

# 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 8th, 2014

# Outline

- 1 Two Problems of Dynamic Resource Allocation
- 2 The Bandit Model
  - Lower Bound for the Regret
  - Optimistic Algorithms
  - The General UCB Algorithm
  - Non-parametric setting : Empirical Likelihood
- 3 Optimal Discovery with Expert Advice
  - The Good-UCB algorithm
  - Optimality results
- 4 Conclusion and perspectives

# Clinical Trials

Imagine you are a doctor :

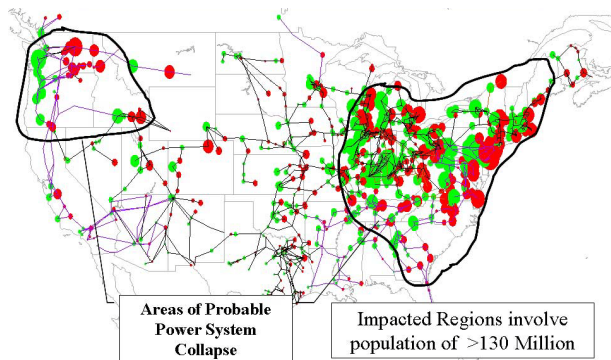
- patients visit you *one after another* for a given disease
- you prescribe one of the (say) *5 treatments* available
- the treatments are *not equally efficient*
- you do not know which one is the best, you *observe the effect* of the prescribed treatment on each patient

⇒ **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*

# Detection of Anomaly in Electrical Systems

**Power system  
security  
assessment**



By Mark MacAlester, Federal Emergency Management Agency [Public domain], via Wikimedia Commons

**Identifying contingencies/scenarios** that could lead to unacceptable operating conditions (dangerous contingencies) if no preventive actions were taken.

# Outline

- 1 Two Problems of Dynamic Resource Allocation
- 2 The Bandit Model
  - Lower Bound for the Regret
  - Optimistic Algorithms
  - The General UCB Algorithm
  - Non-parametric setting : Empirical Likelihood
- 3 Optimal Discovery with Expert Advice
  - The Good-UCB algorithm
  - Optimality results
- 4 Conclusion and perspectives

# 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

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

# Real challenges

- Randomized clinical trials
  - original motivation since the 1930's
  - dynamic strategies can save resources
- Recommender systems :
  - advertisement
  - website optimization
  - news, blog posts, ...
- Computer experiments
  - large systems can be simulated in order to optimize some criterion over a set of parameters
  - but the simulation cost may be high, so that only few choices are possible for the parameters
- Games and planning (tree-structured options)

## 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  $\nu_a \in \mathcal{F}_a$

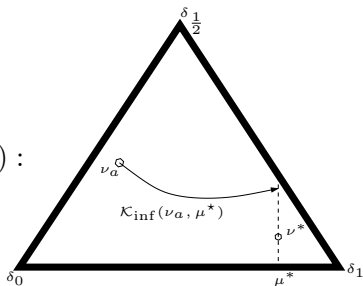
Theorem [Burnetas and Katehakis, '96]

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

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



# Outline

- 1 Two Problems of Dynamic Resource Allocation
- 2 The Bandit Model
  - Lower Bound for the Regret
  - Optimistic Algorithms
  - The General UCB Algorithm
  - Non-parametric setting : Empirical Likelihood
- 3 Optimal Discovery with Expert Advice
  - The Good-UCB algorithm
  - Optimality results
- 4 Conclusion and perspectives

# Optimism in the Face of Uncertainty

**Optimism** is a 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

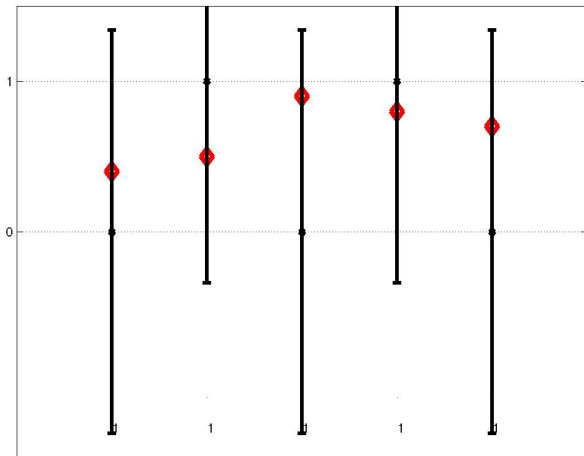
## UCB [Lai&Robins '85; Agrawal '95; Auer&al '02]

