$d(\theta', \theta_*) := \sqrt{T \dot{\mu}(a_*^\top \theta_*)} \|\theta' - \theta_*\|^2$$



- ▶ compromise between (1) and (2):

$$\theta \in \Xi := \{\theta' \in \mathcal{H}_\perp^*, d(\theta', \theta_*) = 1\} .$$
$$\implies \|\theta - \theta_*\|^2 = (T \dot{\mu}(a_*^\top \theta_*))^{-1/2}$$

- ▶ what about regret?

$$\text{Regret}_\theta^\pi(T) \approx \dot{\mu}(a_*(\theta)^\top \theta) \sum_{t=1}^T \|a_t - a_*(\theta)\|^2$$

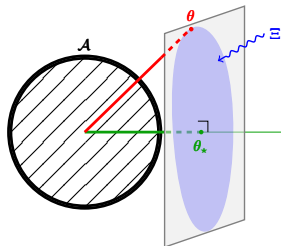
# Optimality (3/3)

- **Proof sketch.** To find a hard nearby instance  $\theta$  for  $\theta_*$ :

- (1)  $\pi$  must behave similarly on  $\theta_*$  and  $\theta$ .
- (2) the best arm differs:  $a_*(\theta) \neq a_*(\theta_*)$ .

- ▶ discrepancy measure: for  $\theta' \in \mathcal{H}_\perp^*$ :

$$d(\theta', \theta_*) := \sqrt{T \dot{\mu}(a_*^\top \theta_*)} \|\theta' - \theta_*\|^2$$



- ▶ compromise between (1) and (2):

$$\theta \in \Xi := \{\theta' \in \mathcal{H}_\perp^*, d(\theta', \theta_*) = 1\}.$$
$$\implies \|\theta - \theta_*\|^2 = (T \dot{\mu}(a_*^\top \theta_*))^{-1/2}$$

- ▶ what about regret?

$$\text{Regret}_\theta^\pi(T) \approx \dot{\mu}(a_*^\top \theta_*) \sum_{t=1}^T \|a_*(\theta_*) - a_*(\theta)\|^2$$

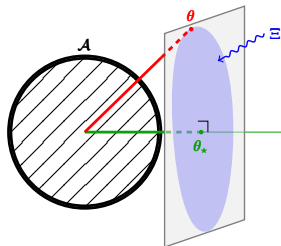
# Optimality (3/3)

- **Proof sketch.** To find a hard nearby instance  $\theta$  for  $\theta_*$ :

- (1)  $\pi$  must behave similarly on  $\theta_*$  and  $\theta$ .
- (2) the best arm differs:  $a_*(\theta) \neq a_*(\theta_*)$ .

- ▶ discrepancy measure: for  $\theta' \in \mathcal{H}_\perp^*$ :

$$d(\theta', \theta_*) := \sqrt{T\dot{\mu}(a_*^\top \theta_*)} \|\theta' - \theta_*\|^2$$



- ▶ compromise between (1) and (2):

$$\theta \in \Xi := \{\theta' \in \mathcal{H}_\perp^*, d(\theta', \theta_*) = 1\} .$$
$$\implies \|\theta - \theta_*\|^2 = (T\dot{\mu}(a_*^\top \theta_*))^{-1/2}$$

- ▶ what about regret?

$$\text{Regret}_\theta^\pi(T) \approx \dot{\mu}(a_*(\theta_*)^\top \theta_*) T \|\theta - \theta_*\|^2$$

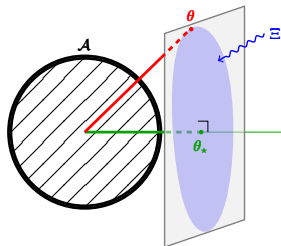
# Optimality (3/3)

- **Proof sketch.** To find a hard nearby instance  $\theta$  for  $\theta_*$ :

- (1)  $\pi$  must behave similarly on  $\theta_*$  and  $\theta$ .
- (2) the best arm differs:  $a_*(\theta) \neq a_*(\theta_*)$ .

- ▶ discrepancy measure: for  $\theta' \in \mathcal{H}_\perp^*$ :

$$d(\theta', \theta_*) := \sqrt{T\dot{\mu}(a_*^\top \theta_*)} \|\theta' - \theta_*\|^2$$



- ▶ compromise between (1) and (2):

$$\theta \in \Xi := \{\theta' \in \mathcal{H}_\perp^*, d(\theta', \theta_*) = 1\} .$$
$$\implies \|\theta - \theta_*\|^2 = (T\dot{\mu}(a_*^\top \theta_*))^{-1/2}$$

- ▶ what about regret?

$$\text{Regret}_\theta^\pi(T) \approx \sqrt{\dot{\mu}(a_*(\theta_*)^\top \theta_*)T}$$

# Extensions

# Contextual bandits

- Reward is also a function of exogenous context  $x_t \in \mathcal{X}$ :

$$\mathbb{E}[r_{t+1} | a_t] = \mu(\phi(a_t, x_t)^\top \theta_\star).$$

for some  $\phi : \mathcal{A} \times \mathcal{X} \mapsto \mathbb{R}^d$ .

