

# Multi-armed bandit models: a tutorial

Emilie Kaufmann



CERMICS seminar,  
March 30th, 2016

# Multi-Armed Bandit model: general setting

$K$  arms:

for  $a \in \{1, \dots, K\}$ ,  $(X_{a,t})_{t \in \mathbb{N}}$  is a **stochastic process**.

(**unknown** distributions)

**Bandit game:** at each round  $t$ , an agent

- chooses an arm  $A_t \in \{1, \dots, K\}$
- receives a reward  $X_t = X_{A_t,t}$

**Goal:** Build a **sequential strategy**

$$A_t = F_t(A_1, X_1, \dots, A_{t-1}, X_{t-1})$$

maximizing

$$\mathbb{E} \left[ \sum_{t=1}^{\infty} \alpha_t X_t \right],$$

where  $(\alpha_t)_{t \in \mathbb{N}}$  is a discount sequence. [Berry and Fristedt, 1985]

# Multi-Armed Bandit model: the i.i.d. case

$K$  independent arms:

for  $a \in \{1, \dots, K\}$ ,  $(X_{a,t})_{t \in \mathbb{N}}$  is i.i.d.  $\sim \nu_a$

(unknown distributions)

**Bandit game:** at each round  $t$ , an agent

- chooses an arm  $A_t \in \{1, \dots, K\}$
- receives a reward  $X_t \sim \nu_{A_t}$

**Goal:** Build a sequential strategy

$$A_t = F_t(A_1, X_1, \dots, A_{t-1}, X_{t-1})$$

maximizing

Finite-horizon MAB	Discounted MAB
$\mathbb{E} \left[ \sum_{t=1}^T X_t \right]$	$\mathbb{E} \left[ \sum_{t=1}^{\infty} \alpha^{t-1} X_t \right]$

# Why MABs?



$V_1$



$V_2$



$V_3$



$V_4$



$V_5$

**Goal:** maximize ones' gains in a casino ?  
(HOPELESS)

# Why MABs? Real motivations

## Clinical trials:



$B(\mu_1)$



$B(\mu_2)$



$B(\mu_3)$



$B(\mu_4)$



$B(\mu_5)$

- choose a **treatment**  $A_t$  for patient  $t$
- observe a **response**  $X_t \in \{0, 1\} : \mathbb{P}(X_t = 1) = \mu_{A_t}$
- Goal: maximize the number of patient healed

## Recommendation tasks:



$\nu_1$



$\nu_2$



$\nu_3$



$\nu_4$



$\nu_5$

- recommend a **movie**  $A_t$  for visitor  $t$
- observe a **rating**  $X_t \sim \nu_{A_t}$  (e.g.  $X_t \in \{1, \dots, 5\}$ )
- Goal: maximize the sum of ratings

# Bernoulli bandit models

$K$  independent **arms**:

for  $a \in \{1, \dots, K\}$ ,  $(X_{a,t})_{t \in \mathbb{N}}$  is **i.i.d**  $\sim \mathcal{B}(\mu_a)$

**Bandit game:** at each round  $t$ , an agent

- chooses an arm  $A_t \in \{1, \dots, K\}$
- receives a reward  $X_t \sim \mathcal{B}(\mu_{A_t}) \in \{0, 1\}$

**Goal:** maximize

Finite-horizon MAB	Discounted MAB
$\mathbb{E} \left[ \sum_{t=1}^T X_t \right]$	$\mathbb{E} \left[ \sum_{t=1}^{\infty} \alpha^{t-1} X_t \right]$

# Bernoulli bandit models

$K$  independent **arms**:

for  $a \in \{1, \dots, K\}$ ,  $(X_{a,t})_{t \in \mathbb{N}}$  is **i.i.d.**  $\sim \mathcal{B}(\mu_a)$

**Bandit game:** at each round  $t$ , an agent

- chooses an arm  $A_t \in \{1, \dots, K\}$
- receives a reward  $X_t \sim \mathcal{B}(\mu_{A_t}) \in \{0, 1\}$

**Goal:** maximize

Finite-horizon MAB	Discounted MAB
$\mathbb{E} \left[ \sum_{t=1}^T X_t \right]$	$\mathbb{E} \left[ \sum_{t=1}^{\infty} \alpha^{t-1} X_t \right]$

<b>Frequentist model</b>	<b>Bayesian model</b>
$\mu_1, \dots, \mu_K$ <b>unknown parameters</b>	$\mu_1, \dots, \mu_K$ drawn from a <b>prior distribution</b> : $\mu_a \sim \pi_a$
arm $a$ : $(X_{a,t})_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{B}(\mu_a)$	arm $a$ : $(X_{a,t})_t   \boldsymbol{\mu} \stackrel{\text{i.i.d.}}{\sim} \mathcal{B}(\mu_a)$

