



BANDIT PROBLEMS

Part I - Stochastic Bandits (2/2)

RLSS, Lille, July 2019

PART I: Solving the stochastic MAB

PART II: Structured Bandits

PART III: Bandit for Optimization

RECAPS

The Stochastic Multi-Armed Bandit Setup

K arms $\leftrightarrow K$ probability distributions : ν_a has mean μ_a



ν_1



ν_2



ν_3



ν_4



ν_5

At round t , an agent:

- ▶ chooses an arm A_t
- ▶ receives a reward $R_t = X_{A_t,t} \sim \nu_{A_t}$

Sequential sampling strategy (**bandit algorithm**):

$$A_{t+1} = F_t(A_1, R_1, \dots, A_t, R_t).$$

Goal: Maximize $\mathbb{E} \left[\sum_{t=1}^T R_t \right]$

Regret of a bandit algorithm

Bandit instance: $\nu = (\nu_1, \nu_2, \dots, \nu_K)$, mean of arm a : $\mu_a = \mathbb{E}_{X \sim \nu_a}[X]$.

$$\mu_\star = \max_{a \in \{1, \dots, K\}} \mu_a \quad a_\star = \operatorname{argmax}_{a \in \{1, \dots, K\}} \mu_a.$$

Maximizing rewards \leftrightarrow selecting a_\star as much as possible
 \leftrightarrow minimizing the **regret** [Robbins, 52]

$$\mathcal{R}_\nu(\mathcal{A}, T) := \underbrace{T\mu_\star}_{\text{sum of rewards of an oracle strategy always selecting } a_\star} - \underbrace{\mathbb{E} \left[\sum_{t=1}^T R_t \right]}_{\text{sum of rewards of the strategy } \mathcal{A}}$$

What regret rate can we achieve?

- $\rightarrow \mathcal{R}_\nu(\mathcal{A}, T) = C_\nu \log(T)$ problem-dependent regret
- $\rightarrow \mathcal{R}_\nu(\mathcal{A}, T) = C\sqrt{KT}$ problem-independent (worse-case) regret

Regret of a bandit algorithm

Bandit instance: $\nu = (\nu_1, \nu_2, \dots, \nu_K)$, mean of arm a : $\mu_a = \mathbb{E}_{X \sim \nu_a}[X]$.

$$\mu_\star = \max_{a \in \{1, \dots, K\}} \mu_a \quad a_\star = \operatorname{argmax}_{a \in \{1, \dots, K\}} \mu_a.$$

Maximizing rewards \leftrightarrow selecting a_\star as much as possible
 \leftrightarrow minimizing the **regret** [Robbins, 52]

$$\mathcal{R}_\nu(\mathcal{A}, T) := \sum_{a=1}^K \underbrace{(\mu_\star - \mu_a)}_{\Delta_a: \text{sub-optimality gap of arm } a} \times \underbrace{\mathbb{E}_\nu[N_a(T)]}_{\text{expected number of selections of arm } a}$$

What regret rate can we achieve?

- $\rightarrow \mathcal{R}_\nu(\mathcal{A}, T) = C_\nu \log(T)$ problem-dependent regret
- $\rightarrow \mathcal{R}_\nu(\mathcal{A}, T) = C\sqrt{KT}$ problem-independent (worse-case) regret

Performance lower bounds

- Problem-dependent for simple parametric model
(Bernoulli, Gaussian with known variance, Exponential, Poisson...)

Theorem [Lai and Robbins, 1985]

For uniformly efficient algorithms, in a regime of large values of T ,

$$\mathcal{R}_\nu(\mathcal{A}, T) \gtrsim \left(\sum_{a: \mu_a < \mu_*} \frac{\Delta_a}{\text{kl}(\mu_a, \mu_*)} \right) \ln(T).$$

- Problem independent (worse-case)

Theorem [Cesa-Bianchi and Lugosi, 06][Bubeck and Cesa-Bianchi, 12]

Fix $T \in \mathbb{N}$. For every bandit algorithm \mathcal{A} , there exists a stochastic bandit model ν with rewards supported in $[0, 1]$ such that

$$\mathcal{R}_\nu(\mathcal{A}, T) \geq \frac{1}{20} \sqrt{KT}$$

Two naive strategies

► Idea 1 : Uniform Exploration

Draw each arm T/K times

► Idea 2 : Follow The Leader (FTL)

$$A_{t+1} = \operatorname{argmax}_{a \in \{1, \dots, K\}} \hat{\mu}_a(t)$$

where $\hat{\mu}_a(t)$ is an estimate of the unknown mean μ_a .

→ Linear regret!

(Sequential) Explore-Then-Commit

For 2 (Gaussian) arms:

- explore uniformly until the **random time**

$$\tau = \inf \left\{ t \in \mathbb{N} : |\hat{\mu}_1(t) - \hat{\mu}_2(t)| > \sqrt{\frac{8\sigma^2 \ln(T/t)}{t}} \right\}$$

- $\hat{a}_\tau = \operatorname{argmax}_a \hat{\mu}_a(\tau)$ and $(A_{t+1} = \hat{a}_\tau)$ for $t \in \{\tau + 1, \dots, T\}$

Logarithmic regret!

$$\mathcal{R}_\nu(\text{S-ETC}, T) \leq \frac{4\sigma^2}{\Delta} \ln(T\Delta^2) + C\sqrt{\ln(T)}.$$

- ➔ this approach can be generalized to more than 2 arms, but cannot be asymptotically optimal (= match Lai and Robbins lower bound)

