

# Non-Asymptotic Pure Exploration by Solving Games

Rémy Degenne, Wouter M. Koolen



and Pierre Ménard



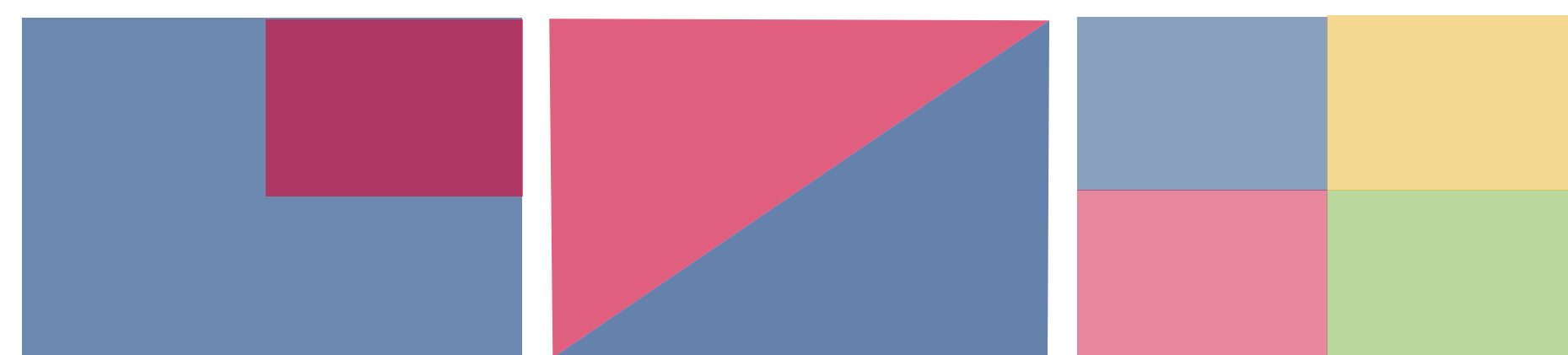
## Topic: Pure Exploration

We want to answer a question about the parameters of a stochastic bandit.

- $K$  arms with unknown distribution parameters (vector  $\mu$  of means).
- A query. Ex: is there an arm with mean  $\mu < 0$ ?  $\rightarrow$  correct answer at  $\mu$  is given by function  $i^*(\mu)$
- At each stage, choose an arm and get an observation from the arm.
- Decide when to stop and return an answer.

Goals:

- answer correctly with probability  $> 1 - \delta$ .
- small sample complexity: stop at  $\tau_\delta$  s.t.  $\mathbb{E}_\mu[\tau_\delta]$  is small.



$\exists k, \mu^k < 0$ ?  $\operatorname{argmax}_k \mu^k$ ? signs of all  $\mu^k$ ?

## Exploration as a Game

Lower Bound, with value of a game:

$$\mathbb{E}_\mu[\tau_\delta] \geq \log\left(\frac{1}{\delta}\right) / \max_{w \in \Delta_K} \inf_{\lambda \in \neg i^*(\mu)} \sum_{k=1}^K w^k d(\mu^k, \lambda^k)$$

Optimal fixed sampling:

$$w^* = \operatorname{argmax}_{w \in \Delta_K} \inf_{\lambda \in \neg i^*(\mu)} \sum_{k=1}^K w^k d(\mu^k, \lambda^k)$$

## Previous work: Track and Stop

- my estimate  $\hat{\mu}_t$  has answer YES.
- the optimal way to sample at  $\hat{\mu}_t$  from the lower bound is  $w_t^*$ .
- I sample to track  $w_t^*$  (+ forced exploration)

Is asymptotically optimal (but sometimes very asymptotically).

Need  $\operatorname{argmax}_{w \in \Delta_K} \inf_{\lambda \in \neg i^*(\hat{\mu}_t)} \sum_{k=1}^K w^k d(\hat{\mu}_t^k, \lambda^k)$  at every time step.

Relies on forced exploration, not adaptive to data.

Solving bandit pure exploration problems with games is computationally efficient and has optimal asymptotic sample complexity.

## Our Strategy

Emulate Nature with a second algorithm.

