

## Week 9 - Lecture 43

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

$\Psi(x, a) \in \mathbb{R}^d$   
 $\forall (x, a)$

Example: Recommendation Engine  $\begin{bmatrix} \text{User: context} \\ \text{Product: arm} \end{bmatrix}$

Products can be categorized : Sports,  
daily use, movie etc.

Users can be categorized: sex, Age, location

young User : Girl	Product : DVD, tennis ball
(sex, Age, location)	(Sport, movie, joy)
(1, 1, 0)	(0, 1, 1) (1, 0, 1)

Depending on user & product feature map  
can be formed.

Such feature maps can be created.

These feature maps can be in millions, but  
these are generated offline.

Thus the problem is to find  $\theta^*$ , with interaction with env.

$\theta^*$  is a const., doesn't depend on context & action.

Fixed for an env., changes with change in env.

Assume dimension of  $\theta^*$  is known, also assume norm is bounded.

$$\|\theta^*\|_2^2 \leq L \quad \leftarrow L \text{ is known}$$

So,

$$\begin{aligned} |R(x, a) - R(x', a')| \\ \leq \|\theta^*\|_2 \|\psi(x, a) - \psi(x', a')\| \end{aligned}$$

Lipschitz function:

$$|f(x) - f(x')| \leq L \|x - x'\|$$

Thus  $|R(x, a) - R(x', a')|$  is Lipschitz with const.  $L$ .

This limits  $\theta^*$  to a bounded space.

Thus, assuming that:

Rewards are linear & parametrized by  $\theta^*$ .

Thus in the game :

- $\underbrace{x_t}_{\text{context}} \rightarrow$
- ①  $\{ \psi(x_t, a) \}_{a \in [K]} = D_t \leftarrow \begin{array}{l} \text{feature set} \\ \text{(Decision set)} \end{array}$
  - ②
  - ③ Choose the feature map with the given context.
  - $\vdots$
  - ④ What we have thus is :

$$R(x, a) = \langle \psi(x, a), \theta^* \rangle$$

where we need to maximize this over the feature vectors.

Thus this can be seen as linear optimization problem over a feature set

$\{ D_t \}_{t \geq 1}$  Now find :

$$d_t^* = \arg \max_{d \in D_t} \langle d, \theta^* \rangle$$

After making the assumption that rewards are linearly parametrized by  $\theta^*$   $\therefore$  arms have no real significance, but we need to maximize the features.

Thus: Stochastic Linear Bandits

$\{ D_t \}_{t \geq 1}$  are sequence of decision sets

Choose  $d_t^*$  s.t.

$$d_t^* = \operatorname{argmax}_{d \in D_t} \langle d, \theta^* \rangle$$

$d_t \in D_t$  in round  $t$ ,

$$r_t = \langle d_t, \theta^* \rangle + \eta_t$$

Regret :

$$\hat{R}_T = \mathbb{E} \left[ \sum_{t=1}^T \max_{d_t \in D_t} \langle d, \theta^* \rangle - \sum_{t=1}^T \langle d_t, \theta^* \rangle \right]$$

---

If  $D_t$  is a set of unit vectors

$$D_t = \{e_1, e_2, e_3, \dots, e_d\}$$

$$\text{where } e_1 = \{1, 0, 0, \dots, 0\}$$

$$e_2 = \{0, 1, 0, \dots, 0\}$$

$$\vdots$$

$$e_d = \{0, 0, 0, \dots, 1\}$$

Assume  $\theta^*$  is a vector with each component being the mean of each arm.

$$\therefore e_i^T \theta^* = \theta_i$$

Thus the stochastic Linear Bandits captures  
the Stochastic D-armed Bandits.

(if  $\theta^*$  has D-dimensions)

But  $D_t$  need not be unit vectors

$$D_t \subset \mathbb{R}^d$$