# Contextual bandits

- Reward is also a function of exogenous context  $\mathbf{x}_t \in \mathcal{X}$ :

$$\mathbb{E}[r_{t+1} | a_t] = \mu(\phi(a_t, \mathbf{x}_t)^\top \theta_\star).$$

for some  $\phi : \mathcal{A} \times \mathcal{X} \mapsto \mathbb{R}^d$ .

- Similar regret upper-bounds:

$$\text{Regret}(T) = \tilde{O} \left( d\sqrt{T} \sqrt{\frac{1}{T} \sum_{t=1}^T \mu(\phi(a_{\star,t}, \mathbf{x}_t)^\top \theta_\star)} \right)$$

where  $a_{\star,t} = \operatorname{argmax}_{a \in \mathcal{A}} \phi(a, \mathbf{x}_t)$  best arm at round  $t$ .



# Contextual bandits

- Reward is also a function of exogenous context  $\mathbf{x}_t \in \mathcal{X}$ :

$$\mathbb{E}[r_{t+1} | a_t] = \mu(\phi(a_t, \mathbf{x}_t)^\top \theta_\star).$$

for some  $\phi : \mathcal{A} \times \mathcal{X} \mapsto \mathbb{R}^d$ .

- Similar regret upper-bounds:

$$\text{Regret}(T) = \tilde{\mathcal{O}} \left( d\sqrt{T} \sqrt{\frac{1}{T} \sum_{t=1}^T \mu(\phi(a_{\star,t}, \mathbf{x}_t)^\top \theta_\star)} \right)$$

where  $a_{\star,t} = \operatorname{argmax}_{a \in \mathcal{A}} \phi(a, \mathbf{x}_t)$  best arm at round  $t$ .

- Same goes for time-varying arm-sets.

# Non-stationary bandits

- Piece-wise stationary environment:

$$\mathbb{E}[r_{t+1} | a_t] = \mu(a_t^\top \theta_\star^t) \quad \text{where} \quad \sum_{t=2}^T \mathbb{1}(\theta_\star^t \neq \theta_\star^{t-1}) = \Gamma_T$$

# Non-stationary bandits

- Piece-wise stationary environment:

$$\mathbb{E}[r_{t+1} | a_t] = \mu(a_t^\top \theta_\star^t) \quad \text{where} \quad \sum_{t=2}^T \mathbb{1}(\theta_\star^t \neq \theta_\star^{t-1}) = \Gamma_T$$

- Change the estimation process to forget the past:

$$\hat{\theta}_t = \operatorname{argmin}_\theta - \sum_{s=1}^t \gamma^{t-s} \log d\mathbb{P}(r_{t+1} | a_t) / d\nu(r) + \lambda \|\theta\|^2 / 2.$$

# Non-stationary bandits

- Piece-wise stationary environment:

$$\mathbb{E}[r_{t+1} | a_t] = \mu(a_t^\top \theta_\star^t) \quad \text{where} \quad \sum_{t=2}^T \mathbb{1}(\theta_\star^t \neq \theta_\star^{t-1}) = \Gamma_T$$

- Change the estimation process to forget the past:

$$\hat{\theta}_t = \operatorname{argmin}_\theta - \sum_{s=1}^t \gamma^{t-s} \log d\mathbb{P}(r_{t+1} | a_t) / d\nu(r) + \lambda \|\theta\|^2 / 2.$$

- Similar conclusion:

Theorem (improves (Russac, F., Cappé and Garivier, 2021))

*There exists an algorithm on the piece-wise stationary GLB problem s.t.:*

$$\text{DynamicRegret}(T) = \tilde{\mathcal{O}} \left( T^{2/3} \Gamma_T^{1/3} \sqrt{\ell_\mu^\star} + \kappa_\mu T^{1/3} \Gamma_T^{2/3} \right).$$

where  $\ell_\mu^\star := \frac{1}{T} \sum_{t=1}^T \dot{\mu}(a_{\star,t}^\top \theta_\star^t)$  is averaged sensitivity at best-arm.

# Towards computationally efficient algorithms

- **Computationally** hungry algorithms;

$$\text{total computational cost} = \tilde{O}(|\mathcal{A}|T^2)$$

# Towards computationally efficient algorithms

- Computationally hungry algorithms;

$$\text{total computational cost} = \tilde{O}(|\mathcal{A}|T^2)$$

- Two computational bottlenecks; at each round:
  - ▶ (*learning*) compute maximum likelihood (up to precision  $\varepsilon = 1/T$ ).
  - ▶ (*planning*) solve  $|\mathcal{A}|$  likelihood-based convex programs.

