

Lecture #6: Lower bounds and Best arm Identification

Lower bound

In Lecture 3, we proposed algorithms with (pseudo)regrets bounded as

$$R_T \leq c \sum_{k, \Delta_k > 0} \frac{\ln T}{\Delta_k} \quad (\text{instance dependent regret})$$

Is it possible to do better?

our goal: show that for "good" algorithms $\mathbb{E}[R_T(\nu, \pi)] \geq f(\nu, \pi)$ for some lower bound f .

Before that, we need to introduce some information theory tools.

Definition

let P, Q be two probability measures over (Ω, \mathcal{F})

$$KL(P, Q) = \begin{cases} +\infty & \text{if } P \text{ is not absolutely continuous wrt } Q \\ \int_{\Omega} \left(\frac{dP}{dQ} \ln \left(\frac{dP}{dQ} \right) \right) dQ = \int_{\Omega} \ln \left(\frac{dP}{dQ} \right) dP & \text{if } P \ll Q. \end{cases}$$

$Q(A) = 0 \Rightarrow P(A) = 0$

First properties:

• $KL(P, Q) \geq 0$ (by concavity of the log)

- Joint convexity: for $\lambda \in [0, 1]$ and P_1, P_2, Q_1, Q_2 .

$$KL(\lambda P_1 + (1-\lambda)P_2, \lambda Q_1 + (1-\lambda)Q_2) \leq \lambda KL(P_1, Q_1) + (1-\lambda)KL(P_2, Q_2)$$

- Not a distance: not symmetric and no triangle inequality

Theorem (data processing inequality)

Let $X, Y \in \mathcal{X}$ be random variables, let U, V be a r.v. independent from X, Y and let

$\varphi: \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{Z}$ be a measurable function. Then

$$KL(\varphi(X, U), \varphi(Y, U)) \leq KL(X, Y)$$

(we write $KL(X, Y)$ for the KL between the distributions of X and Y)

"Processing" random variables can only lose information and make them closer in KL.

Define P^X the law of X under P

$$\text{and } KL(P^{X|Y}, Q^{X|Y}) = E_{y \sim Q} [KL(P^{X|y}, Q^{X|y})]$$

Chain rule for KL

$$KL(P^{(X,Y)}, Q^{(X,Y)}) = KL(P^{X|Y}, Q^{X|Y}) + KL(P^Y, Q^Y)$$

Let $H_t = (U_0, X_{0:t}, U_1, \dots, X_{t:t}, U_t)$ be the history until time t (r.v.s U account for potential randomisation of algorithm)

A policy π is a sequence of measurable functions π_t s.t.

$$\pi_{t+1}(H_t) = a_{t+1}$$

Lemma: (fundamental inequality for stochastic bandits)

For two different bandit instances $\nu = (\nu_k)_{k \in \mathcal{K}}$ and $\nu' = (\nu'_k)_{k \in \mathcal{K}}$,

for all policies and random variables Z taking values in $[0, 1]$ that are $\mathcal{F}(H_t)$ -measurable,

$$\sum_{k=1}^K E_{\nu} [N_k(T)] KL(\nu_k, \nu'_k) = KL(P_{\nu}^{H_T}, P_{\nu'}^{H_T}) \geq KL(B_{\nu}(E_{\nu}[Z]), B_{\nu'}(E_{\nu'}[Z]))$$

Proof: The \geq is a direct consequence of the processing inequality.

with $P(H_t, U) = \mathbb{1}(U \leq Z)$.

The equality comes from the chain rule:

$$\begin{aligned} KL(P_{\nu}^{H_{t+2}}, P_{\nu'}^{H_{t+2}}) &= KL(P_{\nu}^{U_{t+2}, X_{t+2}(t+2)|H_t}, P_{\nu'}^{U_{t+2}, X_{t+2}(t+2)|H_t}) + KL(P_{\nu}^{H_t}, P_{\nu'}^{H_t}) \\ &= \mathbb{E} \left[\sum_{k=1}^K \mathbb{1}(k=a_t) KL(P_{\nu}^{X_k(t+2)}, P_{\nu'}^{X_k(t+2)}) \right] + KL(P_{\nu}^{H_t}, P_{\nu'}^{H_t}) \\ &= \sum_{a=1}^K \mathbb{E}[\mathbb{1}(k=a_t)] KL(\nu_a, \nu'_a) + KL(P_{\nu}^{H_t}, P_{\nu'}^{H_t}) \end{aligned}$$

It then follows by induction. \square

This lemma is key to proving the following lower bound.

