

Instance Dependent Lower Bound

Text (Bandit Algorithms, Chap 16) approach:

- Asymptotic (i.e., time horizon $\rightarrow \infty$) lower bound for any "reasonable" policy (formally consistent policy, see Defn 16.1) for unstructured stochastic bandits.

$$d_{\inf}(P, M, M) = \inf_{P' \in M} \left\{ D(P, P'): \mu(P') > \mu \right\},$$

mean when arm dist
 is P'

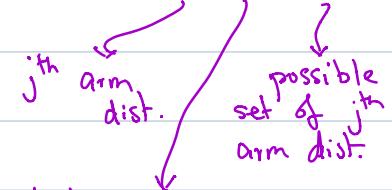
Then,

$$\lim_{n \rightarrow \infty} \frac{R_n}{\ln(n)} \geq \sum_{j: \Delta_j > 0} \frac{\Delta_j}{d_{\inf}(P_j, M_1, M_j)}$$

Roughly,

$$\frac{R_n}{\ln(n)} \geq \sum_{j: \Delta_j > 0} \left(\frac{\Delta_j}{d(P_j, P_1)} \right)$$

asympt.



 best arm mean

- For finite horizon, regret bounds for Gaussian environment with additional conditions on reasonableness of policy,

In this class, I will focus directly on finite horizon bounds. Setting of 2-armed Bernoulli bandit with means μ_1, μ_2 , with $\mu_1 - \mu_2 = \Delta > 0$.

Sources: Thanks to Rajat Sen and Subhashini Krishnasamy (Ph.D. grads from UT) for discussions and notes based on:

- ① S. Bubeck, V. Perchet and P. Rigollet, "Bounded regret in stochastic multi-armed bandits", arXiv:1302.1611
 - ② R. Combes, C. Jiang and R. Srikant, "Bandits with budgets: Regret lower bounds and optimal algorithms", ACM Sigmetrics 2015.
-

Setting: k armed stochastic Bernoulli bandit with unstructured environment $\nu \in \Sigma$,
 $\Sigma = \{ \nu : \nu = (p_1, p_2, \dots, p_k), p_j \sim \text{Bernoulli}(\mu_j) \}$

First, we restrict to "reasonable" policies:

α -consistent Policy: A policy is α -consistent, with $\alpha \in (0, 1)$, if for any $\nu \in \Sigma$, $k \neq k^*$, where $k^* = \operatorname{argmax} \{\mu_1, \mu_2, \dots, \mu_k\}$, we have

$$E \left[\sum_{s=1}^t \mathbb{X}_{\{A_s=k\}} \right] \leq C t^\alpha, \text{ for some } C > 0,$$

w.r.t (\mathcal{D}, π)

and $t \geq T_2$.

Discussion: This condition states that the policies we consider are "reasonable" in the sense that the policy discovers sub-optimal arms quickly enough such that it plays these bad arms only infrequently (sub linear with respect to time).

e.g. Policy: Play arm 1 always. This is a bad policy because for $\nu = (\text{Bernoulli}(0.1) \text{ Bern}(0.3))$, this will incur linear regret. The α -consistency condition disallows such a policy.

* All policies we have considered so far for stochastic bandits are α -consistent.

* Suppose a policy π is NOT α -consistent. Then $\exists \nu \in \mathcal{E}$ s.t. $R_n(\pi, \nu)$ incurs $\geq \Omega(t^\alpha)$ regret.
Henceforth, we restrict to α -consistent policies.

Proof below for a two armed bandit, with arms μ_1, μ_2 with $1 > \mu_1 > \mu_2 > 0$.

