

Dynamic resource allocation

Bandit problems and extensions

Aurélien Garivier

Institut de Mathématiques de Toulouse

Séminaire MIA-T, (INRA) le 19 septembre 2014

Roadmap

- 1 The Bandit Model
- 2 Lower Bound for the Regret
- 3 Optimistic Algorithms
- 4 An Optimistic Algorithm based on Kullback-Leibler Divergence
- 5 Parametric setting: the kl-UCB Algorithm
- 6 Non-parametric setting and Empirical Likelihood
- 7 Extensions

Dynamic resource allocation

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*

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

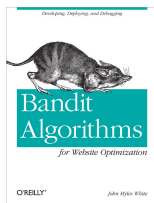
Example Bernoulli rewards: $\theta \in [0, 1]^K$, $\nu_a = \mathcal{B}(\theta_a)$

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

Roadmap

- 1 The Bandit Model
- 2 Lower Bound for the Regret
- 3 Optimistic Algorithms
- 4 An Optimistic Algorithm based on Kullback-Leibler Divergence
- 5 Parametric setting: the kl-UCB Algorithm
- 6 Non-parametric setting and Empirical Likelihood
- 7 Extensions

Asymptotically Optimal Strategies

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

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

- There are efficient strategies and we consider the **best achievable asymptotic performance among 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 π is an 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:

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

The Bound of Burnetas and Katehakis

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

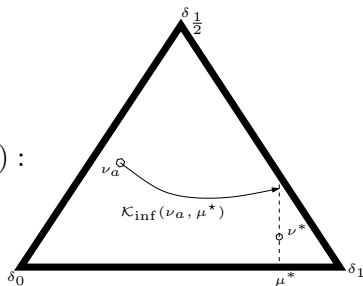
Theorem [Burnetas and Katehakis, '96]

If π 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

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



Intuition

- First assume that μ^* is known and that T is fixed
- How many draws n_a of ν_a are necessary to know that $\mu_a < \mu^*$ with probability at least $1 - 1/T$?
- Test: $H_0 : \mu_a = \mu^*$ against $H_1 : \nu = \nu_a$
- Stein's Lemma: if the first type error $\alpha_{n_a} \leq 1/T$, then

$$\beta_{n_a} \gtrsim \exp(-n_a K_{inf}(\nu_a, \mu^*))$$

\Rightarrow it can be smaller than $1/T$ if

$$n_a \geq \frac{\log(T)}{K_{inf}(\nu_a, \mu^*)}$$

- How to do as well without knowing μ^* and T in advance?
Not asymptotically?

Roadmap

- 1 The Bandit Model
- 2 Lower Bound for the Regret
- 3 Optimistic Algorithms**
- 4 An Optimistic Algorithm based on Kullback-Leibler Divergence
- 5 Parametric setting: the kl-UCB Algorithm
- 6 Non-parametric setting and Empirical Likelihood
- 7 Extensions

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

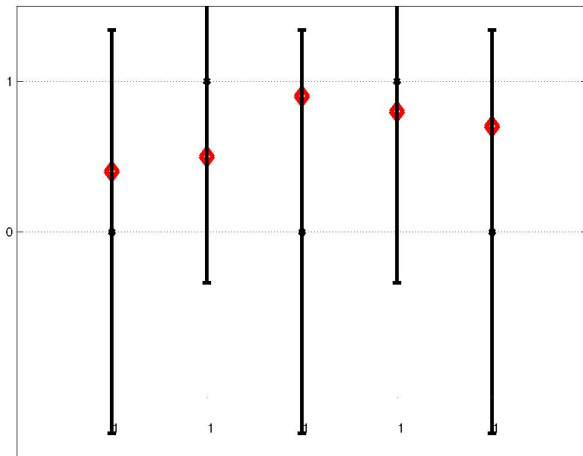
UCB [Lai&Robins '85; Agrawal '95; Auer&al '02]

