

A minimax and asymptotically optimal algorithm for stochastic bandits

Pierre Ménard

Aurélien Garivier

Institut de Mathématiques de Toulouse

October 17, 2017

K-armed bandit problem

$$\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

K-armed bandit problem

$$\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

Game: for each round $1 \leq t \leq T$:

1. Player pulls arm $A_t \in \{1, \dots, K\}$.
2. He gets a reward $Y_t \sim \mathcal{B}(\mu_{A_t})$.

K-armed bandit problem

$$\nu = (\mathcal{B}(\mu_1), \dots, \mathcal{B}(\mu_a), \dots, \mathcal{B}(\mu_K))$$

Game: for each round $1 \leq t \leq T$:

1. Player pulls arm $A_t \in \{1, \dots, K\}$.
2. He gets a reward $Y_t \sim \mathcal{B}(\mu_{A_t})$.

Regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T (\mu^* - \mu_{A_t}) \right] = \sum_{a=1}^K (\mu^* - \mu_a) \mathbb{E}[N_a(T)],$$

where $\mu^* = \max_{a=1, \dots, K} \mu_a$ and $N_a(T) = \sum_{t=1}^T \mathbb{I}_{\{A_t=a\}}$.

UCB algorithm

UCB algorithm:

Play each arm once, then for $K \leq t \leq T - 1$ play:

$$A_{t+1} \in \arg \max_a \hat{\mu}_{a, N_a(t)} + \sqrt{\frac{2 \log(t)}{N_a(t)}}.$$

UCB algorithm

UCB algorithm:

Play each arm once, then for $K \leq t \leq T - 1$ play:

$$A_{t+1} \in \arg \max_a \hat{\mu}_{a, N_a(t)} + \sqrt{\frac{2 \log(t)}{N_a(t)}}.$$

Regret bound: for all a such that $\mu^* - \mu_a > 0$

$$\mathbb{E}[N_a(T)] \leq \frac{8}{(\mu^* - \mu_a)^2} \ln(T) + o(\ln(T)),$$

UCB algorithm

UCB algorithm:

Play each arm once, then for $K \leq t \leq T - 1$ play:

$$A_{t+1} \in \arg \max_a \hat{\mu}_{a, N_a(t)} + \sqrt{\frac{2 \log(t)}{N_a(t)}}.$$

Regret bound: for all a such that $\mu^* - \mu_a > 0$

$$\mathbb{E}[N_a(T)] \leq \frac{8}{(\mu^* - \mu_a)^2} \ln(T) + o(\ln(T)),$$

Is that the best we can do? \Rightarrow Lower bound

Asymptotic lower bound

Kullback-Leibler divergence :

$$\text{kl}(p, q) := \text{KL}(\mathcal{B}(p), \mathcal{B}(q)) = p \ln(p/q) + (1 - p) \ln((1 - p)/(1 - q))$$

Theorem (Asymptotic lower bound from Lai & Robbins)

For all consistent algorithms, for all suboptimal arms a (a such that $\mu^ - \mu_a > 0$),*

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\ln T} \geq \frac{1}{\text{kl}(\mu_a, \mu^*)}.$$

Asymptotic lower bound

Kullback-Leibler divergence :

$$\text{kl}(p, q) := \text{KL}(\mathcal{B}(p), \mathcal{B}(q)) = p \ln(p/q) + (1 - p) \ln((1 - p)/(1 - q))$$

Theorem (Asymptotic lower bound from Lai & Robbins)

For all consistent algorithms, for all suboptimal arms a (a such that $\mu^* - \mu_a > 0$),

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\ln T} \geq \frac{1}{\text{kl}(\mu_a, \mu^*)}.$$

An algorithm is **asymptotically** optimal if (a such that $\mu^* - \mu_a > 0$)

$$\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_a(T)]}{\ln T} \leq \frac{1}{\text{kl}(\mu_a, \mu^*)} \stackrel{\text{Pinsker}}{\leq} \frac{2}{(\mu^* - \mu_a)^2},$$

kl-UCB algorithm

kl-UCB algorithm:

Play each arm once, then for $K \leq t \leq T - 1$ play:

$$A_{t+1} \in \arg \max_a \sup \left\{ \mu \in [0, 1] : \text{kl}(\hat{\mu}_a(t), \mu) \leq \frac{\log(t) + 3 \log \log(t)}{N_a(t)} \right\}.$$

Regret bound For a such that $\mu^* - \mu_a > 0$

$$\mathbb{E}[N_a(T)] \leq \frac{\log(T)}{\text{kl}(\mu_a, \mu^*)} + O\left(\sqrt{\log(T)}\right).$$

⇒ kl-UCB is asymptotically optimal.

