# Multi-Armed Bandits

Given: $K$ arms; for each arm $\alpha$, reward dist $D_\alpha$ w/ mean $\mu_\alpha$ (unknown)
$\in [0,1]$

$T$ rounds $(T \gg K)$

~~Output: Arms~~ $\alpha_1, \alpha_2, \ldots, \alpha_T$ ~~played.~~

## Procedure In each round $t = 1, \ldots, T$,

~~Defn (Regret)~~
1. ALG plays an arm $\alpha_t$
2. Reward $r_t$ is sampled indep from $D_{\alpha_t}$
3. ALG learns $r_t$.

## Defn (Regret)

$$R(T) = \mu^* \cdot T - \sum_{t=1}^{T} \mu_{\alpha_t}, \quad \text{where } \mu^* = \max_{\alpha \in K} \mu_\alpha.$$

Generally, consider $\mathbb{E}[R(T)]$ where the expectation is over the randomness of $D_\alpha$'s & the alg's choices.

## Fact (Hoeffding's Ineq)

Given mutually indep (not necessarily identically dist) $X_1, X_2, \ldots, X_n$, let

$$\bar{X}_n := \frac{X_1 + \cdots + X_n}{n} \quad \text{and} \quad \mu_n := \mathbb{E}[\bar{X}_n] = \frac{\mu_1 + \mu_2 + \cdots + \mu_n}{n}.$$

We then have, for any $T$, $\Pr\left[\,|\bar{X}_n - \mu_n| \leq \sqrt{\frac{2 \log T}{n}}\,\right] \geq 1 - \frac{2}{T^4}.$

## Alg 1 (Uniform exploration)

1. Try each arm $N$ times.
2. Select one w/ highest avg reward.
3. Play the chosen one, in the remaining rounds.
   $\hat{\alpha}$

Obsv At +

F

$\mu_0$

Let "clea

Obsv

By the

Lem Cond

Pf) Alg

condi

Lem $\mathbb{E}[$

pf) $\mathbb{E}$

Ob

Choosing

Obsv Fix

the

rid

See U

**Obsv** At the end of step 4, for any arm $\alpha$,

$$\Pr\left[|\bar{\mu}_\alpha - \mu_\alpha| \leq \text{rad}\right] \geq 1 - \frac{2}{T^4}, \quad \text{where}$$

$\bar{\mu}_\alpha$ denotes the observed avg reward of $\alpha$, and $\text{rad} := \sqrt{\frac{2\lg T}{N}}$.

Let "clean event" be ~~the event that~~ $\bigwedge_{\alpha \in K} (|\bar{\mu}_\alpha - \mu_\alpha| \leq \text{rad})$.

**Obsv**

By the union bound, $\Pr\left[\text{~~clean~~ bad event}\right] \leq \frac{2K}{T^4}$.

**Lem** Cond on "clean event", if ~~the chosen arm~~ $\hat{\alpha} \neq \alpha^*$, ~~we have~~

~~$\mu_{\alpha^*} \leq$~~ $\mu_{\alpha^*} - \mu_{\hat{\alpha}} \leq 2\,\text{rad}$.

**Pf** Alg chose $\hat{\alpha}$ instead of $\alpha^*$ since $\bar{\mu}_{\hat{\alpha}} \geq \bar{\mu}_{\alpha^*}$. Due to the condition, $\mu_{\hat{\alpha}} + \text{rad} \geq \bar{\mu}_{\hat{\alpha}}$ & $\bar{\mu}_{\alpha^*} \geq \mu_{\alpha^*} - \text{rad}$. ∎

**Lem** $\mathbb{E}[R(T)] \leq NK + 2\,\text{rad} \cdot T + o(1)$.

**pf** $\mathbb{E}[R(T)] = \mathbb{E}[R(T) \mid \text{clean}]\Pr[\text{clean}] + \mathbb{E}[R(T) \mid \text{bad}]\Pr[\text{bad}]$

$$\leq \mathbb{E}[R(T) \mid \text{clean}] + T \cdot O\left(\frac{K}{T^4}\right).$$

**Obsv** $\mathbb{E}[R(T) \mid \text{clean}] \leq N(K-1) + 2\,\text{rad}(T - NK)$. ∎

Choosing $N = O\left(\left(\frac{T}{K}\right)^{2/3} \cdot (\lg T)^{1/3}\right)$, we have $\mathbb{E}[R(T)] \leq O\left(T^{\frac{2}{3}} \cdot K^{\frac{1}{3}} \cdot (\lg T)^{\frac{1}{3}}\right)$.

---

**Obsv** Fix round $t$ and arm $\alpha$. If $\alpha$ is played $n_t(\alpha)$ times where $\bar{\mu}_t(\alpha)$ is the observed avg reward, we then have

$$\Pr\left[|\bar{\mu}_t(\alpha) - \mu(\alpha)| \leq \text{rad}_t(\alpha)\right] \geq 1 - \frac{2}{T^4}, \quad \text{where}$$

$$\text{rad}_t(\alpha) := \sqrt{\frac{2\lg T}{n_t(\alpha)}}.$$

Set $\text{UCB}_t(\alpha) := \bar{\mu}_t(\alpha) + \text{rad}_t(\alpha)$ & $\text{LCB}_t(\alpha) := \bar{\mu}_t(\alpha) - \text{rad}_t(\alpha)$.