Get two algorithms playing against each other.

- Alg: my estimate  $\hat{\mu}_t$  has answer YES.
- Nature: but here is  $\lambda_t$  with answer NO which could have generated the same data with relatively high probability.
- Alg: then I sample  $k_t$  s.t. if  $\mu \approx \hat{\mu}_t$  (+optimism), I get maximal evidence for  $\mu \neq \lambda_t$ .

## Why it works

As long as we do not stop:

$$\log \frac{1}{\delta} \geq \inf_{\lambda \in \neg i^*(\hat{\mu}_t)} \sum_{k=1}^K N_t^k d(\mu^k, \lambda^k) \quad (\text{stop rule})$$

$$\approx \inf_{\lambda \in \neg i^*(\hat{\mu}_t)} \sum_{s=1}^t \sum_{k=1}^K w_s^k d(\mu^k, \lambda^k) \quad (\text{tracking})$$

$$\geq \sum_{s=1}^t \sum_{k=1}^K w_s^k \mathbb{E}_{\lambda \sim q_s} d(\mu^k, \lambda^k) - R_t^\lambda \quad (\text{regret } \lambda)$$

$$\geq \max_k \sum_{s=1}^t \mathbb{E}_{\lambda \sim q_s} d(\mu^k, \lambda^k) - R_t^\lambda - R_t^k \quad (\text{regret } k)$$

$$\geq t \inf_{q \in \mathcal{P}(\neg i^*)} \max_k \mathbb{E}_{\lambda \sim q} d(\mu^k, \lambda^k) - O(\sqrt{t})$$

But  $\mu$  unknown  $\rightarrow$  optimism to explore efficiently.

## Results

- Non-asymptotic sample complexity guarantees
- Asymptotically optimal
- Need only  $\operatorname{argmin}_{\lambda \in \neg i^*(\hat{\mu}_t)} \sum_{k=1}^K N_{t-1}^k d(\hat{\mu}_t^k, \lambda^k)$ .  $\rightarrow$  better computational complexity. Up to 100x faster than T-and-S on best arm identification.

