

MULTI-ARM BANDITS

DILEMMA: To **EXPLORE** \rightarrow To explore the settings for parameters or to **EXPLOIT** \rightarrow To choose the best option.

EXAMPLES : ① Online advertising \rightarrow Template Optimisation
 ② Clinical Trials
 ③ Packet Routing in a network
 ④ Game playing and reinforcement learning.

STOCHASTIC MULTI-ARMED BANDITS

$A \rightarrow$ Set of Arms

Arm $a \in A$ has mean reward μ_a

$\mu^* \rightarrow$ Highest mean.

Pick an arm based on History

ALGORITHM: For $t = 0, 1, 2, \dots, T-1$:

DETERMINISTIC ALGORITHM

- Given history $h^t = \{a^0, r^0, a^1, r^1, \dots, a^{t-1}, r^{t-1}\}$
- pick a^t to sample
- obtain reward r^t

$T \rightarrow$ Horizon / Total Sampling Budget

$$P\{h^T\} = \prod_{t=0}^{T-1} P\{a^t | h^t\} P\{r^t | a^t\}$$

Decided by the algorithm

Comes from bandit instance

ϵ - GREEDY ALGORITHMS

EG1: - If $t < \epsilon T$, sample an arm uniformly at random

- At $t = \lfloor \epsilon T \rfloor$, identify a^{best}

- For $t > \epsilon T$, sample a^{best}

Highest empirical mean

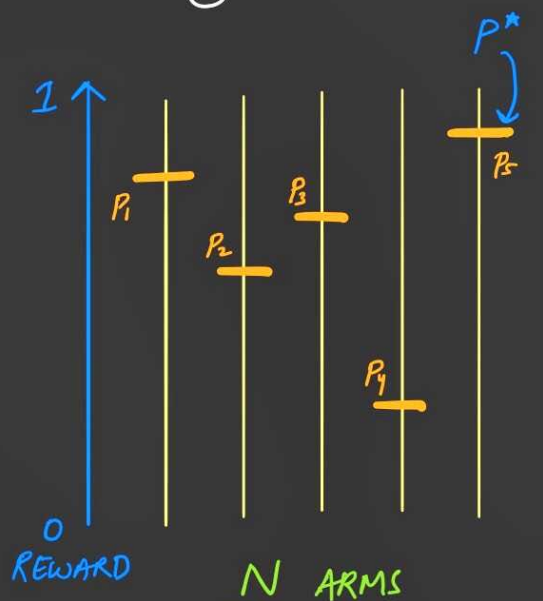
EG2:

- if $t < \epsilon T$, sample uniformly

- $t > \epsilon T$, choose arm with highest empirical mean

EG3:

- with probability ϵ , sample an arm uniformly at random, with $1 - \epsilon$ sample arm with highest mean