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

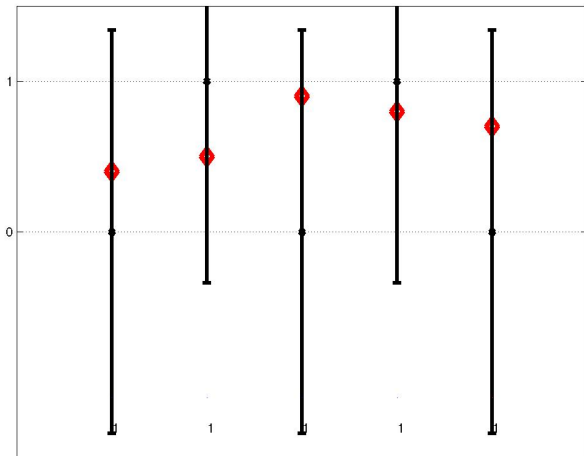
$$\underbrace{\frac{S_a(t)}{N_a(t)}}_{\text{estimated reward}} + \underbrace{\sqrt{\frac{\log(t)}{2N_a(t)}}}_{\text{exploration bonus}}$$

- 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

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

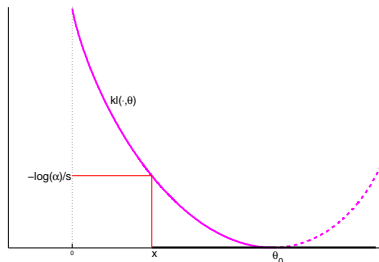
    pick an arm  $A_{t+1} \in \arg \max_{a \in \{1, \dots, K\}} U_a(t)$

**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  $\alpha = \exp(-s \text{kl}(x, \theta_0))$  :

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

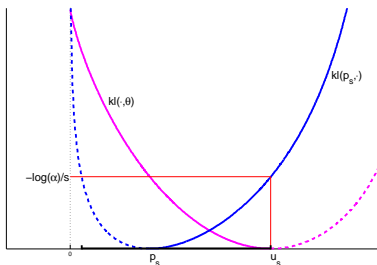
$\implies$  upper confidence bound for  $p$  at risk  $\alpha$  :

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

# 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

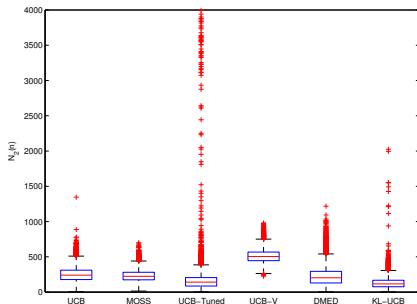
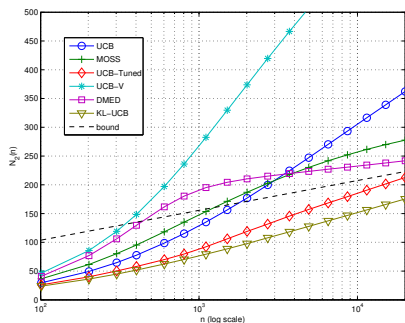
**Theorem :** Assume that all arms belong to a canonical, regular, exponential family  $\mathcal{F} = \{\nu_\theta : \theta \in \Theta\}$  of probability distributions indexed by its natural parameter space  $\Theta \subseteq \mathbb{R}$ . Then, with the choice  $f(t) = \log(t) + 3 \log \log(t)$  for  $t \geq 3$ , 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  $d(\cdot, \mu^*)$ .

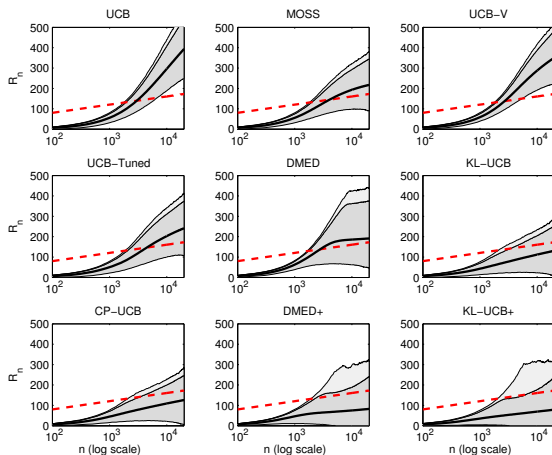


# Results : Two-Arm Scenario



**FIGURE:** Performance of various algorithms when  $\theta = (0.9, 0.8)$ . Left : average number of draws of the sub-optimal arm as a function of time. Right : box-and-whiskers plot for the number of draws of the sub-optimal arm at time  $T = 5,000$ . Results based on 50,000 independent replications

