

# SPADRO kick-off meeting: Optimistic Solutions for Dynamic Resource Allocation

Sébastien Bubeck, Olivier Cappé, Damien Ernst, Aurélien Garivier,  
Sarah Filippi, Emilie Kaufmann, Odalric-Ambrym Maillard, Eric  
Moulines, Rémi Munos, Gilles Stoltz

Institut de Mathématique de Toulouse, Université Paul Sabatier

April 9th, 2014

# Clinical Trials

Idealized situation of clinical trials :

- patients visit you *one after another* for a given disease
- you prescribe one of the (say) *5 treatments* available
- the treatments are *not equally efficient*
- you do not know which one is the best, you *observe the effect* of the prescribed treatment on each patient

⇒ **What do you do ?**

- You must choose each prescription using only the *previous observations*
- Your goal is not to estimate each treatment's efficiency precisely, but to *heal as many patients as possible*

# The (stochastic) Multi-Armed Bandit Model

**Environment**  $K$  arms with parameters  $\theta = (\theta_1, \dots, \theta_K)$  such that for any possible choice of arm  $a_t \in \{1, \dots, K\}$  at time  $t$ , one receives the reward

$$X_t = X_{a_t, t}$$

where, for any  $1 \leq a \leq K$  and  $s \geq 1$ ,  $X_{a, s} \sim \nu_a$ , and the  $(X_{a, s})_{a, s}$  are independent.

**Reward distributions**  $\nu_a \in \mathcal{F}_a$  parametric family, or not. Examples :  
canonical exponential family, general bounded rewards

**Example** Bernoulli rewards :  $\theta \in [0, 1]^K$ ,  $\nu_a = \mathcal{B}(\theta_a)$

**Strategy** The agent's actions follow a dynamical strategy  $\pi = (\pi_1, \pi_2, \dots)$  such that

$$A_t = \pi_t(X_1, \dots, X_{t-1})$$

## Performance Evaluation, Regret

Cumulated Reward  $S_T = \sum_{t=1}^T X_t$

Our goal Choose  $\pi$  so as to maximize

$$\begin{aligned}\mathbb{E}[S_T] &= \sum_{t=1}^T \sum_{a=1}^K \mathbb{E}[\mathbb{E}[X_t \mathbb{1}\{A_t = a\} | X_1, \dots, X_{t-1}]] \\ &= \sum_{a=1}^K \mu_a \mathbb{E}[N_a^\pi(T)]\end{aligned}$$

where  $N_a^\pi(T) = \sum_{t \leq T} \mathbb{1}\{A_t = a\}$  is the number of draws of arm  $a$  up to time  $T$ , and  $\mu_a = E(\nu_a)$ .

Regret Minimization equivalent to minimizing

$$R_T = T\mu^* - \mathbb{E}[S_T] = \sum_{a: \mu_a < \mu^*} (\mu^* - \mu_a) \mathbb{E}[N_a^\pi(T)]$$

where  $\mu^* \in \max\{\mu_a : 1 \leq a \leq K\}$

# Asymptotically Optimal Strategies

- A strategy  $\pi$  is said to be **consistent** if, for any  $(\nu_a)_a \in \mathcal{F}^K$ ,

$$\frac{1}{T} \mathbb{E}[S_T] \rightarrow \mu^*$$

- The strategy is uniformly efficient if for all  $\theta \in [0, 1]^K$  and all  $\alpha > 0$ ,

$$R_T = o(T^\alpha)$$

- There are uniformly efficient strategies and we consider the **best achievable asymptotic performance among uniformly efficient strategies**

# The Bound of Lai and Robbins

One-parameter reward distribution  $\nu_a = \nu_{\theta_a}, \theta_a \in \Theta \subset \mathbb{R}$ .

Theorem [Lai and Robbins, '85]

If  $\pi$  is a uniformly efficient strategy, then for any  $\theta \in \Theta^K$ ,

$$\liminf_{T \rightarrow \infty} \frac{R_T}{\log(T)} \geq \sum_{a: \mu_a < \mu^*} \frac{\mu^* - \mu_a}{\text{KL}(\nu_a, \nu^*)}$$

where  $\text{KL}(\nu, \nu')$  denotes the **Kullback-Leibler divergence**

For example, in the Bernoulli case :

$$\text{KL}(\mathcal{B}(p), \mathcal{B}(q)) = d_{\text{BER}}(p, q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$$

