

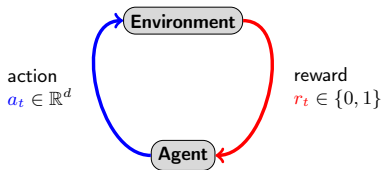
Jointly Efficient and Optimal Algorithms for Logistic Bandits

ML Big Weeks

Marc Abeille, Louis Faury

The Learning Problem

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



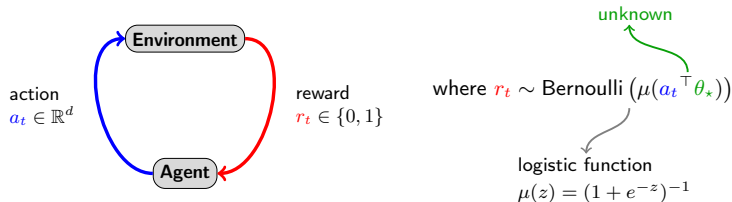
where $r_t \sim \text{Bernoulli}(\mu(a_t^\top \theta_*)$)

logistic function
 $\mu(z) = (1 + e^{-z})^{-1}$

unknown

The Learning Problem

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

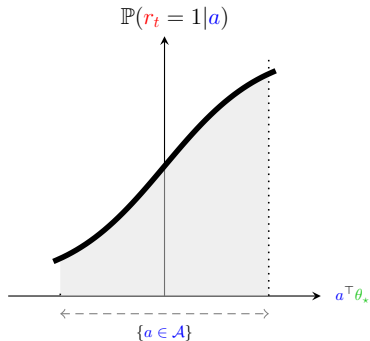


- Goal.** Maximize $\sum_{t=1}^T \mu(a_t^\top \theta_*) =$ performance over time.

θ_* is unknown! \Rightarrow **exploration-exploitation** dilemma.

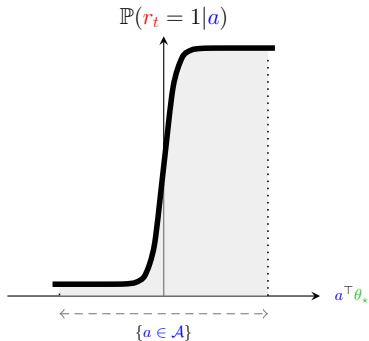
Reward Model: Closer Look

Different regimes



$$\forall a \in \mathcal{A}, \mathbb{P}(r_t = 1|a) \simeq 0.5$$

✓ SOTA!

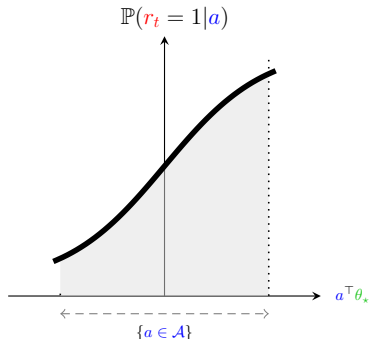


$$\exists a \in \mathcal{A}, \mathbb{P}(r_t = 1|a) \approx 0$$

✗ SOTA!

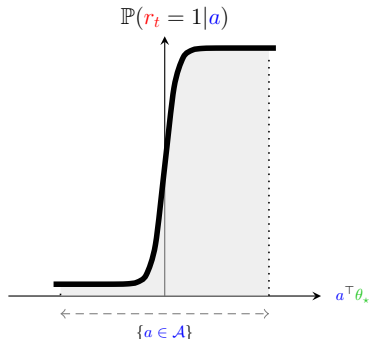
Reward Model: Closer Look

Different regimes



$$\forall a \in \mathcal{A}, \mathbb{P}(r_t = 1|a) \approx 0.5$$

✓ SOTA!



$$\exists a \in \mathcal{A}, \mathbb{P}(r_t = 1|a) \approx 0$$

✗ SOTA!