Minimax optimality

Theorem (Minimax lower bound)

$$\sup_{\nu} R_T \geq C\sqrt{KT}.$$

An algorithm is **minimax** optimal if

$$R_T \leq C'\sqrt{KT}.$$

UCB algorithm not minimax optimal

$$R_T \leq C''\sqrt{KT \log(T)}$$

MOSS-algorithm

Algorithm:

Play each arm once, then for $K \leq t \leq T - 1$ play:

$$A_{t+1} \in \arg \max_a \hat{\mu}_{a, N_a(t)} + \sqrt{\frac{g(t, N_a(t))}{N_a(t)}}.$$

Algorithm	$g(t, n)$	Minimax opt.	Asymptotic opt.
UCB	$2 \log(t)$	✗	✗
MOSS	$\log_+ \left(\frac{T}{Kn} \right)$	✓	✗

kl-UCB⁺⁺ algorithm

Generic algorithm:

Play each arm once, then for $K \leq t \leq T - 1$ play:

$$A_{t+1} \in \arg \max_a \sup \left\{ \mu \in [0, 1] : \text{kl}(\hat{\mu}_a(t), \mu) \leq \frac{g(t, N_a(t))}{N_a(t)} \right\}.$$

Algorithm	$g(T, n)$	Minimax opt.	Asymptotic opt.
kl-UCB	$\log(t) + 3 \log \log(t)$	✗	✓
kl-UCB ⁺⁺	$\log_+ \left(\frac{T}{Kn} \left(\log_+ \left(\frac{T}{Kn} \right)^2 + 1 \right) \right)$	✓	✓

Theorem (Optimality of kl-UCB⁺⁺)

Minimax optimality:

$$R_T \leq 38\sqrt{KT} + K.$$

Asymptotic optimality: For any sub-optimal arm a and any δ such that $\sqrt{11K/(2T)} \leq \delta \leq (\mu^* - \mu_a)/3$,

$$\mathbb{E}[N_a(T)] \leq \frac{\log(T)}{\text{kl}(\mu_a + \delta, \mu^* - \delta)} + O\left(\frac{\log\log(T)}{\delta^2}\right)$$

Sketch of proof: minimax bound

$$U_a(t) := \sup \left\{ \mu \in [0, 1] : \text{kl}(\hat{\mu}_a(t), \mu) \leq \frac{g(N_a(t))}{N_a(t)} \right\}.$$

Decomposition of the regret a^* optimal $\mu_{a^*} = \mu^*$

$$\begin{aligned} R_T &\leq K + \sum_{t=K}^{T-1} \mathbb{E}[\mu^* - U_{A_{t+1}}(t) + U_{A_{t+1}}(t) - \mu_{A_{t+1}}] \\ &\leq K + \underbrace{\sum_{t=K}^{T-1} \mathbb{E}[\mu^* - U_{a^*}(t)]}_A + \underbrace{\sum_{t=K}^{T-1} \mathbb{E}[U_{A_{t+1}}(t) - \mu_{A_{t+1}}]}_B \end{aligned}$$

Sketch of proof: minimax bound

$$U_a(t) := \sup \left\{ \mu \in [0, 1] : \text{kl}(\hat{\mu}_a(t), \mu) \leq \frac{g(N_a(t))}{N_a(t)} \right\}.$$

