

## Minimax Lower Bounds

Source: Chap 15, Bandit Algorithms, L & S.

Recall from Notes (Set 10) — Minimax bound (i.e., showing existence of an environment for which regret is lower bounded by  $\Theta(\sqrt{n})$ , where  $n$  is the time horizon).

This set (Chap 15): Formally deriving the result for a k-arm unstructured bandit environment.

First, divergence decomposition lemma.

Setting: k-armed, environment  $\mathcal{E}$ .

$$v = (P_1, P_2 \dots P_k) \in \mathcal{E} \text{ , and } v' = (P'_1, P'_2 \dots P'_k) \in \mathcal{E}$$

Fix a policy  $\pi$ , and a time-horizon  $n$ .

Let  $P_v \equiv P_{\pi v} \equiv P_{\pi v}^{(n)} \rightarrow$  probability measure induced on  $(A_1, X_1, A_2, X_2, \dots, A_n, X_n)$

under environment  $v \in \mathcal{E}$  and policy  $\pi$   
 (i.e., randomness due to environment  $v$  as well as  
 policy  $\pi$ , and with dependence over time  
 steps induced by  $\pi$ ).

Similarly,  $P_{v'} \equiv P_{v\pi} \equiv P_{v'\pi}^{(n)} \rightarrow$  prob. measure  
 due to  $(v', \pi)$  over  $n$  time-steps.

Let  $E_v[\cdot] = E_{P_{v\pi}^{(n)}}[\cdot]$ ,

and  $E_{v'}[\cdot] = E_{P_{v'\pi}^{(n)}}[\cdot]$ .

Lemma (13.1 in text): Setting as above. Then,

$$D(P_v, P_{v'}) = \sum_{j=1}^k E_v[T_j(n)] D(P_j, P'_j).$$

Proof: Assume  $D(P_j, P'_j) < \infty \quad \forall j=1, 2, \dots, k$ .

$$D(P_v, P_{v'}) = E_v \left[ \ln \left( \frac{P_v}{P_{v'}} \right) \right] \quad \begin{array}{l} \text{(see text for} \\ \text{more careful)} \\ \text{reasoning of existence/finiteness} \end{array}$$

$$P_{\nu\pi}(A_1=a_1, X_1=x_1, \dots, A_n=a_n, X_n=x_n)$$

$$= P_{\nu\pi}(a_1, x_1, a_2, x_2, \dots, a_n, x_n)$$

$$= \prod_{t=1}^n \pi_t(a_t | a_1, x_1, \dots, a_{t-1}, x_{t-1}) \cdot P_{a_t}(x_t),$$

and  $\pi_t(a_t | a_0, x_0) \equiv \pi_t(a_t)$  (slightly misusing notation).

Similarly,

$$P_{\nu'\pi}(a_1, x_1, \dots, a_n, x_n) = \prod_{t=1}^n \pi_t(a_t | a_1, x_1, \dots, a_{t-1}, x_{t-1}) P'_{a_t}(x_t)$$

Note: Policy is same in both environments. Thus, given the same sequence of observations of past actions and rewards, the next action has the same dist. under either  $\nu$  or  $\nu'$  environments.

$$\therefore \frac{dP_\nu}{dP_{\nu'}}(a_1, x_1, \dots, a_n, x_n) = \prod_{t=1}^n \left( \frac{P_{a_t}(x_t)}{P'_{a_t}(x_t)} \right)$$

$$D(P_0, P_0') = E_{\nu} \left[ \ln \frac{dP_0}{dP_0'} \right] = \sum_{t=1}^n E_{\nu} \left[ \ln \left( \frac{P_{A_t}(x_t)}{P'_{A_t}(x_t)} \right) \right]$$

$$= \sum_{t=1}^n E_{\nu} \left[ E_{\nu} \left[ \ln \left( \frac{P_{A_t}(x_t)}{P'_{A_t}(x_t)} \right) \middle| A_t \right] \right]$$