Key Quantity

$$\kappa \simeq \max_{a \in \mathcal{A}} 1/\mathbb{P}(r_t = 1|a)$$

← typically 10^3 !

Objective

- The performance metric is the **regret**:

$$\text{Regret}(T) = T \max_{a \in \mathcal{A}} \mu(a^\top \theta_*) - \sum_{t=1}^T \mu(a_t^\top \theta_*).$$

Objective

- The performance metric is the **regret**:

$$\text{Regret}(T) = T \max_{a \in \mathcal{A}} \mu(a^\top \theta_*) - \sum_{t=1}^T \mu(a_t^\top \theta_*).$$

- Challenges.** Obtain **optimal** regret with **efficient** algorithms.

Lower-bound <i>[Abeille et al. AISTATS21]</i>			
Algorithm	Regret Bound	Minimax	Efficient
GLM-UCB <i>[Filippi et al. NIPS10.]</i>			
OFULog-r <i>[Fauray et al. ICML20]</i>			
OFU-ECOLog (submitted)			

Objective

- The performance metric is the **regret**:

$$\text{Regret}(T) = T \max_{a \in \mathcal{A}} \mu(a^\top \theta_*) - \sum_{t=1}^T \mu(a_t^\top \theta_*).$$

- Challenges.** Obtain **optimal** regret with **efficient** algorithms.

Lower-bound <i>[Abeille et al. AISTATS21]</i>			
Algorithm	Regret Bound	Minimax	Efficient
GLM-UCB <i>[Filippi et al. NIPS10.]</i>	$\mathcal{O}(\kappa d \sqrt{T})$	X	X
OFULog-r <i>[Fauray et al. ICML20]</i>			
OFU-ECOLog (submitted)			

Objective

- The performance metric is the **regret**:

$$\text{Regret}(T) = T \max_{a \in \mathcal{A}} \mu(a^\top \theta_*) - \sum_{t=1}^T \mu(a_t^\top \theta_*).$$

- Challenges.** Obtain **optimal** regret with **efficient** algorithms.

Lower-bound <i>[Abeille et al. AISTATS21]</i>			
Algorithm	Regret Bound	Minimax	Efficient
GLM-UCB <i>[Filippi et al. NIPS10.]</i>	$\mathcal{O}(\kappa d \sqrt{T})$	X	X
OFULog-r <i>[Faury et al. ICML20]</i>	$\mathcal{O}(d \sqrt{T/\kappa})$	✓	X
OFU-ECOLog (submitted)			

Objective

- The performance metric is the **regret**:

$$\text{Regret}(T) = T \max_{a \in \mathcal{A}} \mu(a^\top \theta_*) - \sum_{t=1}^T \mu(a_t^\top \theta_*).$$

- Challenges.** Obtain **optimal** regret with **efficient** algorithms.

Lower-bound <i>[Abeille et al. AISTATS21]</i>	$\Omega(d\sqrt{T/\kappa})$		
Algorithm	Regret Bound	Minimax	Efficient
GLM-UCB <i>[Filippi et al. NIPS10.]</i>	$\mathcal{O}(\kappa d\sqrt{T})$	✗	✗
OFULog-r <i>[Fauray et al. ICML20]</i>	$\mathcal{O}(d\sqrt{T/\kappa})$	✓	✗
OFU-ECOLog (submitted)			

Objective

- The performance metric is the **regret**:

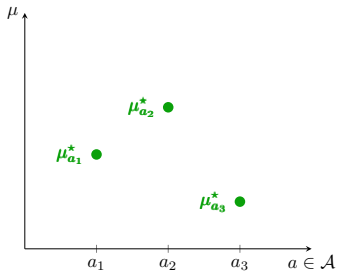
$$\text{Regret}(T) = T \max_{a \in \mathcal{A}} \mu(a^\top \theta_*) - \sum_{t=1}^T \mu(a_t^\top \theta_*).$$

- Challenges.** Obtain **optimal** regret with **efficient** algorithms.