- 1 Bayesian bandits: a planning problem
- 2 Frequentist bandits: asymptotically optimal algorithms
- 3 Non stochastic bandits: minimax algorithms

- 1 Bayesian bandits: a planning problem
- 2 Frequentist bandits: asymptotically optimal algorithms
- 3 Non stochastic bandits: minimax algorithms

# A Markov Decision Process

Bandit model  $(\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_K))$

- prior distribution:  $\mu_a \stackrel{\text{i.i.d}}{\sim} \mathcal{U}([0, 1])$
- posterior distribution:  $\pi_a^t := \mathcal{L}(\mu_a | X_1, \dots, X_t)$

$$\pi_a^t = \text{Beta}\left(\underbrace{S_a(t)}_{\# \text{ones}} + 1, \underbrace{N_a(t) - S_a(t)}_{\# \text{zeros}} + 1\right)$$

$S_a(t)$ : sum of the rewards gathered from  $a$  up to time  $t$

$N_a(t)$ : number of draws of arm  $a$  up to time  $t$



State  $\Pi^t = (\pi_a^t)_{a=1}^K$  that evolves in a MDP.

**An example of transition:**

$$\begin{pmatrix} 1 & 2 \\ 5 & 1 \\ 0 & 2 \end{pmatrix} \xrightarrow{A_t=2} \begin{pmatrix} 1 & 2 \\ 6 & 1 \\ 0 & 2 \end{pmatrix} \text{ if } X_t = 1$$

**Solving a planning problem:** there exists an exact solution to

- The finite-horizon MAB:
- The discounted MAB:

$$\arg \max_{(A_t)} \mathbb{E}_{\mu \sim \pi} \left[ \sum_{t=1}^T X_t \right]$$

$$\arg \max_{(A_t)} \mathbb{E}_{\mu \sim \pi} \left[ \sum_{t=1}^{\infty} \alpha^{t-1} X_t \right]$$

Optimal policy = solution to dynamic programming equations.

**An example of transition:**

$$\begin{pmatrix} 1 & 2 \\ 5 & 1 \\ 0 & 2 \end{pmatrix} \xrightarrow{A_t=2} \begin{pmatrix} 1 & 2 \\ 6 & 1 \\ 0 & 2 \end{pmatrix} \text{ if } X_t = 1$$

**Solving a planning problem:** there exists an exact solution to

- The finite-horizon MAB:
- The discounted MAB:

$$\arg \max_{(A_t)} \mathbb{E}_{\mu \sim \pi} \left[ \sum_{t=1}^T X_t \right]$$

$$\arg \max_{(A_t)} \mathbb{E}_{\mu \sim \pi} \left[ \sum_{t=1}^{\infty} \alpha^{t-1} X_t \right]$$

Optimal policy = solution to dynamic programming equations.

**Problem:** The state space is too large !

# A reduction of the dimension

[Gittins 79]: the solution of the **discounted** MAB reduces to an index policy:

$$A_{t+1} = \operatorname{argmax}_{a=1\dots K} G_{\alpha}(\pi_a^t).$$

- **The Gittins indices:**

$$G_{\alpha}(p) = \sup_{\substack{\text{stopping} \\ \text{times } \tau > 0}} \frac{\mathbb{E}_{\substack{Y_t \text{ i.i.d } \mathcal{B}(\mu) \\ \mu \sim p}} \left[ \sum_{t=1}^{\tau} \alpha^{t-1} Y_t \right]}{\mathbb{E}_{\substack{Y_t \text{ i.i.d } \mathcal{B}(\mu) \\ \mu \sim p}} \left[ \sum_{t=1}^{\tau} \alpha^{t-1} \right]}$$

“instantaneous rewards when committing to arm  $\mu \sim p$ , when rewards are discounted by  $\alpha$ ”

## An alternative formulation:

$$G_\alpha(p) = \inf\{\lambda \in \mathbb{R} : V_\alpha^*(p, \lambda) = 0\},$$

with

$$V_\alpha^*(p, \lambda) = \sup_{\substack{\text{stopping} \\ \text{times } \tau > 0}} \mathbb{E}_{\substack{Y_t \text{ i.i.d. } \mathcal{B}(\mu) \\ \mu \sim p}} \left[ \sum_{t=1}^{\tau} \alpha^{t-1} (Y_t - \lambda) \right].$$

