

What's new at Inria? DELTA Project Meeting



Liège, April 30th, 2019

Permanent researcher:

- Michal Valko (Inria → DeepMind... but still in Delta)
- Emilie Kaufmann (CNRS)

PhD candidate:

- Omar Darwiche-Domingues (since October 2018)

Post-doc:

- Pierre Ménard (since February 2019)

Several collaborators

Task 2.1, “Best Arm Identification Tools for Planning”

- 1 News (BAI) tools that can be useful for Planning
- 2 Keeping Non-Stationarity in Mind
- 3 Recent Work on Planning (and the Simulator)

- 1 News (BAI) tools that can be useful for Planning
- 2 Keeping Non-Stationarity in Mind
- 3 Recent Work on Planning (and the Simulator)

The Power of Zipf Sampling

Abbasi-Yadkori, Bartlett, Gabillon, Malek and Valko. *Best of both worlds: Stochastic & adversarial best-arm identification*, COLT'18

For $t = 1, 2, \dots$

- ▶ Sort and rank the arms by decreasing order of estimated cumulated gain $\hat{G}_k(t-1)$: Rank arm k as $\langle \tilde{k} \rangle_t$
- ▶ Select arm $A_t \in [K]$ at random such that

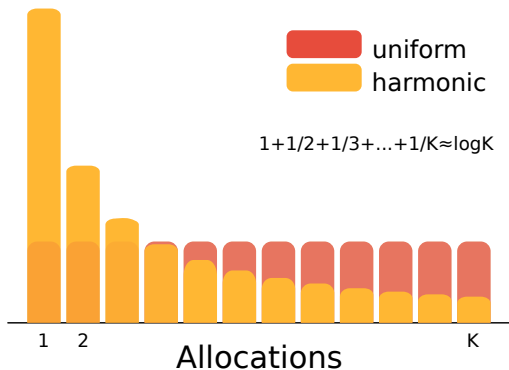
$$\mathbb{P}(A_t = k) = \frac{1}{\langle \tilde{k} \rangle_t \log(K)}$$

Recommend, at any given round t ,

$$J_t \triangleq \arg \max_{k \in \{1, \dots, K\}} \hat{G}_k(t).$$

Figure: The P1 algorithm

The Power of Zipf Sampling

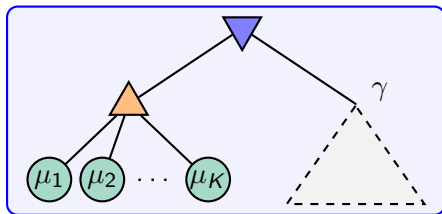


Variants of this simple allocation rule has already been used in different settings:

- black-box optimization (SequOOL, ALT'19)
- planning (Plat γ POOs, ICML' 19)

Some new insights from a toy Active Testing Problem

Kaufmann, Koolen and Garivier, *Sequential Test for the Lowest Mean: from Thompson to Murphy Sampling*, NeurIPS'18



Fix **threshold** γ .

$$\mu^* := \min_i \mu_i \leq \gamma?$$

For $t = 1, \dots, \tau$

- pick an arm A_t
- observe $X_t \sim \mu_{A_t}$

After stopping, recommend $\hat{m} \in \{<, >\}$

Goal: **controlled error** $\mathbb{P}_\mu \{\text{error}\} < \delta$
and small **sample complexity** $\mathbb{E}_\mu[\tau]$



Lower Bound and Oracle Allocation

Generic lower bound [Garivier et al. 16] shows *sample complexity* for any δ -correct algorithm is at least