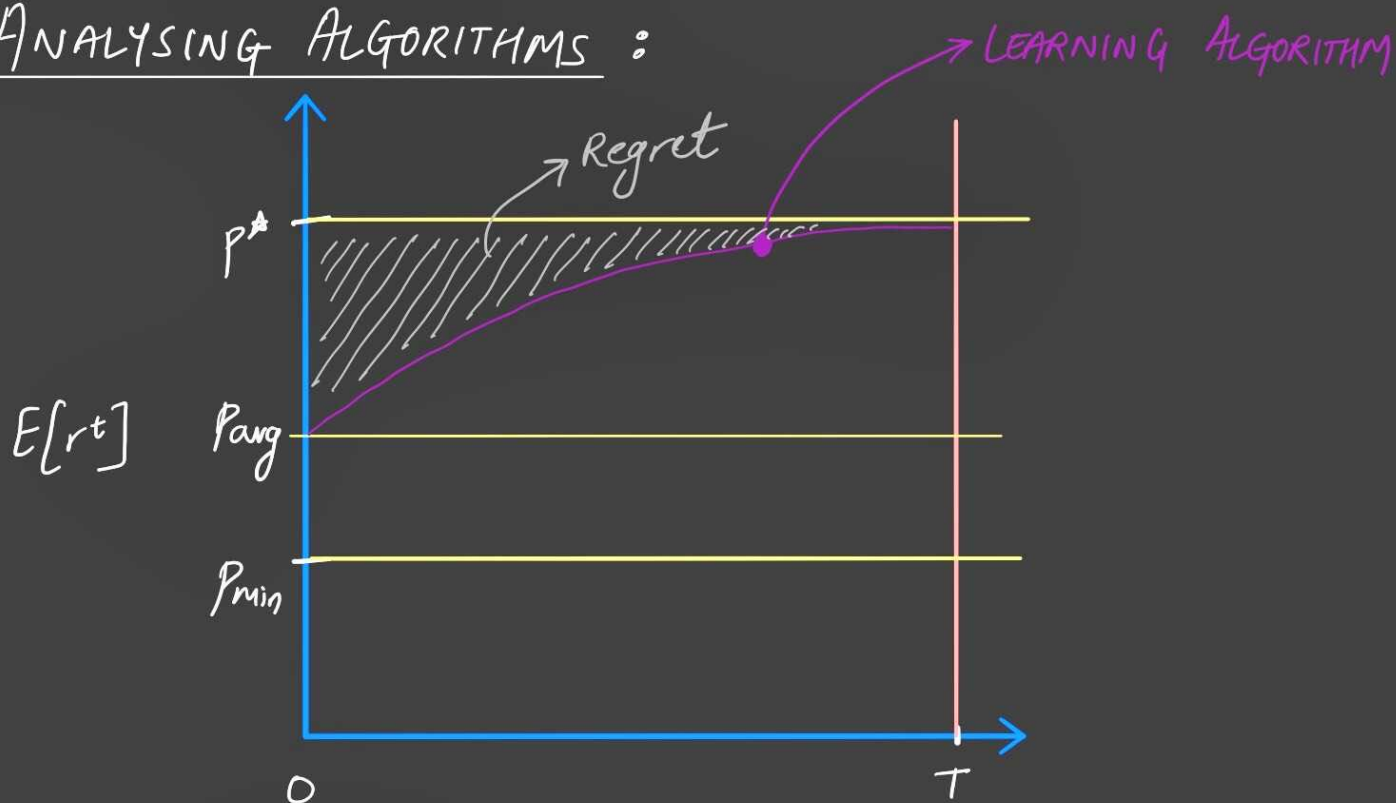
with Bernoulli distribution of rewards (Rewards are 0/1)

(*) NON-DETERMINISTIC ALGORITHM



Parameter $\epsilon \rightarrow$ Controls amount of exploration

ANALYSING ALGORITHMS :



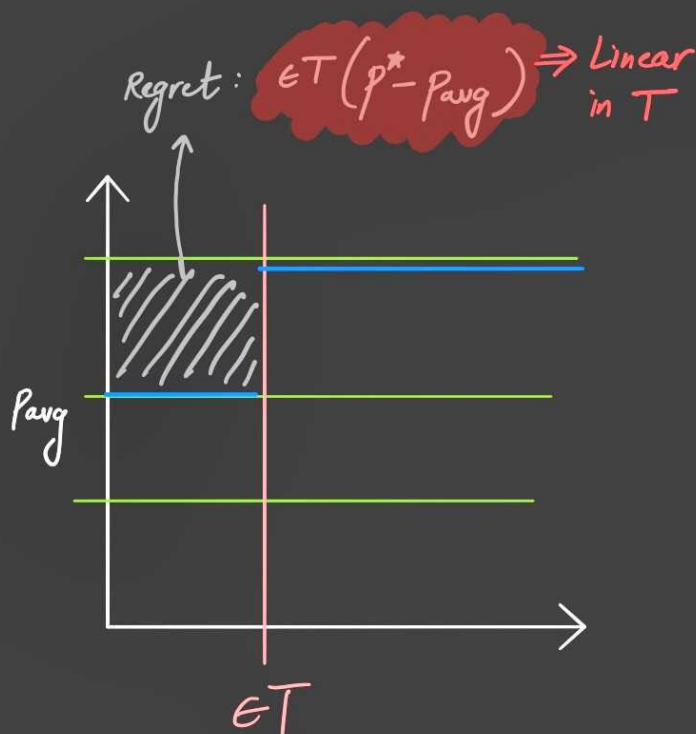
$$\text{Regret } R_T = T p^* - \sum_{t=0}^{T-1} E[r_t]$$

$$\text{Goal} \rightarrow \lim_{T \rightarrow \infty} \left(\frac{R_T}{T} \right) = 0 \rightarrow \text{Sublinear in } T$$