Lower-bound <i>[Abeille et al. AISTATS21]</i>	$\Omega(d\sqrt{T/\kappa})$		
Algorithm	Regret Bound	Minimax	Efficient
GLM-UCB <i>[Filippi et al. NIPS10.]</i>	$\mathcal{O}(\kappa d\sqrt{T})$	✗	✗
OFULog-r <i>[Faury et al. ICML20]</i>	$\mathcal{O}(d\sqrt{T/\kappa})$	✓	✗
OFU-ECOLog (submitted)	$\mathcal{O}(d\sqrt{T/\kappa})$	✓	✓

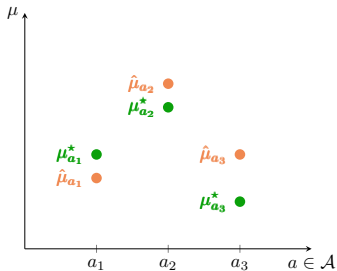
The exploration-exploitation dilemma

High-level idea.



The exploration-exploitation dilemma

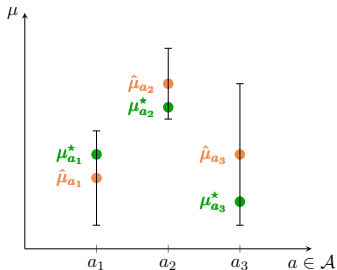
High-level idea.



- **Exploit.** Predict rewards

The exploration-exploitation dilemma

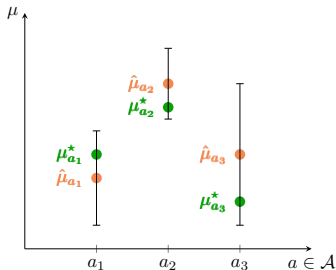
High-level idea.



- **Exploit.** Predict rewards
- **Explore.** Quantify uncertainty

The exploration-exploitation dilemma

High-level idea.



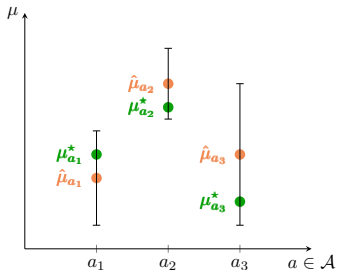
Leveraging the structure.

$$\mu_a^* = \mu(a^\top \theta_*)$$

- **Exploit.** Predict rewards
- **Explore.** Quantify uncertainty

The exploration-exploitation dilemma

High-level idea.



Leveraging the structure.

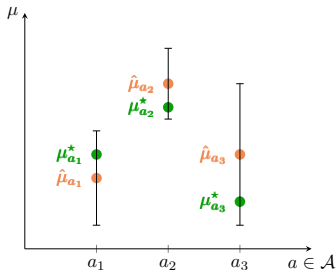
$$\mu_a^* = \mu(a^\top \theta_*)$$



- **Exploit.** Predict rewards
- **Explore.** Quantify uncertainty

The exploration-exploitation dilemma

High-level idea.



Leveraging the structure.

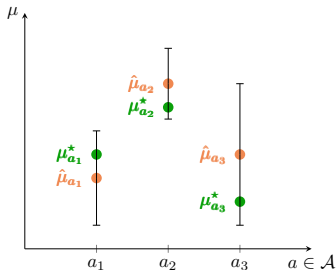
$$\mu_a^* = \mu(a^\top \theta_*)$$



- **Exploit.** Predict rewards $\rightarrow \hat{\mu}_a = \mu(a^\top \hat{\theta})$.
- **Explore.** Quantify uncertainty

The exploration-exploitation dilemma

High-level idea.



Leveraging the structure.

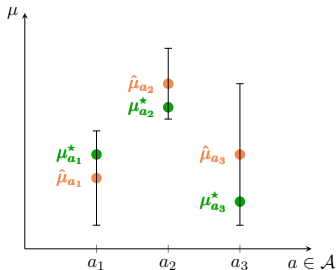
$$\mu_a^* = \mu(a^\top \theta_*)$$



- **Exploit.** Predict rewards $\longrightarrow \hat{\mu}_a = \mu(a^\top \hat{\theta})$.
- **Explore.** Quantify uncertainty \longrightarrow quantify deviation between $\hat{\theta}$ and θ_* .