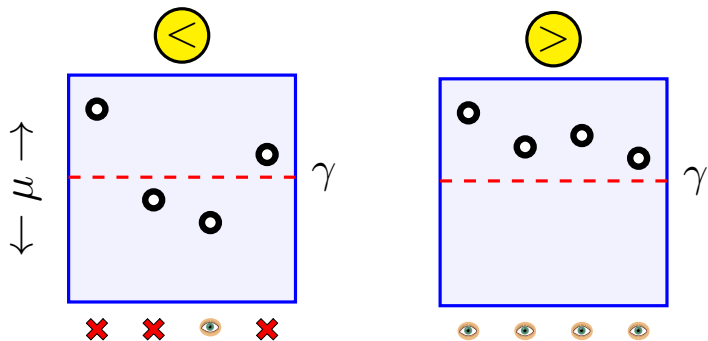
$$\mathbb{E}_{\mu}[\tau] \geq T^*(\mu) \ln\left(\frac{1}{\delta}\right).$$

For our problem the *characteristic time* and *oracle weights* are

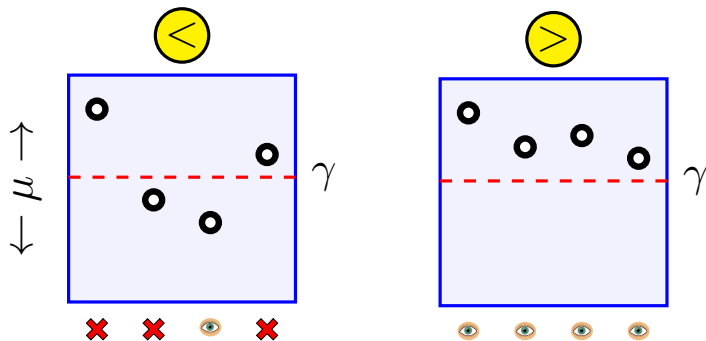
$$T^*(\mu) = \begin{cases} \frac{1}{d(\mu^*, \gamma)} & \mu^* < \gamma, \\ \sum_a \frac{1}{d(\mu_a, \gamma)} & \mu^* > \gamma, \end{cases} \quad w_a^*(\mu) = \begin{cases} \mathbf{1}_{(a=a^*)} & \mu^* < \gamma, \\ \frac{1}{d(\mu_a, \gamma)} & \mu^* > \gamma. \\ \frac{1}{\sum_j \frac{1}{d(\mu_j, \gamma)}} & \mu^* > \gamma. \end{cases}$$

$w_a^*(\mu)$: fraction of selections of arm a under a strategy that would match the lower bound

Dichotomous Oracle Behaviour! Sampling Rule?



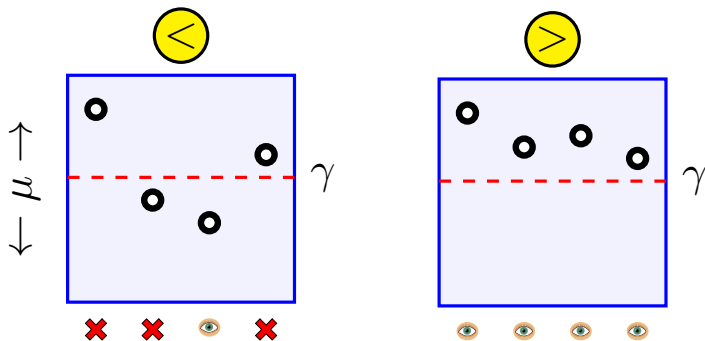
Dichotomous Oracle Behaviour! Sampling Rule?



Two different ideas to get those sampling profiles:

- **Thompson Sampling** (Π_{t-1} is posterior after $t-1$ rounds)
Sample $\theta \sim \Pi_{t-1}$, then play $A_t = \arg \min_a \theta_a$.
- **a Lower Confidence Bound algorithm**
Play $A_t = \arg \min_a \text{LCB}_a(t)$

A Solution: Murphy Sampling!



A more flexible idea:









- **Murphy Sampling** **condition on low minimum mean**

Sample $\theta \sim \Pi_{t-1}(\cdot | \min_a \theta_a < \gamma)$, then play $A_t = \arg \min_a \theta_a$.

→ converges to the optimal allocation in both cases!