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

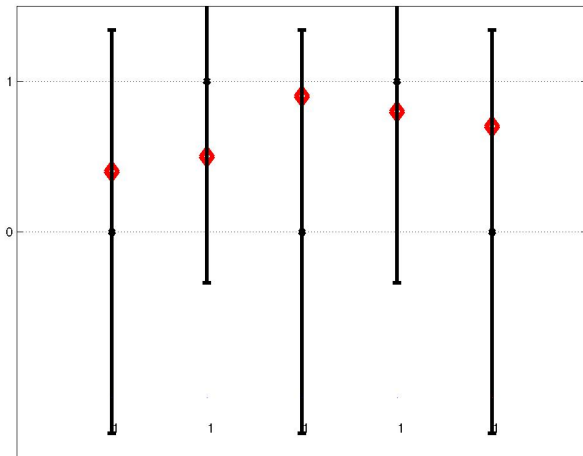
$$\underbrace{\frac{S_a(t)}{N_a(t)}}_{\text{estimated reward}} + \underbrace{\sqrt{\frac{\log(t)}{2N_a(t)}}}_{\text{exploration bonus}}$$

- 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

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

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])

Roadmap

- 1 The Bandit Model
- 2 Lower Bound for the Regret
- 3 Optimistic Algorithms
- 4 An Optimistic Algorithm based on Kullback-Leibler Divergence**
- 5 Parametric setting: the kl-UCB Algorithm
- 6 Non-parametric setting and Empirical Likelihood
- 7 Extensions

The KL-UCB algorithm [Cappé, G. & al '13]

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

Roadmap

- 1 The Bandit Model
- 2 Lower Bound for the Regret
- 3 Optimistic Algorithms
- 4 An Optimistic Algorithm based on Kullback-Leibler Divergence
- 5 Parametric setting: the kl-UCB Algorithm**
- 6 Non-parametric setting and Empirical Likelihood
- 7 Extensions

Exponential Family Rewards

- 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.

Parametric version: 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

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

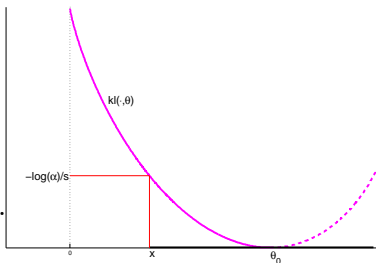
 pick an arm $A_{t+1} \in \arg \max_{a \in \{1, \dots, K\}} U_a(t)$

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 by Chernoff's inequality

$$\mathbb{P}_{\theta_0}(\hat{p}_s \leq x) \leq \exp(-sd_{\text{BER}}(x, \theta_0)).$$



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

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

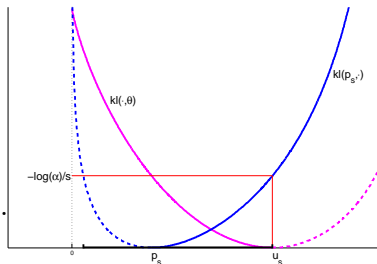
\implies Upper Confidence Bound for p at risk α :

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

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 by Chernoff's inequality

$$\mathbb{P}_{\theta_0}(\hat{p}_s \leq x) \leq \exp(-sd_{\text{BER}}(x, \theta_0)).$$



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

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

\implies Upper Confidence Bound for p at risk α :

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

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

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

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

Results: Two-Arm Scenario

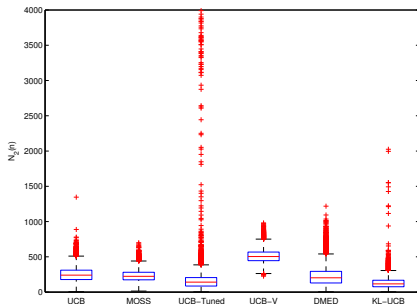
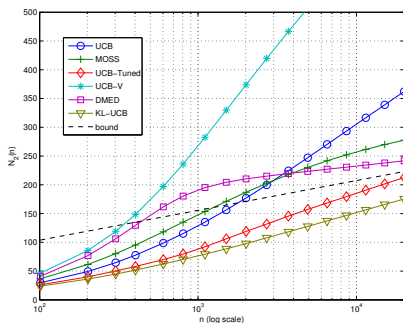


Figure: Performance of various algorithms when $\theta = (0.9, 0.8)$. Left: average number of draws of the sub-optimal arm as a function of time. Right: box-and-whiskers plot for the number of draws of the sub-optimal arm at time $T = 5,000$. Results based on 50,000 independent replications