The optimism principle

- For each arm a , build a confidence interval on the mean μ_k :

$$\mathcal{I}_a(t) = [\text{LCB}_a(t), \text{UCB}_a(t)]$$

LCB = Lower Confidence Bound

UCB = Upper Confidence Bound

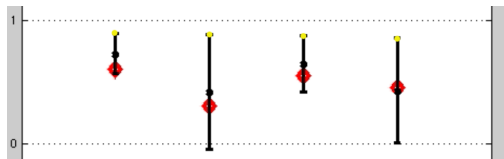


Figure: Confidence intervals on the means after t rounds

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

$$A_{t+1} = \operatorname{argmax}_{a=1,\dots,K} \text{UCB}_a(t).$$

- ▶ UCB for σ^2 -sub Gaussian rewards

$$A_{t+1} = \operatorname{argmax}_{a=1,\dots,K} \hat{\mu}_a(t) + \sqrt{\frac{2\sigma^2 \ln t}{N_a(t)}}$$

- asymptotically optimal for Gaussian distributions, can be used for bounded distribution (with $\sigma^2 = 1/4$).
- $O(\sqrt{KT \ln(T)})$ worse-case regret

- ▶ kl-UCB with divergence $\text{kl}(x, y)$

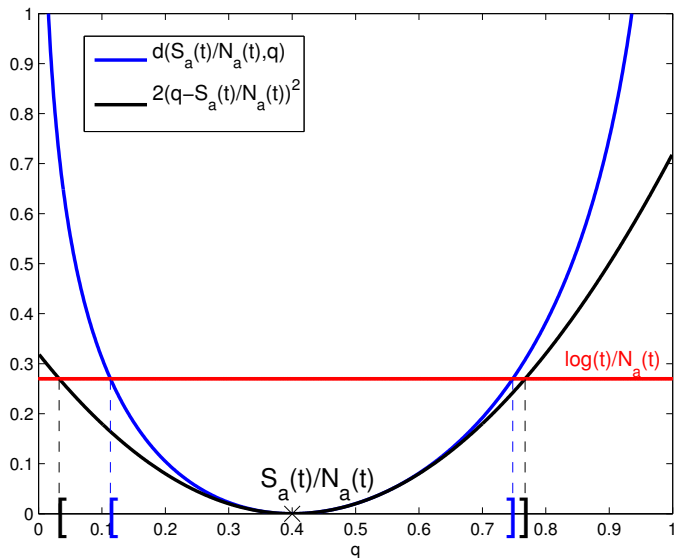
$$A_{t+1} = \underset{a=1, \dots, K}{\operatorname{argmax}} \max \left\{ q : \text{kl}(\hat{\mu}_a(t), q) \leq \frac{\ln(t)}{N_a(t)} \right\}$$

- asymptotically optimal for Bernoulli distribution and can be used for bounded distributions with

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

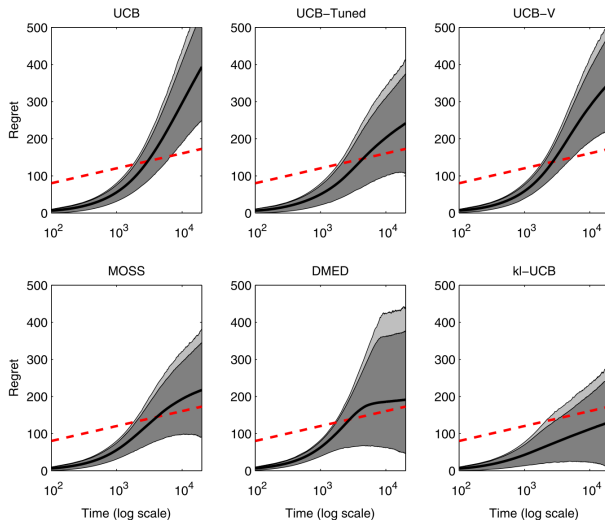
- $O(\sqrt{KT \ln(T)})$ worse-case regret

Comparison of the confidence intervals



UCB versus kl-UCB

$$\mu = [0.1 \ 0.05 \ 0.05 \ 0.05 \ 0.02 \ 0.02 \ 0.02 \ 0.01 \ 0.01 \ 0.01]$$



(Credit: Cappé et al.)

A BAYESIAN LOOK AT THE MULTI-ARMED BANDIT MODEL

Historical perspective

1952 Robbins, formulation of the MAB problem

1985 Lai and Robbins: lower bound, first asymptotically optimal algorithm

1987 Lai, asymptotic regret of kl-UCB

1995 Agrawal, UCB algorithms

1995 Katehakis and Robbins, a UCB algorithm for Gaussian bandits

2002 Auer et al: UCB1 with finite-time regret bound

2009 UCB-V, MOSS...

2011,13 Cappé et al: finite-time regret bound for kl-UCB

Historical perspective

- 1933 Thompson: a Bayesian mechanism for clinical trials
- 1952 Robbins, formulation of the MAB problem
- 1956 Bradt et al, Bellman: optimal solution of a Bayesian MAB problem
- 1979 Gittins: first Bayesian index policy
- 1985 Lai and Robbins: lower bound, first asymptotically optimal algorithm
- 1985 Berry and Fristedt: Bandit Problems, a survey on the Bayesian MAB
- 1987 Lai, asymptotic regret of kl-UCB + study of its Bayesian regret
- 1995 Agrawal, UCB algorithms
- 1995 Katehakis and Robbins, a UCB algorithm for Gaussian bandits
- 2002 Auer et al: UCB1 with finite-time regret bound
- 2009 UCB-V, MOSS...
- 2010 Thompson Sampling is re-discovered
- 2011,13 Cappé et al: finite-time regret bound for kl-UCB
- 2012,13 Thompson Sampling is asymptotically optimal

Frequentist versus Bayesian bandit

$$\nu_{\mu} = (\nu^{\mu_1}, \dots, \nu^{\mu_K}) \in (\mathcal{P})^K.$$

- ▶ Two probabilistic models

Frequentist model	Bayesian model
μ_1, \dots, μ_K unknown parameters	μ_1, \dots, μ_K drawn from a prior distribution : $\mu_a \sim \pi_a$
arm a : $(Y_{a,s})_s \stackrel{\text{i.i.d.}}{\sim} \nu^{\mu_a}$	arm a : $(Y_{a,s})_s \mu \stackrel{\text{i.i.d.}}{\sim} \nu^{\mu_a}$