Theorem

Asymptotic optimality: $N_a(t)/t \rightarrow w_a^*(\mu)$ for all μ

Sampling rule		
Thompson Sampling		
Lower Confidence Bounds		
Murphy Sampling		

Lemma

Any anytime sampling strategy $(A_t)_t$ ensuring $\frac{N_t}{t} \rightarrow w^*(\mu)$ and good stopping rule τ_δ guarantee $\limsup_{\delta \rightarrow 0} \frac{\tau_\delta}{\ln \frac{1}{\delta}} \leq T^*(\mu)$.

→ Murphy Sampling combined with a **good stopping rule** asymptotically attains the optimal sample complexity.

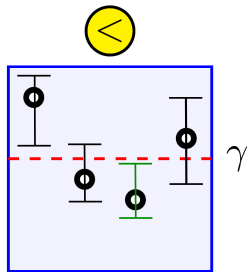
What is a “good stopping rule”?

Example: a stopping rule based on **individual confidence bounds**:

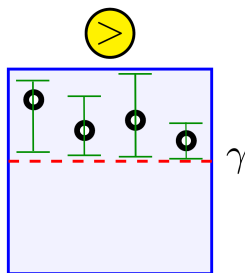
$\tau^{\text{Box}} := \min(\tau_{<}; \tau_{>})$ where

$$\tau_{<} = \inf\{t \in \mathbb{N} : \exists a : \text{UCB}_a(t) < \gamma\}$$

$$\tau_{>} = \inf\{t \in \mathbb{N} : \forall a, \text{LCB}_a(t) > \gamma\}$$



$$\tau = \tau_{<}$$



$$\tau = \tau_{>}$$

What is a “good stopping rule”?

Example: a stopping rule based on **individual confidence bounds**:

$\tau^{\text{Box}} := \min(\tau_{<}; \tau_{>})$ where

$$\tau_{<} = \inf\{t \in \mathbb{N} : \exists a : \text{UCB}_a(t) < \gamma\}$$

$$\tau_{>} = \inf\{t \in \mathbb{N} : \forall a, \text{LCB}_a(t) > \gamma\}$$

→ enough to have the previous (asymptotic) results, but in practice we want to leverage the following:

Multiple low arms
identical or similar \Rightarrow $\left\{ \begin{array}{l} \text{conclude } \mu^* < \gamma \text{ **faster**} \\ \text{**tighter** confidence interval for } \mu^* \end{array} \right.$

Improved Upper Confidence Bound on a Minimum

Given a subset $\mathcal{S} \subseteq \{1, \dots, K\}$, let

- $N_{\mathcal{S}}(t)$ the number of selections of an arm in \mathcal{S}
- $\hat{\mu}_{\mathcal{S}}(t)$ the **aggregated empirical mean**

Theorem

For any prior π , for an appropriate choice of threshold \mathcal{T} ,

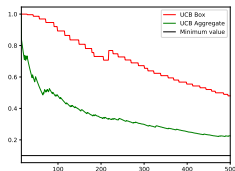
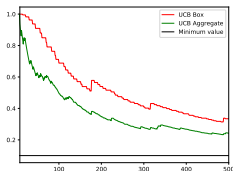
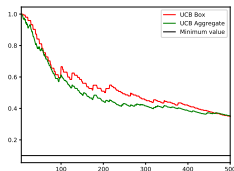
$$\text{UCB}_{\min}^{\pi}(t) := \max \left\{ q : \exists \mathcal{S} \subseteq [K] : \left[N_{\mathcal{S}} d^+(\hat{\mu}_{\mathcal{S}}, q) - \ln \ln N_{\mathcal{S}} \right] \leq \mathcal{T} \left(\ln \frac{1}{\delta \pi(\mathcal{S})} \right) \right\}$$

satisfies $\mathbb{P}(\forall t \in \mathbb{N}, \text{UCB}_{\min}^{\pi}(t) > \mu^*) \geq 1 - \delta$.