Results: Ten-Arm Scenario with Low Rewards

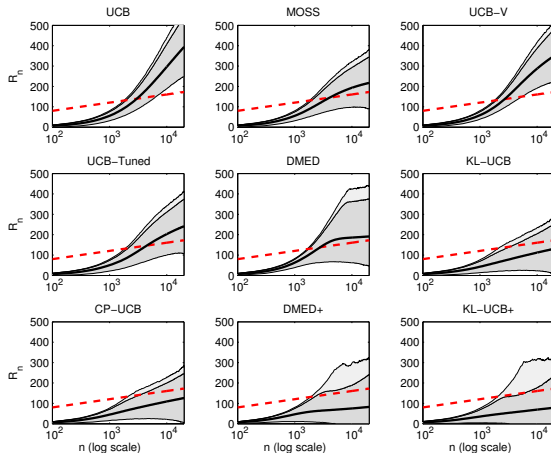


Figure: Average regret as a function of time when $\theta = (0.1, 0.05, 0.05, 0.05, 0.02, 0.02, 0.02, 0.01, 0.01, 0.01)$. Red line: Lai & Robbins lower bound; thick line: average regret; shaded regions: central 99% region and upper 99.95% quantile

Roadmap

- 1 The Bandit Model
- 2 Lower Bound for the Regret
- 3 Optimistic Algorithms
- 4 An Optimistic Algorithm based on Kullback-Leibler Divergence
- 5 Parametric setting: the kl-UCB Algorithm
- 6 Non-parametric setting and Empirical Likelihood**
- 7 Extensions

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_{BER} (Hoeffding)

⇒ Idea 2: Empirical Likelihood [Owen '01]

- Bad idea: use Bernstein / Bennett

First idea: use d_{BER}

Idea: rescale to $[0, 1]$, and take the divergence d_{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_{BER} , for all $T > 3$,

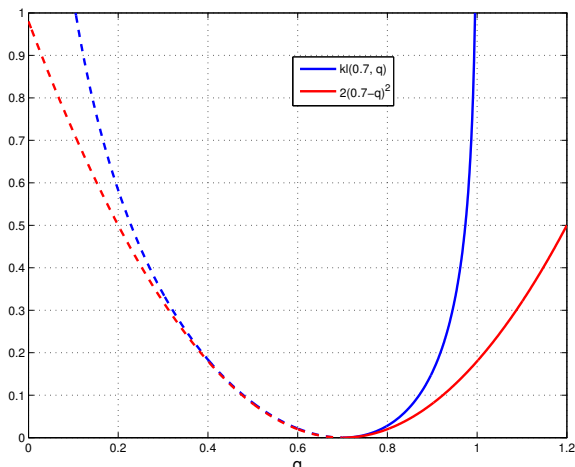
$$\begin{aligned} \mathbb{E}[N_a(T)] \leq & \frac{\log(T)}{d_{\text{BER}}(\mu_a, \mu^*)} + \frac{\sqrt{2\pi} \log\left(\frac{\mu^*(1-\mu_a)}{\mu_a(1-\mu^*)}\right)}{(d_{\text{BER}}(\mu_a, \mu^*))^{3/2}} \sqrt{\log(T) + 3 \log(\log(T))} \\ & + \left(4e + \frac{3}{d_{\text{BER}}(\mu_a, \mu^*)}\right) \log(\log(T))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

- 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

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

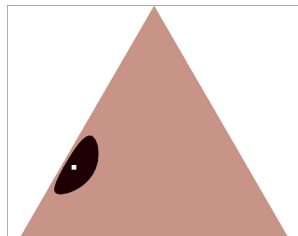
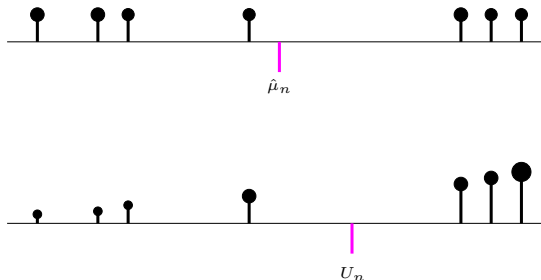


Idea 2: Empirical Likelihood

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

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.