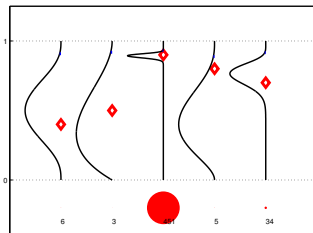
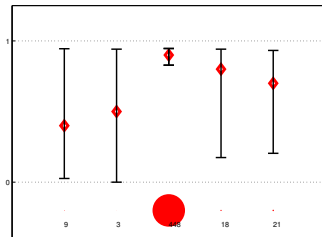
- ▶ The regret can be computed in each case

Frequentist regret (regret)	Bayesian regret (Bayes risk)
$\mathcal{R}_{\mu}(\mathcal{A}, T) = \mathbb{E}_{\mu} \left[\sum_{t=1}^T (\mu_{\star} - \mu_{A_t}) \right]$	$\mathbb{R}^{\pi}(\mathcal{A}, T) = \mathbb{E}_{\mu \sim \pi} \left[\sum_{t=1}^T (\mu_{\star} - \mu_{A_t}) \right]$ $= \int \mathcal{R}_{\mu}(\mathcal{A}, T) d\pi(\mu)$

Frequentist and Bayesian algorithms

- Two types of tools to build bandit algorithms:

Frequentist tools	Bayesian tools
MLE estimators of the means Confidence Intervals	Posterior distributions $\pi_a^t = \mathcal{L}(\mu_a Y_{a,1}, \dots, Y_{a,N_a(t)})$



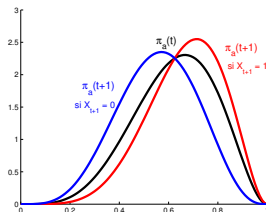
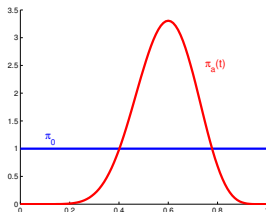
Example: Bernoulli bandits

Bernoulli bandit model $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)$

► **Bayesian view:** μ_1, \dots, μ_K are random variables
prior distribution : $\mu_a \sim \mathcal{U}([0, 1])$

→ posterior distribution:

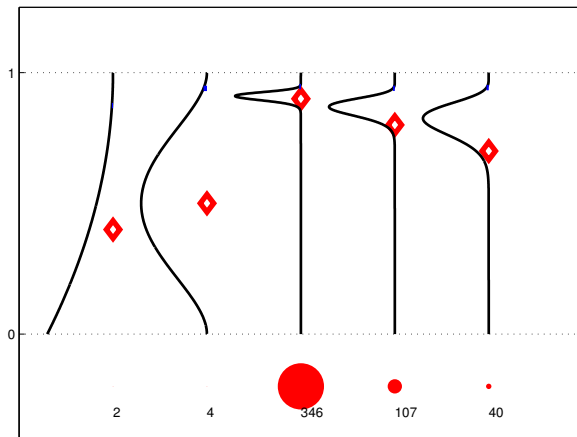
$$\begin{aligned}\pi_a(t) &= \mathcal{L}(\mu_a | R_1, \dots, R_t) \\ &= \text{Beta}\left(\underbrace{S_a(t)+1}_{\text{\#ones}}, \underbrace{N_a(t) - S_a(t) + 1}_{\text{\#zeros}}\right)\end{aligned}$$



$S_a(t) = \sum_{s=1}^t R_s \mathbb{1}_{(A_s=a)}$ sum of the rewards from arm a

Bayesian algorithm

A **Bayesian bandit algorithm** exploits the posterior distributions of the means to decide which arm to select.



Bayesian Bandits

Insights from the Optimal Solution

Bayes-UCB

Thompson Sampling

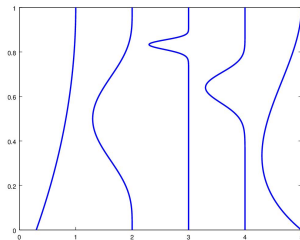
Some insights from the Bayesian solution

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

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

The posterior distribution is fully summarized by a matrix containing the number of ones and zeros observed for each arm.

$$\Pi^t = \begin{pmatrix} 0 & 2 \\ 3 & 3 \\ 13 & 4 \\ 5 & 2 \\ 1 & 3 \end{pmatrix}$$



“State” Π^t that evolves.

A first Markov Decision Process

After each arm selection A_t , we receive a reward R_t such that

$$\mathbb{P}\left(R_t = 1 | \Pi^{t-1} = \Pi, A_t = a\right) = \frac{\Pi^t(a, 1) + 1}{\underbrace{\Pi^t(a, 1) + \Pi^t(a, 2) + 2}_{\text{mean of } \pi_a(t-1)}}$$

and the posterior gets updated:

$$\begin{aligned}\Pi^t(A_t, 1) &= \Pi^{t-1}(A_t, 1) + R_t \\ \Pi^t(A_t, 2) &= \Pi^{t-1}(A_t, 2) + (1 - R_t)\end{aligned}$$

Example of transition:

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

→ Markov Decision Process with state Π^t

A first Markov Decision Process

After each arm selection A_t , we receive a reward R_t such that

$$\mathbb{P}\left(R_t = 1 | \Pi^{t-1} = \Pi, A_t = a\right) = \frac{\Pi^t(a, 1) + 1}{\underbrace{\Pi^t(a, 1) + \Pi^t(a, 2) + 2}_{\text{mean of } \pi_a(t-1)}}$$

and the posterior gets updated:

$$\Pi^t(A_t, 1) = \Pi^{t-1}(A_t, 1) + R_t$$

$$\Pi^t(A_t, 2) = \Pi^{t-1}(A_t, 2) + (1 - R_t)$$

Example of transition:

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

→ Markov Decision Process with state Π^t

An exact solution

Solving the Bayesian bandit \leftrightarrow maximizing rewards in some Markov Decision Process (modern perspective)

There exists an exact solution to

► The finite-horizon MAB:

$$\operatorname{argmax}_{(A_t)} \mathbb{E}_{\mu \sim \pi} \left[\sum_{t=1}^T R_t \right]$$

► The discounted MAB:

$$\operatorname{argmax}_{(A_t)} \mathbb{E}_{\mu \sim \pi} \left[\sum_{t=1}^{\infty} \gamma^{t-1} R_t \right]$$

[Berry and Fristedt, *Bandit Problems*, 1985]

Optimal solution: solution to dynamic programming equations.