Alg2 (Successive elimination)

1. Activate all arms.
2. For each phase:
3.    Play all active arms & update [LCB, UCB].
4.    Deactivate all arms whose CB does not overlap $\varnothing$ "highest" CB.

Clean event? $\bigwedge\limits_{\alpha \in K, t \in T} \left( |\bar{\mu}_t(\alpha) - \mu(\alpha)| \le rad_t(\alpha) \right)$.

$\Leftrightarrow \mu(\alpha) \in [LCB_t(\alpha), UCB_t(\alpha)]$

By the union bound, $Pr[\text{bad event}] \le T \cdot K \cdot \dfrac{2}{T^4} = O\left(\dfrac{1}{T^2}\right)$.

Cond on clean event,

- $\alpha^*$ never deactivated.



- For each arm $\alpha$, the CB of $\alpha$ overlaps CB of $\alpha^*$ at $n_t(\alpha)$'th phase.

$\Rightarrow \mu(\alpha^*) - \mu(\alpha) \le 2 rad_t(\alpha) = 2\sqrt{\dfrac{2 \lg T}{n_t(\alpha)}}$

- Contribution of $\alpha$ to $\mathbb{E}[R(t)] \overset{\text{clean}}{=} n_t(\alpha) \cdot O\left(\sqrt{\dfrac{\lg T}{n_t(\alpha)}}\right) = O\left(\sqrt{n_t(\alpha) \lg T}\right)$

$\Rightarrow \mathbb{E}[R(t) \mid \text{clean}] = \sum\limits_{\alpha \in K} O\left(\sqrt{n_t(\alpha) \lg T}\right) = O(\sqrt{\lg T}) \sum\limits_{\alpha \in K} \sqrt{n_t(\alpha)} = O\left(\sqrt{Kt \lg T}\right)$.

Note: by Jensen's Theorem,

$\dfrac{1}{K} \sum\limits_{\alpha \in K} \sqrt{n_t(\alpha)} \le \sqrt{\dfrac{\sum n_t(\alpha)}{K}} = \sqrt{\dfrac{t}{K}} \Rightarrow \sum\limits_{\alpha \in K} \sqrt{n_t(\alpha)} = \sqrt{Kt}$

Lower bounds

Thm Fix $T$ & $K$. ~~There exists~~ s.t. For any $\text{alg}$, $\exists$ family of instances s.t. $\mathbb{E}[R(t)] \ge \Omega(\sqrt{KT})$.

Today, wts for $K = 2$!

Consider ① $I_1$ — $D_{d_1} = \cancel{BBBB} Ber(\frac{1+\varepsilon}{2})$  ② $I_2$ — $D_{d_1} = Ber(\frac{1}{2})$
            — $D_{d_2} = Ber(\frac{1}{2})$                     — $D_{d_2} = Ber(\frac{1+\varepsilon}{2})$.

Outline: think of MAB as a seq of "best-arm identification".
> In the same setting, the goal is to choose the max-reward arm.
> at round $\cancel{x} \cdot t$

<u>Lem</u> Consider a "best-arm identification" w/ $\cancel{x} \leq \frac{1}{16\varepsilon^2}$. Fix any alg.
$\exists \overset{arm}{\alpha} \in \{d_1, d_2\}$ s.t.
$$Pr[\text{Alg chose } \alpha \mid I_\alpha] < 3/4.$$

~~(Using this lem, we have)~~     $T \leq \frac{1}{16\varepsilon^2} \Rightarrow \boxed{\varepsilon \leq \sqrt{\frac{1}{16T}}}$

<u>Thm</u> Fix T & ~~pg~~ any ~~alg~~ ALG. Choose an arm $\alpha$ uar, & run ALG on $I_\alpha$.
   Then, $\mathbb{E}[R(T)] \geq \Omega(\sqrt{\cancel{\alpha}T})$.                          w/ $\varepsilon = O(\sqrt{\frac{1}{T}})$

pf) By choice of $\varepsilon$, we can use Lem for each round $t \leq T = \frac{1}{16\varepsilon^2}$.

$$Pr[\alpha_t \neq \alpha] = \underbrace{Pr[\alpha_t \neq d_1 \mid I_1]}_{} \underbrace{Pr[I_1]}_{=\frac{1}{2}} + \underbrace{Pr[\alpha_t \neq d_2 \mid I_2]}_{} \underbrace{Pr[I_2]}_{=\frac{1}{2}}.$$
$$\geq \frac{1}{8} \qquad\qquad\qquad \underset{\text{one would be}}{\geq \frac{1}{4}}$$