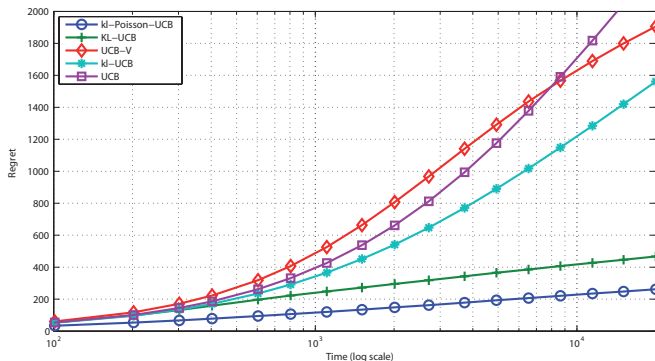
Regret bound

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 $\mu^\star < 1$. There exists a constant $M(\nu_a, \mu^\star) > 0$ only depending on ν_a and μ^\star such that, with the choice $f(t) = \log(t) + \log(\log(t))$ for $t \geq 2$, for all $T \geq 3$:

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

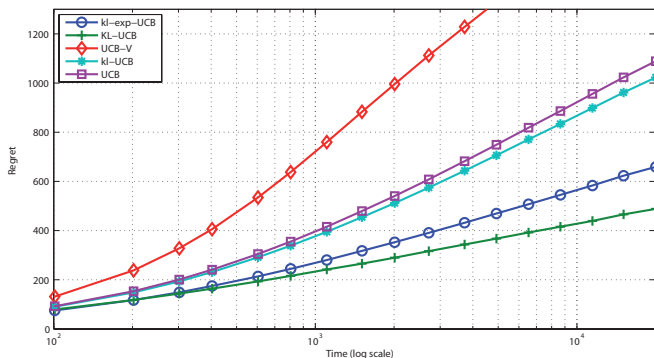
Example: truncated Poisson rewards

- for each arm $1 \leq a \leq 6$ is associated with ν_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.



Take-home message on bandit algorithms

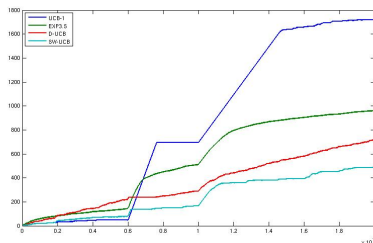
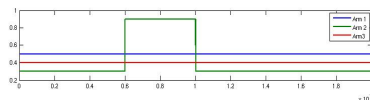
- 1 Use kl-UCB rather than UCB-1 or UCB-2
- 2 Use KL-UCB if speed is not a problem
- 3 todo: improve on the deviation bounds, address general non-parametric families of distributions
- 4 Alternative: Bayesian-flavored methods:
 - Bayes-UCB [Kaufmann, Cappé, G.]
 - Thompson sampling [Kaufmann & al.]

Roadmap

- 1 The Bandit Model
- 2 Lower Bound for the Regret
- 3 Optimistic Algorithms
- 4 An Optimistic Algorithm based on Kullback-Leibler Divergence
- 5 Parametric setting: the kl-UCB Algorithm
- 6 Non-parametric setting and Empirical Likelihood
- 7 Extensions**

Non-stationary Bandits [G. Moulines '11]

- Changepoint : reward distributions change *abruptly*
- Goal : *follow the best arm*
- Application : scanning tunnelling microscope



- Variants D-UCB et SW-UCB including a progressive *discount* of the past
- Bounds $O(\sqrt{n \log n})$ are proved, which is (almost) optimal

(Generalized) Linear Bandits [Filippi, Cappé, G. & Szepesvári '10]

- Bandit with contextual information:

$$\mathbb{E}[X_t|A_t] = \mu(m'_{A_t}\theta_*)$$

where $\theta_* \in \mathbb{R}^d$ is an unknown parameter and $\mu : \mathbb{R} \rightarrow \mathbb{R}$ is a link function

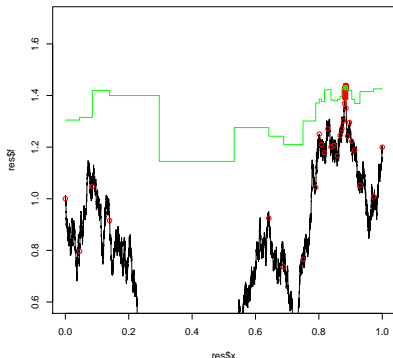
- Example : binary rewards

$$\mu(x) = \frac{\exp(x)}{1 + \exp(x)}$$

- Application : targeted web ads
- GLM-UCB : regret bound depending on dimension d and not on the number of arms