Problem: The state space is very large

\rightsquigarrow often intractable

Gittins indices

[Gittins 79]: the solution of the **discounted** MAB

$$\operatorname{argmax}_{(A_t)} \mathbb{E}_{\mu \sim \pi} \left[\sum_{t=1}^{\infty} \gamma^{t-1} R_t \right]$$

is an **index policy**:

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

► The Gittins indices:

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

with

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

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

Gittins indices for Finite Horizon?

The solution of the **finite horizon** MAB

$$\operatorname{argmax}_{(A_t)} \mathbb{E}_{\mu \sim \pi} \left[\sum_{t=1}^T R_t \right]$$

is NOT an index policy. [Berry and Fristedt 85]

- **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

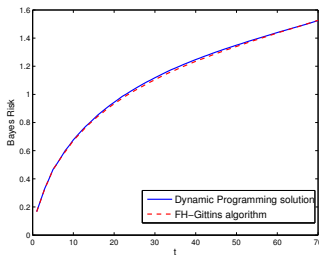
$$V_r^*(p, \lambda) = \sup_{\substack{\text{stopping times} \\ 0 < \tau \leq r}} \mathbb{E}_{\substack{Y_t \text{ i.i.d } \mathcal{B}(\mu) \\ \mu \sim p}} \left[\sum_{t=1}^{\tau} (Y_t - \lambda) \right].$$

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

FH Gittins algorithm:

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

does NOT coincide with the Bayesian optimal solution but is conjectured to be a good approximation!



- ▶ good performance in terms of frequentist regret as well
- ▶ ... with logarithmic regret [Lattimore, 2016]

Approximating the FH-Gittins indices

- ▶ [Burnetas and Katehakis, 03]: when n is large,

$$G(\pi_a(t-1), n) \simeq \max \left\{ q : N_a(t) \times \text{kl}(\hat{\mu}_a(t), q) \leq \ln \left(\frac{n}{N_a(t)} \right) \right\}$$

- ▶ [Lai, 87]: the index policy associated to

$$I_a(t) = \max \left\{ q : N_a(t) \times \text{kl}(\hat{\mu}_a(t), q) \leq \ln \left(\frac{T}{N_a(t)} \right) \right\}$$

is a good approximation of the Bayesian solution for large T .

- looks like the **kl-UCB index**, with a **different exploration rate...**

Bayesian Bandits

Insights from the Optimal Solution

Bayes-UCB

Thompson Sampling

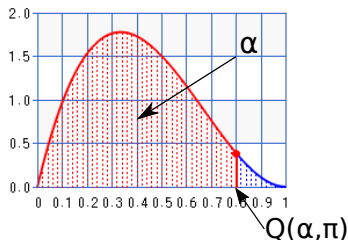
The Bayes-UCB algorithm

- ▶ $\Pi_0 = (\pi_1(0), \dots, \pi_K(0))$ be a prior distribution over (μ_1, \dots, μ_K)
- ▶ $\Pi_t = (\pi_1(t), \dots, \pi_K(t))$ be the posterior distribution over the means (μ_1, \dots, μ_K) after t observations

The **Bayes-UCB algorithm** chooses at time t

$$A_{t+1} = \operatorname{argmax}_{a=1, \dots, K} Q\left(1 - \frac{1}{t(\ln t)^c}, \pi_a(t)\right)$$

where $Q(\alpha, \pi)$ is the quantile of order α of the distribution π .



The Bayes-UCB algorithm

- ▶ $\Pi_0 = (\pi_1(0), \dots, \pi_K(0))$ be a prior distribution over (μ_1, \dots, μ_K)
- ▶ $\Pi_t = (\pi_1(t), \dots, \pi_K(t))$ be the posterior distribution over the means (μ_1, \dots, μ_K) after t observations

The **Bayes-UCB algorithm** chooses at time t

$$A_{t+1} = \operatorname{argmax}_{a=1, \dots, K} Q\left(1 - \frac{1}{t(\ln t)^c}, \pi_a(t)\right)$$

where $Q(\alpha, \pi)$ is the quantile of order α of the distribution π .

Bernoulli reward with uniform prior:

- ▶ $\pi_a(0) \stackrel{i.i.d}{\sim} \mathcal{U}([0, 1]) = \text{Beta}(1, 1)$
- ▶ $\pi_a(t) = \text{Beta}(S_a(t) + 1, N_a(t) - S_a(t) + 1)$

The Bayes-UCB algorithm

- ▶ $\Pi_0 = (\pi_1(0), \dots, \pi_K(0))$ be a prior distribution over (μ_1, \dots, μ_K)
- ▶ $\Pi_t = (\pi_1(t), \dots, \pi_K(t))$ be the posterior distribution over the means (μ_1, \dots, μ_K) after t observations

The **Bayes-UCB algorithm** chooses at time t

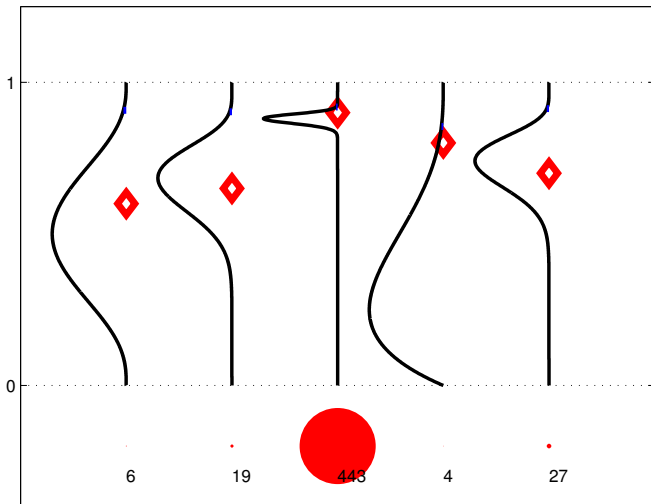
$$A_{t+1} = \operatorname{argmax}_{a=1, \dots, K} Q \left(1 - \frac{1}{t(\ln t)^c}, \pi_a(t) \right)$$

where $Q(\alpha, \pi)$ is the quantile of order α of the distribution π .