# Results : Ten-Arm Scenario with Low Rewards



**FIGURE:** Average regret as a function of time when  $\theta = (0.1, 0.05, 0.05, 0.05, 0.02, 0.02, 0.02, 0.01, 0.01, 0.01)$ . Red line : Lai & Robbins lower bound ; thick line : average regret ; shaded regions : central 99% region and upper 99.95% quantile

# Outline

- 1 Two Problems of Dynamic Resource Allocation
- 2 The Bandit Model
  - Lower Bound for the Regret
  - Optimistic Algorithms
  - The General UCB Algorithm
  - Non-parametric setting : Empirical Likelihood
- 3 Optimal Discovery with Expert Advice
  - The Good-UCB algorithm
  - Optimality results
- 4 Conclusion and perspectives

# 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

## First idea : use $d_{\text{BER}}$

Idea : rescale to  $[0, 1]$ , and take the divergence  $d_{\text{BER}}$ .

→ because Bernoulli distributions **maximize deviations among bounded variables with given expectation** :

### Lemma (Hoeffding '63)

Let  $X$  denote a random variable such that  $0 \leq X \leq 1$  and denote by  $\mu = \mathbb{E}[X]$  its mean. Then, for any  $\lambda \in \mathbb{R}$ ,

$$E [\exp(\lambda X)] \leq 1 - \mu + \mu \exp(\lambda) .$$

This fact is well-known for the variance, but also true for all exponential moments and thus for Cramer-type deviation bounds

# Regret Bound for kl-UCB

## Theorem

With the divergence  $d_{\text{BER}}$ , for all  $T > 3$ ,

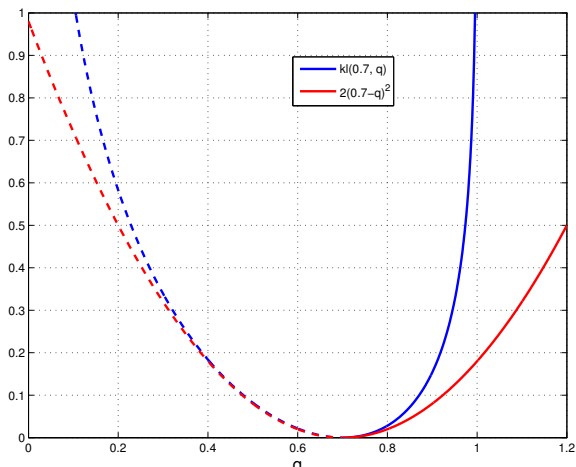
$$\mathbb{E}[N_a(T)] \leq \frac{\log(T)}{d_{\text{BER}}(\mu_a, \mu^*)} + \frac{\sqrt{2\pi} \log\left(\frac{\mu^*(1-\mu_a)}{\mu_a(1-\mu^*)}\right)}{(d_{\text{BER}}(\mu_a, \mu^*))^{3/2}} \sqrt{\log(T) + 3 \log(\log(T))} \\ + \left(4e + \frac{3}{d_{\text{BER}}(\mu_a, \mu^*)}\right) \log(\log(T)) + \frac{2 \left(\log\left(\frac{\mu^*(1-\mu_a)}{\mu_a(1-\mu^*)}\right)\right)^2}{(d_{\text{BER}}(\mu_a, \mu^*))^2} + 6.$$

- kl-UCB satisfies an **improved logarithmic finite-time regret bound**
- Besides, it is **asymptotically optimal in the Bernoulli case**