# The Bound of Burnetas and Katehakis

More general reward distributions  $\nu_a \in \mathcal{F}_a$

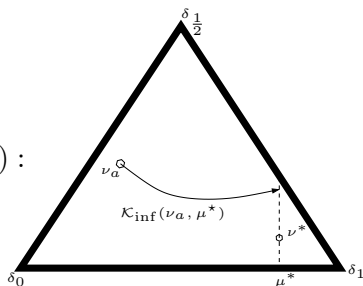
Theorem [Burnetas and Katehakis, '96]

If  $\pi$  is an efficient strategy, then, for any  $\theta \in [0, 1]^K$ ,

$$\liminf_{T \rightarrow \infty} \frac{R_T}{\log(T)} \geq \sum_{a: \mu_a < \mu^*} \frac{\mu^* - \mu_a}{K_{inf}(\nu_a, \mu^*)}$$

where

$$K_{inf}(\nu_a, \mu^*) = \inf \left\{ K(\nu_a, \nu') : \nu' \in \mathcal{F}_a, E(\nu') \geq \mu^* \right\}$$



# Optimism in the Face of Uncertainty

**Optimism** in an heuristic principle popularized by [Lai&Robins '85; Agrawal '95] which consists in letting the agent

play as if the environment was the most favorable among all environments that are sufficiently likely given the observations accumulated so far

Surprisingly, this simple heuristic principle can be instantiated into algorithms that are robust, efficient and easy to implement in many scenarios pertaining to reinforcement learning



# Upper Confidence Bound Strategies

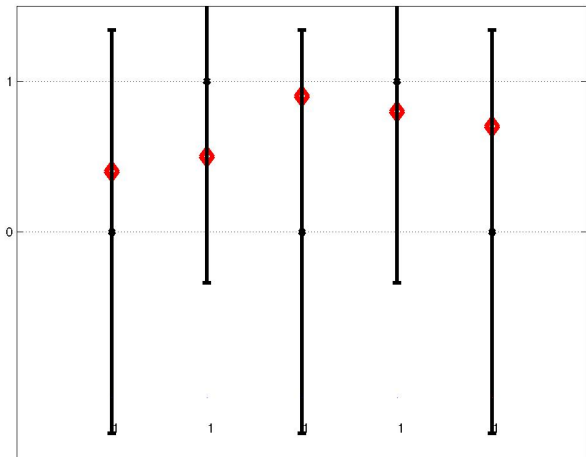
## UCB [Lai&Robins '85; Agrawal '95; Auer&al '02]

- Construct an upper confidence bound for the expected reward of each arm :

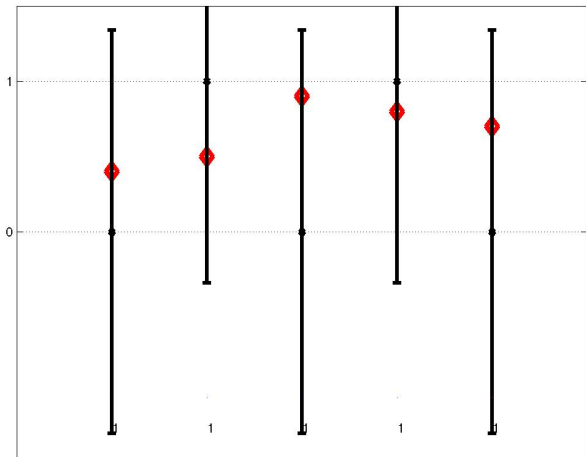
$$\underbrace{\frac{S_a(t)}{N_a(t)}}_{\text{estimated reward}} + \underbrace{\sqrt{\frac{\log(t)}{2N_a(t)}}}_{\text{exploration bonus}}$$

- Choose the arm with the highest UCB
- It is an *index strategy* [Gittins '79]
- Its behavior is easily interpretable and intuitively appealing

# UCB in Action



# UCB in Action



## Performance of UCB

For rewards in  $[0, 1]$ , the regret of UCB is upper-bounded as

$$E[R_T] = O(\log(T))$$

(finite-time regret bound) and

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[R_T]}{\log(T)} \leq \sum_{a: \mu_a < \mu^*} \frac{1}{2(\mu^* - \mu_a)}$$

Yet, in the case of Bernoulli variables, the rhs. is greater than suggested by the bound by Lai & Robbins

Many variants have been suggested to incorporate an estimate of the variance in the exploration bonus (e.g., [Audibert&al '07])

# The KL-UCB algorithm

**Parameters :** An operator  $\Pi_{\mathcal{F}} : \mathfrak{M}_1(\mathcal{S}) \rightarrow \mathcal{F}$ ; a non-decreasing function  $f : \mathbb{N} \rightarrow \mathbb{R}$

