

# Lec 6

Page 1

## CS 59000 Reinforcement Learning (RL)

### - Finite-Armed Stochastic Bandit

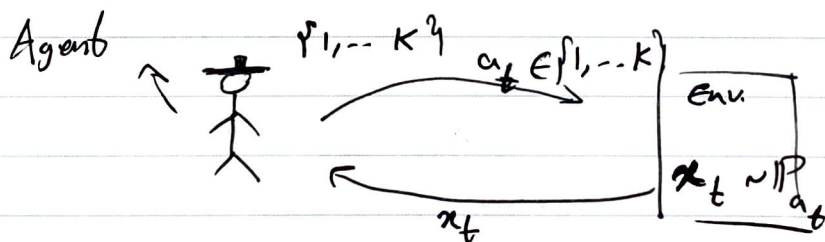
Agenda:

- Problem setup
- Regret definition
- Explore-then-Commit.

### Finite-Armed Stochastic Bandit:

A  $K$ -armed stochastic bandit is a tuple of probability measures  $\mathcal{V} = (P_1, P_2, \dots, P_K)$  where  $P_i$  is a probability measure over reals. for each  $i \in [K]$

The ~~an~~ agent/learner and environment interact sequentially with each others.



In each round  $t$ , the agent chooses an action  $A_t \in [K]$  which is fed to the environment. Then, the environment samples a reward  $X_t \in \mathbb{R}$  from  $P_{A_t}$  and reveals it to the agent.

2)

The agent does not know  $\mathcal{V}$  in advance and is interested in the best arms, the arms with highest  $\mu_i(\mathcal{V}) = \int_{-\infty}^{\infty} x dP_i(x)$ , in other words, arms with highest expected rewards.

In the following, for notation simplicity, we assume one of the best arms with expected reward  $\mu^*$  is the first arm.

---

Let's imagine our ~~agent~~ interacts with the environment  $\mathcal{V}$  for  $n$  timesteps. If the agent knew  $\mathcal{V}$  in advance, it could go after the best arm and receive an expected return of  $n\mu^*$  😊

But, the agent does not have this knowledge. This is the knowledge that Oracle has.  
→ what is Oracle's expected return?

↳  $n\mu^*$

Since the agent does not know the  $\mathcal{V}$ , it needs to learn its property in order to thrive, (what does it mean?)

3)

The agent makes decision  $A_1$  at time 1, and observes  $X_1$ . Then, the agent uses  $(A_1, X_1)$  to decide on  $A_2$ . The agent receives  $X_2$  after committing to  $A_2$  and based on  $(A_1, X_1, A_2, X_2)$  decides on  $A_3, \dots$

Let's denote the agent's strategy,  $\pi$   
It means  $A_t \sim \pi(A_1, \dots, X_{t-1})$

In other words,  $A_t$  is  $\mathcal{F}_t$ -measurable with respect to proper definition of  $\mathcal{F}_t$ .

$$A_t \sim \pi(\mathcal{F}_{t-1})$$

What is our agent's expected return?

$$E\left[\sum_{t=1}^n X_t\right]$$

Def. (Regret): How much more agent could have earned if it knew the environment in advance.

$$R_n(\pi, \nu) = \underbrace{n \mu^*(\nu)}_{\text{How much we could earn (oracle)}} - \underbrace{E\left[\sum_{t=1}^n X_t\right]}_{\text{How much we actually earned}}$$

$$4) \quad R_n(\pi, v) \geq 0 \quad \text{for all } \pi.$$

$$- R_n(\pi, v) = 0 \rightarrow \text{if } \pi \text{ choose } A_t \in \arg \max_i \mu_i$$

$$A_t \in \arg \max_{i \in [K]} \mu_i$$

→ we are interested in strategies that have small regret.

$$- \begin{cases} \text{random regret} = n \mu^*(v) - \sum x_t \\ \text{pseudo regret} = n \mu^*(v) - \sum \mu_{A_t} \end{cases}$$

$$E(x_t | \mathcal{F}_{t-1} \cup A_t)$$

Let  $\Delta_i(v) = \mu^*(v) - \mu_i(v)$  be suboptimality gap of action  $i$  in environment  $v$ .

$$\text{Let, at time } t, \quad T_i(t) = \sum_{s=1}^t \mathbb{I}\{A_s = i\}$$

be a counting random variable for action/arm  $i \in [K]$ .



5)

Theorem: For any policy  $\pi$ , and  $k$  armed stochastic bandit  $\mathcal{V}$ , and horizon  $n \in \mathbb{N}$ , the regret of  $\pi$  in  $\mathcal{V}$  is equal to

$$R_n(\pi, \mathcal{V}) = \sum_{i=1}^K \Delta_i \underbrace{E[T_i(n)]}_{\substack{\text{How many times } i\text{-th} \\ \text{arm has been pulled}}}$$

proof.  $S_n = \sum_{t=1}^n X_t \rightarrow$  return random variable.

For round  $t$ , we have  $\sum_{i=1}^K \mathbb{I}\{A_t = i\} = 1$

$$S_n = \sum_{t=1}^n \sum_{i=1}^K X_t \mathbb{I}\{A_t = i\}$$

$$R_n(\pi, \mathcal{V}) = \underbrace{n\mu^*}_{\sum_{t=1}^n \sum_{i=1}^K \mu^* \mathbb{I}\{A_t = i\}} - E[S_n]$$

$$\sum_{t=1}^n \sum_{i=1}^K \mu^* \mathbb{I}\{A_t = i\}$$

$$\Rightarrow R_n(\pi, \mathcal{V}) = \sum_{t=1}^n \sum_{i=1}^K E[(\mu^* - X_t) \mathbb{I}\{A_t = i\}]$$

$$= \sum_{t=1}^n \sum_{i=1}^K E[E[(\mu^* - X_t) \mathbb{I}\{A_t = i\} | A_t]]$$

$$= \sum_{t=1}^n \sum_{i=1}^K E[\mathbb{I}\{A_t = i\} E[\mu^* - X_t | A_t]]$$

$$= \sum_{t=1}^n \sum_{i=1}^K E[\mathbb{I}\{A_t = i\} \Delta_{A_t}]$$

$$= \sum_{t=1}^n \sum_{i=1}^K E[\mathbb{I}\{A_t = i\} \Delta_{A_t}]$$

6)

$$\begin{aligned}
 \rightarrow R_n(\pi, \nu) &= \sum_{i=1}^k \sum_{t=1}^n E \left[ I(A_t = i) \Delta_{A_t} \right] \\
 &= \sum_{i=1}^k \sum_{t=1}^n E \left[ I(A_t = i) \Delta_i \right] \\
 &= \sum_{i=1}^k \Delta_i \sum E \left[ I(A_t = i) \right] \\
 &= \sum_{i=1}^k \Delta_i E \left[ \sum_{t=1}^n I(A_t = i) \right] \\
 &= \sum_{i=1}^k \Delta_i E[T_i(n)]
 \end{aligned}$$


---