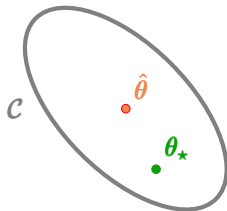
The exploration-exploitation dilemma

High-level idea.



Leveraging the structure.

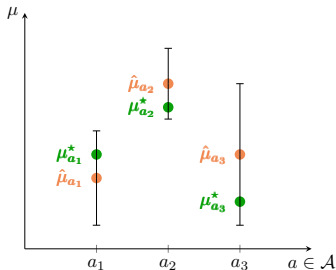
$$\mu_a^* = \mu(a^\top \theta_*)$$



- **Exploit.** Predict rewards $\longrightarrow \hat{\mu}_a = \mu(a^\top \hat{\theta})$.
- **Explore.** Quantify uncertainty \longrightarrow quantify deviation between $\hat{\theta}$ and θ_* .

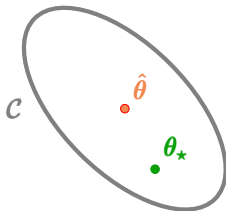
The exploration-exploitation dilemma

High-level idea.



Leveraging the structure.

$$\mu_a^* = \mu(a^\top \theta_*)$$



- **Exploit.** Predict rewards $\rightarrow \hat{\mu}_a = \mu(a^\top \hat{\theta})$.
- **Explore.** Quantify uncertainty \rightarrow quantify deviation between $\hat{\theta}$ and θ_* .

Challenge: design *sharp* and *tractable* confidence set \mathcal{C}

[Filippi et al. NIPS2010]

- Based on the maximum-likelihood estimator and linear design matrix \mathbf{V}_t :

$$\left\{ \begin{array}{l} \hat{\theta}_t = \sum_{s=1}^t \ell_s(\theta) \quad \text{where } \ell_s(\theta) \leftarrow \text{cross-entropy on } (\mathbf{a}_s, \mathbf{r}_s), \\ \mathbf{V}_t = \sum_{s=1}^t \mathbf{a}_s \mathbf{a}_s^\top . \end{array} \right.$$

[Filippi et al. NIPS2010]

- Based on the maximum-likelihood estimator and linear design matrix \mathbf{V}_t :

$$\left\{ \begin{array}{l} \hat{\theta}_t = \sum_{s=1}^t \ell_s(\theta) \quad \text{where } \ell_s(\theta) \leftarrow \text{cross-entropy on } (a_s, r_s), \\ \mathbf{V}_t = \sum_{s=1}^t a_s a_s^\top. \end{array} \right.$$

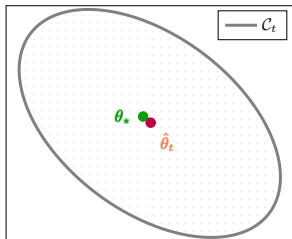
$$\mathcal{C}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{\mathbf{V}_t} \leq \kappa \right\}$$

[Filippi et al. NIPS2010]

- Based on the maximum-likelihood estimator and linear design matrix \mathbf{V}_t :

$$\begin{cases} \hat{\theta}_t = \sum_{s=1}^t \ell_s(\theta) & \text{where } \ell_s(\theta) \leftarrow \text{cross-entropy on } (a_s, r_s), \\ \mathbf{V}_t = \sum_{s=1}^t a_s a_s^\top. \end{cases}$$

$$\mathcal{C}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{\mathbf{V}_t} \leq \kappa \right\}$$

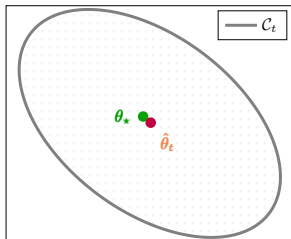


[Filippi et al. NIPS2010]