Improved stopping rule:

$$\tau_{<} = \inf \{ t \in \mathbb{N} : \text{UCB}_{\min}^{\pi}(t) \leq \gamma \}$$

Improved Upper Confidence Bounds on a Minimum



UCB for minimum: **Agg** dominates **Box** with 1, 3 and 10 low arms.

(on can also get a larger LCB on the maximum mean)

- 1 News (BAI) tools that can be useful for Planning
- 2 Keeping Non-Stationarity in Mind
- 3 Recent Work on Planning (and the Simulator)

The Complexity of Rotting Bandits

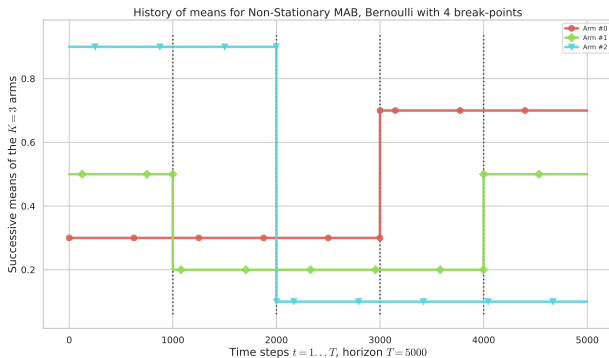
Rotting bandits: each time an arm is played, its mean decreases (a specific form of non-stationnarity)

Seznec, Locatelli, Carpentier, Lazaric, Valko. *Rotting bandits are not harder than stochastic ones*, AISTATS'19 (oral presentation)

(the reason Michal is not here today...)



Piecewise-Stationary Model: one Example



nb of breakpoints: $\Upsilon_T = 4$

(Quick) related work

- Existing guarantees for a variant of EXP3
EXP3.S [Auer et al. 2002]
- Still, many attempts to adapt *stochastic bandit algorithms* to this problem: CUSUM-UCB [Liu et al, 2018], Monitored-UCB [Cao et al, 2019]
- Those attempts require the knowledge of

the number of breakpoints + a lower bound on the minimal magnitude of change

(Quick) related work

- Existing guarantees for a variant of EXP3
EXP3.S [Auer et al. 2002]
- Still, many attempts to adapt *stochastic bandit algorithms* to this problem: CUSUM-UCB [Liu et al, 2018], Monitored-UCB [Cao et al, 2019]
- Those attempts require the knowledge of

the number of breakpoints + a lower bound on the minimal magnitude of change

Our contributions:

- kl-UCB + un efficient adaptive sliding window
- no need to know anything about the size of a change

Context: Piecewise i.i.d. bandit with bounded rewards.

Key tool: an efficient change-point detector to detect a change in the mean of a bounded distribution, the **Bernoulli-GLR**

→ given a stream of samples $(X_s) \in [0, 1]$, detection occurs after n samples if

$$\sup_{s \in [1, n]} \left[s \times \text{kl}(\hat{\mu}_{1:s}, \hat{\mu}_{1:n}) + (n - s) \times \text{kl}(\hat{\mu}_{s+1:n}, \hat{\mu}_{1:n}) \right] \geq \beta(n, \delta)$$

where $\hat{\mu}_{s:s'} = (\sum_{k=s}^{s'} X_k) / (s' - s + 1)$ and

$$\text{kl}(x, y) = x \ln(x/y) + (1 - x) \ln((1 - x)/(1 - y)).$$

Context: Piecewise i.i.d. bandit with bounded rewards.

Key tool: an efficient change-point detector to detect a change in the mean of a bounded distribution, the **Bernoulli-GLR**

→ given a stream of samples $(X_s) \in [0, 1]$, detection occurs after n samples if

$$\sup_{s \in [1, n]} \left[s \times \text{kl}(\hat{\mu}_{1:s}, \hat{\mu}_{1:n}) + (n - s) \times \text{kl}(\hat{\mu}_{s+1:n}, \hat{\mu}_{1:n}) \right] \geq \beta(n, \delta)$$

Lemma (false alarm probability)

For $\beta(n, \delta) \simeq \ln(3n\sqrt{n}/\delta)$ the probability that a detection occurs on a i.i.d. stream is at most δ .

kl-UCB meets the GLRT

Parameters: $\alpha \in (0, 1)$, $\delta > 0$.