“price worth paying for committing to arm  $\mu \sim p$  when rewards are discounted by  $\alpha$ ”

# Gittins indices for finite horizon

**The Finite-Horizon Gittins indices:** depend on the **remaining time to play**  $r$

$$G(p, r) = \inf\{\lambda \in \mathbb{R} : V_r^*(p, \lambda) = 0\},$$

with

$$V_r^*(p, \lambda) = \sup_{\substack{\text{stopping times} \\ 0 < \tau \leq r}} \mathbb{E}_{\substack{Y_t \text{ i.i.d. } \mathcal{B}(\mu) \\ \mu \sim p}} \left[ \sum_{t=1}^{\tau} (Y_t - \lambda) \right].$$

“price worth paying for playing arm  $\mu \sim p$  for at most  $r$  rounds”

## The Finite-Horizon Gittins algorithm

$$A_{t+1} = \operatorname{argmax}_{a=1 \dots K} G(\pi_a^t, T - t)$$

does NOT coincide with the optimal solution [Berry and Fristedt 85]... but is conjectured to be a good approximation !

- 1 Bayesian bandits: a planning problem
- 2 Frequentist bandits: asymptotically optimal algorithms
- 3 Non stochastic bandits: minimax algorithms

# Regret minimization

$\mu = (\mu_1, \dots, \mu_K)$  unknown parameters,  $\mu^* = \max_a \mu_a$ .

- The **regret** of a strategy  $\mathcal{A} = (A_t)$  is defined as

$$R_\mu(\mathcal{A}, T) = \mathbb{E}_\mu \left[ \mu^* T - \sum_{t=1}^T X_t \right]$$

and can be rewritten

$$R_\mu(\mathcal{A}, T) = \sum_{a=1}^K (\mu^* - \mu_a) \mathbb{E}_\mu [N_a(T)].$$

$N_a(t)$  : number of draws of arm  $a$  up to time  $t$

Maximizing rewards  $\Leftrightarrow$  Minimizing regret

**Goal:** Design strategies that have small regret for all  $\mu$ .

# Optimal algorithms for regret minimization

All the arms should be drawn infinitely often !

- [Lai and Robbins, 1985]: a uniformly efficient strategy ( $\forall \mu, \forall \alpha \in ]0, 1[$ ,  $R_\mu(\mathcal{A}, T) = o(T^\alpha)$ ) satisfies

$$\mu_a < \mu^* \Rightarrow \liminf_{T \rightarrow \infty} \frac{\mathbb{E}_\mu[N_a(T)]}{\log T} \geq \frac{1}{d(\mu_a, \mu^*)},$$

where

$$\begin{aligned} d(\mu, \mu') &= \text{KL}(\mathcal{B}(\mu), \mathcal{B}(\mu')) \\ &= \mu \log \frac{\mu}{\mu'} + (1 - \mu) \log \frac{1 - \mu}{1 - \mu'}. \end{aligned}$$

## Definition

A bandit algorithm is **asymptotically optimal** if, for every  $\mu$ ,

$$\mu_a < \mu^* \Rightarrow \limsup_{T \rightarrow \infty} \frac{\mathbb{E}_\mu[N_a(T)]}{\log T} \leq \frac{1}{d(\mu_a, \mu^*)}$$

- **Idea 1** : Draw each arm  $T/K$  times

⇒ EXPLORATION

- **Idea 1** : Draw each arm  $T/K$  times

⇒ EXPLORATION

- **Idea 2**: Always choose the **empirical best arm**:

$$A_{t+1} = \operatorname{argmax}_a \hat{\mu}_a(t)$$

⇒ EXPLOITATION

- **Idea 1** : Draw each arm  $T/K$  times

⇒ EXPLORATION

- **Idea 2**: Always choose the empirical best arm:

$$A_{t+1} = \operatorname{argmax}_a \hat{\mu}_a(t)$$

⇒ EXPLOITATION

- **Idea 3** : Draw the arms uniformly during  $T/2$  rounds, then draw the empirical best until the end

⇒ EXPLORATION followed EXPLOITATION

- **Idea 1** : Draw each arm  $T/K$  times

⇒ EXPLORATION

- **Idea 2**: Always choose the **empirical best arm**:

$$A_{t+1} = \operatorname{argmax}_a \hat{\mu}_a(t)$$

⇒ EXPLOITATION

- **Idea 3** : Draw the arms uniformly during  $T/2$  rounds, then draw the empirical best until the end