$D(P_{A_t}, P'_{A_t})$

$$= \sum_{t=1}^n E_{\nu} \left[ D(P_{A_t}, P'_{A_t}) \right]$$

$$= \sum_{t=1}^n E_{\nu} \left[ \sum_{j=1}^k \chi_{\{A_t=j\}} D(P_j, P'_j) \right]$$

$$= \sum_{j=1}^k E_{\nu} \left[ \sum_{t=1}^n \chi_{\{A_t=j\}} D(P_j, P'_j) \right]$$

$$= \sum_{j=1}^k D(P_j, P'_j) E_{\nu} \left[ \underbrace{\sum_{t=1}^n \chi_{\{A_t=j\}}}_{T_j(n)} \right]$$

$$= \sum_{j=1}^k E_{\nu} [T_j(n)] D(P_j, P'_j)$$
⊗

Rest of this section: Assume the arms are 1-Gaussian (i.e.,  $\Sigma = \Sigma_N^k(1)$ ).

Then, if  $P_j = N(\mu_j, 1)$ ,  $P'_j = N(\mu'_j, 1)$ ,

$$D(P_j, P'_j) = \frac{(\mu_j - \mu'_j)^2}{2}$$

(Recall KL divergence for Gaussians in Set 10).

Theorem (15.2 in text): For any policy  $\pi$ ,  
 $\exists \mu = (\mu_1, \dots, \mu_k) \in [0, 1]^k$  with  
 $P_{\nu_\pi} = (P_1, P_2, \dots, P_k)$ ,  $P_j \sim N(\mu_j, 1)$

such that  $R_n(\pi, \nu_\pi) \geq \frac{1}{27} \sqrt{(k-1)n}$

Proof: Fix a policy  $\pi$ .

Let  $\Delta = \sqrt{\frac{k-1}{4n}}$  (observe  $\Delta < \frac{1}{2}$   
 $\forall k \geq 1, n \geq 1$ ).

<u>System 1</u> ( $\nu$ )		<u>System 2</u> ( $\nu'$ )
arm 1 2 3 ... k		1 2 3 ... $i^*$ ... k
• • • - - - -		• • • - - - -
$\Delta$ 0 0 ... 0		$\Delta$ 0 0 ... $2\Delta$ 0
$\downarrow$ $M_1$		$\downarrow$ $M_1$
$\downarrow$ $M_k$		$\downarrow$ $M_{i^*}$
		$\downarrow$ $M_k$
$i^* = \arg \min \mathbb{E}_\nu [T_j(n)]$		
$\therefore j > 1$		
- - - - - - - - -		

Choose the arm that is EXPECTED to be played least frequently under  $(\pi, \nu)$  (System 1) to be the changed arm in System 2.

$$\nu: (p_1, \dots, p_k) = (N(\Delta, 1), N(0, 1), \dots, N(0, 1))$$

$$\nu': (p'_1, \dots, p'_k) = \underbrace{(N(\Delta, 1), N(0, 1), \dots, N(2\Delta, 1), \dots, N(0, 1))}_{1 \dots i^* \dots k}$$

$$\text{Now, } R_n(\pi, \nu) = \sum_{j=1}^k \Delta_j \mathbb{E}_\nu [T_j(n)]$$

$$\geq \Delta \left(\frac{n}{2}\right) P_\nu(T_1(n) \leq n/2).$$

Similarly,

$$R_n(\pi, \nu') \geq P_{\nu'}(T_1(n) > \frac{n}{2}) \cdot \binom{n}{2} \cdot \Delta.$$

$\overbrace{\quad \quad \quad}$  strict lower bound on expected  
number of times arm 1  
is played.  $\overbrace{\quad \quad \quad}$  regret  
incurred  
whenever  
arm 1  
played

Now, let  $P = P_\nu$ ,  $Q = P_{\nu'}$ ,

$$A = \left\{ T_1(n) \leq \frac{n}{2} \right\}.$$

As discussed in Set 10,  $A$  is a "bad event" under environment  $\nu$  (the optimal arm is played for less than half the time steps), and  $A^c$  is a bad event under  $Q$  (the first arm, which is sub-optimal under  $\nu'$  is played more than half the time steps).