Arm selection: at round t ,

- if $\alpha > 0$ and $t \bmod \lfloor K/\alpha \rfloor \in \{1, \dots, K\}$,

$$\text{(forced exploration)} \quad A_t \leftarrow t \bmod \lfloor K/\alpha \rfloor$$

- else, select

$$\text{(kl-UCB)} \quad A_t \leftarrow \arg \max_a \text{UCB}_a(t)$$

$\tau_a(t)$: instant of the last **restart**

$n_a(t)$: number of selection of arm a since the last restart

$\hat{\mu}_a(t)$: empirical mean of samples from arm a since last restart

$$\text{UCB}_a(t) := \max\{q \in [0, 1] : n_a(t) \times \text{kl}(\hat{\mu}_a(t), q) \leq f(t - \tau_a(t))\}.$$

Restarts: **Local** or **Global** after a change is detected by the Bernoulli-GLRT on the mean of the selected arm

- a unified analysis of Local and Global changes
- a tuning of the algorithm that ensures $O(\Upsilon_T \sqrt{T})$ when Υ_T is unknown and $O(\sqrt{\Upsilon_T T})$ regret if Υ_T is known
- good practical performance !

work in progress with Lilian Besson (CentraleSupélec Rennes) and Odalric Maillard (Inria)

- 1 News (BAI) tools that can be useful for Planning
- 2 Keeping Non-Stationarity in Mind
- 3 Recent Work on Planning (and the Simulator)

Planning in Regularized MDPs and Games

- Planning problem: given a generative model, estimate the value function at a state s ;
- K actions, state space of any cardinality;
- We study value functions of entropy regularized MDPs and games.
- Example: Bellman equations for MDPs with entropy regularization

$$V(s) = \max_{\pi(\cdot|s) \in \mathcal{P}(\mathcal{A})} \mathbb{E}[r(s, a) + \gamma V(z)] + \lambda \underbrace{\mathcal{H}(\pi(\cdot|s))}_{\text{entropy}} \quad (1)$$

$$= \lambda \log \sum_{a \in \mathcal{A}} \exp \left(\frac{1}{\lambda} \mathbb{E}[r(s, a) + \gamma V(z)] \right), \quad z \sim P(\cdot|s, a) \quad (2)$$

Planning in Regularized MDPs and Games

General case:

$$V(s) = F_s(Q_s), \quad \text{with} \quad Q_s(a) = \mathbb{E}[r(s, a) + \gamma V(z)], \quad z \sim P(\cdot | s, a) \quad (3)$$

- $F_s = \max$: Bellman equations for MDPs;
- $F_s = \max$ or \min according to the player: value function for turn-based two-player game (discounted).
- Replace \max and \min by smooth approximations with $\text{LogSumExp} =$ entropy regularization on the policy.
- Main assumptions:
 - (smoothness) $|F_s(x) - F_s(x_0) - (x - x_0)^T \nabla F_s(x_0)| \leq L \|x - x_0\|_2^2$;
 - F_s is 1-Lipschitz, nonnegative gradient.