⇒ EXPLORATION followed EXPLOITATION

**Linear regret...**

# Optimistic algorithms

- For each arm  $a$ , build a confidence interval on  $\mu_a$  :

$$\mu_a \leq \text{UCB}_a(t) \text{ w.h.p}$$



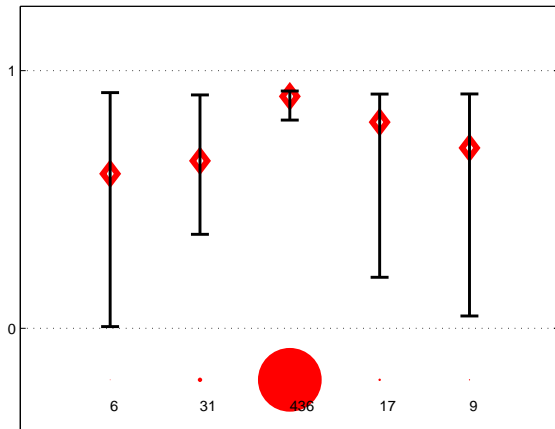
Figure : Confidence intervals on the arms at round  $t$

- Optimism principle:

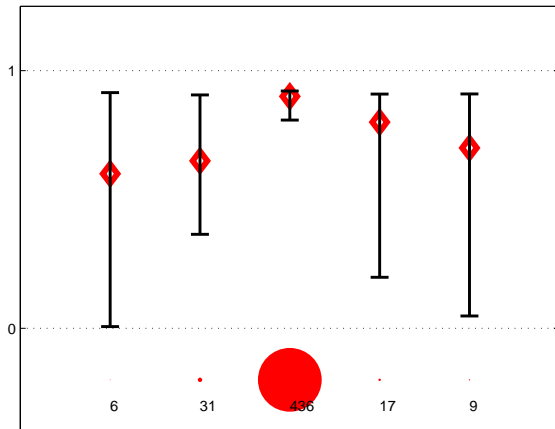
“act as if the best possible model were the true model”

$$A_{t+1} = \arg \max_a \text{UCB}_a(t)$$

# A UCB algorithm in action !



# A UCB algorithm in action !



# The UCB1 algorithm

UCB1 [Auer et al. 02] is based on the index

$$\text{UCB}_a(t) = \hat{\mu}_a(t) + \sqrt{\frac{\alpha \log(t)}{2N_a(t)}}$$

- Hoeffding's inequality:

$$\mathbb{P}\left(\hat{\mu}_{a,s} + \sqrt{\frac{\alpha \log(t)}{2s}} \leq \mu_a\right) \leq \exp\left(-2s \left(\frac{\alpha \log(t)}{2s}\right)\right) = \frac{1}{t^\alpha}.$$

- Union bound:

$$\begin{aligned}\mathbb{P}(\text{UCB}_a(t) \leq \mu_a) &\leq \mathbb{P}\left(\exists s \leq t : \hat{\mu}_{a,s} + \sqrt{\frac{\alpha \log(t)}{2s}} \leq \mu_a\right) \\ &\leq \sum_{s=1}^t \frac{1}{t^\alpha} = \frac{1}{t^{\alpha-1}}.\end{aligned}$$

## Theorem

For every  $\alpha > 2$  and every sub-optimal arm  $a$ , there exists a constant  $C_\alpha > 0$  such that

$$\mathbb{E}_\mu[N_a(T)] \leq \frac{2\alpha}{(\mu^* - \mu_a)^2} \log(T) + C_\alpha.$$

It follows that

$$R_T \leq 2\alpha \left( \sum_{a \neq a^*} \frac{1}{(\mu^* - \mu_a)^2} \right) \log(T) + KC_\alpha.$$

Assume  $\mu^* = \mu_1$  and  $\mu_2 < \mu_1$ .

$$\begin{aligned}
 N_2(T) &= \sum_{t=0}^{T-1} \mathbb{1}_{(A_{t+1}=2)} \\
 &= \sum_{t=0}^{T-1} \mathbb{1}_{(A_{t+1}=2) \cap (\text{UCB}_1(t) \leq \mu_1)} + \sum_{t=0}^{T-1} \mathbb{1}_{(A_{t+1}=2) \cap (\text{UCB}_1(t) > \mu_1)} \\
 &\leq \sum_{t=0}^{T-1} \mathbb{1}_{(\text{UCB}_1(t) \leq \mu_1)} + \sum_{t=0}^{T-1} \mathbb{1}_{(A_{t+1}=2) \cap (\text{UCB}_2(t) > \mu_1)}
 \end{aligned}$$

Assume  $\mu^* = \mu_1$  and  $\mu_2 < \mu_1$ .