**Initialization :** Pull each arm of  $\{1, \dots, K\}$  once

**for**  $t = K$  to  $T - 1$  **do**

    compute for each arm  $a$  the quantity

$$U_a(t) = \sup \left\{ E(\nu) : \nu \in \mathcal{F} \text{ and } KL\left(\Pi_{\mathcal{F}}(\hat{\nu}_a(t)), \nu\right) \leq \frac{f(t)}{N_a(t)} \right\}$$

    pick an arm  $A_{t+1} \in \arg \max_{a \in \{1, \dots, K\}} U_a(t)$

**end for**

## Parametric setting : Exponential Families

- Assume that  $\mathcal{F}_a = \mathcal{F} = \text{canonical exponential family}$ , i.e. such that the pdf of the rewards is given by

$$p_{\theta_a}(x) = \exp(x\theta_a - b(\theta_a) + c(x)), \quad 1 \leq a \leq K$$

for a parameter  $\theta \in \mathbb{R}^K$ , expectation  $\mu_a = \dot{b}(\theta_a)$

- The KL-UCB is simply :

$$U_a(t) = \sup \left\{ \mu \in \bar{I} : d(\hat{\mu}_a(t), \mu) \leq \frac{f(t)}{N_a(t)} \right\}$$

- For instance,
  - for Bernoulli rewards :

$$d_{\text{BER}}(p, q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$$

- for exponential rewards  $p_{\theta_a}(x) = \theta_a e^{-\theta_a x}$  :

$$d_{\text{EXP}}(u, v) = u - v + u \log \frac{u}{v}$$

- The analysis is generic and yields a non-asymptotic regret bound optimal in the sense of Lai and Robbins.

## The kl-UCB algorithm

**Parameters** :  $\mathcal{F}$  parameterized by the expectation  $\mu \in I \subset \mathbb{R}$  with divergence  $d$ , a non-decreasing function  $f : \mathbb{N} \rightarrow \mathbb{R}$

**Initialization** : Pull each arm of  $\{1, \dots, K\}$  once

**for**  $t = K$  to  $T - 1$  **do**

    compute for each arm  $a$  the quantity

$$U_a(t) = \sup \left\{ \mu \in \bar{I} : d(\hat{\mu}_a(t), \mu) \leq \frac{f(t)}{N_a(t)} \right\}$$

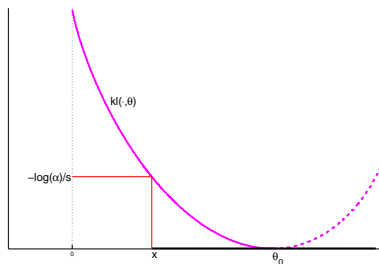
    pick an arm  $A_{t+1} \in \arg \max_{a \in \{1, \dots, K\}} U_a(t)$

**end for**

## The kl Upper Confidence Bound in Picture

If  $Z_1, \dots, Z_s \stackrel{iid}{\sim} \mathcal{B}(\theta_0)$ ,  $x < \theta_0$   
and if  $\hat{p}_s = (Z_1 + \dots + Z_s)/s$ ,  
then

$$\mathbb{P}_{\theta_0}(\hat{p}_s \leq x) \leq \exp(-s \text{kl}(x, \theta_0))$$



In other words, if  $\alpha = \exp(-s \text{kl}(x, \theta_0))$  :

$$\mathbb{P}_{\theta_0}(\hat{p}_s \leq x) = \mathbb{P}_{\theta_0} \left( \text{kl}(\hat{p}_s, \theta_0) \leq -\frac{\log(\alpha)}{s}, \hat{p}_s < \theta_0 \right) \leq \alpha$$

$\implies$  upper confidence bound for  $p$  at risk  $\alpha$  :

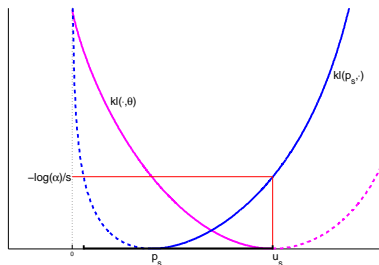
$$u_s = \sup \left\{ \theta > \hat{p}_s : \text{kl}(\hat{p}_s, \theta) \leq -\frac{\log(\alpha)}{s} \right\}$$



