

Week 6 - Lecture 29

R.V.s. that satisfy subgaussianity

Example:

- 1) If X has zero mean and $|X| \leq B$ for some B , then X is B -subgaussian almost surely
- 2) If X has zero mean and $X \in [a, b]$ almost surely, then X is $(b-a)/2$ -subgaussian

Henceforth focus will be on subgaussian r.v.s.

$$\therefore \mathcal{E} = \left\{ \underset{\substack{\uparrow \\ \text{Environment} \\ \text{class}}}{V} = (V_a)_{a \in [K]} : \left(V_a - \mathbb{E}[X] \right)_{\substack{x \sim V_a \\ \forall a \in [K]}} \text{ is } \sigma\text{-subgaussian} \right\}$$

Stochastic Bandits

Input : No. of arms K

for $t = 1, 2, \dots, n$

learner selects $I_t \in [K]$

Environment $X_{I_t t} \sim V_{I_t}$

Policy $\pi_t(\cdot | X_{I_1 1}, I_1, X_{I_2 2}, I_2, \dots, X_{I_{t-1} t-1}, I_{t-1})$

$$\pi = (\pi_1, \pi_2, \dots, \pi_n)$$

$$\text{Regret}(\pi, n) = \mathbb{E} \left[\max_i \sum_{t=1}^n X_{it} - \sum_{t=1}^n X_{I_t t} \right]$$

$$= \underbrace{\mathbb{E} \left[\max_i \sum_{t=1}^n X_{it} \right]}_{\text{This is harder to deal with, thus will relax it}} - \mathbb{E} \left[\sum_{t=1}^n X_{I_t t} \right]$$

I_t has randomness
and X_{I_t} also
has randomness

$$\geq \max_i \mathbb{E} \left[\sum_{t=1}^n X_{it} \right] - \mathbb{E} \left[\sum_{t=1}^n \mu_{I_t t} \right]$$

$$= \max_i \sum_{t=1}^n \mu_i - \mathbb{E} \left[\sum_{t=1}^n \mu_{I_t} \right]$$

$$= \max_i n \mu_i - \mathbb{E} \left[\sum_{t=1}^n \mu_{I_t} \right]$$

V_i has a mean μ_i
 $\therefore X_{it} \sim V_i$

Let $\mu^* = \max_i \mu_i$

$$= n \mu^* - \mathbb{E} \left[\sum_{t=1}^n \mu_{I_t} \right]$$

$= \bar{R}(\pi, n) \rightarrow$ Pseudo-regret is a lower bound on actual regret

$$R(\pi, n) \geq \bar{R}(\pi, n)$$

We will upper bound $\bar{R}(\pi, n)$, but it will not be an upper bound on $R(\pi, n)$; but we will ignore this

We are interested in making this learnable, which means

$$\frac{\bar{R}(\pi, n)}{n} \rightarrow 0$$

Decomposition of regret

$$\bar{R}(\pi, n) = n\mu^* - \mathbb{E} \left[\sum_{t=1}^n \mu_{I_t} \right]$$

μ_{I_t} depends on π

Defining no. of pulls :

$$T_i(n) = \sum_{t=1}^n \mathbb{1}_{\{I_t = i\}}$$

$T_i(n)$ is a random variable

$$\sum_{i=1}^K \mathbb{E}[T_i(n)] = n$$

$$\therefore \bar{R}(\pi, n) = \sum_{i=1}^K \mathbb{E}[T_i(n)] \mu^*$$

$$- \mathbb{E} \left[\sum_{t=1}^n \sum_{i=1}^K \mu_i \mathbb{1}_{\{I_t = i\}} \right]$$

□

$$\begin{aligned}
&= \sum_{i=1}^K \mathbb{E}[T_i(n)] \mu^* - \mathbb{E} \left[\sum_{i=1}^K \sum_{t=1}^n \mu_i 1_{\{I_t=i\}} \right] \\
&= \sum_{i=1}^K \mathbb{E}[T_i(n)] \mu^* - \mathbb{E} \left[\sum_{i=1}^K \mu_i T_i(n) \right] \\
&= \sum_{i=1}^K \mathbb{E}[T_i(n)] \mu^* - \sum_{i=1}^K \mu_i \mathbb{E}[T_i(n)] \\
&= \sum_{i=1}^K \mathbb{E}[T_i(n)] (\mu^* - \mu_i)
\end{aligned}$$

Let $\mu^* - \mu_i = \Delta_i \leftarrow$ gap between best arm and i th arm

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

Regret decomposition formula

If i is optimal then no Δ_i . But for sub-optimal arm Δ_i will have a value. This formula calculates the gap for all arms w.r.t. the optimal arm, and the gap with optimal arm would be zero.

Let $i^* = \arg\max \mu_i$

So our policy should be such that T_i for i^* should be higher than T_i for sub-optimal arms.

The bound for $\bar{R}(\pi, n)$ would thus be decided by no. of pulls of the sub-optimal arms.

Thus we need a policy that will help us identify the arm with the highest mean quickly.