

Lecture #6: Lower bound

Recall bandit setting

- to each arm k is associated a probability distribution $\nu_k \in \mathcal{D}$.

- \mathcal{D} is the bandit model ($\mathcal{D} \subset \mathcal{P}_1(\mathbb{R})$)

- A bandit instance is denoted by $\nu = (\nu_k)_{k \in [K]}$

- Goal: minimise the regret, which can be rewritten as:

$$R_T = \sum_{k=2}^K \Delta_k \mathbb{E}[N_k(T)]$$

Bounding the regret \Leftrightarrow bounding $\mathbb{E}[N_k(T)]$.

what are the best possible (by an algorithm) bounds?

- what is a randomised strategy π ?

a sequence of measurable functions $(\pi_t)_{t \geq 1}$ with

$$\pi_{t+1}: \underbrace{H_t = (U_0, X_{0_1}(1), U_1, \dots, X_{0_t}(t), U_t)}_{\text{history of observations + randomisation for first } t \text{ rounds}} \mapsto \underbrace{\pi_{t+1}(H_t) = a_{t+1}}_{\text{arm picked at } t+1}$$

- a strategy is consistent w.r.t a model \mathcal{D} if,

for all bandit instances $\nu \in \mathcal{D}^K$, $\forall \alpha \in (0, 1]$, $\forall k$ s.t. $\Delta_k > 0$,

$$\mathbb{E}[N_k(T)] = o(T^\alpha)$$

Defn 1

for well behaved models, there exist consistent strategies
eg UCB with $\mathcal{D} = \mathcal{P}([0,1])$.

(asymptotic)
- typical bounds for good strategies

$$\forall v \in \mathcal{D}^K, \forall k \text{ s.t. } \Delta_k > 0, \limsup_{T \rightarrow \infty} \frac{\mathbb{E}[N_k(T)]}{\ln T} \leq \underbrace{C_k(v)}_{\text{problem dependent term}}$$

- optimal such term: $C_k(v) = \frac{1}{K_{\text{inf}}(v_k, \mu^*, \mathcal{D})}$

problem dependent
term

$$\text{where } K_{\text{inf}}(v_k, \mu^*, \mathcal{D}) = \inf \left\{ KL(v_k \| v_k') \mid \begin{array}{l} v_k' \in \mathcal{D} \\ \mathbb{E}(v_k') > \mu^* \end{array} \right\}$$

we will now prove one part of this optimality: a lower bound on $C_k(v)$.

Theorem (Lai and Robbins, 1985,
Burnetas and Katehakis, 1996)

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

for any consistent strategy w.r.t. \mathcal{D} ,

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

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

Conclusion

for all bandit models \mathcal{D} , any consistent strategy w.r.t \mathcal{D} , all bandit instances $v \in \mathcal{D}^K$:

$$\liminf_{T \rightarrow \infty} \frac{R_T}{\ln T} \geq \sum_{\substack{a, \\ \Delta a > 0}} \frac{\Delta a}{\text{Kull}(v_a, \mu^*, \mathcal{D})}$$

To prove this theorem (and other lower bounds), we need the following fundamental inequality.

Notation: for a strategy π , we note $H_t = (U_0, X_{a_1}(1), U_1, X_{a_2}(2), \dots, X_{a_t}(t), U_t)$

Recall that a_{t+1} is $\sigma(H_t)$ -measurable.
depends on π

Lemma: (fundamental inequality for stochastic bandits)

For all bandit problems $v = (v_k)_{k \in [K]}$ and $v' = (v'_k)_{k \in [K]}$ in \mathcal{D}^K with $v_k \ll v'_k$ for all k ,

for all strategies and random variables Z taking values in $[0, 1]$ that are $\sigma(H_t)$ -measurable,

$$\sum_{k=1}^K \mathbb{E}_v[N_k(T)] \text{KL}(v_k, v'_k) = \text{KL}(\mathbb{P}_v^{H_T}, \mathbb{P}_{v'}^{H_T}) \geq \text{KL}(\text{Bu}(\mathbb{E}_v[Z]), \text{Bu}(\mathbb{E}_{v'}[Z]))$$

law of H_t under v (and π)

dependence in strategy π hidden everywhere here.