Decomposition of the regret a^* optimal $\mu_{a^*} = \mu^*$

$$\begin{aligned} R_T &\leq K + \sum_{t=K}^{T-1} \mathbb{E}[\mu^* - U_{A_{t+1}}(t) + U_{A_{t+1}}(t) - \mu_{A_{t+1}}] \\ &\leq K + \underbrace{\sum_{t=K}^{T-1} \mathbb{E}[\mu^* - U_{a^*}(t)]}_A + \underbrace{\sum_{t=K}^{T-1} \mathbb{E}[U_{A_{t+1}}(t) - \mu_{A_{t+1}}]}_B \end{aligned}$$

B term: go back to MOSS-Index!

$$U_a(t) \leq B_a(t) = \hat{\mu}_{a, N_a(t)} + \sqrt{\frac{g(N_a(t))}{2N_a(t)}}.$$

Sketch of proof: minimax bound

$$U_a(t) := \sup \left\{ \mu \in [0, 1] : \text{kl}(\widehat{\mu}_a(t), \mu) \leq \frac{g(N_a(t))}{N_a(t)} \right\}.$$

Decomposition of the regret a^* optimal $\mu_{a^*} = \mu^*$

$$\begin{aligned} R_T &\leq K + \sum_{t=K}^{T-1} \mathbb{E}[\mu^* - U_{A_{t+1}}(t) + U_{A_{t+1}}(t) - \mu_{A_{t+1}}] \\ &\leq K + \underbrace{\sum_{t=K}^{T-1} \mathbb{E}[\mu^* - U_{a^*}(t)]}_A + \underbrace{\sum_{t=K}^{T-1} \mathbb{E}[U_{A_{t+1}}(t) - \mu_{A_{t+1}}]}_B \end{aligned}$$

B term: go back to MOSS-Index!

$$B \leq C_B \sqrt{KT}$$

Sketch of proof 2

A term: peeling trick.

Integrate the deviations $\delta_0 \sim \sqrt{K/T}$

$$E[\mu^* - U_{a^*}(t)] \leq \delta_0 + \int_{\delta_0}^{+\infty} \mathbb{P}(U_{a^*}(t) \leq \mu^* - u) du,$$

Split the probability

$$\mathbb{P}(U_{a^*}(t) \leq \mu^* - u) \leq$$

$$\mathbb{P}(\exists 1 \leq n \leq T, \text{kl}_+(\hat{\mu}_{a^*,n}, \mu^* - u) \geq g(n)/n) \leq$$

$$\underbrace{\mathbb{P}(\exists 1 \leq n \leq n_u, \text{kl}_+(\hat{\mu}_{a^*,n}, \mu^*) \geq g(n)/n)}_{A_1} + \underbrace{\mathbb{P}(\exists n_u \leq n \leq T, \hat{\mu}_{a^*,n} \leq \mu^* - u)}_{A_2}$$

Sketch of proof 2

A term: peeling trick.

Integrate the deviations $\delta_0 \sim \sqrt{K/T}$

$$E[\mu^* - U_{a^*}(t)] \leq \delta_0 + \int_{\delta_0}^{+\infty} \mathbb{P}(U_{a^*}(t) \leq \mu^* - u) du,$$

Split the probability

$$\mathbb{P}(U_{a^*}(t) \leq \mu^* - u) \leq$$

$$\mathbb{P}(\exists 1 \leq n \leq T, \text{kl}_+(\hat{\mu}_{a^*,n}, \mu^* - u) \geq g(n)/n) \leq$$

$$\underbrace{\mathbb{P}(\exists 1 \leq n \leq n_u, \text{kl}_+(\hat{\mu}_{a^*,n}, \mu^*) \geq g(n)/n)}_{A_1} + \underbrace{\mathbb{P}(\exists n_u \leq n \leq T, \hat{\mu}_{a^*,n} \leq \mu^* - u)}_{A_2}$$

With

$$g(n_u)/n_u \sim \frac{u^2}{2}$$

since

$$\text{kl}_+(\hat{\mu}_{a^*,n}, \mu^* - u) \geq g(n)/n \Rightarrow \text{kl}_+(\hat{\mu}_{a^*,n}, \mu^*) \geq g(n)/n + u^2/2$$

Sketch of proof 2

A term: peeling trick.