# The kl Upper Confidence Bound in Picture

If  $Z_1, \dots, Z_s \stackrel{iid}{\sim} \mathcal{B}(\theta_0)$ ,  $x < \theta_0$   
and if  $\hat{p}_s = (Z_1 + \dots + Z_s)/s$ ,  
then

$$\mathbb{P}_{\theta_0}(\hat{p}_s \leq x) \leq \exp(-s \text{kl}(x, \theta_0))$$



In other words, if  $\alpha = \exp(-s \text{kl}(x, \theta_0))$  :

$$\mathbb{P}_{\theta_0}(\hat{p}_s \leq x) = \mathbb{P}_{\theta_0} \left( \text{kl}(\hat{p}_s, \theta_0) \leq -\frac{\log(\alpha)}{s}, \hat{p}_s < \theta_0 \right) \leq \alpha$$

$\implies$  upper confidence bound for  $p$  at risk  $\alpha$  :

$$u_s = \sup \left\{ \theta > \hat{p}_s : \text{kl}(\hat{p}_s, \theta) \leq -\frac{\log(\alpha)}{s} \right\}$$

## Key Tool : Deviation Inequality for Self-Normalized Sums

- Problem : random number of summands
- Solution : peeling trick (as in the proof of the LIL)

**Theorem** For all  $\epsilon > 1$ ,

$$\mathbb{P}(\mu_a > \hat{\mu}_a(t) \quad \text{and} \quad N_a(t) d(\hat{\mu}_a(t), \mu_a) \geq \epsilon) \leq e \lceil \epsilon \log(t) \rceil e^{-\epsilon}.$$

Thus,

$$P(U_a(t) < \mu_a) \leq e \lceil f(t) \log(t) \rceil e^{-f(t)}$$

## Regret bound

**Theorem** : Assume that all arms belong to a canonical, regular, exponential family  $\mathcal{F} = \{\nu_\theta : \theta \in \Theta\}$  of probability distributions indexed by its natural parameter space  $\Theta \subseteq \mathbb{R}$ . Then, with the choice  $f(t) = \log(t) + 3 \log \log(t)$  for  $t \geq 3$ , the number of draws of any suboptimal arm  $a$  is upper bounded for any horizon  $T \geq 3$  as

$$\mathbb{E}[N_a(T)] \leq \frac{\log(T)}{d(\mu_a, \mu^*)} + 2 \sqrt{\frac{2\pi\sigma_{a,\star}^2 (d'(\mu_a, \mu^*))^2}{(d(\mu_a, \mu^*))^3}} \sqrt{\log(T) + 3 \log(\log(T))} \\ + \left(4e + \frac{3}{d(\mu_a, \mu^*)}\right) \log(\log(T)) + 8\sigma_{a,\star}^2 \left(\frac{d'(\mu_a, \mu^*)}{d(\mu_a, \mu^*)}\right)^2 + 6,$$

where  $\sigma_{a,\star}^2 = \max \{ \text{Var}(\nu_\theta) : \mu_a \leq E(\nu_\theta) \leq \mu^* \}$  and where  $d'(\cdot, \mu^*)$  denotes the derivative of  $d(\cdot, \mu^*)$ .

# Non-parametric setting

- Rewards are only assumed to be bounded (say in  $[0, 1]$ )
- Need for an estimation procedure
  - with non-asymptotic guarantees
  - efficient in the sense of Stein / Bahadur

⇒ Idea 1 : use  $d_{\text{BER}}$  (Hoeffding)

⇒ Idea 2 : Empirical Likelihood [Owen '01]

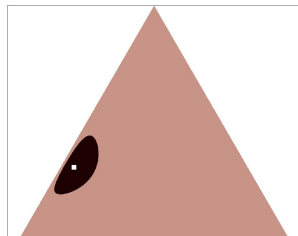
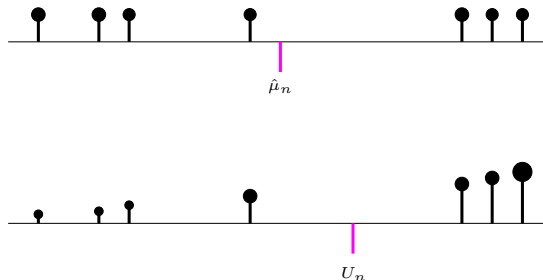
- Bad idea : use Bernstein / Bennett