Gaussian rewards with Gaussian prior:

- ▶ $\pi_a(0) \stackrel{i.i.d}{\sim} \mathcal{N}(0, \kappa^2)$
- ▶ $\pi_a(t) = \mathcal{N} \left(\frac{S_a(t)}{N_a(t) + \sigma^2/\kappa^2}, \frac{\sigma^2}{N_a(t) + \sigma^2/\kappa^2} \right)$

Bayes UCB in action



- Bayes-UCB is **asymptotically optimal** for Bernoulli rewards

Theorem [K., Cappé, Garivier 2012]

Let $\epsilon > 0$. The Bayes-UCB algorithm using a uniform prior over the arms and parameter $c \geq 5$ satisfies

$$\mathbb{E}_{\mu}[N_a(T)] \leq \frac{1 + \epsilon}{\text{kl}(\mu_a, \mu_{\star})} \ln(T) + o_{\epsilon, c}(\ln(T)).$$

Lemma [K. et al., 12]

The index $q_a(t)$ used by Bayes-UCB satisfies

$$\tilde{u}_a(t) \leq q_a(t) \leq u_a(t)$$

where

$$u_a(t) = \max \left\{ q : \text{kl} \left(\frac{S_a(t)}{N_a(t)}, q \right) \leq \frac{\ln(t) + c \ln(\ln(t))}{N_a(t)} \right\}$$
$$\tilde{u}_a(t) = \max \left\{ q : \text{kl} \left(\frac{S_a(t)}{N_a(t) + 1}, q \right) \leq \frac{\ln \left(\frac{t}{N_a(t) + 2} \right) + c \ln(\ln(t))}{(N_a(t) + 1)} \right\}$$

Proof: rely on the [Beta-Binomial trick](#) :

$$F_{\text{Beta}(a,b)}(x) = 1 - F_{\text{Bin}(a+b-a,x)}(a-1)$$

[Agrawal and Goyal, 12]

- For **one-dimensional exponential families**, Bayes-UCB rewrites

$$A_{t+1} = \operatorname{argmax}_a Q \left(1 - \frac{1}{t(\ln t)^c}, \pi_{a, N_a(t), \hat{\mu}_a(t)} \right)$$

Extra assumption: there exists μ^-, μ^+ such that for all a , $\mu_a \in [\mu^-, \mu^+]$

Theorem [K. 17]

Let $\bar{\mu}_a(t) = (\hat{\mu}_a(t) \vee \mu^-) \wedge \mu^+$. The index policy

$$A_{t+1} = \operatorname{argmax}_a Q \left(1 - \frac{1}{t(\ln t)^c}, \pi_{a, N_a(t), \bar{\mu}_a(t)} \right)$$

with parameter $c \geq 7$ is such that, for all $\epsilon > 0$,

$$\mathbb{E}_\mu[N_a(T)] \leq \frac{1 + \epsilon}{\text{kl}(\mu_a, \mu_*)} \ln(T) + O_\epsilon(\sqrt{\ln(T)}).$$

An interesting by-product

- Tools from the analysis of Bayes-UCB can be used to analyze two variants of kl-UCB

kl-UCB-H⁺

$$u_a^{H,+}(t) = \max \left\{ q : N_a(t) \times \text{kl}(\hat{\mu}_a(t), q) \leq \ln \left(\frac{T \ln^c T}{N_a(t)} \right) \right\}$$

kl-UCB⁺

$$u_a^+(t) = \max \left\{ q : N_a(t) \times \text{kl}(\hat{\mu}_a(t), q) \leq \ln \left(\frac{t \ln^c t}{N_a(t)} \right) \right\}$$

The index policy associated to $u_a^{H,+}(t)$ and $u_a^+(t)$ satisfy, for all $\epsilon > 0$,

$$\mathbb{E}_\mu[N_a(T)] \leq \frac{1 + \epsilon}{\text{kl}(\mu_a, \mu_\star)} \ln(T) + O_\epsilon(\sqrt{\ln(T)}).$$

Bayesian Bandits

Insights from the Optimal Solution

Bayes-UCB

Thompson Sampling

Historical perspective

- 1933 Thompson: in the context of clinical trial, the allocation of a treatment should be some increasing function of its **posterior probability to be optimal**
- 2010 Thompson Sampling rediscovered under different names
 - Bayesian Learning Automaton [Granmo, 2010]
 - Randomized probability matching [Scott, 2010]
- 2011 An empirical evaluation of Thompson Sampling: **an efficient algorithm**, beyond simple bandit models
 - [Li and Chapelle, 2011]
- 2012 First (logarithmic) **regret bound** for Thompson Sampling
 - [Agrawal and Goyal, 2012]
- 2012 Thompson Sampling is **asymptotically optimal for Bernoulli bandits**
 - [K., Korda and Munos, 2012][Agrawal and Goyal, 2013]
- 2013- Many **successful uses of Thompson Sampling** beyond Bernoulli bandits (contextual bandits, reinforcement learning)