$$\begin{aligned}
 N_2(T) &= \sum_{t=0}^{T-1} \mathbb{1}_{(A_{t+1}=2)} \\
 &= \sum_{t=0}^{T-1} \mathbb{1}_{(A_{t+1}=2) \cap (\text{UCB}_1(t) \leq \mu_1)} + \sum_{t=0}^{T-1} \mathbb{1}_{(A_{t+1}=2) \cap (\text{UCB}_1(t) > \mu_1)} \\
 &\leq \sum_{t=0}^{T-1} \mathbb{1}_{(\text{UCB}_1(t) \leq \mu_1)} + \sum_{t=0}^{T-1} \mathbb{1}_{(A_{t+1}=2) \cap (\text{UCB}_2(t) > \mu_1)} \\
 \\
 \mathbb{E}[N_2(T)] &\leq \underbrace{\sum_{t=0}^{T-1} \mathbb{P}(\text{UCB}_1(t) \leq \mu_1)}_A + \underbrace{\sum_{t=0}^{T-1} \mathbb{P}(A_{t+1} = 2, \text{UCB}_2(t) > \mu_1)}_B
 \end{aligned}$$

$$\mathbb{E}[N_2(T)] \leq \underbrace{\sum_{t=0}^{T-1} \mathbb{P}(\text{UCB}_1(t) \leq \mu_1)}_A + \underbrace{\sum_{t=0}^{T-1} \mathbb{P}(A_{t+1} = 2, \text{UCB}_2(t) > \mu_1)}_B$$

- **Term A:** if  $\alpha > 2$ ,

$$\begin{aligned} \sum_{t=0}^{T-1} \mathbb{P}(\text{UCB}_1(t) \leq \mu_1) &\leq 1 + \sum_{t=1}^{T-1} \frac{1}{t^{\alpha-1}} \\ &\leq 1 + \zeta(\alpha - 1) := C_\alpha/2. \end{aligned}$$

- Term B:

$$\begin{aligned}
 (B) &= \sum_{t=0}^{T-1} \mathbb{P}(A_{t+1} = 2, \text{UCB}_2(t) > \mu_1) \\
 &\leq \sum_{t=0}^{T-1} \mathbb{P}(A_{t+1} = 2, \text{UCB}_2(t) > \mu_1, \text{LCB}_2(t) \leq \mu_2) + C_\alpha/2
 \end{aligned}$$

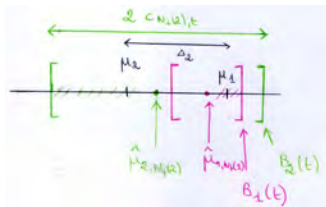
with

$$\text{LCB}_2(t) = \hat{\mu}_2(t) - \sqrt{\frac{\alpha \log t}{2N_2(t)}}$$

$$(\text{LCB}_2(t) < \mu_2 < \mu_1 \leq \text{UCB}_2(t))$$

$$\Rightarrow (\mu_1 - \mu_2) \leq 2\sqrt{\frac{\alpha \log(T)}{2N_2(t)}}$$

$$\Rightarrow N_2(t) \leq \frac{2\alpha}{(\mu_1 - \mu_2)^2} \log(T)$$



- **Term B:** (continued)

$$\begin{aligned}
 (B) &\leq \sum_{t=0}^{T-1} \mathbb{P}(A_{t+1} = 2, \text{UCB}_2(t) > \mu_1, \text{LCB}_2(t) \leq \mu_2) + C_\alpha/2 \\
 &\leq \sum_{t=0}^{T-1} \mathbb{P}\left(A_{t+1} = 2, N_2(t) \leq \frac{2\alpha}{(\mu_1 - \mu_2)^2} \log(T)\right) + C_\alpha/2 \\
 &\leq \frac{2\alpha}{(\mu_1 - \mu_2)^2} \log(T) + C_\alpha/2
 \end{aligned}$$

- **Conclusion:**

$$\mathbb{E}[N_2(T)] \leq \frac{2\alpha}{(\mu_1 - \mu_2)^2} \log(T) + C_\alpha.$$

## Theorem

For every  $\alpha > 2$  and every sub-optimal arm  $a$ , there exists a constant  $C_\alpha > 0$  such that

$$\mathbb{E}_\mu[N_a(T)] \leq \frac{2\alpha}{(\mu^* - \mu_a)^2} \log(T) + C_\alpha.$$

## Theorem

For every  $\alpha > 2$  and every sub-optimal arm  $a$ , there exists a constant  $C_\alpha > 0$  such that

$$\mathbb{E}_\mu[N_a(T)] \leq \frac{2\alpha}{(\mu^* - \mu_a)^2} \log(T) + C_\alpha.$$

