

Week 9 - Lecture 40

$$P(A) + Q(A^c) \geq \frac{1}{2} \exp(-D(P||Q))$$

Let $\mathcal{V} = (P_i)_{i \in [K]}$

π - Policy

Let $P_{\mathcal{V}\pi} \sim (A_1, X_1, A_2, X_2, A_3, X_3, \dots, A_T, X_T)$

Also \nearrow a_1, a_2, \dots, a_t are discrete
while x_1, x_2, \dots, x_t come from continuous distribution

$$P_{\mathcal{V}\pi}(a_1, x_1, a_2, x_2, \dots, a_T, x_T) =$$

$$\prod_{t=1}^T \underbrace{\pi(a_t | a_1, x_1, \dots, a_{t-1}, x_{t-1})}_{\text{Policy}} \underbrace{P_{a_t}(x_t)}_{\text{reward from env. for action } a_t}$$

(1)

Similarly: (diff. env. same policy) $\mathcal{V}' = (P'_i)_{i \in [K]}$

$$P_{\mathcal{V}'\pi}(a_1, x_1, a_2, x_2, \dots, a_t, x_t)$$

$$\prod_{t=1}^t \pi(a_t | a_1, x_1, \dots, a_{t-1}, x_{t-1}) P'_{a_t}(x_t)$$

(2)

Lemma:

$$D(P_V \| P_{V'}) = \sum_{i=1}^K \mathbb{E}_V[N_i(T)] D(P_i \| P'_i)$$

Proof:

$$\log \frac{dP_V}{dP_{V'}} (a_1 x_1, a_2 x_2, \dots, a_t x_t)$$

Putting ① & ② and simplifying

$$= \sum_{t=1}^T \log \frac{P_{A_t}(x_t)}{P'_{A_t}(x_t)} \quad \leftarrow \text{This was for one realization}$$

$$\mathbb{E}_V \left[\log \frac{dP_V}{dP_{V'}} (A_1 X_1, \dots, A_T X_T) \right] = D(P_V \| P_{V'})$$

*

$$= \mathbb{E}_V \left[\sum_{t=1}^T \log \frac{P_{A_t}(x_t)}{P'_{A_t}(x_t)} \right]$$

* $D(P \| Q) = \int \log \frac{dP(\omega)}{dQ} dP(\omega)$

$$= \mathbb{E}_V \left[\sum_{t=1}^T \mathbb{E}_V \left[\log \frac{P_{A_t}(x_t)}{P'_{A_t}(x_t)} \mid A_t \right] \right] \quad \leftarrow \begin{array}{l} \text{conditioned} \\ \text{on } A_t \end{array}$$

$$= \mathbb{E}_V \left[\sum_{t=1}^T D(P_{A_t} \| P'_{A_t}) \right] \quad \leftarrow \text{Now the}$$

$\sim \lfloor t^{\frac{1}{2}} \rfloor$

expectation is
only over the
randomness of arm pulls

$$= \sum_{k=1}^K \mathbb{E}_{\mathcal{V}} \left[\sum_{t=1}^T \mathbb{1}\{A_t = k\} D(P_{A_t} \| P'_{A_t}) \right]$$

$$= \sum_{k=1}^K \mathbb{E}_{\mathcal{V}} [N_k(T)] D(P_{A_t} \| P'_{A_t})$$

Hence proved

Divergence between the two environments
can be decomposed into divergence between
arms

Theorem: $K > 1$, $T \geq K-1$

$\forall \pi \exists$ a mean vector $\mu \in [0, 1]^K$ s.t.

$$R_T(\pi, \mathcal{V}_{\mu}) \geq \frac{1}{27} \sqrt{(K-1)T}$$

$\mathcal{V}_{\mu} \rightarrow$ env.
whose rewards
come from μ

$$\Rightarrow R_T^*(\varepsilon^K) \geq \frac{1}{27} \sqrt{(K-1)T}$$

\uparrow
minmax regret