Review of $\epsilon G1$, $\epsilon G2$:

ϵ -first **Explore** for ϵT pulls and thereafter exploit

UNIFORM SAMPLING

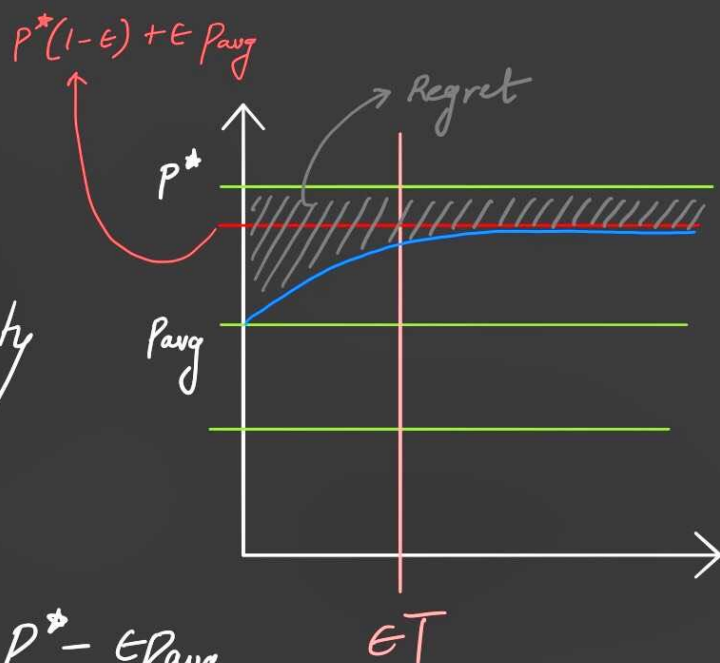


Mathematically :

$$\begin{aligned} R_T &= T p^* - \sum_{t=0}^{T-1} E[r^t] = T p^* - \sum_{t=0}^{\epsilon T - 1} E[r^t] - \sum_{t=\epsilon T}^{T-1} E[r^t] \\ &= T p^* - \epsilon T p_{\text{avg}} - \sum_{t=\epsilon T}^{T-1} E[r^t] \geq T p^* - \epsilon T p_{\text{avg}} - (T - \epsilon T) p^* \\ &= \epsilon (p^* - p_{\text{avg}}) T = \Omega(T) \rightarrow \text{Linear Regret} \end{aligned}$$

Review of ϵ -Greedy :

EXPLORE with probability ϵ , EXPLOIT with probability $1 - \epsilon$.



$E[r^t]$ can never exceed $p^* - \epsilon p_{\text{avg}}$

Mathematically :

$$\begin{aligned} R_T &= T p^* - \sum_{t=0}^{T-1} E[r^t] \geq T p^* - \sum_{t=0}^{T-1} (\epsilon p_{\text{avg}} + (1 - \epsilon) p^*) \\ &= \epsilon (p^* - p_{\text{avg}}) T = \Omega(T) \rightarrow \text{Linear Regret} \end{aligned}$$

CONDITIONS FOR SUB-LINEAR REGRET:

C1. INFINITE EXPLORATION: As $T \rightarrow \infty$, each arm must be pulled infinite number of times

REASON with probability $(1-p^*)^U$, optimal arm will not be chosen and hence regret will be linear. This is because non-optimal arm will get pulled forever.

ϵ -Greedy satisfy C1

C2. GREED IN THE LIMIT: Let $\text{exploit}(T)$ denote the number of pulls that are greedy with respect to empirical mean upto horizon T .

For sub-linear regret: $\lim_{T \rightarrow \infty} \frac{E[\text{exploit}(T)]}{T} = 1$

ϵ -Greedy do not satisfy C2.