Thompson Sampling

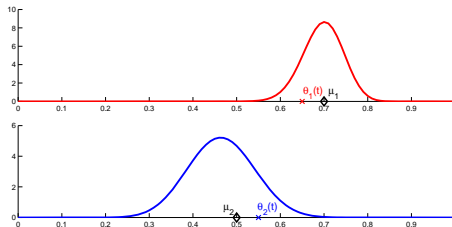
Two equivalent interpretations:

- ▶ “select an arm at random according to its probability of being the best”
- ▶ “draw a possible bandit model from the posterior distribution and act optimally in this sampled model”

≠ optimistic

Thompson Sampling: a randomized Bayesian algorithm

$$\begin{cases} \forall a \in \{1..K\}, \theta_a(t) \sim \pi_a(t) \\ A_{t+1} = \operatorname{argmax}_{a=1..K} \theta_a(t). \end{cases}$$



Thompson Sampling is asymptotically optimal

Problem-dependent regret

$$\forall \epsilon > 0, \quad \mathbb{E}_{\mu}[N_a(T)] \leq (1 + \epsilon) \frac{1}{\text{kl}(\mu_a, \mu_{\star})} \ln(T) + o_{\mu, \epsilon}(\ln(T)).$$

This results holds:

- ▶ for **Bernoulli bandits**, with a **uniform prior**
[K. Korda, Munos 12][Agrawal and Goyal 13]
- ▶ for **Gaussian bandits**, with **Gaussian prior** [Agrawal and Goyal 17]
- ▶ for **exponential family bandits**, with **Jeffrey's prior** [Korda et al. 13]

Problem-independent regret [Agrawal and Goyal 13]

For Bernoulli and Gaussian bandits, Thompson Sampling satisfies

$$\mathcal{R}_{\mu}(\text{TS}, T) = O\left(\sqrt{KT \ln(T)}\right).$$

- ▶ Thompson Sampling is also **asymptotically optimal for Gaussian with unknown mean and variance** [Honda and Takemura, 14]

Understanding Thompson Sampling

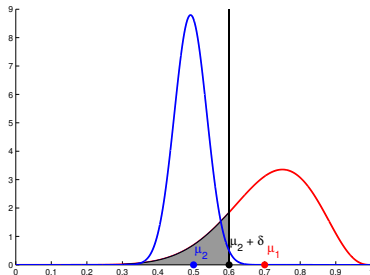
- ▶ a key ingredient in the analysis of [K. Korda and Munos 12]

Proposition

There exists constants $b = b(\mu) \in (0, 1)$ and $C_b < \infty$ such that

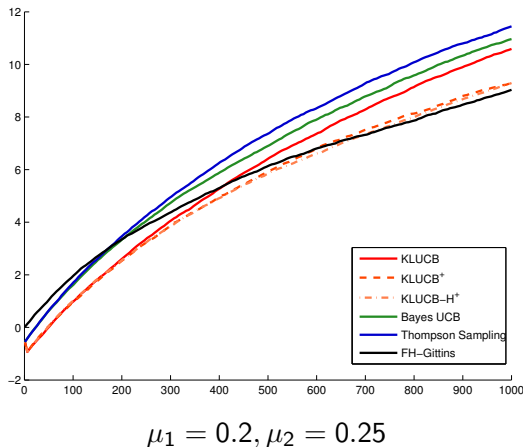
$$\sum_{t=1}^{\infty} \mathbb{P} \left(N_1(t) \leq t^b \right) \leq C_b.$$

$\{N_1(t) \leq t^b\} = \{\text{there exists a time range of length at least } t^{1-b} - 1$
with no draw of arm 1 } $\}$



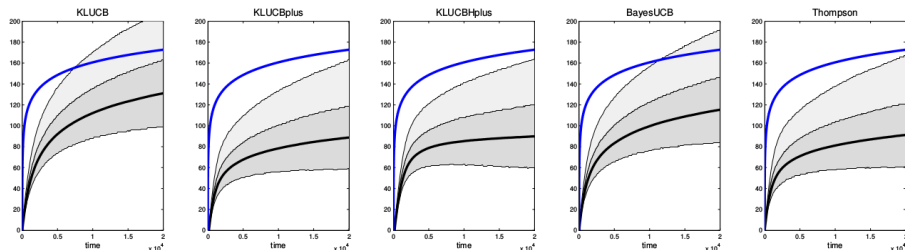
Bayesian versus Frequentist algorithms

- Short horizon, $T = 1000$ (average over $N = 10000$ runs)



Bayesian versus Frequentist algorithms

- Long horizon, $T = 20000$ (average over $N = 50000$ runs)



10 arms bandit problem

$$\mu = [0.1 \ 0.05 \ 0.05 \ 0.05 \ 0.02 \ 0.02 \ 0.02 \ 0.01 \ 0.01 \ 0.01]$$

OTHER RANDOMIZED ALGORITHMS

Limitation of existing approaches

Two families of asymptotically optimal algorithms

- ▶ Confidence bound algorithms
 - ▶ Thompson Sampling
-
- ▶ Provably optimal finite-time regret under the assumption that the rewards distribution belong to **some class \mathcal{D}**
 - ▶ A **different algorithm for each \mathcal{D}** : TS or kl-UCB for Bernoulli, Poisson, for Exponential, etc.

Can we build a **universal algorithm** that would be asymptotically optimal over different classes \mathcal{D} ?

Best Empirical Sub-sampling Average

"Sub-sampling for multi-armed bandits",
Baransi, Maillard, Mannor *ECML*, 2014.

BESA

- ▶ Competitive regret against state-of-the-art for various \mathcal{D} .
- ▶ Same algorithm for all \mathcal{D} .
- ▶ Not relying on upper confidence bounds, not Bayesian...
- ▶ ...and extremely simple to implement.

→ How? Optimality? For which distributions ?