Note: this lemma is our key to perform an implicit change of measures

in the proof of the theorem.

Proof of the theorem (based on the lemma)

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

convention if $\phi = +\infty$

This is why we will:

- fix \mathcal{D} , strategy π , v and k s.t. $\Delta_k > 0$ (π is consistent w.r.t. \mathcal{D})

- fix an alternative model v' with

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

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

- take $Z = \frac{N_k(T)}{T}$ which is $[0, 1]$ -valued $\sigma(H_T)$ -measurable

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

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

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

$$= p \ln\left(\frac{1}{q}\right) + (1-p) \ln\left(\frac{1}{1-q}\right) + (p \ln(p) + (1-p) \ln(1-p))$$

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

π is consistent, so:

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

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

so: $\frac{1}{1 - \mathbb{E}_{v'}\left[\frac{N_a(T)}{T}\right]} = \frac{T}{T - \mathbb{E}_{v'}[N_a(T)]} = \frac{T}{o(T^\alpha)}$

$\geq T^{1-\alpha}$ 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 usually > 0))

for any $\alpha \in (0, 1]$, so $\liminf_{T \rightarrow +\infty} \frac{\mathbb{E}_v[N_a(T)]}{\ln T} \geq \frac{1}{KL(v_a, v'_a)}$

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

$$\liminf_{T \rightarrow \infty} \frac{\mathbb{E}_v[N_a(T)]}{L_n T} \geq \frac{1}{\text{Kinf}(v_a, \mathcal{D})_{\mu^*}} \quad \square$$

Proof of the Lemma

• The inequality \geq is a direct application of the data processing inequality with expectations.

• For the equality: and same for v' : $P_v'^{HT} = K_T'(K_{T-1}' \dots)$

(1) we show by induction that $P_v^{HT} = K_T(K_{T-1}(\dots(K_2 \lambda_0)))$

we check below that \downarrow
is regular

$$\left\{ \begin{array}{l} \text{where } K_T \text{ is the transition kernel:} \\ h \in [0, 1] \times (\mathbb{R} \times [0, 1])^{T-1} \mapsto K_T(h, \cdot) = \nu_{\frac{h}{T}} \otimes \lambda_0 \end{array} \right. \quad \begin{array}{l} \text{with } \nu_{\frac{h}{T}} \text{ (E.O. 2.5)} \\ \text{with } \nu_T \sim \lambda_0 \end{array}$$

prob. measure on $\mathbb{R} \times [0, 1]$

$$\underline{T=0}: \quad H_0 = U_0 \sim \lambda_0; \quad P_v^{U_0} = \lambda_0$$

$$\underline{T \rightarrow T+1} \quad \forall A \in \mathcal{B}([0, 1] \times (\mathbb{R} \times [0, 1])^T), \quad \forall B' \in \mathcal{B}(\mathbb{R}), \quad \forall B \in \mathcal{B}([0, 1]);$$

$$P_v^{H_{T+1}}(A \times B' \times B) = P_v(H_T \in A \text{ and } X_{\frac{H_T}{T+1}} \in B' \text{ and } U_{T+1} \in B)$$

$$= \mathbb{E}_v [\mathbb{1}_A(H_T) P_v[X_{\frac{H_T}{T+1}} \in B' \text{ and } U_{T+1} \in B | H_T]]$$

tower rule \downarrow

$$= E_{\nu} [\mathbb{1}_A(H_T) \cdot \nu_{\pi_{T+1}(H_T)}(B) \cdot \lambda_0(B)]$$

← defn of the model and strategy

$$= E_{\nu} [\mathbb{1}_A(H_T) K_{T+1}(H_T, B' \times B)]$$

↓ defn of K_{T+1}

$$= \int \mathbb{1}_A(h) K_{T+1}(h, B' \times B) dP_{\nu}^{H_T}(h)$$

↓ rewriting

↓ defn of $K_{T+1} P_{\nu}^{H_T}$

$$= K_{T+1} P_{\nu}^{H_T}(A \times B' \times B)$$

→ we've shown the induction

(2) we check that the assumptions of the chain rule are satisfied.