RESULT: An algorithm L achieves sub-linear regret on all instances $I \in \tilde{\mathcal{I}}$ if and only if C1 and C2 are satisfied on all $I \in \tilde{\mathcal{I}}$.

GLIE \iff Sub-linear regret

Greedy Limit \swarrow Infinite Exploration \searrow

GLIE-fying ϵ -GREEDY STRATEGIES:

* ϵ_T -first with $\epsilon_T = \frac{1}{\sqrt{T}}$

* ϵ_t -greedy with $\epsilon_t = \frac{1}{t+1}$

At t^{th} step,
EXPLORE = $\frac{1}{t+1}$
EXPLOIT = $\frac{t}{t+1}$

QUESTION: What if $\epsilon_t = \frac{1}{(t+1)^2}$? ANSWER: C1 violated, no infinite exploration.

LOWER BOUND ON REGRET :

* We desire low regret on all instances

LAI and ROBBINS

THEOREM 2 Let L be an algorithm such that for every bandit instance $I \in \tilde{\mathcal{I}}$ and for every $\alpha > 0$, as $T \rightarrow \infty$:

$$R_T(L, I) = O(T^\alpha)$$

Then for every bandit instance $I \in \tilde{\mathcal{I}}$, as $T \rightarrow \infty$:

$$\frac{R_T(L, I)}{\ln(T)} \geq \sum_{a: p(a) \neq p^*(I)} \frac{p^*(I) - p_a(I)}{KL(p_a(I), p^*(I))}$$

KL Divergence on Bernoulli Distribution of these probabilities

* $KL(x, y) \stackrel{\text{def}}{=} x \ln\left(\frac{x}{y}\right) + (1-x) \ln\left(\frac{1-x}{1-y}\right)$ [$0 \ln 0 \stackrel{\text{def}}{=} 0$]

ALGORITHMS :

① UCB \rightarrow Upper Confidence Bounds

- At time t , for every arm a , define $ucb_a^t = \hat{p}_a^t + \sqrt{\frac{2 \ln(t)}{u_a^t}}$
- Pull arm having maximal ucb_a^t

REGRET: $O(\log(T))$

② KL-UCB

- At time t , for every arm, define $ucb\text{-}kl_a^t = \max\{q \in [\hat{p}_a^t, 1] \text{ such that } u_a^t KL(\hat{p}_a^t, q) \leq \ln(t) + c \ln(\ln(t))\}$
- * $ucb\text{-}kl_a^t \leq ucb_a^t \rightarrow$ Tighter bound

REGRET: Matches Lai and Robbins lower bound asymptotically

$$\boxed{C \geq 3}$$

BETA DISTRIBUTION :

$$f(x; \alpha, \beta) = \text{Constant} \cdot x^{\alpha-1} (1-x)^{\beta-1}$$

$$\text{Mean} = \frac{\alpha}{\alpha+\beta}$$

$$\text{Variance} = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$$

③ THOMPSON SAMPLING

- At a time t , let arm a have s_a^t success and f_a^t failures
- $\text{Beta}(s_a^t+1, f_a^t+1) \Rightarrow$ Represents belief about the true mean

$$\text{Mean} = \frac{s_a^t+1}{s_a^t+f_a^t+2}$$

$$\text{Variance} = \frac{(s_a^t+1)(f_a^t+1)}{(s_a^t+f_a^t+2)^2(s_a^t+f_a^t+3)}$$

① Computation: For every arm, draw $x_a^t \sim \text{Beta}(s_a^t+1, f_a^t+1)$

② Sampling: Sample arm a for which x_a^t is maximal

REGRET: Matches Lai and Robbins lower bound

CONCENTRATION BOUNDS

① Hoeffding's Inequality: $x \in [0,1]$ $E[x] = \mu$

$$P\{\bar{x} \geq \mu + \epsilon\} \leq e^{-2n\epsilon^2}$$

$$P\{\bar{x} \leq \mu - \epsilon\} \leq e^{-2n\epsilon^2}$$

② KL Inequality: $x \in [0,1]$ $E[x] = \mu$

Tighter \downarrow

$$P\{\bar{x} \geq \mu + \epsilon\} \leq e^{-n \text{KL}(\mu+\epsilon, \mu)}$$

$$P\{\bar{x} \leq \mu - \epsilon\} \leq e^{-n \text{KL}(\mu-\epsilon, \mu)}$$

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

ANALYSIS OF UCB ALGORITHM :

