

## Week 9 - Lecture 42

### Stochastic Contextual Bandits

$$t = 1, 2, \dots, T$$

Actions:

[1]

For every user have a context at every time step

[2]

$\therefore x_t = \text{context or side-information}$

$\vdots$

$$x_t \in \mathcal{C} \subset \mathbb{R}^d$$

[K]

$$r_t = R(x_t, I_t) \leftarrow \text{reward}$$

where,  $I_t \in [K]$  is the action chosen based on  $x_t$ .

Reward depends on context & arm-played

Actually reward also has a noise component.

$$r_t = R(x_t, I_t) + \eta_t \leftarrow \eta_t \text{ is conditionally sub-gaussian}$$

$$R : \mathcal{C} \times [K] \rightarrow \mathbb{R} \quad \underline{R \text{ is reward function}}$$

$\hookrightarrow$  is a function from context, arm giving reward

For the learner the reward is unknown.

Thus goal is to figure out the function  $R$ .

If  $\eta_t$  is  $\sigma$ -subgaussian, then:

$$\mathbb{E} \left[ e^{\lambda \eta_t} \right] \leq e^{\lambda^2 \sigma^2 / 2}$$

But  $\eta_t$  is conditionally  $\sigma$ -subgaussian on  $\mathcal{F}_t$

$$\mathbb{E} \left[ e^{\lambda \eta_t} \mid \mathcal{F}_t \right] \leq e^{\lambda^2 \sigma^2 / 2}$$

$$\mathcal{F}_t = \sigma(x_1, I_1, x_2, I_2, \dots, x_t, I_t)$$

$\uparrow$   $\sigma$ -algebra generated by observations (context, action pair)

$\sigma$ -subgaussian will have zero mean

$\therefore$

$$\mathbb{E} [r_t \mid \mathcal{F}_t] = R(x_t, I_t)$$

The policy would be mapping from context to arm

$$\pi : \mathcal{C} \rightarrow [K]$$

The context are generated here (assume) stochastically in an i.i.d. fashion according to some common distribution, that would be revealed to the learner.

$\therefore X_1, X_2, \dots, X_T \leftarrow$  randomly generated

Say a particular realization is:

$$x_1, x_2, \dots, x_T$$

$$\sum_{t=1}^T R(x_t, a_t^*) \leftarrow \text{Best total reward}$$

where

$$a_t^* = \operatorname{argmax}_{a \in [K]} R(x_t, a)$$

The learner does not know  $R \therefore$  Regret of a policy would be:

$$R_T(\pi) = \mathbb{E} \left[ \sum_{t=1}^T R(x_t, a_t^*) - \sum_{t=1}^T R(x_t, I_t) \right]$$

$\uparrow$   
randomness over contexts & arms played

Assume that:

$$R(x, a) = \left\langle \underbrace{\Psi(x, a)}_{\substack{\text{known} \\ \text{feature map} \\ \Psi(x, a) \in \mathbb{R}^d \\ \forall (x, a)}}, \underbrace{\theta^*}_{\substack{\text{unknown} \\ \theta^* \in \mathbb{R}^d}} \right\rangle$$

$\left. \begin{array}{l} \text{known} \\ \text{feature map} \\ \Psi(x, a) \in \mathbb{R}^d \\ \forall (x, a) \end{array} \right\} \text{ does not depend on context}$