Defn

A strategy is **consistent** w.r.t a model \mathcal{D} if, for all bandit instances $\nu \in \mathcal{D}^K$, $\forall \alpha \in (0, 1]$, $\forall \Delta_a > 0$, $\mathbb{E}[N_a(T)] = o(T^\alpha)$.

- typical bounds for good strategies: $\forall \nu \in \mathcal{D}^K$, $\forall \Delta_a > 0$, $\limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_k(T)]}{\ln T} \leq C_k(\nu)$

(remember UCB)

- optimal such term is $C_k(\nu) = \frac{1}{K_{\inf}(\nu_k, \mu^*, \mathcal{D})}$

where $K_{\inf}(\nu_k, \mu^*, \mathcal{D}) = \inf \left\{ KL(\nu_k, \nu'_k) \mid \nu'_k \in \mathcal{D}, \mathbb{E}(\nu'_k) > \mu^* \right\}$

We will only prove the lower bound for this term.

Theorem (Lai and Robbins, 1985,
Bennett and Kiefer, 1996)

For all bandit models $\mathcal{D} \in \mathcal{P}_2(\mathbb{R})$,

for any consistent strategy wrt \mathcal{D} ,

for any bandit instance $v \in \mathcal{D}^K$,

for all suboptimal arms k (i.e. $\Delta_k > 0$), $\liminf_{T \rightarrow \infty} \frac{\mathbb{E}[N_k(T)]}{\ln T} \geq \frac{1}{\text{Kinf}(v_k, \mu^*, \mathcal{D})}$.

Corollary

for all bandit models \mathcal{D} , any consistent strategy wrt \mathcal{D} , all bandit instances

$$v \in \mathcal{D}^K: \liminf_{T \rightarrow \infty} \frac{R_T}{\ln T} \geq \sum_{\substack{k \\ \Delta_k > 0}} \frac{\Delta_k}{\text{Kinf}(v_k, \mu^*, \mathcal{D})}$$

Proof:

$$\text{Kinf}(v_k, \mathcal{D}, \mu^*) = \inf \left\{ \text{KL}(v_k, v'_k) \mid v'_k \in \mathcal{D}, v_k \ll v'_k \text{ and } \mathbb{E}(v'_k) > \mu^* \right\}.$$

For the proof, we

- fix \mathcal{D} , strategy π , v and k s.t. $\Delta_k > 0$

- fix an alternative model v' with

$$\begin{cases} v'_i = v_i & \text{for all } i \neq k \\ v_k \text{ s.t. } v_k \in \mathcal{D}, v_k \ll v'_k \text{ and } \mathbb{E}(v'_k) > \mu^* \end{cases}$$

convention if $\phi = +\infty$

That is v and v' only differ at k , the unique optimal arm in v' .

We are using the fundamental inequality with

$$z = \frac{N_k(T)}{T} \quad \text{which is } [0, 1] \text{-valued}$$

$\mathcal{F}(H_T)$ -measurable

The fundamental inequality (lemma) yields, since v and v' only differ at k :

$$\begin{aligned} \mathbb{E}_v[N_k(T)] \text{KL}(v_k, v'_k) &\geq \text{KL}(\text{Ber}(\mathbb{E}_v[\frac{N_k(T)}{T}], \text{Ber}(\mathbb{E}_{v'}[\frac{N_k(T)}{T}])) \\ &\geq -\ln(2) + (1 - \mathbb{E}_v[\frac{N_k(T)}{T}]) \ln\left(\frac{1}{1 - \mathbb{E}_{v'}[\frac{N_k(T)}{T}]} \right) \end{aligned}$$