Stochastic Optimization

- Goal : Find the maximum of a function $f : C \subset \mathbb{R}^d \rightarrow \mathbb{R}$ (possibly) observed in noise
- Application : DAS



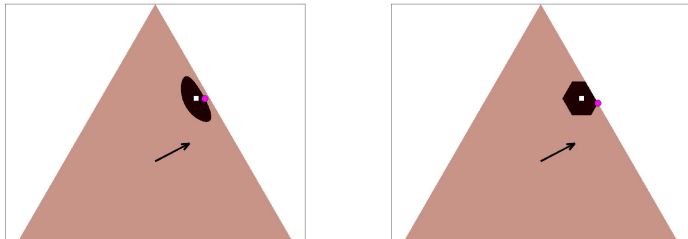
- Model : f is the realization of a Gaussian Process (or has a small norm in some RKHS)
- GP-UCB : evaluate f at the point $x \in C$ where the confidence interval for $f(x)$ has the highest upper-bound

Markov Decision Processes (MDP) [Filippi, Cappé & G. '10]

The system is in state S_t which evolves as a Markov Chain:

$$S_{t+1} \sim P(\cdot; S_t, A_t) \text{ et } R_t = r(S_t, A_t) + \varepsilon_t$$

Optimistic algorithm: search the best transition matrix in a neighborhood of the ML estimate



The use of Kullback-Leibler neighborhoods leads to better performance and has desirable properties.

Optimal Exploration with Probabilistic Expert Advice

Search space : $B \subset \Omega$ discrete set

Probabilistic experts : $P_a \in \mathfrak{M}_1(\Omega)$ for $a \in \mathcal{A}$

Requests : at time t , calling expert A_t yields a realization of $X_t = X_{A_t, t}$ independent with law P_a

Goal : find as many distinct elements of B as possible with few requests :

$$F_n = \text{Card} (B \cap \{X_1, \dots, X_n\})$$

\neq bandit : finding the same element twice is no use !

Oracle : selects the expert with highest 'missing mass'

$$A_{t+1}^* = \arg \max_{a \in \mathcal{A}} P_a (B \setminus \{X_1, \dots, X_t\})$$

Estimating the missing mass

Notation :

- $X_t \stackrel{iid}{\sim} P \in \mathfrak{M}_1(\Omega)$, $O_n(\omega) = \sum_{t=1}^n \mathbb{1}\{X_t = \omega\}$
- $Z_n(x) = \mathbb{1}\{O_n(\omega) = 0\}$
- $H_n(\omega) = \mathbb{1}\{O_n(\omega) = 1\}$, $H_n = \sum_{\omega \in B} H_n(\omega)$

Problem : estimate the missing mass

$$R_n = \sum_{\omega \in B} P(\omega) Z_n(\omega)$$

Good-Turing : 'estimator' $\hat{R}_n = H_n/n$ st. $\mathbb{E}[\hat{R}_n - R_n] \in [0, 1/n]$.

Concentration : by McDiarmid's inequality, with probability
 $\geq 1 - \delta$

$$\left| \hat{R}_n - E[\hat{R}_n] \right| \leq \sqrt{\frac{(2/n + p_{\max})^2 n \log(2/\delta)}{2}}$$

The Good-UCB algorithm [Bubeck, Ernst & G.]

Optimistic algorithm based on Good-Turing's estimator :

$$A_{t+1} = \arg \max_{a \in \mathcal{A}} \left\{ \frac{H_a(t)}{N_a(t)} + c \sqrt{\frac{\log(t)}{N_a(t)}} \right\}$$

- $N_a(t)$ = number of draws of P_a up to time t
- $H_a(t)$ = number of elements of B seen exactly once thanks to P_a
- c = tuning parameter

Good-UCB en action

Macroscopic optimality

Hypotheses :

- $\Omega = \mathcal{A} \times \{1, \dots, N\}$
- $\forall a \in \mathcal{A}, \forall j \in \{1, \dots, N\}, P_a(\{(a, j)\}) = 1/N$

