

## Lec 9

Page 1

CS59009 - RL

Agenda:

- MOSS & Bernoulli bandit, Adversarial bandit
- contextual and linear bandit
- Linear regression.

Minimax Optimal Strategy is Stochastic  
multi-armed bandit (MOSS)

With a different analysis, we can set

MOSS

$$\rightarrow UCB_i(\delta, t) = \sqrt{\frac{4}{T_i(t-1)} \log^+ \left( \frac{n}{K T_i(t-1)} \right)} + \hat{\mu}_{i,t-1}$$

where  $\log^+(x) = \log(\max\{1, x\})$

Theorem: For any 1-sub-Gaussian  $K$ -armed bandit,  
the regret of MOSS satisfies

$$R_n \leq 38\sqrt{Kn} + \sum_{i=1}^K \Delta_i$$

proof: Chapter 9. Bandit Algorithm.

You are given some ~~a~~ side information.

Bernoulli bandit:

Consider a  $K$ -armed bandit with  $\mu_i \in [0, 1]$  for all  $i$  and  $x_t \sim B(\mu_{A_t})$

Def: Relative entropy between Bernoulli distributions with parameters  $p, q \in [0, 1]$  is

$$d(p, q) = p \log \left( \frac{p}{q} \right) + (1-p) \log \left( \frac{1-p}{1-q} \right)$$

Relative entropy is also known as Kullback-Leibler (KL) divergence.

KL-UCB

$$A_t = \arg \max_i \max_{\mu \in [0, 1]} d(\hat{\mu}_{i,t-1}, \mu)$$

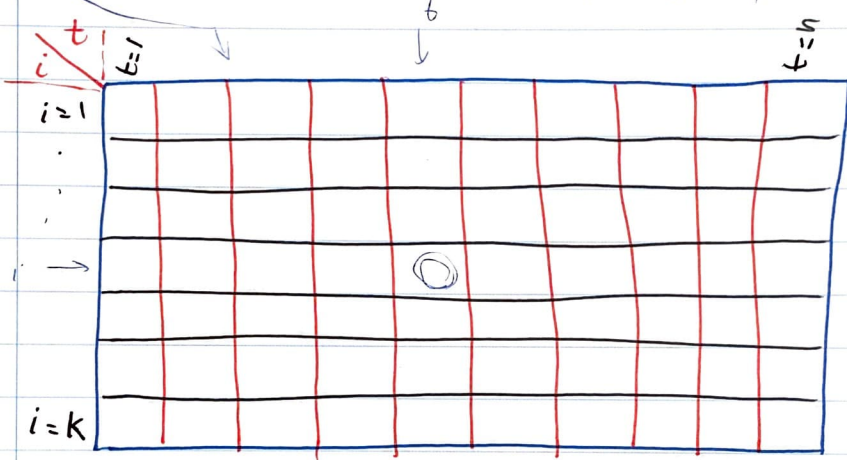
$$\left\{ \frac{\log(1 + t \log^2 t)}{T_i(t-1)} \right\}$$

chapter 10 Bandit Algorithms

## Adversarial bandit:

- A  $K$ -armed adversarial bandit is an arbitrary sequence of reward vectors

$\vec{r} = (r_1, \dots, r_n)$  where  $r_t \in [0, 1]^K$  for each  $t \in [n]$ .



$\hookrightarrow (i, t)$  entry represent the reward at time  $t$  if arm  $i$  is pulled.

In each round ~~at~~ the learner chooses an action  $A_t \in [K]$  and then observes  $X_t = r_{t, A_t}$ .

Let's look at the worst-case regret:

$$R_n^*(\pi) = \sup_{\vec{r} \in [0, 1]^{K \times n}} R_n(\pi, \vec{r})$$

For deterministic policies  $\rightarrow R_n^* \geq n(1 - \frac{1}{K})$

There is an algorithm called

Exponential-weight algorithm for

Exploration and Exploitation (EXP3)

which achieves regret of  $\leq 2\sqrt{NK \log(K)}$   
(chapter 11 - Bandit Algorithm)

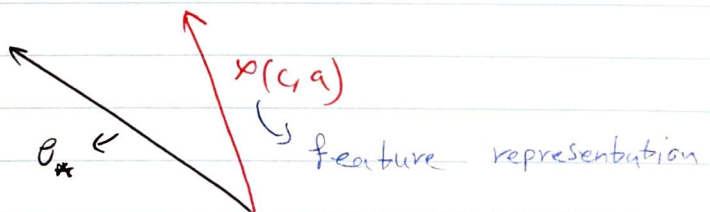
Contextual bandit:

- Consider a bandit setting where at each time step  $t$ , we have a set of actions  $A_t$ , and, before making a decision, we are also given a context  $C_t$ . When we observe a context  $C$ , and choose an action  $a$ , we receive a reward  $X$ , which depends on  $a$  and  $C$ .

Linear bandit

Consider a setting of contextual bandit where

$$X_t = \langle \phi(C, a), \theta_* \rangle + \eta_t$$



Def: Let  $X_0, \dots$  be a sequence of random variables, on  $(\Omega, \mathcal{F}, \mathbb{P})$ , and  $\mathbb{F} = \{\mathcal{F}_t\}_{t \geq 0}$  a filtration. An  $\mathbb{F}$ -adapted sequence of random variable  $(X_t)_{t \geq 0}$  is a  $\mathbb{F}$ -adapted supermartingale sequence if

$$a) E[X_t | \mathcal{F}_{t-1}] \leq X_{t-1} \quad \text{a.s.} \quad \forall t > 0$$

$$b) X_t \text{ is integrable} \quad t \geq 0$$

Note that a martingale sequence is also supermartingale (similarly, submartingale is defined when  $E[X_t | \mathcal{F}_{t-1}] \geq X_{t-1}$ )

The general case of maximal inequality:

Theorem: (Maximal inequality): Let  $(X_t)_{t \in \mathbb{N}}$  be a supermartingale with  $X_t \geq 0$  a.s. for all  $t$ . Then for any  $S > 0$

$$\mathbb{P}\left(\sup_{t \in \mathbb{N}} X_t > S\right) \leq \frac{E[X_0]}{S}$$


---