$\therefore \mathbb{E}[R(T)] = \sum_{t=1}^{T} Pr[\alpha_t \neq \alpha] \cdot \frac{\varepsilon}{2} \geq \frac{\varepsilon T}{16} = \Omega(\sqrt{T})$.  ∎
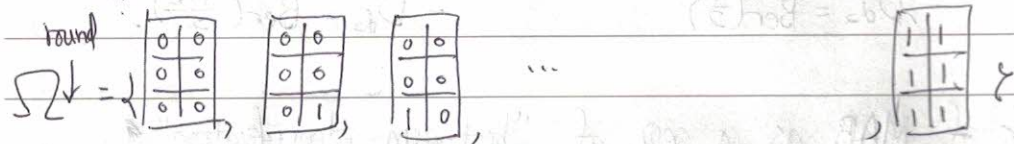
<u>KL divergence</u>

$$KL(p, z) = \sum_{\lambda \in \Omega} p(\lambda) \cdot \ln \frac{p(\lambda)}{z(\lambda)}.$$

<u>Fact</u> (Chain rule) For $P = P_1 \times P_2 \times \cdots \times P_n$ & $z = z_1 \times z_2 \times \cdots \times z_n$ where
$p_i, z_i$ on the same sample space, $KL(p, z) = \sum_{i=1}^{n} KL(p_i, z_i)$.

<u>Fact</u> (Pinsker's ineq) For any event $A \subset \Omega$, $2(p(A) - z(A))^2 \leq KL(p, z)$!

<u>Fact</u> $KL(Ber(\frac{1+\varepsilon}{2}), Ber(\frac{1}{2})) \leq 2\varepsilon^2$ & $KL(Ber(\frac{1}{2}), Ber(\frac{1+\varepsilon}{2})) < \varepsilon^2$ $\forall \varepsilon < \frac{1}{2}$.

pf of Lem) Consider all possible outcomes. For example, $t=3$,

arm →

round

$$\Omega^t = \left\{ \begin{array}{|cc|} \hline 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ \hline \end{array} , \begin{array}{|cc|} \hline 0 & 0 \\ 0 & 0 \\ 0 & 1 \\ \hline \end{array} , \begin{array}{|cc|} \hline 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ \hline \end{array} , \dots , \begin{array}{|cc|} \hline 1 & 1 \\ 1 & 1 \\ 1 & 1 \\ \hline \end{array} \right\}$$

Let $p_1(\omega)$ & $p_2(\omega)$ be the prob of outcome $\omega$ in $J_1$ & $J_2$, resp.

e.g.,

$$p_1\left( \begin{array}{|cc|} \hline 0 & 0 \\ 0 & 0 \\ 0 & 1 \\ \hline \end{array} \right) = \left( \frac{1-\varepsilon}{2} \right)^3 \cdot \left( \frac{1}{2} \right)^3 \qquad p_2 \left( \begin{array}{|cc|} \hline 0 & 0 \\ 0 & 0 \\ 0 & 1 \\ \hline \end{array} \right) = \left( \frac{1}{2} \right)^3 \left( \frac{1-\varepsilon}{2} \right)^2 \left( \frac{1+\varepsilon}{2} \right)$$

As ALG is not aware of the real dist, the output of ALG depends on $\omega$. the outcome, i.e., $ALG : \Omega \to \{d_1, d_2\}$. ~~keep same of ALG~~

Let $\mathcal{E} := \{\omega : ALG(\omega) = d_1 \}$. Then $\overline{\mathcal{E}} = \{\omega : ALG(\omega) = d_2 \}$.

\* Spse t.c. that $p_1(\mathcal{E}) \geq 3/4$ & $p_2(\overline{\mathcal{E}}) \geq 3/4$. We then have

$$p_1(\mathcal{E}) - p_2(\mathcal{E}) \geq 3/4 - 1/4 = 1/2. \quad \cdots (a)$$

However, observe that

$$2(p_1(\mathcal{E}) - p_2(\mathcal{E}))^2 \leq KL(p_1, p_2) = \sum_{\text{each cell}} KL(p_1^{(cell)}, p_2^{(cell)})$$

$$\leq 2 t \cdot 2\varepsilon^2$$

$$\Rightarrow p_1(\mathcal{E}) - p_2(\mathcal{E}) \leq \varepsilon \sqrt{2t} \leq \frac{1}{2\sqrt{2}} < \frac{1}{2} \quad \text{since } t \leq \frac{1}{16\varepsilon^2}. \quad \text{[a]}$$

~~For general K,~~

For general $K$, $I_d : \mathcal{D}_d = Ber\left( \frac{1+\varepsilon}{2} \right)$, the others $Ber\left( \frac{1}{2} \right)$.

Lem Consider a "best arm identification" w/ $t \leq \frac{cK}{\varepsilon^2}$ (for a small enough absolute constant $c$). Fix any $Alg$, $\exists$ at least $\lceil \frac{K}{8} \rceil$ arms $d$ s.t.

$$Pr[Alg \text{ chose } d \mid I_d] < 3/4.$$