Integrate the deviations $\delta_0 \sim \sqrt{K/T}$

$$E[\mu^* - U_{a^*}(t)] \leq \delta_0 + \int_{\delta_0}^{+\infty} \mathbb{P}(U_{a^*}(t) \leq \mu^* - u) du,$$

Split the probability

$$\mathbb{P}(U_{a^*}(t) \leq \mu^* - u) \leq$$

$$\mathbb{P}(\exists 1 \leq n \leq T, \text{kl}_+(\hat{\mu}_{a^*,n}, \mu^* - u) \geq g(n)/n) \leq$$

$$\underbrace{\mathbb{P}(\exists 1 \leq n \leq n_u, \text{kl}_+(\hat{\mu}_{a^*,n}, \mu^*) \geq g(n)/n)}_{A_1} + \underbrace{\mathbb{P}(\exists n_u \leq n \leq T, \hat{\mu}_{a^*,n} \leq \mu^* - u)}_{A_2}$$

$$\left. \begin{array}{l} A_1 : \quad \text{Peeling trick} \\ A_2 : \quad \text{Maximal inequality} \end{array} \right\} \Rightarrow A \leq C_A \sqrt{KT}$$

One-dimensional exponential family

Bandit problem:

$$\nu = (\nu_{\theta_1}, \dots, \nu_{\theta_a}, \dots, \nu_{\theta_K}),$$

one exponential family

$$\frac{d\nu_\theta}{d\rho}(x) = \exp(x\theta - b(\theta)), \quad \text{where } b(\theta) = \log \int_{\mathbb{R}} e^{x\theta} d\rho(x).$$

mean parametrization $\mu = b'(\theta)$, Pinsker-like inequality,

$$\text{kl}(\mu, \mu') := \text{KL}(\nu_\theta, \nu_{\theta'}) \geq \frac{1}{2V} (\mu - \mu')^2.$$

One-dimensional exponential family

Bandit problem:

$$\nu = (\nu_{\theta_1}, \dots, \nu_{\theta_a}, \dots, \nu_{\theta_K}),$$

one exponential family

$$\frac{d\nu_\theta}{d\rho}(x) = \exp(x\theta - b(\theta)), \quad \text{where } b(\theta) = \log \int_{\mathbb{R}} e^{x\theta} d\rho(x).$$

mean parametrization $\mu = b'(\theta)$, Pinsker-like inequality,

$$\text{kl}(\mu, \mu') := \text{KL}(\nu_\theta, \nu_{\theta'}) \geq \frac{1}{2V}(\mu - \mu')^2.$$

Regret bound for kl-UCB⁺⁺:

$$\mathbb{E}[N_a(T)] \leq \frac{\log(T)}{\text{kl}(\mu_a + \delta, \mu^* - \delta)} + O\left(\frac{\log \log(T)}{\delta^2}\right)$$

$$R_T \leq 76\sqrt{VKT} + (\mu^+ - \mu^-)K.$$

Open question: $\mathbb{P}[0, 1]$

Bandit problem with $\nu_a \in \mathbb{P}[0, 1]$:

$$\nu = (\nu_1, \dots, \nu_a, \dots, \nu_K),$$

Open question: $\mathbb{P}[0, 1]$

Bandit problem with $\nu_a \in \mathbb{P}[0, 1]$:

$$\nu = (\nu_1, \dots, \nu_a, \dots, \nu_K),$$

Empirical measure: $\hat{\nu}_a(t) = 1/N_a(t) \sum_{s=1}^t \delta_{Y_s} \mathbb{I}_{A_s=a}$

Algorithm:

Play each arm once, then for $K \leq t \leq T - 1$ play:

$$A_{t+1} \in \arg \max_a \sup \left\{ E \nu'_a : \nu'_a \in \mathbb{P}[0, 1], \text{KL}(\hat{\nu}_a(t), \nu'_a) \leq \frac{g(t, N_a(t))}{N_a(t)} \right\}.$$

Algorithm	$g(t, n)$	Minimax opt.	Asymptotic opt.
KL-UCB	$\log(t) + \log \log(t)$	✗	✓
KL-UCB+	$\log_+ \left(\frac{T}{Kn} \right)$?	?