

Minimax Lower Bound - Sketch.

Source: Chap. 13, Bandit Algorithms, L & S.

Running example: 2 armed bandit, unstructured environment, Gaussian(1) reward,

i.e.

Arm 1

Arm 2

$$\{x_{t1}, t=1, 2, \dots, n\}$$

$$\{x_{t2}, t=1, 2, \dots, n\}$$

$$X_{t1} \sim N(\Delta, 1)$$

$$X_{t2} \sim N(0, 1)$$

$$\mu_1 = \Delta$$

$$\mu_2 = 0$$

Fact: With m samples, $X_{\tau_2} \sim N(0, 1)$

$$\frac{\sqrt{s/\pi}}{\sqrt{m\Delta^2 + \sqrt{m\Delta^2 + 16}}} e^{-m\Delta^2/8} \leq P\left(\frac{1}{m} \sum_{s=1}^m X_{\tau_2} > \frac{\Delta}{2}\right)$$

$$\leq \frac{\sqrt{s/\pi}}{\sqrt{m\Delta^2 + \sqrt{m\Delta^2 + 32/\pi}}} e^{-m\Delta^2/8}$$

Suppose $\Delta = \frac{1}{\sqrt{m}}$.

$$\text{Then, } P\left(\frac{1}{m} \sum_{r=1}^m x_{r2} > \frac{\Delta}{2}\right) \approx e^{-1/8} > 0$$

i.e., with a non-vanishing probability, we cannot resolve if samples from arm 2 might have been generated by arm 1.

Similarly,

$$P\left(\frac{1}{m} \sum_{r=1}^m x_{r1} < \frac{\Delta}{2}\right) \approx e^{-1/8} > 0$$

and similar reasoning as above.

This means that with constant prob (of $e^{-1/8}$), it is plausible that we play the wrong arm for the entire duration of m time steps

$$\text{i.e., Regret} \approx e^{-1/8} \cdot \Delta \cdot m = e^{-1/8} \cdot \frac{1}{\sqrt{m}} \cdot \sqrt{m}$$
$$\approx O(\sqrt{m}).$$

The idea for constructing bandit lower bounds:

Construct two alternative models of the bandit system that are "difficult" to distinguish between, while satisfying two criteria:

(a) The best arm in one system is a sub-optimal arm in the other system.

(b) The systems are statistically as "close" as possible (while satisfying (a) above), so that given the samples from arm plays, it is hard to distinguish between them.

----- time horizon

Back to example: choose small $\Delta = \frac{1}{\sqrt{n}}$

$$\boxed{\text{System 1}} \Rightarrow$$

arm 1



$$N(\Delta, 1)$$

best arm

arm 2



$$N(0, 1)$$

$$v' = \boxed{\text{System 2}}$$

arm 1



$$N(\Delta, 1)$$

best arm



$$N(2\Delta, 1)$$

Fix any policy π , i.e., π is a mapping from past samples and actions to current play decision.

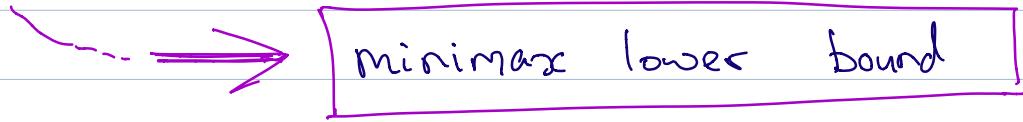
* For either system, π is the same.

\Rightarrow If the player sees the same (action, reward) sequence in either system, the player makes the same action decision at the next play.

Sketch below to show that:

$$\max\{R_n(\pi, v), R_n(\pi, v')\} \geq \frac{1}{2}\sqrt{n}$$

i.e., with any fixed policy π , the player suffers a regret of $O(\sqrt{n})$ in one of these two environments.

 minimax lower bound

Fix policy π , and let $T_1(n), T_2(n)$ be the number of arm 1 and arm 2 plays.

Consider the environment v . (System 1)

$$E[\cdot] = E_{P_v}[\cdot]$$

$$(\Delta, 0)$$

$$\text{Then, } R_n(\pi, v) = \underbrace{(n - E[T_1(n)])}_{E[T_2(n)]} \cdot \Delta \xrightarrow{\Delta, 0} \left(\frac{1}{\sqrt{n}}\right)$$

Now, in environment v' (System 2)

$$E'[\cdot] = E_{P_{v'}}[\cdot]$$

$$(\Delta, 2\Delta)$$

$$R_n(\pi, v') = \Delta E'[T_1(n)]$$

Suppose: $E[T_1(n)] \approx E'[T_1(n)]$, i.e.,

This is key, and
so far unjustified
in any formal
sense ---- but
seems reasonable;
based on Gaussian
intuition earlier.

the policy π is unable to
distinguish between v and v' ,
and thus plays arm 1 the
same number of times under
either scenario.

Then $R_n(\pi, v) \approx \left(\sqrt{n} - \frac{a}{\sqrt{n}} \right)$

$$\Delta = \frac{1}{\sqrt{n}}$$

$$R_n(\pi, v') = \frac{a}{\sqrt{n}} \quad \begin{matrix} \xrightarrow{\text{a} = E'[\tau_{i(n)}]} \\ \approx E[\tau_{i(n)}] \end{matrix}$$

$$1 \leq a \leq n$$

$$\Rightarrow \max \left\{ R_n(\pi, v), R_n(\pi, v') \right\} \geq \frac{\sqrt{n}}{2}$$

\downarrow
 $a = \eta_2$

Note: This is not a sketch of proof because we have assumed $E[\tau_{i(n)}] \approx E'[\tau_{i(n)}]$ under any policy π , and scenarios:
 $v = (\Delta, 0)$, $v' = (\Delta, 2\Delta)$, $\Delta = \frac{1}{\sqrt{n}}$.

In subsequent lectures:

k-armed, unstructured environment $v \in \mathcal{E}$, stochastic setting.

(2) Minimax bound — Showing existence of environment $v' \in \mathcal{E}$ s.t. $R_n(\pi, v') \geq \Theta(\sqrt{nk})$.

(b) Instance dependent bound - For any environment v ,
 showing $R_n(\pi, v) \geq c^* \cdot \ln(n)$, for reasonable π .

Note that c^* depends on $\{\Delta_i\}$, the arm gaps.
 Specifically, we will see that roughly,

$$c^* \approx \sum_{\{i : \Delta_i > 0\}} \frac{\Delta_i}{d(\mu_i, \mu_1)}$$

(c) High probability bound - Bounds for both
stochastic and adversarial environments, that hold
 with high prob., i.e., results of the form

stochastic: $P(R_n(\pi, v) > c \sqrt{k_n \ln(1/\delta)}) \geq \delta$

adversarial: \exists rewards $\{x_1, x_2, \dots, x_n\}$ s.t.

$$P(R_n(\pi, x) > c \sqrt{k_n \ln(1/\delta)}) \geq \delta$$

for $\delta \in (0, 1)$.