**Remark:**

$$\frac{2\alpha}{(\mu^* - \mu_a)^2} > 4\alpha \frac{1}{d(\mu_a, \mu^*)}$$

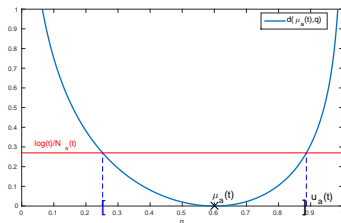
(UCB1 not asymptotically optimal)

# The KL-UCB algorithm

- A UCB-type algorithm:  $A_{t+1} = \arg \max_a u_a(t)$
- ... associated to **the right upper confidence bounds**:

$$u_a(t) = \max \left\{ q \geq \hat{\mu}_a(t) : d(\hat{\mu}_a(t), q) \leq \frac{\log(t)}{N_a(t)} \right\},$$

$\hat{\mu}_a(t)$ : empirical mean of rewards from arm  $a$  up to time  $t$ .

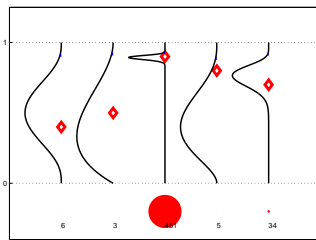
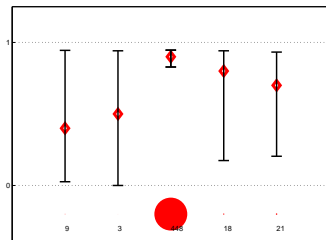


[Cappé et al. 13]: KL-UCB satisfies

$$\mathbb{E}_{\mu} [N_a(T)] \leq \frac{1}{d(\mu_a, \mu^*)} \log T + O(\sqrt{\log(T)}).$$

# Bayesian algorithms for regret minimization?

Algorithms based on **Bayesian tools**  
can be good to solve (frequentist) **regret minimization**



## Ideas:

- use the Finite-Horizon Gittins
- use posterior quantiles
- use **posterior samples**

# Thompson Sampling

$(\pi_a^t, \dots, \pi_K^t)$  posterior distribution on  $(\mu_1, \dots, \mu_K)$  at round  $t$ .

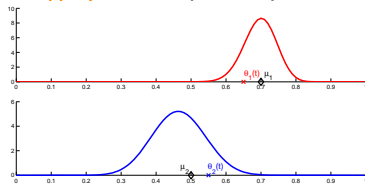
## Algorithm: Thompson Sampling

**Thompson Sampling** is a randomized Bayesian algorithm:

$$\forall a \in \{1..K\}, \theta_a(t) \sim \pi_a^t$$

$$A_{t+1} = \operatorname{argmax}_a \theta_a(t)$$

“Draw each arm according to its posterior probability of being optimal” [Thompson 1933]



Thompson Sampling is asymptotically optimal.  
[K.,Korda,Munos 2012]

# Bayesian algorithms in contextual linear bandit models

At time  $t$ , a set of 'contexts'  $\mathcal{D}_t \subset \mathbb{R}^d$  is revealed.

= characteristics of the items to recommend

## The model:

- if the context  $x_t \in \mathcal{D}_t$  is selected
- a reward  $r_t = x_t^T \theta + \epsilon_t$  is received

$\theta \in \mathbb{R}^d$  = underlying preference vector

## A Bayesian model: (with Gaussian prior)

$$r_t = x_t^T \theta + \epsilon_t, \quad \theta \sim \mathcal{N}(0, \kappa^2 I_d), \quad \epsilon_t \sim \mathcal{N}(0, \sigma^2).$$

Explicit posterior:  $p(\theta | x_1, r_1, \dots, x_t, r_t) = \mathcal{N}(\hat{\theta}(t), \Sigma_t)$ .