# Towards computationally efficient algorithms

- Computationally hungry algorithms;

$$\text{total computational cost} = \tilde{O}(|\mathcal{A}|T^2)$$

- Two computational bottlenecks; at each round:
  - ▶ (*learning*) compute maximum likelihood (up to precision  $\varepsilon = 1/T$ ).
  - ▶ (*planning*) solve  $|\mathcal{A}|$  likelihood-based convex programs.
- Simultaneously computationally and statistically efficient GLB algorithms?
  - ▶ (*learning*) confidence sets with  $\tilde{O}(1)$  sufficient statistic
  - ▶ (*planning*) Thompson Sampling alternative.

# Towards computationally efficient algorithms

- Computationally hungry algorithms;

$$\text{total computational cost} = \tilde{O}(|\mathcal{A}|T^2)$$

- Two computational bottlenecks; at each round:
  - ▶ (*learning*) compute maximum likelihood (up to precision  $\varepsilon = 1/T$ ).
  - ▶ (*planning*) solve  $|\mathcal{A}|$  likelihood-based convex programs.
- Simultaneously computationally and statistically efficient GLB algorithms?
  - ▶ (*learning*) confidence sets with  $\tilde{O}(1)$  sufficient statistic
  - ▶ (*planning*) Thompson Sampling alternative.
- ✔ tools from online convex optimization literature ([Jézéquel et al. 2020]).



# Towards computationally efficient algorithms

- Computationally hungry algorithms;

$$\text{total computational cost} = \tilde{O}(|\mathcal{A}|T^2)$$

- Two computational bottlenecks; at each round:
  - ▶ (*learning*) compute maximum likelihood (up to precision  $\varepsilon = 1/T$ ).
  - ▶ (*planning*) solve  $|\mathcal{A}|$  likelihood-based convex programs.
- Simultaneously computationally and statistically efficient GLB algorithms?
  - ▶ (*learning*) confidence sets with  $\tilde{O}(1)$  sufficient statistic
  - ▶ (*planning*) Thompson Sampling alternative.

✓ tools from online convex optimization literature ([Jézéquel et al. 2020]).

⇒ same regret guarantees and computationally efficient algorithm.

# Conclusion

## Key Take-Aways.

- **Generalized Linear Bandits:**
  - ▶ Flexible yet simple model for many real-world situations.

# Conclusion

## Key Take-Aways.

- **Generalized Linear Bandits:**
  - ▶ Flexible yet simple model for many real-world situations.
  - ▶ Neat study of non-linearity in parametric bandits.

# Conclusion

## Key Take-Aways.

- **Generalized Linear Bandits:**
  - ▶ Flexible yet simple model for many real-world situations.
  - ▶ Neat study of non-linearity in parametric bandits.
- **Contributions:**
  - ▶ Improved algorithms (much smaller regret).

# Conclusion

## Key Take-Aways.

- **Generalized Linear Bandits:**
  - ▶ Flexible yet simple model for many real-world situations.
  - ▶ Neat study of non-linearity in parametric bandits.
  
- **Contributions:**
  - ▶ Improved algorithms (much smaller regret).
  - ▶ Refined analysis tool for **local** treatment.

# Conclusion

## Key Take-Aways.

- **Generalized Linear Bandits:**
  - ▶ Flexible yet simple model for many real-world situations.
  - ▶ Neat study of non-linearity in parametric bandits.
  
- **Contributions:**
  - ▶ Improved algorithms (much smaller regret).
  - ▶ Refined analysis tool for **local** treatment.
  - ▶ Not harder to solve than Linear Bandit!

# Conclusion

## Key Take-Aways.

- **Generalized Linear Bandits:**
  - ▶ Flexible yet simple model for many real-world situations.
  - ▶ Neat study of non-linearity in parametric bandits.
  
- **Contributions:**
  - ▶ Improved algorithms (much smaller regret).
  - ▶ Refined analysis tool for **local** treatment.
  - ▶ Not harder to solve than Linear Bandit!
  
- **Limitations and Perspectives.**
  - ▶ Towards richer reward models?
  - ▶ Adversarial bandits.

Thank you!



# Bibliography

- Sarah Filippi, Olivier Cappé, Aurélien Garivier, Csaba Szepesvári. *Parametric Bandits: The Generalized Linear Case*, NeurIPS, 2010.
- Francis Bach. *Self-Concordant Analysis for Logistic Regression*, EJS, 2010.
- Yasin Abbasi-Yadkori, Csaba Szepesvári, David Pal. Improved Algorithms for Linear Stochastic Bandits, NeurIPS, 2011.
- Lihong Li, Yu Lu, Dengyong Zhou. Provably Optimal Algorithms for Generalized Linear Contextual Bandits, ICML, 2017.
- Kwang-Sung Jun, Aniruddha Bhargava, Robert Nowak, Rebecca Willett. Scalable Generalized Linear Bandits: Online Computation and Hashing, NeurIPS, 2017.
- Shi Dong, Tengyu Ma, Benjamin Van Roy. *On the Performance of Thompson Sampling on Logistic Bandits*, COLT, 2019.