From B-H inequality,  $P(A) + Q(A^c) \geq e^{-D(P, Q)} \frac{2}{2}$ .

i.e.,  $P_\nu(T_1(n) \leq \frac{n}{2}) + P_{\nu'}(T_1(n) > \frac{n}{2}) \geq e^{-D(P_\nu, P_{\nu'})} \frac{2}{2}$

Multiply by  $\frac{n\Delta}{2}$ :

$$\leq R_n(\pi, \nu) \quad \leftarrow R_n(\pi, \nu')$$

$$\underbrace{\frac{n\Delta}{2} P_\nu(T_1(n) \leq \frac{n}{2})}_{\text{Left side}} + \underbrace{\frac{n\Delta}{2} P_\nu(T_1(n) > \frac{n}{2})}_{\text{Right side}}$$

$$\geq \frac{n\Delta}{4} e^{-D(P_\nu, P_{\nu'})}.$$

$$\Rightarrow R_n(\pi, \nu) + R_n(\pi, \nu') \geq \frac{n\Delta}{4} e^{-D(P_\nu, P_{\nu'})}$$

$\downarrow \rightarrow \textcircled{X}$ .

Next

$$D(P_\nu, P_{\nu'}) = \sum_{j=1}^n E_\nu[T_j(n)] D(P_j, P'_j)$$

$$= E_\nu[T_{i^*}(n)] (2\Delta^2)$$

$$= \begin{cases} 0 & j \neq i^* \\ \frac{(2\Delta)^2}{2} & j = i^* \end{cases}$$

Finally, observe that

$$\sum_{j=2}^n E_{\nu} [\tau_j(n)] \leq n.$$

Also,  $i^* = \arg \min_{j \geq 2} E_{\nu} [\tau_j(n)]$

Pigeon hole principle  $\Rightarrow E_{\nu} [\tau_{i^*}(n)] \leq \frac{n}{k-1}$ .

$$\therefore D(P_{\nu}, P_{\nu'}) \leq (2\Delta^2) \cdot \frac{n}{k-1} = \frac{1}{2}.$$

From  $\textcircled{1}$  and  $\textcircled{**}$ ,

$\rightarrow \textcircled{**}$ .

$$R_n(\pi, v) + R_n(\pi, v') \geq \frac{n\Delta}{4} \cdot e^{-1/2} \quad \underline{0.0379}$$

$$= \sqrt{(k-1)n} \left( \frac{e^{-1/2}}{8} \right)$$

(Pigeon hole principle)

$\Rightarrow$

$$\max \{ R_n(\pi, v), R_n(\pi, v') \} \geq \frac{1}{2} \left( \sqrt{(k-1)n} \cdot \frac{e^{-1/2}}{8} \right)$$

i.e., either for  $\tilde{v} = v$  or  $\tilde{v} = v'$ , we have

$$\begin{aligned} R_n(\pi, \tilde{v}) &\geq \sqrt{n(k-1)} \cdot \left( \frac{e^{-1/2}}{16} \right) \\ &\geq \frac{\sqrt{n(k-1)}}{27} \end{aligned}$$



Note: The only place in the proof where Gaussian noise was used was in the step:

$$D(P_{i^*}, P'_{i^*}) = \frac{(2\Delta)^2}{2}$$

Recall from the KL-UCB discussion (Set 7) that for  $P_{i^*} \sim \text{Bernoulli}(\mu)$ ,  $P'_{i^*} \sim \text{Bernoulli}(\mu')$ .

$$d(\mu, \mu') = \frac{\Delta^2}{2\mu(1-\mu)} + o(\Delta^2),$$

and thus we again have quadratic scaling in  $\Delta$ . Thus, we can construct a minimax LB with Bernoulli rvs. Please read 15.3, Note 1 for additional discussion on more general rvs.