## Thompson Sampling:

$$\tilde{\theta}(t) \sim \mathcal{N}(\hat{\theta}(t), \Sigma_t), \quad \text{and} \quad x_{t+1} = \underset{x \in \mathcal{D}_{t+1}}{\operatorname{argmax}} x^T \tilde{\theta}(t).$$

- 1 Bayesian bandits: a planning problem
- 2 Frequentist bandits: asymptotically optimal algorithms
- 3 Non stochastic bandits: minimax algorithms

# Minimax regret

- In stochastic bandits, we exhibited algorithm such that

$$\forall \mu, \mathbb{E}_{\mu}[N_a(T)] \leq \log T / d(\mu_a, \mu^*) + o(\log T).$$

Their regret satisfy

$$R_{\mu}(\mathcal{A}, T) = \sum_{a=2}^K \min \left[ \underbrace{\frac{(\mu^* - \mu_a)}{d(\mu_a, \mu^*)} \log(T)}_{\text{large when } \mu_a \rightarrow \mu^*}, \underbrace{(\mu^* - \mu_a) T}_{\text{small when } \mu_a \rightarrow \mu^*} \right] + o(\log(T)).$$

→ There exist some constant  $C$  such that

$$\forall \mu, R_{\mu}(\mathcal{A}, T) \leq \underbrace{C \sqrt{KT \log(T)}}_{\text{problem-independent bound}}.$$

## Minimax rate of the regret

$$\inf_{\mathcal{A}} \sup_{\mu} R_{\mu}(\mathcal{A}, T) = O\left(\sqrt{KT}\right)$$

**A new bandit game:** at round  $t$

- the player chooses arm  $A_t$
- simultaneously, an **adversary** chooses the vector of rewards

$$(x_{t,1}, \dots, x_{t,K})$$

- the player receives the reward  $x_t = x_{A_t,t}$

**Goal:** maximize rewards, or minimize **regret**

$$R(\mathcal{A}, T) = \max_a \mathbb{E} \left[ \sum_{t=1}^T x_{a,t} \right] - \mathbb{E} \left[ \sum_{t=1}^T x_t \right].$$

# Full information: Exponential Weighted Forecaster

**The full-information game:** at round  $t$

- the player chooses arm  $A_t$
- simultaneously, an **adversary** chooses the vector of rewards

$$(x_{t,1}, \dots, x_{t,K})$$

- the player receives the reward  $x_t = x_{A_t,t}$
- and he observes the reward vector  $(x_{t,1}, \dots, x_{t,K})$

The EWF algorithm [Littellstone, Warmuth 1994]

With  $\hat{p}_t$  the probability distribution

$$\hat{p}_{a,t} \propto e^{\eta(\sum_{s=1}^{t-1} x_{a,s})}$$

at round  $t$ , choose

$$A_t \sim \hat{p}_t$$

We don't have access to the  $(x_{a,t})$  for all  $a...$

$$\hat{x}_{a,t} = \frac{x_{a,t}}{\hat{p}_{a,t}} \mathbb{1}_{(A_t=a)}$$

satisfies  $\mathbb{E}[\hat{x}_{a,t}] = x_{a,t}$ .

The EXP3 strategy [Auer et al. 2003]

With  $\hat{p}_t$  the probability distribution

$$\hat{p}_{a,t} \propto e^{\eta(\sum_{s=1}^{t-1} \hat{x}_{a,s})}$$

at round  $t$ , choose

$$A_t \sim \hat{p}_t$$

## The EXP3 strategy [Auer et al. 2003]

With  $\hat{p}_t$  the probability distribution

$$\hat{p}_{a,t} \propto e^{\eta(\sum_{s=1}^{t-1} \hat{x}_{a,s})}$$

at round  $t$ , choose

$$A_t \sim \hat{p}_t$$

[Bubeck and Cesa-Bianchi 12] EXP3 with

$$\eta = \sqrt{\frac{\log(K)}{KT}}$$

satisfies

$$R(\text{EXP3}, T) \leq \sqrt{2 \log K} \sqrt{KT}$$

### Remarks:

- almost the same guarantees for  $\eta_t = \sqrt{\frac{\log(K)}{Kt}}$
- extra exploration is needed to have high probability results

Under different assumptions, different types of strategies to achieve an exploration-exploitation tradeoff in bandit models:

## **Index policies:**

- Gittins indices
- UCB-type algorithms

## **Randomized algorithms:**

- Thompson Sampling
- Exponential weights

More complex bandit models not covered today:  
restless bandits, contextual bandits, combinatorial bandits...

Under different assumptions, different types of strategies to achieve an **exploration-exploitation tradeoff** in bandit models:

## **Index policies:**

- Gittins indices
- UCB-type algorithms

## **Randomized algorithms:**

- Thompson Sampling
- Exponential weights

More complex bandit models not covered today:  
restless bandits, contextual bandits, combinatorial bandits...