# Empirical Likelihood

$$U(\hat{\nu}_n, \epsilon) = \sup \left\{ E(\nu') : \nu' \in \mathfrak{M}_1(\text{Supp}(\hat{\nu}_n)) \text{ and } \text{KL}(\hat{\nu}_n, \nu') \leq \epsilon \right\}$$

or, rather, *modified Empirical Likelihood* :

$$U(\hat{\nu}_n, \epsilon) = \sup \left\{ E(\nu') : \nu' \in \mathfrak{M}_1(\text{Supp}(\hat{\nu}_n) \cup \{1\}) \text{ and } \text{KL}(\hat{\nu}_n, \nu') \leq \epsilon \right\}$$



## Coverage properties of the modified EL confidence bound

**Proposition :** Let  $\nu_0 \in \mathfrak{M}_1([0, 1])$  with  $E(\nu_0) \in (0, 1)$  and let  $X_1, \dots, X_n$  be independent random variables with common distribution  $\nu_0 \in \mathfrak{M}_1([0, 1])$ , not necessarily with finite support. Then, for all  $\epsilon > 0$ ,

$$\begin{aligned} \mathbb{P}\{U(\hat{\nu}_n, \epsilon) \leq E(\nu_0)\} &\leq \mathbb{P}\{K_{inf}(\hat{\nu}_n, E(\nu_0)) \geq \epsilon\} \\ &\leq e(n+2) \exp(-n\epsilon) . \end{aligned}$$

**Remark :** For  $\{0, 1\}$ -valued observations, it is readily seen that  $U(\hat{\nu}_n, \epsilon)$  boils down to the upper-confidence bound above.

$\implies$  This proposition is at least not always optimal : the presence of the factor  $n$  in front of the exponential  $\exp(-n\epsilon)$  term is questionable.

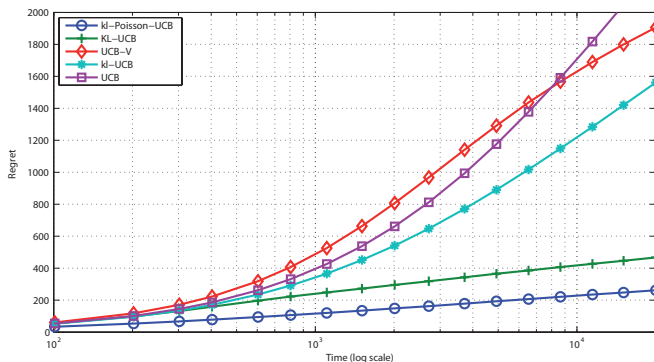
## Regret bound

**Theorem :** Assume that  $\mathcal{F}$  is the set of finitely supported probability distributions over  $\mathcal{S} = [0, 1]$ , that  $\mu_a > 0$  for all arms  $a$  and that  $\mu^* < 1$ . There exists a constant  $M(\nu_a, \mu^*) > 0$  only depending on  $\nu_a$  and  $\mu^*$  such that, with the choice  $f(t) = \log(t) + \log(\log(t))$  for  $t \geq 2$ , for all  $T \geq 3$  :

$$\begin{aligned} \mathbb{E}[N_a(T)] &\leq \frac{\log(T)}{K_{inf}(\nu_a, \mu^*)} + \frac{36}{(\mu^*)^4} (\log(T))^{4/5} \log(\log(T)) \\ &\quad + \left( \frac{72}{(\mu^*)^4} + \frac{2\mu^*}{(1 - \mu^*) K_{inf}(\nu_a, \mu^*)^2} \right) (\log(T))^{4/5} \\ &\quad + \frac{(1 - \mu^*)^2 M(\nu_a, \mu^*)}{2(\mu^*)^2} (\log(T))^{2/5} \\ &\quad + \frac{\log(\log(T))}{K_{inf}(\nu_a, \mu^*)} + \frac{2\mu^*}{(1 - \mu^*) K_{inf}(\nu_a, \mu^*)^2} + 4. \end{aligned}$$

## Example : truncated Poisson rewards

- for each arm  $1 \leq a \leq 6$  is associated with  $\nu_a$ , a Poisson distribution with expectation  $(2 + a)/4$ , truncated at 10.
- $N = 10,000$  Monte-Carlo replications on an horizon of  $T = 20,000$  steps.





## Example : truncated Exponential rewards

- exponential rewards with respective parameters  $1/5$ ,  $1/4$ ,  $1/3$ ,  $1/2$  and  $1$ , truncated at  $x_{\max} = 10$ ;
- kl-UCB uses the divergence  $d(x, y) = x/y - 1 - \log(x/y)$  prescribed for genuine exponential distributions, but it ignores the fact that the rewards are truncated.