## Comparison to UCB

KL-UCB addresses **exactly the same problem** as UCB, with the same generality, but it has always a **smaller regret** as can be seen from Pinsker's inequality

$$d_{\text{BER}}(\mu_1, \mu_2) \geq 2(\mu_1 - \mu_2)^2$$

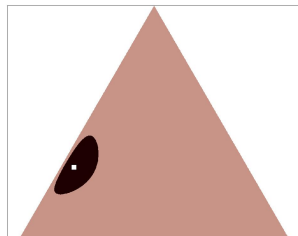
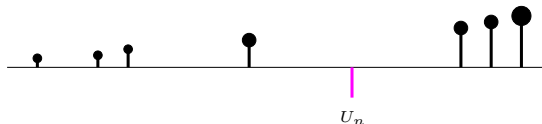
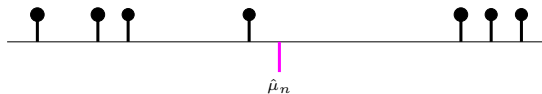


## Idea 2 : Empirical Likelihood

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

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.

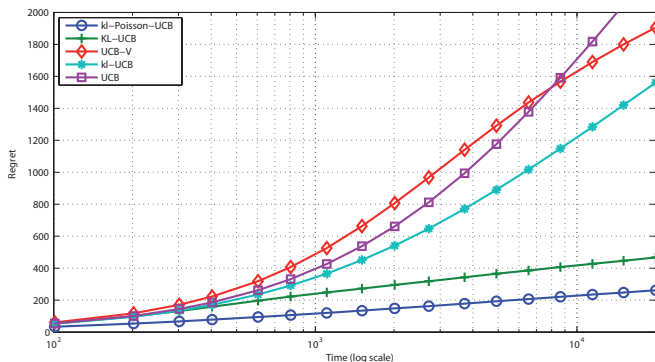
## Regret bound

**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  $T \geq 3$  :

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

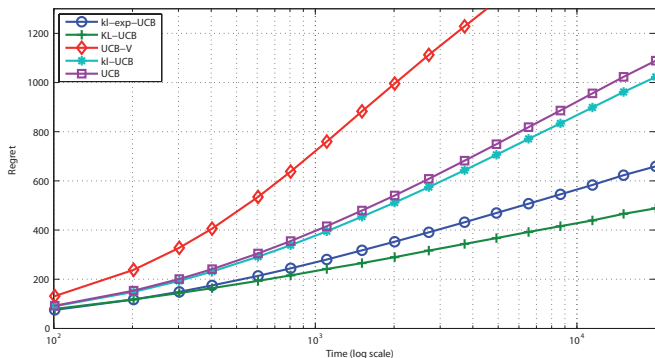
## Example : truncated Poisson rewards

- 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.

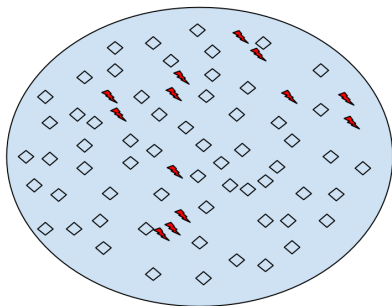


# Outline

- 1 Two Problems of Dynamic Resource Allocation
- 2 The Bandit Model
  - Lower Bound for the Regret
  - Optimistic Algorithms
  - The General UCB Algorithm
  - Non-parametric setting : Empirical Likelihood
- 3 Optimal Discovery with Expert Advice
  - The Good-UCB algorithm
  - Optimality results
- 4 Conclusion and perspectives

# The model

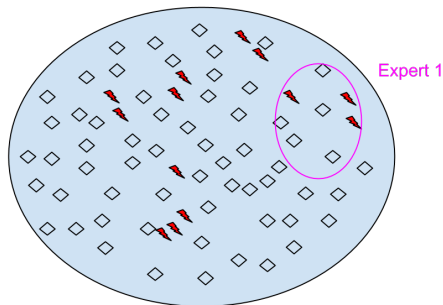
- Subset  $A \subset \mathcal{X}$  of important items
- $|\mathcal{X}| \gg 1$ ,  $|A| \ll |\mathcal{X}|$
- Access to  $\mathcal{X}$  only by probabilistic experts  $(P_i)_{1 \leq i \leq K}$  : sequential independent draws



**Goal : discover rapidly the elements of  $A$**

# The model

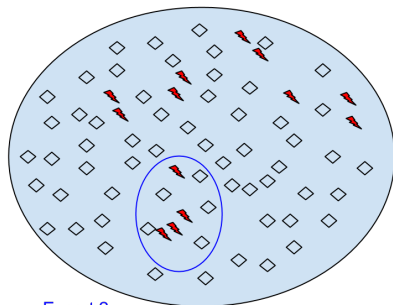
- Subset  $A \subset \mathcal{X}$  of important items
- $|\mathcal{X}| \gg 1$ ,  $|A| \ll |\mathcal{X}|$
- Access to  $\mathcal{X}$  only by probabilistic experts  $(P_i)_{1 \leq i \leq K}$  : sequential independent draws