• the K_t are regular transition kernels: $\forall E \in \mathcal{B}(\mathbb{R}) \otimes \mathcal{B}([0,1])$,

$$h \mapsto K_t(h, E) = \sum_{k=1}^K \mathbb{1}_{\{\pi_t(h)=k\}} (\nu_k \otimes \lambda_0)(E) \quad \text{is measurable as}$$

π_t is measurable (with respect to considered spaces)

• Assumption (i): $\forall h, K_t(h, \cdot) \ll K'_t(h, \cdot)$ as $\forall k, \nu_k \ll \nu'_k$ by ass.

• Assumption (ii): $(h, (y, u)) \mapsto \frac{dK_t(h, \cdot)}{dK'_t(h, \cdot)}(y, u) = \sum_{k=1}^K \mathbb{1}_{\{\pi_t(h)=k\}} \frac{d\nu_k}{d\nu'_k}(y)$

is indeed bi-measurable (product of measurable functions)

(3) we then may apply the chain rule and show by induction the desired result based on:

$$- KL(P_{\nu}^{H_0}, P_{\nu'}^{H_0}) = KL(\lambda_0, \lambda_0) = 0$$

$$- \text{for } t \geq 0, KL(P_{\nu}^{H_{t+1}}, P_{\nu'}^{H_{t+1}}) = KL(K_{t+1} P_{\nu}^{H_t}, K'_{t+1} P_{\nu'}^{H_t})$$

↓ chain rule

$$= KL(P_{\nu}^{H_t}, P_{\nu'}^{H_t}) + \int KL(K_{t+1}(h, \cdot), K'_{t+1}(h, \cdot)) dP_{\nu}^{H_t}(h)$$

$$= \text{KL}(P_{\nu}^{H_T}, P_{\nu'}^{H_T}) + \int \text{KL}(v_{T_{r,h}(h)} \otimes \lambda_0, v'_{T_{r,h}(h)} \otimes \lambda_0) dP_{\nu}^{H_T}(h)$$

$$= \text{KL}(P_{\nu}^{H_T}, P_{\nu'}^{H_T}) + \sum_{k=1}^K \text{KL}(v_k, v'_k) \cdot \underbrace{\int \mathbb{1}_{T_{r,h}(h)=k} dP_{\nu}^{H_T}(h)}_{\mathbb{E}[\mathbb{1}_{(a_{r+1}(H_t)=k)}]} = \mathbb{E}[\mathbb{1}_{(a_{r+1}=k)}]$$

$$= \text{KL}(P_{\nu}^{H_T}, P_{\nu'}^{H_T}) + \sum_{k=1}^K \text{KL}(v_k, v'_k) \mathbb{E}[\mathbb{1}_{(a_{r+1}=k)}]$$

by induction

$$\text{KL}(P_{\nu}^{H_T}, P_{\nu'}^{H_T}) = \sum_{k=1}^T \sum_{h=1}^K \text{KL}(v_h, v'_h) \mathbb{E}[\mathbb{1}_{a_r=k}]$$

$$= \sum_k \text{KL}(v_k, v'_k) \mathbb{E}[N_k(T)]$$

□

Comments on the lower bound

- see a comparison with our upper bounds in exercise session #4.
- 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 ? \rightarrow see exercise session #4

what if we fix T and choose arbitrarily the bandit instance v ?

Theorem (minimax lower bound)

Let $\mathcal{D} = \{ \mathcal{N}(\mu, 1) \mid \mu \in \mathbb{R} \}$, $K \geq 2$ and $T \geq K-1$. Then, there exists a universal constant $c > 0$ such that,

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

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

Proof in exercise session #4.

$$\text{minimax} \iff \min_{\pi} \max_{v \in \mathcal{D}^K} R_T(\pi, v) \geq c \sqrt{KT}$$