# Multi-Armed Bandit Problem

# Key Problem

Exploration vs exploitation dilemma

- inspect new arms with possibly better rewards.
- use existing information to select best arm.

# Stochastic Bandit

- K arms: for each arm $i \in \{1, ..., K\}$.
  - reward distribution $P_i$ .
  - reward mean $\mu_i$.
  - gap to best: $\Delta_i = \mu^* - \mu_i$, where $\mu^* = \max_{i \in [1,K]} \mu_i$.

- Bandit Setting: For $t = 1$ to T do
  - players selects action $I_t \in \{1, \cdots, K\}$(*randomized*)
  - player receives reward $X_t \sim P_{I_t}$

## Objectives

1. Expected Regret

$$\mathrm{E}[\mathrm{R}_T] = \mathrm{E}\Big[\max_{i\in[1,K]}\sum_{t=1}^{T}\mathrm{X}_{i,t} - \sum_{t=1}^{T}\mathrm{X}_{\mathrm{I}_t,t}\Big]$$

2. Pseudo-regret

$$\begin{aligned}
\overline{\mathrm{R}_T} &= \max_{i\in[1,K]}\mathrm{E}\Big[\sum_{t=1}^{T}\mathrm{X}_{i,t} - \sum_{t=1}^{T}\mathrm{X}_{\mathrm{I}_t,t}\Big]\\
&= \max_{i\in[1,K]}\mathrm{E}\Big[\sum_{t=1}^{T}\mathrm{X}_{i,t} - \mathrm{E}\Big[\sum_{t=1}^{T}\mathrm{X}_{\mathrm{I}_t,t}\Big]\\
&= \max_{i\in[1,K]}\sum_{t=1}^{T}\mathrm{E}\big[\mathrm{X}_{i,t}\big] - \mathrm{E}\Big[\sum_{t=1}^{T}\mathrm{X}_{\mathrm{I}_t,t}\Big]\\
&= \mu^* T - \mathrm{E}\Big[\sum_{t=1}^{T}\mathrm{X}_{\mathrm{I}_t,t}\Big]
\end{aligned}$$

3. By Jensen's inequality, $\overline{\mathrm{R}_T} \le \mathrm{E}[\mathrm{R}_T]$

# Pseudo Regret

1. Expression in terms of $\Delta_i$s:

$$\overline{R_T} = \sum_{t=1}^{K} E[T_i(T)]\Delta_i$$

$T_i(T)$: number of times arm $i$ was pulled up to time $t$,
$T_i(t) = \sum_{s=1}^{t} 1_{I_s=i}$

# Pseudo Regret

Proof.

$$\overline{R_T} = \mu^* T - E\big[\sum_{t=1}^{T} X_{I_t,t}\big] = E\big[\sum_{t=1}^{T}(\mu^* - X_{I_t,t})\big]$$

$$= E\big[\sum_{t=1}^{T}\sum_{i=1}^{K}(\mu^* - X_{I_t,t})1_{I_t=i}\big] = \sum_{t=1}^{T}\sum_{i=1}^{K} E[(\mu^* - X_{I_t,t})]\, E[1_{I_t=i}]$$

$$= \sum_{i=1}^{K}(\mu^* - \mu_i)\, E[\sum_{t=1}^{T} 1_{I_t=i}] = \sum_{i=1}^{K} E[T_i(T)]\Delta_i$$

$\square$

# $\epsilon$-Greedy Strategy

---

Input: $\epsilon_t \in (0,1]$
1: for $t \leftarrow 1$ to $T$ do
2:      $I_t = \text{argmax}_{j=1,\cdots,K} \hat{\mu}_j$
3:      Draw u uniformly from [0,1]
4:      if $u > \epsilon_t$ then
5:         Play arm $I_t$
6:      else
7:         Play a random arm
8:      end if
9: end for

---

# Thompson Sampling

For the Bernoulli bandit, $X_t$ follows a Beta distribution, as $X_t$ is essentially the success probability $\theta$ in Bernoulli distribution. The value of $Beta(\alpha, \beta)$ is within the interval [0, 1]; $\alpha$ and $\beta$ correspond to the counts when we succeeded or failed to get a reward respectively.

$$\alpha_i = \alpha_i + X_{I_t, t} 1_{I_t = i}$$
$$\beta_i = \beta_i + (1 - X_{I_t, t}) 1_{I_t = i}$$

# UCB Strategy

Intuition: we are drawn to the bandits that are paying out large rewards and those that we know little about. We want to upper bound the number of times we will pull arm $i$, so we will attempt to compute $E[T_i(T)]$.

- For each $t \in [1, T]$, compute an upper confidence bound estimate on the mean of each arm at some fixed confidence level.

- select arm with largest UCB.

# UCB Strategy

- With the rewards distributions in [0,1]. from Hoeffding's inequality we have:

$$\log \mathrm{E}[e^{t(X-\mathrm{E}[X])}] \leq \Psi(t)$$

- Then

$$
\begin{aligned}
Pr[X - \mathrm{E}[X] > \epsilon] &= Pr[e^{t(X-\mathrm{E}[X])} > e^{t\epsilon}] \\
&\leq \inf_{t>0} e^{-t\epsilon} \, \mathrm{E}[e^{t(X-\mathrm{E}[X])}] \\
&\leq \inf_{t>0} e^{-t\epsilon} e^{\Psi(t)} \\
&= e^{-\sup_{t>0}(t\epsilon - \Psi(t))} \\
&= e^{-\Psi^*(\epsilon)}
\end{aligned}
$$

# UCB Strategy

1. Average reward estimate for arm $i$ by time $t$:

$$\widehat{\mu_{i,t}} = \frac{1}{T_i(t)} \sum_{s=1}^{t} X_{i,s} 1_{I_s=i}$$

2. Concentration inequality:

$$Pr[\mu_i - \frac{1}{t} \sum_{s=1}^{t} X_{i,s} > \epsilon] \leq e^{-\Psi^*(\epsilon)}$$

3. Thus, for any $\delta > 0$, with probability at least $1 - \delta$:

$$\mu_i < \frac{1}{t} \sum_{s=1}^{t} X_{i,s} + \Psi^{*-1}\left(\frac{1}{t} \log \frac{1}{\delta}\right)$$

# $(\alpha, \Psi)$-UCB Strategy

- Parameter $\alpha > 0$; $(\alpha, \Psi)$-UCB Strategy: at time t, select:

$$I_t \in \operatorname{argmax} \left[ \mu_{i,\hat{t}-1} + {\Psi^*}^{-1} \left( \frac{\alpha \log t}{T_i(t-1)} \right) \right]$$

- For $\Psi(\lambda) = \frac{\lambda^2}{8}$, then ${\Psi^*}^{-1} = 2\epsilon^2$ and substituting for $\alpha = 4$, we obtain that at time t, plays the arm (UCB1):

$$I_t \in \operatorname{argmax} \left[ \mu_{i,\hat{t}-1} + \sqrt{\frac{2 \log t}{T_i(t-1)}} \right]$$