**Goal : discover rapidly the elements of  $A$**

# The model

- Subset  $A \subset \mathcal{X}$  of important items
- $|\mathcal{X}| \gg 1$ ,  $|A| \ll |\mathcal{X}|$
- Access to  $\mathcal{X}$  only by probabilistic experts  $(P_i)_{1 \leq i \leq K}$  : sequential independent draws

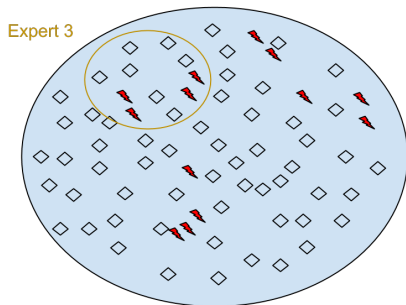


**Goal : discover rapidly the elements of  $A$**



# The model

- Subset  $A \subset \mathcal{X}$  of important items
- $|\mathcal{X}| \gg 1$ ,  $|A| \ll |\mathcal{X}|$
- Access to  $\mathcal{X}$  only by probabilistic experts  $(P_i)_{1 \leq i \leq K}$  : sequential independent draws



**Goal : discover rapidly the elements of  $A$**

## Goal

At each time step  $t = 1, 2, \dots$  :

- pick an index  $I_t = \pi_t(I_1, Y_1, \dots, I_{s-1}, Y_{s-1}) \in \{1, \dots, K\}$  according to past observations
- observe  $Y_t = X_{I_t, n_{I_t, t}} \sim P_{I_t}$ , where

$$n_{i,t} = \sum_{s \leq t} \mathbb{1}\{I_s = i\}$$

**Goal** : design the strategy  $\pi = (\pi_t)_t$  so as to **maximize the number of important items found** after  $t$  requests

$$F^\pi(t) = \left| A \cap \{Y_1, \dots, Y_t\} \right|$$

**Assumption** : non-intersecting supports

$$A \cap \text{supp}(P_i) \cap \text{supp}(P_j) = \emptyset \text{ for } i \neq j$$

# Is it a Bandit Problem ?

It looks like a bandit problem. . .

- sequential choices among  $K$  options
- want to maximize cumulative rewards
- exploration vs exploitation dilemma

# Is it a Bandit Problem ?

It looks like a bandit problem. . .

- sequential choices among  $K$  options
- want to maximize cumulative rewards
- exploration vs exploitation dilemma

. . . but it is **not a bandit problem!**

- rewards are not i.i.d.
- **destructive rewards** : no interest to observe twice the same important item
- all strategies eventually equivalent

## The oracle strategy

**Proposition :** Under the non-intersecting support hypothesis, the greedy oracle strategy

$$I_t^* \in \arg \max_{1 \leq i \leq K} P_i(A \setminus \{Y_1, \dots, Y_t\})$$

is optimal : for every possible strategy  $\pi$ ,  $\mathbb{E}[F^\pi(t)] \leq \mathbb{E}[F^*(t)]$ .

**Remark :** the proposition is false if the supports may intersect

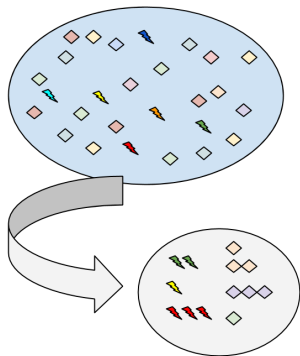
$\implies$  estimate the “**missing mass of important items**” !