- Based on the maximum-likelihood estimator and linear design matrix \mathbf{V}_t :

$$\begin{cases} \hat{\theta}_t = \sum_{s=1}^t \ell_s(\theta) & \text{where } \ell_s(\theta) \leftarrow \text{cross-entropy on } (a_s, r_s), \\ \mathbf{V}_t = \sum_{s=1}^t a_s a_s^\top. \end{cases}$$

$$\mathcal{C}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{\mathbf{V}_t} \leq \kappa \right\}$$



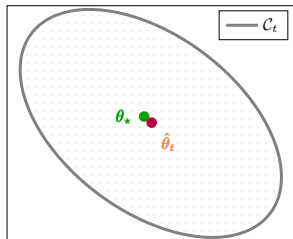
We play $a \approx 10^4$ times.

[Filippi et al. NIPS2010]

- Based on the maximum-likelihood estimator and linear design matrix V_t :

$$\begin{cases} \hat{\theta}_t = \sum_{s=1}^t \ell_s(\theta) & \text{where } \ell_s(\theta) \leftarrow \text{cross-entropy on } (a_s, r_s), \\ V_t = \sum_{s=1}^t a_s a_s^\top. \end{cases}$$

$$\mathcal{C}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{V_t} \leq \kappa \right\}$$



We play $a \approx 10^4$ times.

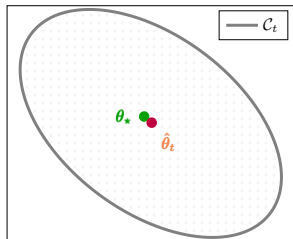
$$\begin{aligned} \hat{\mu}_a - \mu_a^* &\approx \kappa \|a\|_{V_t^{-1}} \\ &\approx \frac{\kappa}{\sqrt{\text{number of times } a \text{ was played}}} \\ &\approx 10^2 / \sqrt{10^4} \\ &\approx 1 \quad \leftarrow \text{trivial bound!} \end{aligned}$$

[Filippi et al. NIPS2010]

- Based on the maximum-likelihood estimator and linear design matrix \mathbf{V}_t :

$$\begin{cases} \hat{\theta}_t = \sum_{s=1}^t \ell_s(\theta) & \text{where } \ell_s(\theta) \leftarrow \text{cross-entropy on } (a_s, r_s), \\ \mathbf{V}_t = \sum_{s=1}^t a_s a_s^\top. \end{cases}$$

$$\mathcal{C}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{\mathbf{V}_t} \leq \kappa \right\}$$



We play $a \approx 10^4$ times.

$$\begin{aligned} \hat{\mu}_a - \mu_a^* &\approx \kappa \|a\|_{\mathbf{V}_t^{-1}} \\ &\approx \frac{\kappa}{\sqrt{\text{number of times } a \text{ was played}}} \\ &\approx 10^2 / \sqrt{10^4} \\ &\approx 1 \quad \leftarrow \text{trivial bound!} \end{aligned}$$

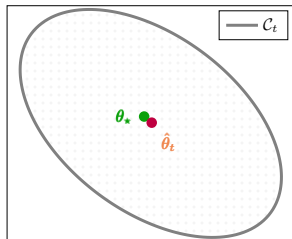
✗ Aggressive exploration

[Filippi et al. NIPS2010]

- Based on the maximum-likelihood estimator and linear design matrix V_t :

$$\begin{cases} \hat{\theta}_t = \sum_{s=1}^t \ell_s(\theta) & \text{where } \ell_s(\theta) \leftarrow \text{cross-entropy on } (a_s, r_s), \\ V_t = \sum_{s=1}^t a_s a_s^\top. \end{cases}$$

$$\mathcal{C}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{V_t} \leq \kappa \right\}$$



We play $a \approx 10^4$ times.