FTL

- 1 Play each arm once.
- 2 At time t , define $\tilde{\mu}_a(t) = \hat{\mu}(R_{1:N_a(t)}^a)$ for all $a \in \mathcal{A}$.
 - ▶ $\hat{\mu}(\mathcal{X})$: empirical average of population \mathcal{X} .
 - ▶ $R_{1:N_a(t)}^a = \{R_s : A_s = a, s \leq t\}$
- 3 Choose (break ties in favor of the smallest $N_a(t)$)

$$A_{t+1} = \operatorname{argmax}_{a' \in \{a, b\}} \tilde{\mu}_{a'}(t).$$

Properties

- ▶ Generally bad: linear regret.
- ▶ A variant (ϵ -greedy) performs ok if well-tuned [Auer et al, 2002].

Follow the FAIR leader (aka BESA)

Idea: Compare two arms based on "equal opportunity"
i.e. same number of observations.

BESA at time t for two arms a, b :

- ① Sample two sets of indices $\mathcal{I}_a(t) \sim \text{Wr}(N_a(t); N_b(t))$ and $\mathcal{I}_b(t) \sim \text{Wr}(N_b(t); N_a(t))$.
 - ▶ $\text{Wr}(n, N)$: sample n points from $\{1, \dots, N\}$ without replacement (return all the set if $n \geq N$).
- ② Define $\tilde{\mu}_a(t) = \hat{\mu}(R_{1:N_a(t)}^a(\mathcal{I}_a(t)))$ and $\tilde{\mu}_b(t) = \hat{\mu}(R_{1:N_{t,b}}^b(\mathcal{I}_b(t)))$.
- ③ Choose (break ties in favor of the smallest $N_{a'}(t)$)

$$A_{t+1} = \underset{a' \in \{a, b\}}{\operatorname{argmax}} \tilde{\mu}_{a'}(t).$$

- ▶ more than two arms? tournament.

Example

- ▶ $\mathcal{X} = (x_1, \dots, x_N)$, a finite population of N real points.

x_1	x_2	x_3	x_4	x_5	\dots	x_{N-2}	x_{N-1}	x_N
-------	-------	-------	-------	-------	---------	-----------	-----------	-------

- ▶ Sub-sample of size $n \leq N$ from \mathcal{X} : X_1, \dots, X_n picked uniformly randomly without replacement from \mathcal{X} .

x_1	X_{n-1}	X_1	x_4	X_2	\dots	x_{N-2}	X_n	x_N
-------	-----------	-------	-------	-------	---------	-----------	-------	-------

- ▶ Example: $N_a(t) = 3$ and $N_b(t) = 10$:

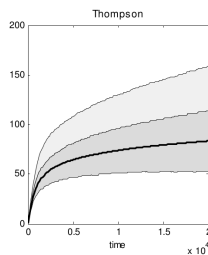
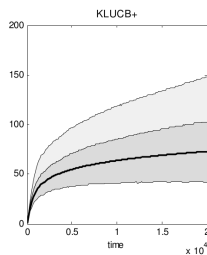
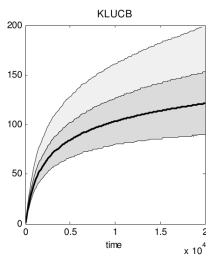
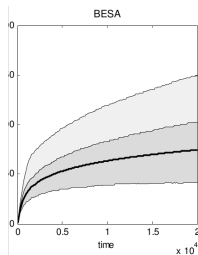
$$\mathcal{I}_a(t) = \{1, 2, 3\},$$

$$|\mathcal{I}_b(t)| = 3, \text{ sampled without replacement from } \{1, \dots, 10\}.$$

Good practical performance ($T = 20,000$, $N = 50,000$)

► 10 Bernoulli($0.1, 3\{0.05\}, 3\{0.02\}, 3\{0.01\}$)

	BESA	kl-UCB	kl-UCB+	TS	Others
Regret	74.4	121.2	72.8	83.4	100-400
Beat BESA	-	1.6%	35.4%	3.1%	
Run Rime	13.9X	2.8X	3.1X	X	



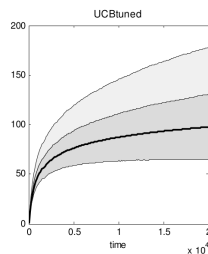
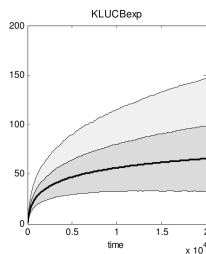
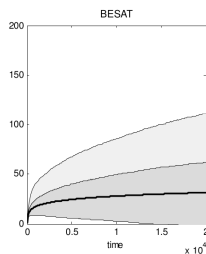
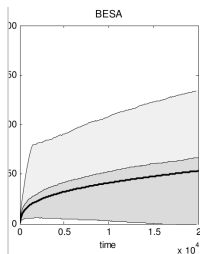
Others: UCB, Moss, UCB-Tunes, DMED, UCB-V.

(Credit: Akram Baransi)

Good practical performance ($T = 20,000$, $N = 50,000$)

► Exponential($\frac{1}{5}, \frac{1}{4}, \frac{1}{3}, \frac{1}{2}, 1$)

	BESA	KL-UCB-exp	UCB-tuned	FTL 10	Others
Regret	53.3	65.7	97.6	306.5	60-110,120+
Beat BESA	-	5.7%	4.3%	-	
Run Rime	6X	2.8X	X	-	



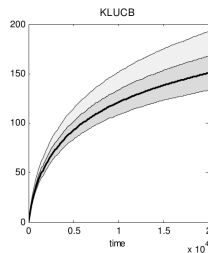
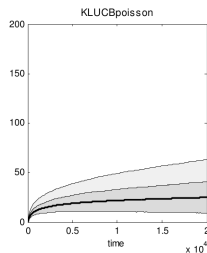
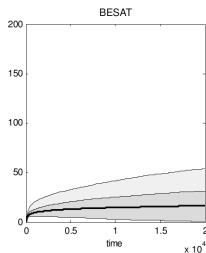
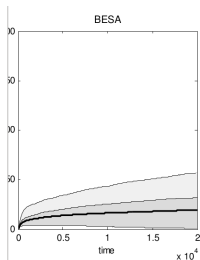
Others: UCB, Moss, kl-UCB, UCB-V.