## Missing mass estimation

Let us first focus on one expert  $i : P = P_i, X_n = X_{i,n}$

$X_1, \dots, X_n$  independent draws of  $P$

$$O_n(x) = \sum_{m=1}^n \mathbb{1}\{X_m = x\}$$



How to 'estimate' the **total mass of the *unseen*** important items

$$R_n = \sum_{x \in A} P(x) \mathbb{1}\{O_n(x) = 0\} ?$$

## The Good-Turing Estimator

Idea : use the **hapaxes** = items seen only once (linguistic)

$$\hat{R}_n = \frac{U_n}{n}, \quad \text{where } U_n = \sum_{x \in A} \mathbb{1}\{O_n(x) = 1\}$$

**Lemma [Good '53]** : For every distribution  $P$ ,

$$0 \leq \mathbb{E}[\hat{R}_n] - \mathbb{E}[R_n] \leq \frac{1}{n}$$

**Proposition** : With probability at least  $1 - \delta$  for every  $P$ ,

$$\hat{R}_n - \frac{1}{n} - (1 + \sqrt{2})\sqrt{\frac{\log(4/\delta)}{n}} \leq R_n \leq \hat{R}_n + (1 + \sqrt{2})\sqrt{\frac{\log(4/\delta)}{n}}$$

See [McAllester and Schapire '00, McAllester and Ortiz '03] :

- deviations of  $\hat{R}_n$  : McDiarmid's inequality
- deviations of  $R_n$  : negative association

# The Good-UCB algorithm

Estimator of the missing important mass for expert  $i$  :

$$\hat{R}_{i,n_i,t-1} = \frac{1}{n_{i,t-1}} \sum_{x \in A} \mathbb{1} \left\{ \sum_{s=1}^{n_{i,t-1}} \mathbb{1}\{X_{i,s} = x\} = 1 \right.$$

$$\left. \text{and } \sum_{j=1}^K \sum_{s=1}^{n_{j,t-1}} \mathbb{1}\{X_{j,s} = x\} = 1 \right\}$$

**Good-UCB algorithm :**

- 1: For  $1 \leq t \leq K$  choose  $I_t = t$ .
- 2: **for**  $t \geq K + 1$  **do**
- 3:   Choose  $I_t = \arg \max_{1 \leq i \leq K} \left\{ \hat{R}_{i,n_i,t-1} + C \sqrt{\frac{\log(4t)}{n_{i,t-1}}} \right\}$
- 4:   Observe  $Y_t$  distributed as  $P_{I_t}$
- 5:   Update the missing mass estimates accordingly
- 6: **end for**



# Classical analysis

**Theorem** : For any  $t \geq 1$ , under the non-intersecting support assumption, Good-UCB (with constant  $C = (1 + \sqrt{2})\sqrt{3}$ ) satisfies

$$\mathbb{E} [F^*(t) - F^{UCB}(t)] \leq 17\sqrt{Kt \log(t)} + 20\sqrt{Kt} + K + K \log(t/K)$$

Remark : Usual result for bandit problem, but not-so-simple analysis

## Sketch of proof

- 1 On a set  $\tilde{\Omega}$  of probability at least  $1 - \sqrt{\frac{K}{t}}$ , the “confidence intervals” hold true simultaneously all  $u \geq \sqrt{Kt}$
- 2 Let  $\bar{I}_u = \arg \max_{1 \leq i \leq K} R_{i, n_{i, u-1}}$ . On  $\tilde{\Omega}$ ,

$$R_{I_u, n_{I_u, u-1}} \geq R_{\bar{I}_u, n_{\bar{I}_u, u-1}} - \frac{1}{n_{I_u, u-1}} - 2(1 + \sqrt{2}) \sqrt{\frac{3 \log(4u)}{n_{I_u, u-1}}}$$

- 3 But one shows that  $\mathbb{E}F^*(t) \leq \sum_{u=1}^t \mathbb{E}R_{\bar{I}_u, n_{\bar{I}_u, u-1}}^\pi$
- 4 Thus

$$\begin{aligned} & \mathbb{E} [F^*(t) - F^{UCB}(t)] \\ & \leq \sqrt{Kt} + \mathbb{E} \left[ \sum_{u=1}^t \frac{1}{n_{I_u, u-1}} + 2(1 + \sqrt{2}) \sqrt{\frac{3 \log(4t)}{n_{I_u, u-1}}} \right] \\ & \leq \sqrt{Kt} + K + K \log(t/K) + 4(1 + \sqrt{2}) \sqrt{3Kt \log(4t)} \end{aligned}$$