$$\begin{aligned} \hat{\mu}_a - \mu_a^* &\approx \kappa \|a\|_{V_t^{-1}} \\ &\approx \frac{\kappa}{\sqrt{\text{number of times } a \text{ was played}}} \\ &\approx 10^2 / \sqrt{10^4} \\ &\approx 1 \quad \leftarrow \text{trivial bound!} \end{aligned}$$

- ✗ Aggressive exploration
- ✗ Computationally Expensive

- Same estimator $\hat{\theta}_t$ but new metric of deviation:

$$H_t = \sum_{s=1}^t \dot{\mu}(a_s^\top \hat{\theta}_t) a_s a_s^\top .$$

- Same estimator $\hat{\theta}_t$ but new metric of deviation:

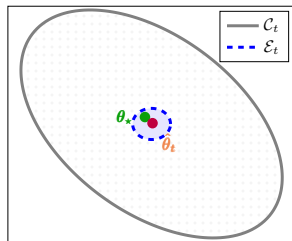
$$\mathbf{H}_t = \sum_{s=1}^t \dot{\mu}(a_s^\top \hat{\theta}_t) a_s a_s^\top .$$

$$\mathcal{E}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{\mathbf{H}_t} \leq 1 \right\}$$

- Same estimator $\hat{\theta}_t$ but new metric of deviation:

$$H_t = \sum_{s=1}^t \dot{\mu}(a_s^\top \hat{\theta}_t) a_s a_s^\top .$$

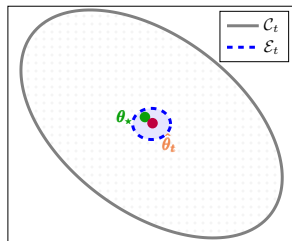
$$\mathcal{E}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{H_t} \leq 1 \right\}$$



- Same estimator $\hat{\theta}_t$ but new metric of deviation:

$$H_t = \sum_{s=1}^t \dot{\mu}(a_s^\top \hat{\theta}_t) a_s a_s^\top .$$

$$\mathcal{E}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{H_t} \leq 1 \right\}$$

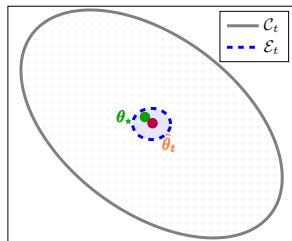


New concentration tools.

- Same estimator $\hat{\theta}_t$ but new metric of deviation:

$$H_t = \sum_{s=1}^t \dot{\mu}(a_s^\top \hat{\theta}_t) a_s a_s^\top .$$

$$\mathcal{E}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{H_t} \leq 1 \right\}$$



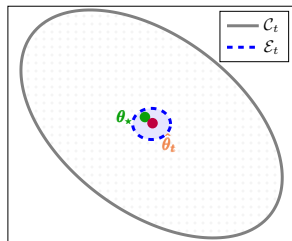
New concentration tools.

We play $a \approx 10^4$ times.

- Same estimator $\hat{\theta}_t$ but new metric of deviation:

$$\mathbf{H}_t = \sum_{s=1}^t \dot{\mu}(a_s^\top \hat{\theta}_t) a_s a_s^\top .$$

$$\mathcal{E}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{\mathbf{H}_t} \leq 1 \right\}$$



New concentration tools.

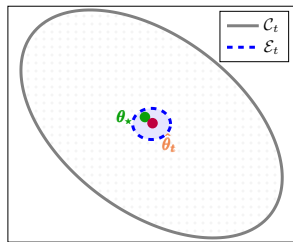
We play $a \approx 10^4$ times.

$$\begin{aligned} \hat{\mu}_a - \mu_a^* &\approx \dot{\mu}(a^\top \hat{\theta}_t) \|a\|_{\mathbf{H}_t}^{-1} \\ &\approx \frac{1}{\sqrt{\text{number of times } a \text{ was played}}} \\ &\approx 10^{-2} \end{aligned}$$

- Same estimator $\hat{\theta}_t$ but new metric of deviation:

$$\mathbf{H}_t = \sum_{s=1}^t \dot{\mu}(a_s^\top \hat{\theta}_t) a_s a_s^\top .$$

$$\mathcal{E}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{\mathbf{H}_t} \leq 1 \right\}$$



New concentration tools.