Theorem: Suppose $\nu = (P_1, P_2)$, with $P_i \sim \text{Bernoulli}(\mu_i)$, with $1 > \mu_1 > \mu_2 > 0$, and π is any α -consistent policy, then for any $\delta \in (0, 1 - \mu_1)$, $n \geq \tau_2$,

$$\begin{aligned} R_n(\pi, \nu) &= \Delta_2 E[\tau_2(n)] && \text{time horizon.} \\ &\geq \frac{\Delta}{d(\mu_2, \mu_1 + \delta)} \left((1-\alpha) \ln(n) - \ln(8C) \right) \\ &\quad \text{from } \alpha\text{-consistent policy defn.} \\ &\quad \text{KL divergence for Bernoulli r.v.s} \end{aligned}$$

Proof:

System 1: ν <table style="margin-left: 20px; border-collapse: collapse;"> <tr> <td style="border-bottom: 1px solid black;">arm 1</td> <td style="border-bottom: 1px solid black;">arm 2</td> </tr> <tr> <td style="text-align: center;">•</td> <td style="text-align: center;">•</td> </tr> <tr> <td style="text-align: center;">μ_1</td> <td style="text-align: center;">μ_2</td> </tr> </table>	arm 1	arm 2	•	•	μ_1	μ_2	System 2: ν' <table style="margin-left: 20px; border-collapse: collapse;"> <tr> <td style="border-bottom: 1px solid black;">arm 1</td> <td style="border-bottom: 1px solid black;">arm 2</td> </tr> <tr> <td style="text-align: center;">•</td> <td style="text-align: center;">•</td> </tr> <tr> <td style="text-align: center;">μ_1</td> <td style="text-align: center;">$\mu_1 + \delta$</td> </tr> </table>	arm 1	arm 2	•	•	μ_1	$\mu_1 + \delta$
arm 1	arm 2												
•	•												
μ_1	μ_2												
arm 1	arm 2												
•	•												
μ_1	$\mu_1 + \delta$												

$$E_{\nu, \pi} [\tau_2(n)] \leq C n^\alpha$$

$\# n \geq \tau_2$

$$E_{\nu', \pi} [\tau_1(n)] \leq C n^\alpha$$

$\# n \geq \tau_2$

Above bounds follows from α -consistency of policy π .

Let $P = \text{prob. measure induced by } (\mathcal{V}, \pi)$ over time horizon n . (i.e., dist. over $A_1, X_1, \dots, A_n, X_n$)

$$P \equiv P_{\mathcal{V}, \pi} (A_1=a_1, X_1=x_1, A_2=a_2, \dots, A_n=a_n, X_n=x_n)$$

$Q = \text{prob. measure induced by } (\mathcal{V}', \pi)$ over time horizon n .

$$A = \left\{ T_2(n) > \frac{n}{2} \right\}. \quad \begin{cases} A \text{ is a bad event} \\ \text{under } P; A^c \text{ is} \\ \text{bad under } Q \end{cases}$$

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

Markov's inequality

Now,

$$P(A) = P\left(T_2(n) > \frac{n}{2}\right) \leq \frac{E_{\mathcal{V}, \pi}[T_2(n)]}{\left(\frac{n}{2}\right)}$$

$$\leq \frac{2Cn^\alpha}{n} = \frac{2C}{n^{1-\alpha}}$$

Similarly,

$$Q(A^c) = Q\left(T_2(n) \leq \frac{n}{2}\right) = Q\left(T_1(n) > \frac{n}{2}\right)$$

$$\leq 2C/n^{1-\alpha}.$$

$$\therefore \frac{4C}{n^{-\alpha}} \geq \frac{e^{-D(P,Q)}}{2}$$

$$\Rightarrow D(P,Q) \geq (\alpha) \ln(n) - \ln(8C)$$

$\forall n \geq \tau_2.$

$\hookrightarrow \oplus$

Claim 1: $D(P,Q) = d(\mu_2, \mu_1 + \delta) E_{v,\pi}[\tau_2(t)].$

\hookrightarrow Divergence Decomposition (See Set 10).

Then, $\oplus + \text{Claim 1} \Rightarrow \forall n \geq \tau_2,$

$$E_{v,\pi}[\tau_2(t)] \geq \frac{1}{d(\mu_2, \mu_1 + \delta)} \left((\alpha) \ln(n) - \ln(8C) \right)$$

$$\therefore R_n(\pi, v) = D_2 E_{v,\pi}[\tau_2(t)]$$

$$\geq \frac{D_2}{d(\mu_2, \mu_1 + \delta)} \left((\alpha) \ln(n) - \ln(8C) \right).$$