

Week 6 - Lecture 30

Stochastic Bandits

$V = (V_a)_{a \in [K]} \leftarrow \text{Bandit Instance}$

$$X \sim V_a$$

$$\mathbb{E}[X] = \mu_a$$

$X \sim V_a$

Goal: Identify the arm with largest μ_a , quickly.

$$\begin{aligned}\bar{R}(\pi, n) &= (\max_a \mu_a) n - \mathbb{E} \left[\sum_{t=1}^n \mu_{I_t} \right] \\ &= \sum_{i=1}^K \mathbb{E} [T_i(n)] \Delta_i\end{aligned}$$

where $\Delta_i = (\max_a \mu_a) - \mu_i$

We will assume that arm 1 is optimal, algorithm doesn't know this. \downarrow

For ease in calculations

$$\therefore \Delta_1 = 0$$

$$\Delta_2 = +ve$$

$$\Delta_3 = +ve$$

Also we will assume that optimal arm is unique

$$\therefore \mu_1 > \mu_i \quad \forall i \neq 1$$

Explore then Commit (ETC) Algorithm

Play all arms to get the estimates of the arms and then play the arm with highest mean.

But key is, how many rounds to explore

Input : m - no. of rounds to explore each time
 K - no. of arms

For $t = 1, 2, \dots$

$$I_t = \begin{cases} t \bmod K + 1 & \text{if } t \leq mk \\ \arg \max_i \hat{\mu}_i(mk) & \text{if } t > mk \end{cases} \leftarrow \begin{array}{l} \text{(play arm} \\ \text{in a} \\ \text{round} \\ \text{robin} \\ \text{fashion)} \end{array}$$

$$\hat{\mu}_i(t) = \frac{\sum_{s=1}^t \mathbb{1}\{I_s = i\} X_{I_s}}{\sum_{s=1}^t \mathbb{1}\{I_s = i\}} \leftarrow \begin{array}{l} \text{(rewards of arm } i \\ \text{till round } t) \\ \text{(counting the number} \\ \text{of plays of arm } i \text{ till} \\ \text{round } t) \end{array}$$

$$T_i(n) = m + \sum_{s=mk+1}^n \mathbb{1}\{I_s = i\}$$

$$T_i(n) = m + (n - mk) \mathbb{1}\left\{\hat{\mu}_i(mk) \geq \max_{j \neq i} \hat{\mu}_j(mk)\right\}$$

$$\mathbb{E}[T_i(n)] \leq m + (n - mk) \times$$

$$\Pr \left\{ \hat{\mu}_i(mk) \geq \max_{j \neq i} \hat{\mu}_j(mk) \right\}$$

Now need to figure this \uparrow

$$\Pr \left\{ \hat{\mu}_i(mk) \geq \max_{j \neq i} \mu_j(mk) \right\} \quad \text{--- (A)} \quad \left[\begin{array}{l} \text{Assume } i \neq 1, \text{ i.e.} \\ \text{it is not the optimal} \\ \text{arm} \end{array} \right]$$

Thus, [an event that shouldn't happen, but checking now what is the probability of that happening]

$$\Pr \left\{ \hat{\mu}_i(mk) \geq \hat{\mu}_1(mk) \right\} \quad \text{--- (B)}$$

$$\hat{\mu}_1(mk) < \max_{j \neq i} \mu_j(mk), \quad \hat{\mu}_1 \text{ is contained in } \max \mu_j(mk)$$

thus A is more stringent than B \therefore

$$\Pr \left\{ \hat{\mu}_i(mk) \geq \max_{j \neq i} \mu_j(mk) \right\} \leq \Pr \left\{ \hat{\mu}_i(mk) \geq \hat{\mu}_1(mk) \right\}$$

$$= \Pr \left\{ \hat{\mu}_i(mk) - \mu_i - (\hat{\mu}_1(mk) - \mu_1) \geq \mu_1 - \mu_i \right\}$$

$$\Rightarrow \Pr \left\{ \hat{\mu}_i(mk) \geq \max_{j \neq i} \hat{\mu}_j(mk) \right\}$$

$$\leq \Pr \left\{ \hat{\mu}_i(mk) - \mu_i - (\hat{\mu}_1(mk) - \mu_1) \geq \Delta_i \right\}$$

$\hat{\mu}_i(mk)$ is mean of arm i played m times;

$$\therefore \hat{\mu}_i(mK) = \frac{\sum_{s=1}^m X_{is}}{m} \quad \leftarrow \begin{array}{l} \text{all samples from arm } i, \\ \text{all samples are iid} \end{array}$$

$$\therefore \hat{\mu}_i(mK) - \mu_i = \frac{\sum_{s=1}^m X_{is}}{m} - \mu_i$$

$$= \frac{\sum_{s=1}^m (X_{is} - \mu_i)}{m} \quad \leftarrow \begin{array}{l} \text{Each quantity is} \\ 1\text{-subgaussian} \end{array}$$

$$\begin{array}{l} \text{is } \frac{1}{\sqrt{m}}\text{-subgaussian} \\ \hline \end{array} \left[\begin{array}{l} \sqrt{\sigma_1^2 + \sigma_2^2 + \sigma_3^2 + \dots + \sigma_m^2} \\ = \sqrt{\frac{1}{m^2} + \frac{1}{m^2} + \dots + \frac{1}{m^2}} \\ = \frac{1}{\sqrt{m}} \end{array} \right]$$

$(\hat{\mu}_i(mK) - \mu_i)$ will also be $\frac{1}{\sqrt{m}}$ -subgaussian

as all arms - μ is 1-subgaussian

Thus:

$$\Pr \left\{ \underbrace{\hat{\mu}_i(mK) - \mu_i}_{\frac{1}{\sqrt{m}}\text{-subgaussian}} - \underbrace{(\hat{\mu}_1(mK) - \mu_1)}_{\frac{1}{\sqrt{m}}\text{-subgaussian}} \geq \Delta_i \right\}$$

Also these two are independent

$$\therefore \Pr \left\{ \underbrace{\hat{\mu}_i(mK) - \mu_i - (\hat{\mu}_1(mK) - \mu_1)}_{\sqrt{2}\text{-subgaussian}} \geq \Delta_i \right\}$$

□

$$\left[\begin{array}{l} \text{Lemma of Subgaussian :} \\ \Pr \{x \geq \varepsilon\} \leq \exp \left\{ -\varepsilon^2 / 2\sigma^2 \right\} \end{array} \right] \text{inputting}$$

$$\begin{aligned} \therefore \Pr \{ \hat{\mu}_i(mk) - \mu_i - (\hat{\mu}_i(mk) - \mu_i) \geq \Delta_i \} \\ \leq \exp \left\{ -\frac{\Delta_i^2 m}{4} \right\} \end{aligned}$$

Thus :

$$\mathbb{E} [T_i(n)] \leq m + (n - mk) \exp \left\{ -\frac{\Delta_i^2 m}{4} \right\}$$

$$\therefore R(\text{ETC}, n) = \sum_{i=1}^K \mathbb{E} [T_i(n)] \Delta_i$$

$$= m \sum_{i=1}^K \Delta_i + (n - mk) \sum_{i=1}^K \Delta_i \exp \left\{ -\frac{\Delta_i^2 m}{4} \right\}$$