(Credit: Akram Baransi)

Good practical performance ($T = 20,000$, $N = 50,000$)

► **Poisson**($\{\frac{1}{2} + \frac{i}{3}\}_{i=1,\dots,6}$)

	BESA	KL-UCB-Poisson	kl-UCB	FTL 10
Regret	19.4	25.1	150.6	144.6
Beat BESA	-	4.1%	0.7%	-
Run Rime	3.5X	1.2X	X	-



(Credit: Akram Baransi)

Regret bound (slightly simplified statement)

With two arms $\{\star, a\}$, define

$$\alpha(M, n) = \mathbb{E}_{Z^\star \sim \nu_{\star, n}} \left[\left(\mathbb{P}_{Z \sim \nu_{a, n}}(Z > Z^\star) + \frac{1}{2} \mathbb{P}_{Z \sim \nu_{a, n}}(Z = Z^\star) \right)^M \right].$$

Theorem [Baransi et al. 14]

If $\exists \alpha \in (0, 1), c > 0$ such that $\alpha(M, 1) \leq c\alpha^M$, then

$$\mathcal{R}_\nu(\text{BESA}, T) \leq \frac{11 \ln(T)}{\mu_\star - \mu_a} + C_\nu + O(1).$$

Example

► Bernoulli μ_a, μ_\star : $\alpha(M, 1) = O\left(\left(\frac{\mu_a \vee (1 - \mu_a)}{2}\right)^M\right)$

Future work: understand when BESA fails, and whether it can be asymptotically optimal in some cases...

Another class of (randomized) bandit algorithms that do not exploit any assumption on \mathcal{D} is that of **adversarial bandit algorithms**.

[Auer, Cesa-Bianchi, Freund, Shapire,
The non-stochastic multi-armed bandit, 2002]

Can we achieve $O(\sqrt{KT})$ regret with respect to the best static action if the rewards are arbitrarily generated?

Some answers in the next classes and practical sessions!

SUMMARY

Take-home messages

Now you are aware of:

- ▶ several methods for facing an exploration/exploitation dilemma
- ▶ notably two powerful classes of methods
 - ▶ optimistic “UCB” algorithms
 - ▶ Bayesian approaches, mostly Thompson Sampling

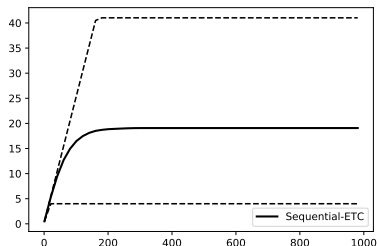
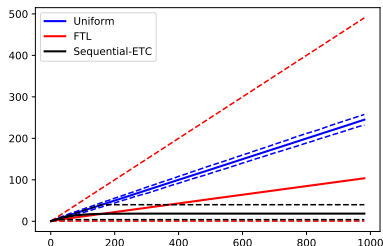
And you are therefore ready to apply them for solving more complex (structured) bandit problems and for Reinforcement Learning!

You also saw a bunch of important tools:

- ▶ performance lower bounds, guiding the design of algorithms
- ▶ Kullback-Leibler divergence to measure deviations
- ▶ self-normalized concentration inequalities
- ▶ Bayesian tools

First practical session

Objective: run UCB, kl-UCB, Thompson Sampling and some tweaks of those algorithms, and see what performs best (on simulated data).



- visualize expected regret *averaged over multiple runs* / distribution of the regret

Files: link on my webpage

- ▶ W.R. Thompson (1933). On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*.
- ▶ H. Robbins (1952). Some aspects of the sequential design of experiments. *Bulletin of the American Mathematical Society*.
- ▶ Bradt, R., Johnson, S., and Karlin, S. (1956). On sequential designs for maximizing the sum of n observations. *Annals of Mathematical Statistics*.
- ▶ R. Bellman (1956). A problem in the sequential design of experiments. *The indian journal of statistics*.
- ▶ Gittins, J. (1979). Bandit processes and dynamic allocation indices. *Journal of the Royal Statistical Society*.
- ▶ Berry, D. and Fristedt, B. Bandit Problems (1985). Sequential allocation of experiments. *Chapman and Hall*.
- ▶ Lai, T. and Robbins, H. (1985). Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*.
- ▶ Lai, T. (1987). Adaptive treatment allocation and the multi-armed bandit problem. *Annals of Statistics*.

- ▶ Agrawal, R. (1995). Sample mean based index policies with $O(\log n)$ regret for the multi-armed bandit problem. *Advances in Applied Probability*.
- ▶ Katehakis, M. and Robbins, H. (1995). Sequential choice from several populations. *Proceedings of the National Academy of Science*.
- ▶ Burnetas, A. and Katehakis, M. (1996). Optimal adaptive policies for sequential allocation problems. *Advances in Applied Mathematics*.
- ▶ Auer, P., Cesa-Bianchi, N., and Fischer, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine Learning*.
- ▶ Auer, P., Cesa-Bianchi, N., Freund, Y., and Schapire, R. (2002). The nonstochastic multiarmed bandit problem. *SIAM Journal of Computing*.
- ▶ Burnetas, A. and Katehakis, M. (2003). Asymptotic Bayes Analysis for the finite horizon one armed bandit problem. *Probability in the Engineering and Informational Sciences*.
- ▶ Cesa-Bianchi, N. and Lugosi, G. (2006). Prediction, Learning and Games. *Cambridge University Press*.

- ▶ Audibert, J.-Y., Munos, R. and Szepesvari, C. (2009). Exploration-exploitation trade-off using variance estimates in multi-armed bandits. *Theoretical Computer Science*.
- ▶ Audibert, J.-Y. and Bubeck, S. (2010). Regret Bounds and Minimax Policies under Partial Monitoring. *Journal of Machine Learning Research*.

to be completed!