# Relevant Publications

- L.F, Marc Abeille, Clément Calauzènes, Olivier Fercoq. *Improved Optimistic Algorithms for Logistic Bandits*, ICML, 2020.
- Marc Abeille, L.F, Clément Calauzènes. *Instance-Wise Minimax-Optimal Algorithms for Logistic Bandits*, AISTATS, 2021.
- Yoan Russac, L.F, Olivier Cappé, Aurélien Garivier. *Self-Concordant Analysis of Generalized Linear Bandits with Forgetting*, AISTATS, 2021.
- L.F, Yoan Russac, Marc Abeille, Clément Calauzènes. *A Technical Note on Non-Stationary Parametric Bandits: Existing Mistakes and Preliminary Solutions*, ALT, 2021.
- L.F, Yoan Russac, Marc Abeille, Clément Calauzènes. *Regret Bounds for Generalized Linear Bandits under Parameter Drift*, preprint, 2021.

# What about self-concordance?

- Mostly used for the **learning** process. Actually, concentration is given by:

$$\left\| \sum_{s=1}^{t-1} \left[ \mu(a_s^\top \hat{\theta}_t) - \mu(a_s^\top \theta_\star) \right] a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)} = \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)} = \mathcal{O} \sqrt{d \log(t/\delta)}.$$

# What about self-concordance?

- Mostly used for the **learning** process. Actually, concentration is given by:

$$\left\| \sum_{s=1}^{t-1} \left[ \mu(a_s^\top \hat{\theta}_t) - \mu(a_s^\top \theta_*) \right] a_s \right\|_{\mathbf{H}_t^{-1}(\theta_*)} = \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_*)} = \mathcal{O} \sqrt{d \log(t/\delta)}.$$

- ▶ To be used for bounding regret, need to be tied to  $\left\| \hat{\theta}_t - \theta_* \right\|_{\mathbf{H}_t^{-1}(\theta_*)}$

# What about self-concordance?

- Mostly used for the **learning** process. Actually, concentration is given by:

$$\left\| \sum_{s=1}^{t-1} \left[ \mu(a_s^\top \hat{\theta}_t) - \mu(a_s^\top \theta_\star) \right] a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)} = \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)} = \mathcal{O} \sqrt{d \log(t/\delta)}.$$

- ▶ To be used for bounding regret, need to be tied to  $\left\| \hat{\theta}_t - \theta_\star \right\|_{\mathbf{H}_t^{-1}(\theta_\star)}$
- By the mean-value theorem:

$$\sum_{s=1}^{t-1} \left[ \mu(a_s^\top \hat{\theta}_t) - \mu(a_s^\top \theta_\star) \right] a_s = \mathbf{G}_t(\hat{\theta}_t, \theta_\star) (\hat{\theta}_t - \theta_\star)$$

where  $\mathbf{G}_t(\hat{\theta}_t, \theta_\star) = \sum_{s=1}^t \left[ \int_{v=0}^1 \dot{\mu}(a_s^\top \theta_\star + v a_s^\top (\hat{\theta}_t - \theta_\star)) dv \right] a_s a_s^\top + \lambda \mathbf{I}_d$ .

# What about self-concordance?

- Mostly used for the **learning** process. Actually, concentration is given by:

$$\left\| \sum_{s=1}^{t-1} \left[ \mu(a_s^\top \hat{\theta}_t) - \mu(a_s^\top \theta_\star) \right] a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)} = \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)} = \mathcal{O} \sqrt{d \log(t/\delta)}.$$

- ▶ To be used for bounding regret, need to be tied to  $\left\| \hat{\theta}_t - \theta_\star \right\|_{\mathbf{H}_t^{-1}(\theta_\star)}$

- By the mean-value theorem:

$$\sum_{s=1}^{t-1} \left[ \mu(a_s^\top \hat{\theta}_t) - \mu(a_s^\top \theta_\star) \right] a_s = \mathbf{G}_t(\hat{\theta}_t, \theta_\star) (\hat{\theta}_t - \theta_\star)$$

where  $\mathbf{G}_t(\hat{\theta}_t, \theta_\star) = \sum_{s=1}^t \left[ \int_{v=0}^1 \dot{\mu}(a_s^\top \theta_\star + v a_s^\top (\hat{\theta}_t - \theta_\star)) dv \right] a_s a_s^\top + \lambda \mathbf{I}_d$ .

- Self-concordance to the rescue:

$$\int_{v=0}^1 \dot{\mu}(a_s^\top \theta_\star + v a_s^\top (\hat{\theta}_t - \theta_\star)) dv \geq (1 + 2S)^{-1} \dot{\mu}(a_s^\top \theta_\star)$$

so  $\mathbf{G}_t(\hat{\theta}_t, \theta_\star) \geq (1 + 2S)^{-1} \mathbf{H}_t(\theta_\star)$ .