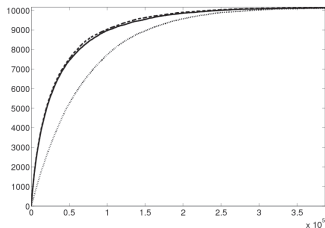
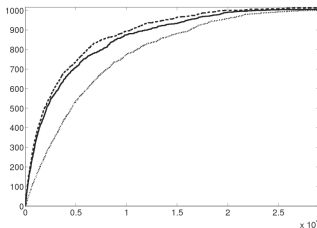
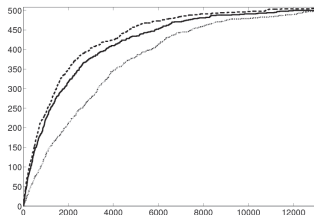
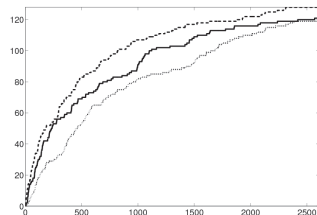
Macroscopic limit :

- $N \rightarrow \infty$
- $\forall a \in \mathcal{A}, \text{Card}(B \cap \{a\} \times \{1, \dots, N\}) / N \rightarrow q_a \in]0, 1[$

Macroscopic optimality

When N goes to infinity, the performance of the Good-UCB algorithm during the discovery $t \mapsto F([Nt])$ uniformly converges to that of the oracle $t \mapsto F^*([Nt])$ on \mathbb{R}^+ .

Simulation



Number of items found by Good-UCB (line), the oracle (bold dashed), and by uniform sampling (light dotted) as a function of time, for sample sizes $N = 128$, $N = 500$, $N = 1000$ et $N = 10000$, in an environment with 7 experts.

Bibliography

- 1 **[Abbasi-Yadkori&al '11]** Yasin Abbasi-Yadkori, Dávid Pál, Csaba Szepesvári: Online Least Squares Estimation with Self-Normalized Processes: An Application to Bandit Problems CoRR abs/1102.2670: (2011)
- 2 **[Agrawal '95]** R. Agrawal. Sample mean based index policies with $O(\log n)$ regret for the multi-armed bandit problem. *Advances in Applied Probability*, 27(4) :1054-1078, 1995.
- 3 **[Audibert&al '09]** J-Y. Audibert, R. Munos, and Cs. Szepesvári. Exploration-exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19), 2009
- 4 **[Auer&al '02]** P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2) :235-256, 2002.
- 5 **[Bubeck, Ernst&G. '11]** S. Bubeck, D. Ernst, A. Garivier, Optimal discovery with probabilistic expert advice: finite time analysis and macroscopic optimality *Journal of Machine Learning Research* vol. 14 Feb. 2013 p.601-623
- 6 **[De La Pena&al '04]** V.H. De La Pena, M.J. Klass, and T.L. Lai. Self-normalized processes : exponential inequalities, moment bounds and iterated logarithm laws. *Annals of Probability*, 32(3) :1902-1933, 2004.
- 7 **[Filippi, Cappé&Garivier '10]** S. Filippi, O. Cappé, and A. Garivier. Optimism in reinforcement learning and Kullback-Leibler divergence. In *Allerton Conf. on Communication, Control, and Computing*, Monticello, US, 2010.
- 8 **[Filippi, Cappé, G.& Szepesvari '10]** S. Filippi, O. Cappé, A. Garivier, and C. Szepesvari. Parametric bandits : The generalized linear case. In *Neural Information Processing Systems (NIPS)*, 2010.
- 9 **[G.&Cappé '11]** A. Garivier and O. Cappé. The KL-UCB algorithm for bounded stochastic bandits and beyond. In *23rd Conf. Learning Theory (COLT)*, Budapest, Hungary, 2011.
- 10 **[Cappé,G.&al '13]** O. Cappé, A. Garivier, O-A. Maillard, R. Munos, G. Stoltz, Kullback-Leibler Upper Confidence Bounds for Optimal Sequential Allocation, *Annals of Statistics*, accepted
- 11 **[G.&Leonardi '11]** A. Garivier and F. Leonardi. Context tree selection : A unifying view. *Stochastic Processes and their Applications*, 121(11) :2488-2506, Nov. 2011.
- 12 **[G.&Moulines '11]** A. Garivier and E. Moulines. On upper-confidence bound policies for non-stationary bandit problems. In *Algorithmic Learning Theory (ALT)*, volume 6925 of *Lecture Notes in Computer Science*, 2011.
- 13 **[Lai&Robins '85]** T.L. Lai and H. Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, 6(1) :4-22, 1985.