indeed $\text{KL}(\text{Ber}(p), \text{Ber}(q)) = p \ln\left(\frac{p}{q}\right) + (1-p) \ln\left(\frac{1-p}{1-q}\right)$

$$= \underbrace{p \ln\left(\frac{1}{q}\right)}_{> 0} + (1-p) \ln\left(\frac{1}{1-q}\right) + \underbrace{(p \ln(p) + (1-p) \ln(1-p))}_{\geq -\ln 2}$$

$$\geq -\ln 2 + (1-p) \ln\left(\frac{1}{1-q}\right) \quad \text{for all } (p, q) \in [0, 1] \quad (\text{and even for } p, q \in [0, 1])$$

π is consistent, so:

- instance $v \rightarrow k$ is suboptimal $\mathbb{E}_v\left[\frac{N_k(T)}{T}\right] \xrightarrow{T \rightarrow \infty} 0$

- instance $v' \rightarrow$ all $i \neq k$ are suboptimal:

for any $\alpha \in [0, 1]$, $\mathbb{E}_{v'}[N_i(T)] = o(T^\alpha)$

In particular: $T - \mathbb{E}_{v'}[N_k(T)] = \sum_{i \neq k} \mathbb{E}_{v'}[N_i(T)] = o(T^\alpha)$

so:

$$\frac{1}{1 - \mathbb{E}_v \left[\frac{N_a(T)}{T} \right]} = \frac{T}{T \cdot \mathbb{E}[N_a(T)]} = \frac{T}{o(T^\alpha)}$$

$\geq T^{1-\alpha}$ for T large enough.

Substituting back and dividing by $\ln T$: for any $\alpha \in (0, 1]$ and T large enough

$$\frac{\mathbb{E}_v[N_a(T)]}{\ln T} \geq \frac{KL(v_a, v_a)}{\ln T} \geq -\frac{\ln 2}{\ln T} + \left(1 - \mathbb{E}_v \left[\frac{N_a(T)}{T} \right]\right) \frac{\ln(T^{1-\alpha})}{\ln T}$$

thus $\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_v[N_a(T)]}{\ln T} \geq \frac{(1-\alpha)}{KL(v_a, v_a)}$ (true whether the KL is $< +\infty$ or $= +\infty$ (it is necessarily > 0))

for any $\alpha \in (0, 1]$, v

$$\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_v[N_a(T)]}{\ln T} \geq \frac{1}{KL(v_a, v_a)}$$

Holds for any $v_a \in \mathcal{D}$ s.t. $v_a \ll v_a'$ and $\mathbb{E}(v_a) > \mu^*$, so that taking the supremum of the right hand side on these v_a' yields the lower bound:

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_v[N_a(T)]}{\ln T} \geq \frac{1}{\sup_{(v_a, v_a') \in \mathcal{D}} KL(v_a, v_a')}$$

□

Comments on the Lower bound

- algorithms with optimal instance dependent bounds are known (e.g. KL-UCB, Thompson sampling) but require a long and technical analysis.
- this is an asymptotic lower bound for $T \rightarrow \infty$.
what about small T ?

Theorem (minimax lower bound)

Let $\mathcal{D} = \{ \mathcal{N}(\mu, \sigma^2) \mid \mu \in \mathbb{R} \}$, $K \geq 2$ and $T \geq K-1$

there exists a universal constant $c > 0$ such that,

for any policy π , there exists $v \in \mathcal{D}^K$ s.t.

$$R_T(\pi, v) \geq c \sqrt{KT}$$

• Case 1: $\mathcal{D} = \{ \mathcal{N}(\mu, \sigma^2) \mid \mu \in \mathbb{R} \}$

then

$$\text{King}(v_a, \mu^*, \mathcal{D}) = \frac{\Delta_a^2}{2\sigma^2}$$

Best possible regret of order $2\sigma^2 \sum_{a, \Delta_a > 0} \frac{\ln T}{\Delta_a}$

$$\text{UCB has regret} \leq 32\sqrt{2} \sum_{k, \Delta_k > 0} \frac{\ln T}{\Delta_k}$$

↳ optimal up to constant factor
can be made optimal with finer version

• Case 2: $\mathcal{D} = \{ \text{Bern}(p) \mid p \in [0, 1] \}$

then

$$\text{Kiv}(\nu_k, \mu^*, \mathcal{D}) = \mu_k \ln \frac{\mu_k}{\mu^*} + (1 - \mu_k) \ln \frac{1 - \mu_k}{1 - \mu^*}$$

Best arm identification

Until now: maximise cumulative reward

→ exploration/exploitation trade-off.