# Bonus vs parameter based optimism

- For the Linear Bandit:

▶ Bonus: play  $a_t = \operatorname{argmax}_{a \in \mathcal{A}} a^\top \hat{\theta}_t + \sqrt{d \log(t)} \|a\|_{V_t^{-1}}$ .

▶ Parameter: play  $a_t = \operatorname{argmax}_{a \in \mathcal{A}} \max_{\theta \in \mathcal{C}_t(\delta)} a^\top \theta$

are exactly **equivalent**.

# Bonus vs parameter based optimism

- For the Linear Bandit:
  - ▶ Bonus: play  $a_t = \operatorname{argmax}_{a \in \mathcal{A}} a^\top \hat{\theta}_t + \sqrt{d \log(t)} \|a\|_{V_t^{-1}}$ .
  - ▶ Parameter: play  $a_t = \operatorname{argmax}_{a \in \mathcal{A}} \max_{\theta \in \mathcal{C}_t(\delta)} a^\top \theta$   
are exactly **equivalent**.
- No longer true with **non-ellipsoidal** confidence sets.



# Bonus vs parameter based optimism

- For the Linear Bandit:
  - ▶ Bonus: play  $a_t = \operatorname{argmax}_{a \in \mathcal{A}} a^\top \hat{\theta}_t + \sqrt{d \log(t)} \|a\|_{V_t^{-1}}$ .
  - ▶ Parameter: play  $a_t = \operatorname{argmax}_{a \in \mathcal{A}} \max_{\theta \in \mathcal{C}_t(\delta)} a^\top \theta$   
are exactly **equivalent**.
- No longer true with **non-ellipsoidal** confidence sets.
- Bonus-based exploration:
  - ▶ much more complicated bonus function.
  - ▶ requires additional projection.
  - ▶ non-tight analysis  $\Rightarrow$  non-tight design.
  - ▶ typically much less performant for GLBs.

# Bonus vs parameter based optimism

- For the Linear Bandit:
  - ▶ Bonus: play  $a_t = \operatorname{argmax}_{a \in \mathcal{A}} a^\top \hat{\theta}_t + \sqrt{d \log(t)} \|a\|_{V_t^{-1}}$ .
  - ▶ Parameter: play  $a_t = \operatorname{argmax}_{a \in \mathcal{A}} \max_{\theta \in \mathcal{C}_t(\delta)} a^\top \theta$   
are exactly **equivalent**.
- No longer true with **non-ellipsoidal** confidence sets.
- Bonus-based exploration:
  - ▶ much more complicated bonus function.
  - ▶ requires additional projection.
  - ▶ non-tight analysis  $\Rightarrow$  non-tight design.
  - ▶ typically much less performant for GLBs.
- Parameter-based:
  - ▶ non-tight analysis remains at analysis time
  - ▶ more adaptive algorithms (e.g second-order term)

# Fast yet optimal algorithms (1/2)

- Replace  $\hat{\theta}_t$  by:

$$\theta_t = \operatorname{argmin}_{\Theta} \|\theta - \theta_{t-1}\|_{\mathbf{W}_{t-1}}^2 + \ell(a_t^\top \theta, r_{t+1}).$$

where:

- ▶  $\Theta$  is a “small” convex set around  $\theta_*$  (forced-exploration)
- ▶ and  $\mathbf{W}_t = \sum_{s=1}^{t-1} \dot{\mu}(a_s^\top \theta_{s+1}) a_s a_s$ .

# Fast yet optimal algorithms (1/2)

- Replace  $\hat{\theta}_t$  by:

$$\theta_t = \operatorname{argmin}_{\Theta} \|\theta - \theta_{t-1}\|_{\mathbf{W}_{t-1}}^2 + \ell(a_t^\top \theta, r_{t+1}) .$$

where:

- ▶  $\Theta$  is a “small” convex set around  $\theta_*$  (forced-exploration)
  - ▶ and  $\mathbf{W}_t = \sum_{s=1}^{t-1} \dot{\mu}(a_s^\top \theta_{s+1}) a_s a_s$ .
- Yields the confidence set:

$$\left\{ \theta, \|\theta - \theta_t\|_{\mathbf{W}_t} \leq \sqrt{d \log(t/\delta)} \right\} .$$

- ▶  $\mathbf{W}_t$  mimics  $\mathbf{H}_t(\theta) \implies$  reward sensitivity.
- ▶ Sufficient statistics can be maintained at  $\tilde{O}(1)$  cost.