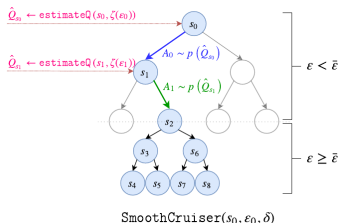
Our algorithm

Algorithm 1 sampleV

```

1: Input:  $(s, \varepsilon) \in \mathcal{S} \times \mathbb{R}_+$ 
2: if  $\varepsilon \geq 1/(1-\gamma)$  then
3:   Output: 0
4: else if  $\varepsilon \geq \bar{\varepsilon}$  then
5:    $\hat{Q}_s \leftarrow \text{estimateQ}(s, \varepsilon)$ 
6:   Output:  $F_s(\hat{Q}_s)$ 
7: else if  $\varepsilon < \bar{\varepsilon}$  then
8:    $\hat{Q}_s \leftarrow \text{estimateQ}(s, \sqrt{\varepsilon\varepsilon})$ 
9:    $A \leftarrow \text{action drawn from } \frac{\nabla F_s(\hat{Q}_s)}{\|\nabla F_s(\hat{Q}_s)\|_1}$ 
10:   $(R, Z) \leftarrow \text{oracle}(s, A)$ 
11:   $\hat{V} \leftarrow \text{sampleV}(Z, \varepsilon/\sqrt{\gamma})$ 
12: end if
13: Output:  $F_s(\hat{Q}_s) - \hat{Q}_s^\top \nabla F_s(\hat{Q}_s) + (R + \gamma \hat{V}) \|\nabla F_s(\hat{Q}_s)\|_1$ 

```



Algorithm 2 estimateQ

```

1: Input:  $(s, \varepsilon)$ 
2: // Compute the value of N
3: if  $\varepsilon \geq 1/(1-\gamma)$  then
4:   Output:  $(0, \dots, 0)$ 
5: else if  $\varepsilon \geq \bar{\varepsilon}$  then
6:    $N \leftarrow \frac{2}{(1-\gamma)^4(1-\sqrt{\gamma})^2} \frac{\log(2K/\delta)}{\varepsilon^2}$ 
7: else if  $\varepsilon < \bar{\varepsilon}$  then
8:    $C \leftarrow \frac{1}{1-\gamma} \left( 4\bar{\varepsilon} + \frac{4}{1-\gamma} + 1 \right)$ 
9:    $N \leftarrow \frac{C^2}{2(1-\sqrt{\gamma})^2} \frac{\log(2K/\delta)}{\varepsilon^2}$ 
10: end if
11: // Average to estimate Q function
12: for  $a \in \mathcal{A}$  do
13:    $q_i \leftarrow 0$  for  $i \in 1, \dots, N$ 
14:   for  $i \in 1, \dots, N$  do
15:      $(R, Z) \leftarrow \text{oracle}(s, a)$ 
16:      $\hat{V} \leftarrow \text{sampleV}(Z, \varepsilon/\sqrt{\gamma})$ 
17:      $q_i \leftarrow R + \gamma \hat{V}$ 
18:   end for
19:    $\hat{Q}_s(a) \leftarrow \text{mean}(q_1, \dots, q_N)$ 
20: end for
21: Output:  $\hat{Q}_s$ 

```

Algorithm 3 SmoothCruiser

```

Input:  $(s, \varepsilon, \delta) \in \mathcal{S} \times \mathbb{R}_+ \times \mathbb{R}_+$ 
 $\bar{\varepsilon} \leftarrow (1-\sqrt{\gamma})/KL$ 
Set  $\delta$  and  $\bar{\varepsilon}$  as a global parameters
 $\hat{Q}_s \leftarrow \text{estimateQ}(s, \varepsilon)$ 
Output:  $F_s(\hat{Q}_s)$ 

```

Theorem

Let $n(\epsilon, \delta)$ be the number of calls to the generative model (oracle) before the algorithm terminates. For any state s and $\epsilon, \delta > 0$,

$$n(\epsilon, \delta) \leq \frac{c_1}{\epsilon^4} \log\left(\frac{c_2}{\delta}\right) \left[c_3 \log\left(\frac{c_4}{\epsilon}\right) \right]^{\log_2(c_5(\log(\frac{c_2}{\delta})))} = \mathcal{O}\left(\frac{1}{\epsilon^{4+c}}\right), \forall c > 0$$

where c_1, c_2, c_3, c_4 and c_5 are constants that depend only on K, L and γ .

Theorem

For any state s and $\delta, \epsilon > 0$,

$$\mathbb{P}\left[\left|\hat{V}(s) - V(s)\right| > \epsilon\right] \leq \delta n(\epsilon, \delta).$$