In some applications, there is no price for exploring.

Think for example of a researcher testing drugs on mice/artificial human cells

or testing products on some people before commercialisation.

Share similarities with regret minimisation, but good algorithms are actually different.

Example: simple regret minimisation:

$$\Delta_{T+1} = \mu^* - \mu_{\alpha_{T+1}}$$

Setting 1: BAI, fixed confidence

At each round $t=1, \dots, \infty$:

- agent picks an arm $\alpha_t \in [K]$ (based on previous observations)
- observes $X_{\alpha_t}(t) \sim \nu_{\alpha_t} \in \mathcal{D}$
- decides whether to continue sampling or stop

If stop: return a final choice $\Psi \in [K]$

The (random) stopping time is called τ

new

Question: which arm has the highest mean?

Goal: 1) Have a δ -correct strategy: $\mathbb{P}(\tau < \infty \text{ and } \mu_\Psi < \mu^*) \leq \delta$ (for any $\delta \in (0, 1)$)
with confidence level $\delta \in (0, 1)$ confidence level

2) minimize the exploration time $\mathbb{E}[\tau]$

Theorem (lower bound)

Let (π, τ, Ψ) be a δ -correct strategy for the bandit model $\mathcal{D} = \{ \mathcal{N}(\mu, 1) \mid \mu \in \mathbb{R} \}$ and let $\nu \in \mathcal{D}^K$. Then:

$$\mathbb{E}[\tau] \geq c^*(\nu) \ln\left(\frac{1}{4\delta}\right) \quad \text{where}$$

$$c^*(\nu)^{-1} = \sup_{\alpha \in [K]} \left(\inf_{\mu' \in \text{ALT}(\mu)} \sum_{k=1}^K \alpha_k \frac{(\mu_k - \mu')^2}{2} \right)$$

$$\text{where } \text{ALT}(\mu) = \left\{ \mu' \in \mathbb{R}^K \mid \arg\max_k \mu_k \cap \arg\max_k \mu'_k = \emptyset \right\}$$

i.e. no arm is optimal for both ν and ν'

Proof is similar to regret lower bound, but with v' given by $\mu' \in \text{ALT}(\mu)$

$$Z = \mathbb{1}_{\left\{ \tau < \infty \text{ and } \psi \notin \arg\max_k \mathbb{E}(v_k') \right\}}$$

Remarks

- α_k represents the "optimal" fraction of pulls on arm k .

It indeed appears in the proof that $c^+(v) \geq \sum_{k=1}^K \mathbb{E} \left[\frac{N_k(\tau)}{\tau} \right] \text{KL}(v_k, v_k')$

- This suggests an optimal non-adaptive sampling allocation:

$$N_k(\tau) = \tau \alpha_k^* \quad \text{with } \alpha^* = \arg\max_{\alpha \in \mathcal{P}_K} \left(\inf_{\mu' \in \text{ALT}(\mu)} \sum_{k=1}^K \alpha_k \frac{(\mu_k - \mu_k')^2}{2} \right)$$

i.e. we should stop when

$$\underbrace{\inf_{\mu' \in \text{ALT}(\mu)} \sum_{k=1}^K N_k(t) \frac{(\mu_k - \mu_k')^2}{2}}_{Z_t} \geq \ln\left(\frac{1}{\delta}\right)$$

Problem μ is unknown, but can be estimated by $\hat{\mu}$

In that case, we can approximate Z_t^* by

$$Z_t := \frac{1}{2} \inf_{\mu' \in \text{ALT}(\hat{\mu})} \sum_{k=1}^K N_k(t) (\hat{\mu}_k(t) - \mu_k')^2$$

