

Ans-33.1  $R_n^{\text{Sample}}(\pi, \nu) = \frac{R_n(\pi, \nu)}{n}$

where  $R_n(\pi, \nu)$  is the cumulative regret of policy  $\pi$  on environment  $\nu$ .

Proof :-

Let  $(\pi)_{t=1}^n$  be a policy & let

$$\pi(i | a_1, x_1, \dots, a_n, x_n) = \frac{1}{n} \sum_{t=1}^n \mathbb{1}_{\{a_t = i\}}$$

Thus by regret decomposition

$$\begin{aligned} R_n(\pi, \nu) &= E \left[ \sum_{i=1}^n \Delta_i N_i(n) \right] \\ &= n E \left[ \frac{\sum \Delta_i N_i(n)}{n} \right] \\ &= n E \left[ \Delta_{A_n} \right] \\ &= n R_n^{\text{Sample}} \left( (\pi_t)_{t=1}^n, \nu \right) \end{aligned}$$

$$\Rightarrow R_n^{\text{Sample}}(\pi, \nu) = \frac{R_n(\pi, \nu)}{n}$$

We know that minimax bound on  $R_n$  is

$$R_n(\pi, \nu) \geq c \sqrt{(K-1)n} \geq c \sqrt{Kn}$$

$$R_n^{\text{Sample}}(\pi, \nu) \geq c \sqrt{\frac{K}{n}}$$

1-33-2. As proved above,

$$R_n^{\text{sample}}(\pi, \nu) = \frac{R_n(\pi, \nu)}{n}$$

For a UE policy,

For any bandit, no. of optimal pulls  $\geq (K-1) \left\lceil \frac{n}{K} \right\rceil$

$$R_n(\pi, \nu) \geq (K-1) \left\lceil \frac{n}{K} \right\rceil \sum_{\substack{i=1 \\ i \neq i^*}}^K \Delta_i$$

Assuming 1 subgaussian rewards,

The term on RHS  $\geq c \sqrt{n K \log K}$

$$R_n^{\text{sample}} \geq c \sqrt{\frac{K \log K}{n}}$$