# Experiment : restoring property

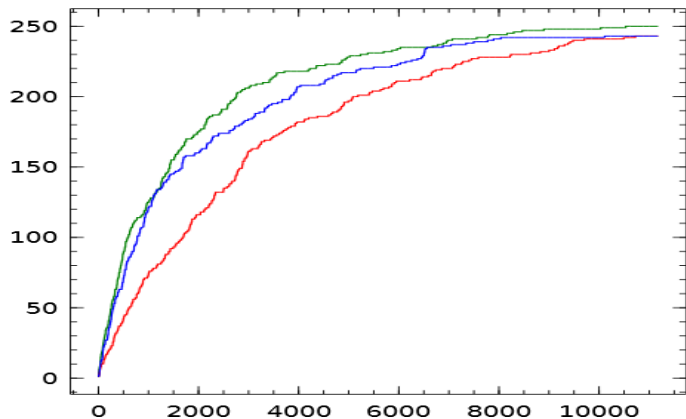


FIGURE: green : oracle, blue : Good-UCB, red : uniform sampling

## Another analysis of Good-UCB

For  $\lambda \in (0, 1)$ ,  $T(\lambda)$  = time at which missing mass of important items is smaller than  $\lambda$  on all experts :

$$T(\lambda) = \inf \left\{ t : \forall i \in \{1, \dots, K\}, P_i(A \setminus \{Y_1, \dots, Y_t\}) \leq \lambda \right\}$$

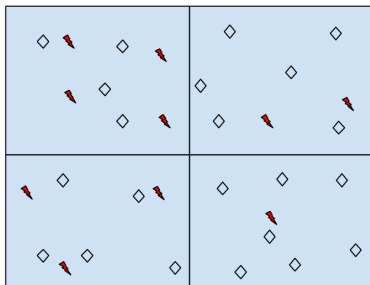
**Theorem :** Let  $c > 0$  and  $S \geq 1$ . Under the non-intersecting support assumption, for Good-UCB with  $C = (1 + \sqrt{2})\sqrt{c+2}$ , with probability at least  $1 - \frac{K}{cS^c}$ , for any  $\lambda \in (0, 1)$ ,

$$T_{UCB}(\lambda) \leq T^* + KS \log(8T^* + 16KS \log(KS)),$$

$$\text{where } T^* = T^* \left( \lambda - \frac{3}{S} - 2(1 + \sqrt{2})\sqrt{\frac{c+2}{S}} \right)$$

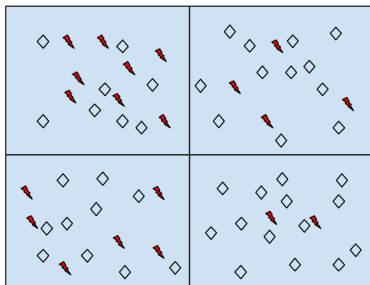
# The macroscopic limit

- Restricted framework :  $P_i = \mathcal{U}\{1, \dots, N\}$
- $N \rightarrow \infty$
- $|A \cap \text{supp}(P_i)|/N \rightarrow q_i \in (0, 1)$ ,  $q = \sum_i q_i$



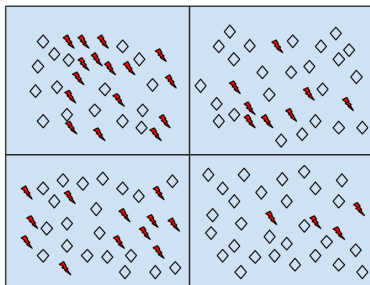
# The macroscopic limit

- Restricted framework :  $P_i = \mathcal{U}\{1, \dots, N\}$
- $N \rightarrow \infty$
- $|A \cap \text{supp}(P_i)|/N \rightarrow q_i \in (0, 1)$ ,  $q = \sum_i q_i$



# The macroscopic limit

- Restricted framework :  $P_i = \mathcal{U}\{1, \dots, N\}$
- $N \rightarrow \infty$
- $|A \cap \text{supp}(P_i)|/N \rightarrow q_i \in (0, 1)$ ,  $q = \sum_i q_i$



# The Oracle behaviour

The limiting discovery process of the Oracle strategy is *deterministic*

**Proposition :** For every  $\lambda \in (0, q_1)$ , for every sequence  $(\lambda^N)_N$  converging to  $\lambda$  as  $N$  goes to infinity, almost surely

$$\lim_{N \rightarrow \infty} \frac{T_*^N(\lambda^N)}{N} = \sum_i \left( \log \frac{q_i}{\lambda} \right)_+$$



## Oracle vs. uniform sampling

**Oracle :** The proportion of important items not found after  $Nt$  draws tends to

$$q - F^*(t) = I(t) \underline{q}_{I(t)} \exp(-t/I(t)) \leq K \underline{q}_K \exp(-t/K)$$

with  $\underline{q}_K = \left( \prod_{i=1}^K q_i \right)^{1/K}$  the geometric mean of the  $(q_i)_i$ .