Track-and-stop algorithm

Input δ and $\beta_T(\delta)$

Pull all arms once

While $Z_t < \beta_T(\delta)$

if $\min_k N_k(t) \leq \sqrt{t}$ then pull $a_{t+1} \in \arg\min_k N_k(t)$ forced exploration

else choose $a_{t+1} \in \arg\max_k \hat{F}_k(t) - N_k(t)$ track

stop and return $\psi \in \arg\max_k \hat{\mu}_k(t)$ stop.

Theorem

If $\mathcal{Z} = \{N(\mu, \sigma) \mid \mu \in \mathbb{R}\}$, Track-and-Stop is δ -correct and

asymptotically optimal for $\beta_T(\delta) \approx K \ln(t) + \ln\left(\frac{1}{\delta}\right)$:

$$\limsup_{\delta \rightarrow 0} \frac{\mathbb{E}_\mu[\mathbb{E}]}{\ln\left(\frac{1}{\delta}\right)} \leq \frac{1}{\sup_{\alpha \in \mathbb{P}_K} \inf_{\mu \in \text{EAB}(\mu)} \sum_k \alpha_k \frac{(\mu_k - \mu_k^*)^2}{z}}$$

Omitted proof relies on heavy concentration bounds

Remarks

- Track-and-Stop is computationally intensive due to the oracle allocation computation (\hat{a})
- Forced exploration is wasteful in practice
(we could use optimism instead)
- still no good bound for fixed $\delta > 0$ (e.g. $\delta = 0.05$)

Setting 2: BAI, fixed budget

At each round $t=1, \dots, T$:

- agent picks an arm $a_t \in [K]$ (based on previous observations)
- observes $X_{a_t}(t) \sim \nu_{a_t} \in \mathcal{J}$

At time T : return a final choice $\Psi \in [K]$

Goal: minimize $P(\mu^\star > \mu_\Psi)$

Complexities: $H_1 = \sum_{\Delta_k > 0} \frac{1}{\Delta_k}$, $H_2 = \max_{k, \Delta_k > 0} \frac{k}{\Delta_k}$
ordered values: $\Delta_{(1)} \leq \Delta_{(2)} \leq \dots \leq \Delta_{(K)}$

$$H_2 \leq H_1 \leq \ln(2K) H_2$$

Lower bound of order $\begin{cases} \exp(-\frac{T}{\ln(K) H_2}) & \text{if } H_2 \text{ unknown} \\ \exp(-\frac{T}{H_2}) & \text{if known} \end{cases}$

First approach: uniform exploration of the arms. Good baseline, but not very good when arms have very different means.

Sequential Halving:

Set $L = \lceil \log_2(K) \rceil$ and $A_1 = [K]$

For $l=1, \dots, L$:

Pull each arm in A_l $T_l = \lfloor \frac{T}{|A_l|} \rfloor$ times

Let $\hat{\mu}_i^l$ be the empirical mean of arm i based only on these last T_l samples

Let A_{l+1} contain the top $\lfloor \frac{|A_l|}{2} \rfloor$ arms in A_l

Return Ψ as the only arm in A_{L+1}

Thm:

If the distributions are 1-sub-Gaussian, then Sequential Halving satisfies:

$$P(\mu^* > \mu_\Psi) \leq 3 \log_2(K) \exp\left(-\frac{T}{16 H_2 \log_2(K)}\right).$$

Remarks

- close to lower bound
- For uniform exploration, we can bound this probability by

$$\sum_{k: \Delta_k > 0} \exp\left(-\frac{\lfloor K \rfloor \Delta_k^2}{4}\right).$$

- VE slightly better than SH when $\Delta_k = \Delta$ for any k ($H_2 = \frac{K}{\Delta^2}$)

but SH much better than VE when $\Delta_2 = \Delta \ll 1$ ($H_2 = \frac{1}{\Delta^2}$)
 $D_k = 1$ for $k \geq 2$