TO SHOW: $R_T = O\left(\sum_{a: p_a \neq p^*} \frac{1}{p^* - p_a} \log(T)\right)$

NOTATION:

① $\Delta_a \stackrel{\text{def}}{=} p^* - p_a$

② $Z_a^t \rightarrow$ Event that arm a is pulled at t

③ $Z_a^t \rightarrow$ Random value that has value 1 if arm a is pulled

$E[Z_a^t] = P\{Z_a^t\}(1) + (1 - P\{Z_a^t\})(0) = P\{Z_a^t\}$

④ Instance specific constant \rightarrow

$\bar{u}_a^T \stackrel{\text{def}}{=} \left\lceil \frac{B}{(\Delta_a)^2} \ln(T) \right\rceil$

STEP 1: Show that $R_T = \sum_{a: p_a \neq p^*} E[u_a^T] \Delta_a$

$$\begin{aligned} R_T &= T p^* - \sum_{t=0}^{T-1} E[r^t] = T p^* - \sum_{t=0}^{T-1} \sum_{a \in A} P\{Z_a^t\} E[r^t | Z_a^t] = T p^* - \sum_{t=0}^{T-1} \sum_{a \in A} E[Z_a^t] p_a \\ &= \left(\sum_{a \in A} E[u_a^T] \right) p^* - \sum_{a \in A} E[u_a^T] p_a = \sum_{a \in A} E[u_a^T] (p^* - p_a) = \sum_{a \in A} E[u_a^T] \Delta_a \end{aligned}$$

STEP 2: TWO REGIMES FOR SUBOPTIMAL PULL. TO SHOW: $E[u_a^T] \leq \bar{u}_a^T + \underbrace{C}_{\text{constant}}$

$$E[u_a^T] = \sum_{t=0}^{T-1} E[Z_a^t] = \underbrace{\sum_{t=0}^{T-1} P\{Z_a^t \wedge (u_a^t < \bar{u}_a^T)\}}_A + \underbrace{\sum_{t=0}^{T-1} P\{Z_a^t \wedge (u_a^t \geq \bar{u}_a^T)\}}_B$$

STEP 3: BOUNDING A

$$A = \sum_{t=0}^{T-1} P\{Z_a^t \wedge (u_a^t < \bar{u}_a^T)\} = \sum_{t=0}^{T-1} \sum_{m=0}^{\bar{u}_a^T-1} P\{Z_a^t \wedge (u_a^t = m)\} \leq \sum_{m=0}^{\bar{u}_a^T-1} 1 = \bar{u}_a^T$$

← FLIP SUMMATION

STEP 4: BOUNDING B

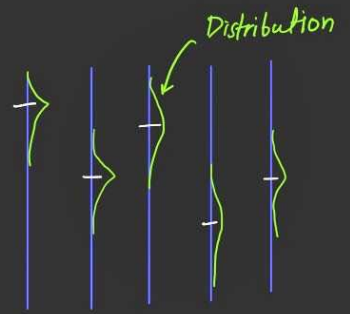
$$\begin{aligned} B &\leq \sum_{t=0}^{T-1} P\left\{ \left(\hat{p}_a^T + \sqrt{\frac{2}{u_a^t} \ln(T)} \geq \hat{p}_a^t + \sqrt{\frac{2}{u_a^t} \ln(t)} \right) \wedge (u_a^t \geq \bar{u}_a^T) \right\} \\ &\leq \sum_{t=0}^{T-1} \sum_{x=\bar{u}_a^T}^t \sum_{y=1}^x P\left\{ \hat{p}_a(x) + \sqrt{\frac{2}{x} \ln(t)} \geq \hat{p}_a^t + \sqrt{\frac{2}{y} \ln(t)} \right\} \\ &\leq \sum_{t=0}^{T-1} \sum_{x=\bar{u}_a^T}^t \sum_{y=1}^x \left(e^{-2x\left(\frac{\Delta_a}{2}\right)^2} + e^{-2y\left(\frac{\sqrt{\frac{2}{y} \ln(t)}}{y}\right)^2} \right) \leq \sum_{t=0}^{T-1} t^2 \left(\frac{2}{t^4}\right) \leq \sum_{t=0}^{\infty} \frac{2}{t^2} = \frac{\pi^2}{3} \end{aligned}$$