# Fast yet optimal algorithms (1/2)

- Replace  $\hat{\theta}_t$  by:

$$\theta_t = \operatorname{argmin}_{\Theta} \|\theta - \theta_{t-1}\|_{\mathbf{W}_{t-1}}^2 + \ell(a_t^\top \theta, r_{t+1}) .$$

where:

- ▶  $\Theta$  is a “small” convex set around  $\theta_*$  (forced-exploration)
  - ▶ and  $\mathbf{W}_t = \sum_{s=1}^{t-1} \dot{\mu}(a_s^\top \theta_{s+1}) a_s a_s$ .
- Yields the confidence set:

$$\left\{ \theta, \|\theta - \theta_t\|_{\mathbf{W}_t} \leq \sqrt{d \log(t/\delta)} \right\} .$$

- ▶  $\mathbf{W}_t$  mimics  $\mathbf{H}_t(\theta) \implies$  reward sensitivity.
  - ▶ Sufficient statistics can be maintained at  $\tilde{O}(1)$  cost.
- Same regret bounds!

# Fast yet optimal algorithms (1/2)

- Replace  $\hat{\theta}_t$  by:

$$\theta_t = \operatorname{argmin}_{\Theta} \|\theta - \theta_{t-1}\|_{\mathbf{W}_{t-1}}^2 + \ell(a_t^\top \theta, r_{t+1}) .$$

where:

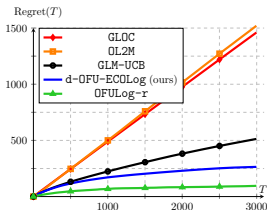
- ▶  $\Theta$  is a “small” convex set around  $\theta_*$  (forced-exploration)
  - ▶ and  $\mathbf{W}_t = \sum_{s=1}^{t-1} \dot{\mu}(a_s^\top \theta_{s+1}) a_s a_s$ .
- Yields the confidence set:

$$\left\{ \theta, \|\theta - \theta_t\|_{\mathbf{W}_t} \leq \sqrt{d \log(t/\delta)} \right\} .$$

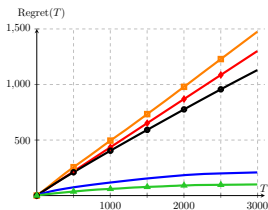
- ▶  $\mathbf{W}_t$  mimics  $\mathbf{H}_t(\theta) \implies$  reward sensitivity.
  - ▶ Sufficient statistics can be maintained at  $\tilde{\mathcal{O}}(1)$  cost.
- Same regret bounds!
  - Forced exploration can be dropped through at data-dependent approach.

# Fast yet optimal algorithms (2/2)

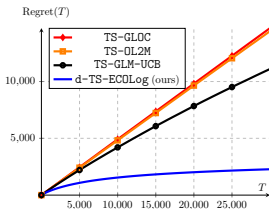
- Some experimental results:



$d = 2$ ,  $|\mathcal{A}| = 20$ ,  $\kappa = 150$



$d = 2$ ,  $|\mathcal{A}| = 20$ ,  $\kappa = 400$



$d = 5$ ,  $\mathcal{A} = \mathcal{B}_5$ ,  $\kappa = 400$ .

# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_\star - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_\star)} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)}$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^\top \theta_\star) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^\top \theta) a_s a_s^\top + \lambda_t \mathbf{I}_d$$



# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_{\star} - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_{\star})} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_{\star})}$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^{\top} \theta_{\star}) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^{\top} \theta) a_s a_s^{\top} + \lambda_t \mathbf{I}_d$$

- Define  $\mathcal{F}_s = \sigma(a_1, r_2, \dots, r_s, a_s)$ ; Exponential family distribution:

# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_\star - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_\star)} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_\star)}$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^\top \theta_\star) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^\top \theta) a_s a_s^\top + \lambda_t \mathbf{I}_d$$

- Define  $\mathcal{F}_s = \sigma(a_1, r_2, \dots, r_s, a_s)$ ; Exponential family distribution:
  - ▶  $\mathbb{E}[\eta_{s+1} | \mathcal{F}_s] = \mathbb{E}[r_{s+1} - \mu(a_s^\top \theta_\star)]$

# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_{\star} - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_{\star})} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_{\star})}$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^{\top} \theta_{\star}) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^{\top} \theta) a_s a_s^{\top} + \lambda_t \mathbf{I}_d$$

- Define  $\mathcal{F}_s = \sigma(a_1, r_2, \dots, r_s, a_s)$ ; Exponential family distribution:
  - ▶  $\mathbb{E}[\eta_{s+1} | \mathcal{F}_s] = 0$ .

# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_{\star} - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_{\star})} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_{\star})}$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^{\top} \theta_{\star}) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^{\top} \theta) a_s a_s^{\top} + \lambda_t \mathbf{I}_d$$

- Define  $\mathcal{F}_s = \sigma(a_1, r_2, \dots, r_s, a_s)$ ; Exponential family distribution:
  - ▶  $\mathbb{E}[\eta_{s+1} | \mathcal{F}_s] = 0$ .
  - ▶  $\text{Var}(\eta_{s+1} | \mathcal{F}_s) = \text{Var}(r_{s+1} | \mathcal{F}_s)$

# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_{\star} - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_{\star})} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_{\star})}$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^{\top} \theta_{\star}) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^{\top} \theta) a_s a_s^{\top} + \lambda_t \mathbf{I}_d$$

- Define  $\mathcal{F}_s = \sigma(a_1, r_2, \dots, r_s, a_s)$ ; Exponential family distribution:
  - ▶  $\mathbb{E}[\eta_{s+1} | \mathcal{F}_s] = 0$ .
  - ▶  $\text{Var}(\eta_{s+1} | \mathcal{F}_s) = \dot{\mu}(a_s^{\top} \theta_{\star})$

# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_{\star} - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_{\star})} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_{\star})}$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^{\top} \theta_{\star}) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^{\top} \theta) a_s a_s^{\top} + \lambda_t \mathbf{I}_d$$

- Define  $\mathcal{F}_s = \sigma(a_1, r_2, \dots, r_s, a_s)$ ; Exponential family distribution:
  - ▶  $\mathbb{E}[\eta_{s+1} | \mathcal{F}_s] = 0$ .
  - ▶  $\text{Var}(\eta_{s+1} | \mathcal{F}_s) = \dot{\mu}(a_s^{\top} \theta_{\star})$  so  $\mathbf{H}_t(\theta_{\star}) = \sum_{s=1}^t \text{Var}(\eta_{s+1} | \mathcal{F}_s) a_s a_s^{\top} + \lambda \mathbf{I}_d$ .

# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_{\star} - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_{\star})} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_{\star})}$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^{\top} \theta_{\star}) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^{\top} \theta) a_s a_s^{\top} + \lambda_t \mathbf{I}_d$$

- Define  $\mathcal{F}_s = \sigma(a_1, r_2, \dots, r_s, a_s)$ ; Exponential family distribution:
  - ▶  $\mathbb{E}[\eta_{s+1} | \mathcal{F}_s] = 0$ .
  - ▶  $\text{Var}(\eta_{s+1} | \mathcal{F}_s) = \dot{\mu}(a_s^{\top} \theta_{\star})$  so  $\mathbf{H}_t(\theta_{\star}) = \sum_{s=1}^t \text{Var}(\eta_{s+1} | \mathcal{F}_s) a_s a_s^{\top} + \lambda \mathbf{I}_d$ .

# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_{\star} - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_{\star})} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_{\star})} = \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right)$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^{\top} \theta_{\star}) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^{\top} \theta) a_s a_s^{\top} + \lambda_t \mathbf{I}_d$$

- Define  $\mathcal{F}_s = \sigma(a_1, r_2, \dots, r_s, a_s)$ ; Exponential family distribution:
  - ▶  $\mathbb{E}[\eta_{s+1} | \mathcal{F}_s] = 0$ .
  - ▶  $\text{Var}(\eta_{s+1} | \mathcal{F}_s) = \dot{\mu}(a_s^{\top} \theta_{\star})$  so  $\mathbf{H}_t(\theta_{\star}) = \sum_{s=1}^t \text{Var}(\eta_{s+1} | \mathcal{F}_s) a_s a_s^{\top} + \lambda \mathbf{I}_d$ .



# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_{\star} - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_{\star})} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_{\star})} = \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right)$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^{\top} \theta_{\star}) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^{\top} \theta) a_s a_s^{\top} + \lambda_t \mathbf{I}_d$$

# Application to GLBs

- Using optimality of  $\hat{\theta}_t$  for the regularized log-loss  $\mathcal{L}_t^{\lambda_t}$ :

$$\forall t \geq 1, \quad \left\| \theta_{\star} - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta_{\star})} \leq \left\| \sum_{s=1}^{t-1} \eta_{s+1} a_s \right\|_{\mathbf{H}_t^{-1}(\theta_{\star})} = \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right)$$

where:

$$\eta_{s+1} = r_{s+1} - \mu(a_s^{\top} \theta_{\star}) \quad \text{and} \quad \mathbf{H}_t(\theta) = \sum_{s=1}^t \dot{\mu}(a_s^{\top} \theta) a_s a_s^{\top} + \lambda_t \mathbf{I}_d$$

- New confidence set:

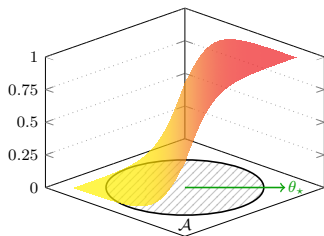
Proposition (F., Abeille, Calauzènes and Fercoq, 2020.)

For  $\delta \in (0, 1]$  let:

$$\mathcal{C}_t(\delta) := \left\{ \|\theta\| \leq S, \left\| \theta - \hat{\theta}_t \right\|_{\mathbf{H}_t(\theta)} \leq \mathcal{O} \left( \sqrt{d \log(t/\delta)} \right) \right\} .$$

Then  $\mathbb{P}(\forall t \geq 1, \theta_{\star} \in \mathcal{C}_t(\delta)) \geq 1 - \delta$ .