**Uniform :** The proportion of important items not found after  $Nt$  draws tends to  $K \bar{q}_K \exp(-t/K)$

$\implies$  Asymptotic ratio of efficiency

$$\rho(q) = \frac{\bar{q}_K}{\underline{q}_K} = \frac{\frac{1}{K} \sum_{i=1}^k q_i}{\left( \prod_{i=1}^k q_i \right)^{1/K}} \geq 1$$

larger if the  $(q_i)_i$  are unbalanced

## Macroscopic optimality

**Theorem :** Take  $C = (1 + \sqrt{2})\sqrt{c + 2}$  with  $c > 3/2$  in the Good-UCB algorithm.

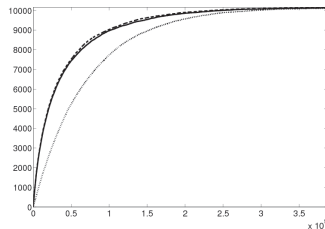
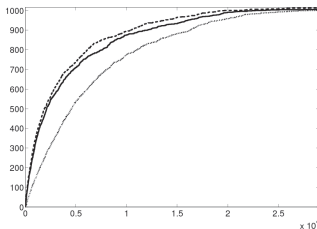
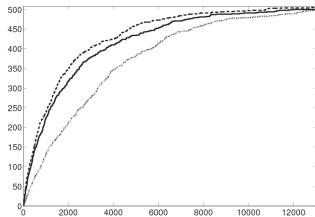
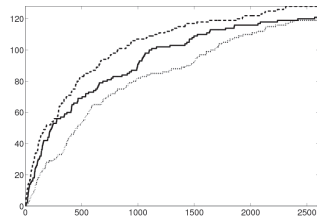
- For every sequence  $(\lambda^N)_N$  converging to  $\lambda$  as  $N$  goes to infinity, almost surely

$$\limsup_{N \rightarrow +\infty} \frac{T_{UCB}^N(\lambda^N)}{N} \leq \sum_i \left( \log \frac{q_i}{\lambda} \right)_+$$

- The proportion of items found after  $Nt$  steps  $F^{GUCB}$  converges *uniformly* to  $F^*$  as  $N$  goes to infinity

## Experiment

Number of items found by Good-UCB (solid), the OCL (dashed), and uniform sampling (dotted) as a function of time for sizes  $N = 128$ ,  $N = 500$ ,  $N = 1000$  and  $N = 10000$  in a 7-experts setting.



# Outline

- 1 Two Problems of Dynamic Resource Allocation
- 2 The Bandit Model
  - Lower Bound for the Regret
  - Optimistic Algorithms
  - The General UCB Algorithm
  - Non-parametric setting : Empirical Likelihood
- 3 Optimal Discovery with Expert Advice
  - The Good-UCB algorithm
  - Optimality results
- 4 Conclusion and perspectives

## For True Bandit Problems :

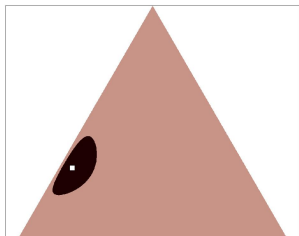
- Use kl-UCB (not UCB), or KL-UCB if speed is not a problem
- To do : improve on the deviation bounds
- address more general non-parametric families of distributions

## Otherwise :

- Ideas can be adapted to specific needs
- For the optimal discovery with probabilistic expert advice, we give a standard regret analysis under the only assumption that the supports of the experts are non-overlapping
- We propose a different optimality result, which permits a macroscopic analysis in the uniform case
- Another interesting limit to consider is when the number of important items to find is fixed, but the total number of items tends to infinity (Poisson regime)
- Then, the behavior of the algorithm is not very good : need tighter deviation bounds

For model-based Reinforcement Learning in Markov Decision Processes, see :

[Filippi *et al.*, *Optimism in Reinforcement Learning and Kullback-Leibler Divergence*, Allerton Conference, 2010]



Thank you for your attention !