We play $a \approx 10^4$ times.

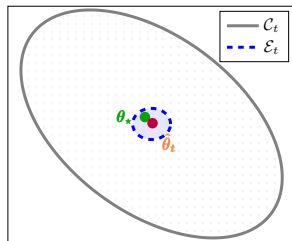
$$\begin{aligned} \hat{\mu}_a - \mu_a^* &\approx \dot{\mu}(a^\top \hat{\theta}_t) \|a\|_{\mathbf{H}_t^{-1}} \\ &\approx \frac{1}{\sqrt{\text{number of times } a \text{ was played}}} \\ &\approx 10^{-2} \end{aligned}$$

- ✓ Efficient exploration
- ✓ Optimal (Cramer-Rao)

- Same estimator $\hat{\theta}_t$ but new metric of deviation:

$$\mathbf{H}_t = \sum_{s=1}^t \dot{\mu}(a_s^\top \hat{\theta}_t) a_s a_s^\top .$$

$$\mathcal{E}_t = \left\{ \theta \in \mathbb{R}^d, \|\theta - \hat{\theta}_t\|_{\mathbf{H}_t} \leq 1 \right\}$$



New concentration tools.

We play $a \approx 10^4$ times.

$$\begin{aligned} \hat{\mu}_a - \mu_a^* &\approx \dot{\mu}(a^\top \hat{\theta}_t) \|a\|_{\mathbf{H}_t^{-1}} \\ &\approx \frac{1}{\sqrt{\text{number of times } a \text{ was played}}} \\ &\approx 10^{-2} \end{aligned}$$

- ✓ Efficient exploration
- ✓ Optimal (Cramer-Rao)
- ✗ Computationally expensive!

Fast and Optimal Approach (submitted)

- Main bottleneck: computation of $\hat{\theta}_t$ and H_t .
 - ▶ $\Omega(t)$ operations at each round!
 - ▶ In practice, **very** slow.

Fast and Optimal Approach (submitted)

- Main bottleneck: computation of $\hat{\theta}_t$ and H_t .
 - ▶ $\Omega(t)$ operations at each round!
 - ▶ In practice, **very** slow.
- Efficient alternative through recursive-least-squares-like operations:

$$\text{ECOLog procedure: } \begin{cases} \theta_t = \operatorname{argmin}_{\theta} \|\theta - \theta_{t-1}\|_{\mathbf{W}_{t-1}}^2 + \ell_t(\theta) , \\ \mathbf{W}_t = \mathbf{W}_{t-1} + \mu(\mathbf{a}_t^\top \theta_t) \mathbf{a}_t \mathbf{a}_t^\top . \end{cases}$$

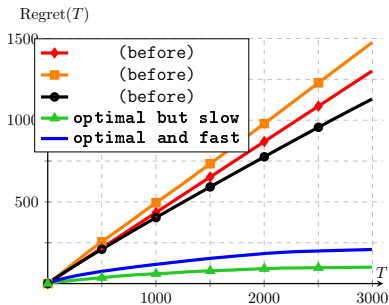
Fast and Optimal Approach (submitted)

- Main bottleneck: computation of $\hat{\theta}_t$ and H_t .
 - ▶ $\Omega(t)$ operations at each round!
 - ▶ In practice, **very** slow.
- Efficient alternative through recursive-least-squares-like operations:

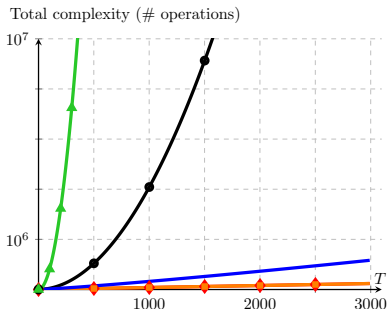
$$\text{ECOLog procedure: } \begin{cases} \theta_t = \operatorname{argmin}_{\theta} \|\theta - \theta_{t-1}\|_{\mathbf{W}_{t-1}}^2 + \ell_t(\theta), \\ \mathbf{W}_t = \mathbf{W}_{t-1} + \mu(a_t^\top \theta_t) a_t a_t^\top. \end{cases}$$

- Best of both world:
 - ✓ **Online computations:** $\tilde{O}(1)$ operations!
 - ✓ **Statistical tightness:** $\mathcal{E}'_t = \left\{ \theta, \|\theta - \theta_t\|_{\mathbf{W}_t} \leq 1 \right\} \approx \mathcal{E}_t$.

Can we see some curves?

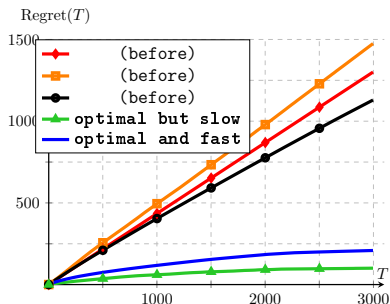


Regret: the smaller the better.

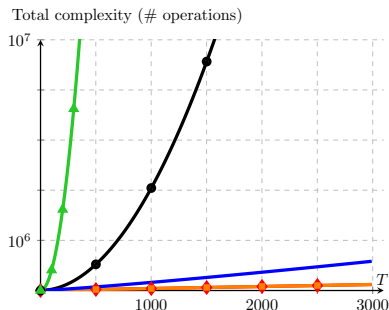


Complexity: the smaller the better

Can we see some curves?



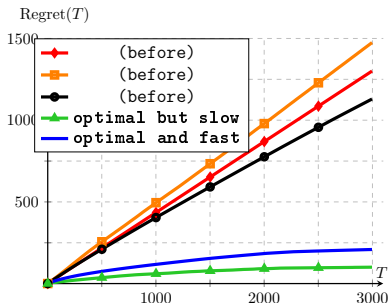
Regret: the smaller the better.



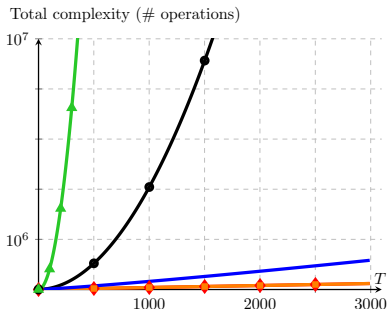
Complexity: the smaller the better

- **Blue curve:** Best of both world behavior.

Can we see some curves?



Regret: the smaller the better.



Complexity: the smaller the better

- **Blue curve**: Best of both world behavior.
- In short: **mature** for deployment in real-life situations.

What's in it for you?

- **Learning.** Principled and efficient **estimation** procedure:
 - ▶ Same convergence guarantee as MLE.
 - ▶ Fully online.
 - ▶ Compatible with non-stationary environments.

What's in it for you?

- **Learning.** Principled and efficient **estimation** procedure:
 - ▶ Same convergence guarantee as MLE.
 - ▶ Fully online.
 - ▶ Compatible with non-stationary environments.

- **Planning.** Principled **exploration**:
 - ▶ Readily usable confidence sets / prediction errors.
 - ▶ Long-term optimal without burning cash.
 - ▶ Compatible with randomized exploration.

What's in it for you?

- **Learning.** Principled and efficient **estimation** procedure:
 - ▶ Same convergence guarantee as MLE.
 - ▶ Fully online.
 - ▶ Compatible with non-stationary environments.

- **Planning.** Principled **exploration**:
 - ▶ Readily usable confidence sets / prediction errors.
 - ▶ Long-term optimal without burning cash.
 - ▶ Compatible with randomized exploration.

Thank you! Questions?