A $\hat{p}_a(x) + \sqrt{\frac{2}{x} \ln(t)} \geq p^* \Rightarrow \hat{p}_a(x) \geq p^* + \frac{\Delta_a}{2}$
OR
B $\hat{p}_a(y) + \sqrt{\frac{2}{y} \ln(t)} < p^*$

THOMPSON SAMPLING:

- At a time t , arm has s_a^t successes and f_a^t failures
- $\text{Beta}(s_a^t+1, f_a^t+1)$ represents belief about p_a

COMPUTATIONAL STEP: For every arm, draw a sample
 $x_a^t \sim \text{Beta}(s_a^t+1, f_a^t+1)$



SAMPLING STEP: Pull arm which has maximum x_a^t

* BAYESIAN INFERENCE: $P\{A|B\} = P\{B|A\} P\{A\} / P\{B\}$ $\text{Belief}_m = P\{w|e_1, e_2 \dots e_m\}$

$$\text{Belief}_{m+1}(w) = P\{w|e_1, e_2 \dots e_{m+1}\} = \frac{P\{e_1, e_2 \dots e_{m+1}|w\} P\{w\}}{P\{e_1, e_2 \dots e_{m+1}\}}$$

$$= \frac{\overset{\text{CONDITIONAL INDEPENDENCE}}{\uparrow} P\{e_1, e_2 \dots, e_{m+1}\}}{P\{e_1, e_2 \dots e_{m+1}\}} P\{w\} = \frac{P\{e_1, e_2 \dots e_m, w\} P\{e_{m+1}|w\}}{P\{e_1, e_2 \dots e_{m+1}\}}$$

$$= \frac{P\{w|e_1, e_2 \dots e_m\} P\{e_1, e_2, \dots e_m\} P\{e_{m+1}|w\}}{P\{e_1, e_2, \dots e_{m+1}\}}$$

$$= \frac{\text{Belief}_m(w) P\{e_{m+1}|w\}}{\sum_{w' \in W} \text{Belief}_m(w') P\{e_{m+1}|w'\}} \quad \xrightarrow{\text{Normalization Constant}}$$

CASE 1: $e_{m+1} = 1$ reward

$$\text{Belief}_{m+1}(x) = \frac{\text{Belief}_m(x) \cdot x}{\int_{y=0}^1 \text{Belief}_m(y) y dy}$$

CASE 2: $e_{m+1} = 0$ reward

$$\text{Belief}_{m+1}(x) = \frac{\text{Belief}_m(x) \cdot (1-x)}{\int_{y=0}^1 \text{Belief}_m(y) (1-y) dy}$$

$$\text{Belief}_m(x) = \text{Beta}_{s_H, f_H}(x) dx$$

PRINCIPLE OF SELECTING ARM TO PULL:

- We sample a bandit instance I from joint belief distribution and act optimally with respect to I

ALTERNATIVE EXPLANATION: Probability of picking an arm is belief that it is optimal.

OTHER BANDIT PROBLEM :

- Incorporating risk/variance in the objective
 - * RISK MINIMIZATION
- What if true means vary over time. * Eg. ONLINE ADS
 - * → Might take recent data from an interval.
- Pure Exploration
 - * PAC Formulation
 - * Simple Regret formulation.
- Limited number of Feedback
 - * $S < T$ Times
- Large number of arms
 - * Quantile Regret \Rightarrow Look for good and not optimal arms
- Interaction with many bandits simultaneously.
 - * Contextual Bandits
- Rewards are not from fixed random processes.
 - * Adversarial Bandits
- Necessary to use randomised algorithms