# A pure-exploration objective

Regret minimization:

maximize the number of patients healed during the trial



Alternative goal: identify as quickly as possible the best treatment  
(no focus on curing patients during the study)

# A pure-exploration objective

Regret minimization:

maximize the number of patients healed during the trial



Alternative goal: identify as quickly as possible the best treatment  
(no focus on curing patients during the study)

Additionally to the **sampling strategy**  $(A_t)$ , one needs

- a **stopping rule**  $\tau$  (stopping time)
- a **recommendation rule**  $\hat{a}_\tau$

such that, for some risk parameter  $\delta \in ]0, 1[$ ,

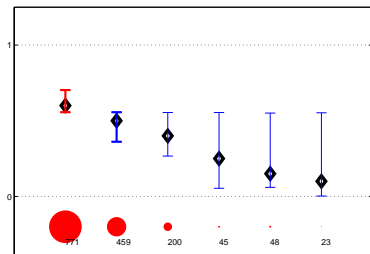
$$\mathbb{P}(\hat{a}_\tau \neq a^*) \leq \delta \quad \text{and} \quad \mathbb{E}[\tau] \text{ is as small as possible.}$$

# An algorithm: KL-LUCB [K., Kalyanakrishnan 13]

An algorithm based on **Upper and Lower** confidence bounds

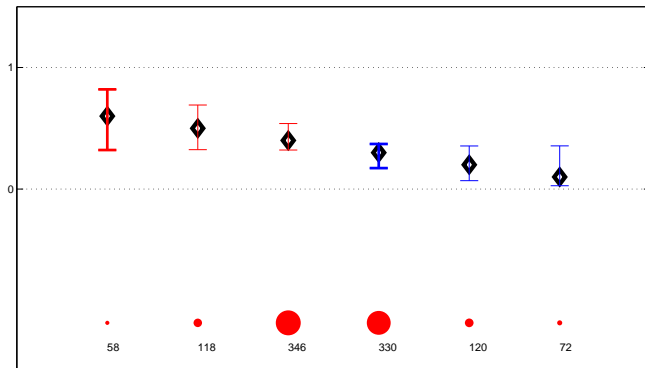
$$u_a(t) = \max \{q : N_a(t)d(\hat{\mu}_a(t), q) \leq \log(Kt/\delta)\}$$

$$l_a(t) = \min \{q : N_a(t)d(\hat{\mu}_a(t), q) \leq \log(Kt/\delta)\}$$



- sampling rule:  $A_{t+1} = \operatorname{argmax}_a \hat{\mu}_a(t)$ ,  $B_{t+1} = \operatorname{argmax}_{b \neq A_{t+1}} u_b(t)$
- stopping rule:  $\tau = \inf \{t \in \mathbb{N} : \ell_{A_t}(t) > u_{B_t}(t)\}$
- recommendation rule:  $\hat{a}_\tau = \operatorname{argmax}_{a=1, \dots, K} \hat{\mu}_a(\tau)$

# KL-LUCB for finding the $m$ best arms



# The complexity of best-arm identification

## Theorem [K. and Garivier, 16]

For any  $\delta$ -PAC algorithm,

$$\mathbb{E}_{\mu}[\tau] \geq T^*(\mu) \log\left(\frac{1}{2.4\delta}\right),$$

where

$$T^*(\mu)^{-1} = \sup_{w \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\mu)} \left( \sum_{a=1}^K w_a d(\mu_a, \lambda_a) \right).$$

→ an optimal strategy satisfies  $\frac{\mathbb{E}_{\mu}[N_a(\tau)]}{\mathbb{E}_{\mu}[\tau]} \simeq w_a^*(\mu)$  with

$$w^*(\mu) = \operatorname{argmax}_{w \in \Sigma_K} \inf_{\lambda \in \text{Alt}(\mu)} \left( \sum_{a=1}^K w_a d(\mu_a, \lambda_a) \right)$$

→ tracking these optimal proportions yield a  $\delta$ -PAC algorithm

$$\text{such that } \limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_{\mu}[\tau_{\delta}]}{\log(1/\delta)} = T^*(\mu).$$

## Bayesian bandits:

- *Bandit Problems. Sequential allocation of experiments.*  
Berry and Fristedt. Chapman and Hall, 1985.
- *Multi-armed bandit allocation indices.*  
Gittins, Glazebrook, Weber. Wiley, 2011.

## Stochastic and non-stochastic bandits:

- *Regret analysis of Stochastic and Non-stochastic Bandit Problems.* Bubeck and Cesa-Bianchi.  
Foundations and Trends in Machine Learning, 2012.
- *Prediction, Learning and Games.*  
Cesa-Bianchi and Lugosi. Cambridge University Press, 2006.

## Best arm identification:

- *Optimal Best Arm Identification with Fixed Confidence.*  
A. Garivier and E. Kaufmann. Preprint, 2016.