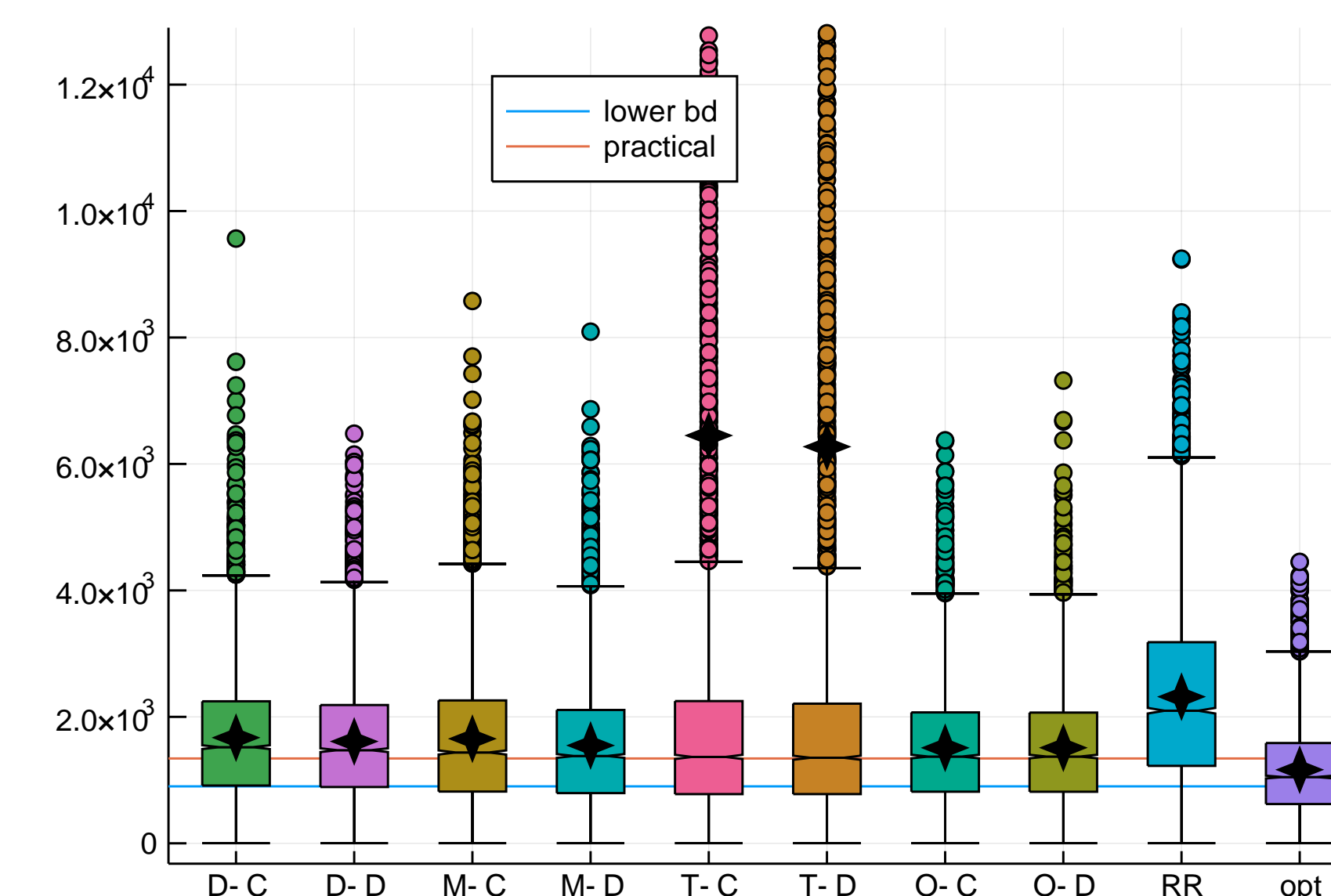


Figure 1: Track and Stop can fail, even for  $\delta = e^{-10}$ . Our algorithms are (provably) good.

## Algorithm

Inputs:

- Algorithms  $\mathcal{A}^k$  and  $\mathcal{A}^\lambda$ , full information adversarial regret minimization algorithms.
- stopping threshold  $\beta(t, \delta) \approx \log \frac{\log t}{\delta}$ , exploration bonus  $f(t) \approx \log t$ .

Algorithm:

- Sample each arm once and form estimate  $\hat{\mu}_K$ . For  $t = K + 1, \dots$ 
  - For  $k \in [K]$ , let  $[\alpha_t^k, \beta_t^k] = \{\xi : N_{t-1}^k d(\hat{\mu}_{t-1}^k, \xi) \leq f(t-1)\}$ . (KL confidence intervals)
  - Let  $i_t = i^*(\hat{\mu}_{t-1})$ .
  - Stop and output  $\hat{i} = i_t$  if  $\inf_{\lambda \in \neg i_t} \sum_k N_{t-1}^k d(\hat{\mu}_{t-1}^k, \lambda^k) > \beta(t, \delta)$ . (GLRT Stopping rule)
  - Get  $w_t$  and  $q_t$  from  $\mathcal{A}_{i_t}^k$  and  $\mathcal{A}_{i_t}^\lambda$ .
  - For  $k \in [K]$ , let  $U_t^k = \max_{\xi \in \{\alpha_t^k, \beta_t^k\}} \mathbb{E}_{\lambda \sim q_t} d(\xi, \lambda^k)$ . (Optimism)
  - Feed  $\mathcal{A}_{i_t}^k$  the loss  $\ell_t^w(w) = -\sum_{k=1}^K w^k U_t^k$ .
  - Feed  $\mathcal{A}_{i_t}^\lambda$  the loss  $\ell_t^\lambda(q) = \mathbb{E}_{\lambda \sim q} \sum_{k=1}^K w_t^k d(\hat{\mu}_{t-1}^k, \lambda^k)$ .
  - Pick arm  $k_t = \operatorname{argmin}_k N_{t-1}^k - \sum_{s=1}^t w_s^k$ . (Cumulative tracking)
  - Observe sample  $X_t \sim \nu_{k_t}$ . Update  $\hat{\mu}_t$ .

I want to know more!



Take a picture  
to download the full paper.