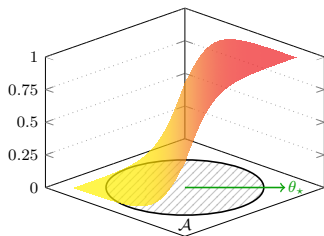
# GLBs: information vs. regret



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

- Varying reward sensitivity:

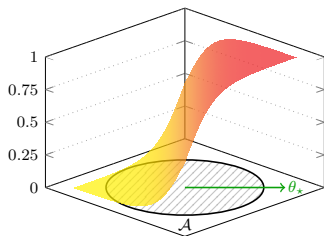
# GLBs: information vs. regret



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

- Varying reward sensitivity:
  - ▶ low-sensitivity:
    - information is hard to get
    - small regret

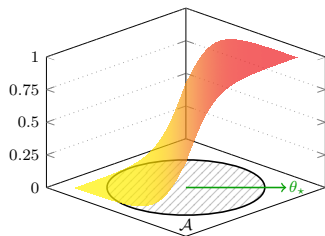
# GLBs: information vs. regret



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

- Varying reward sensitivity:
  - ▶ low-sensitivity:
    - information is hard to get
    - small regret
  - ▶ high-sensitivity:
    - information is easy to get
    - large regret

# GLBs: information vs. regret



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

- Varying reward sensitivity:

- ▶ low-sensitivity:

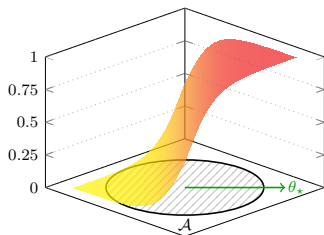
- information is hard to get
- small regret

- ▶ high-sensitivity:

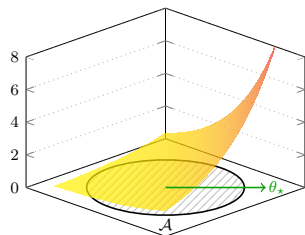
- information is easy to get
- large regret

⇒ linearization: worst of both world.

# GLBs: information vs. regret



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



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

- Varying reward sensitivity:

- ▶ low-sensitivity:

- information is hard to get
- small regret

- ▶ high-sensitivity:

- information is easy to get
- large regret

⇒ linearization: worst of both world.

# Minimax rates in general non-stationary settings

- Beyond piece-wise stationarity thanks to *variation-budget*:

$$B_T := \sum_{t=2}^T \|\theta_\star^t - \theta_\star^{t-1}\| .$$

- ▶ describe broader non-stationary environments.



# Minimax rates in general non-stationary settings

- Beyond piece-wise stationarity thanks to *variation-budget*:

$$B_T := \sum_{t=2}^T \|\theta_\star^t - \theta_\star^{t-1}\| .$$

- ▶ describe broader non-stationary environments.
- 
- Forgetting mechanisms to the rescue?
    - ▶ minimax-optimal for the MAB setting

# Minimax rates in general non-stationary settings

- Beyond piece-wise stationarity thanks to *variation-budget*:

$$B_T := \sum_{t=2}^T \|\theta_\star^t - \theta_\star^{t-1}\| .$$

- ▶ describe broader non-stationary environments.
- 
- Forgetting mechanisms to the rescue?
    - ▶ minimax-optimal for the MAB setting
    - ▶ not so well understood in LB! [F. et al, 2021a]

# Minimax rates in general non-stationary settings

- Beyond piece-wise stationarity thanks to *variation-budget*:

$$B_T := \sum_{t=2}^T \|\theta_\star^t - \theta_\star^{t-1}\| .$$

- ▶ describe broader non-stationary environments.
- 
- Forgetting mechanisms to the rescue?
    - ▶ minimax-optimal for the MAB setting
    - ▶ not so well understood in LB! [F. et al, 2021a]
    - ▶ best know regret bound for GLBs [F. et al. 2021b]:

$$\text{DynamicRegret}(T) = \tilde{O} \left( \kappa_\mu B_T^{1/5} T^{4/5} \right) .$$

# Minimax rates in general non-stationary settings

- Beyond piece-wise stationarity thanks to *variation-budget*:

$$B_T := \sum_{t=2}^T \|\theta_\star^t - \theta_\star^{t-1}\| .$$

- ▶ describe broader non-stationary environments.
- 
- Forgetting mechanisms to the rescue?
    - ▶ minimax-optimal for the MAB setting
    - ▶ not so well understood in LB! [F. et al, 2021a]
    - ▶ best know regret bound for GLBs [F. et al. 2021b]:

$$\text{DynamicRegret}(T) = \tilde{O} \left( \kappa_\mu B_T^{1/5} T^{4/5} \right) .$$

- Room for improvement!  $\rightarrow$  [Wei and Luo, 2021]