

Lecture 22

CS 59000 - RL

MDP

Agenda:

- Finish UCRL2 proof
- other models.

$$R_k \leq D + \sum_{t=r_k}^{T_{k+1}-1} \underbrace{E[V_k(x_{t+1}) | \mathcal{F}_t] - V_k(x_{t+1})}_{\text{mean zero}} \leftarrow A_1$$

$$+ \frac{D\sqrt{L|\mathcal{X}|}}{2} \sum_{x_{1:a} \in \mathcal{X}_{1:A}} \frac{T_k(x_{1:a})}{\sqrt{\max(1, T_{k-1}(x_{1:a}))}} \leftarrow A_2$$

→ number of  $x_{1:a}$  samples  $\leq$  epoch  $k$

$$\Rightarrow \hat{R}_T \leq KD + \sum_{t=1}^T E[V_{k_b}(x_{t+1}) | \mathcal{F}_t] - V_{k_b}(x_{t+1})$$

↳ mean epoch number  $\leq$  time  $t$

$$+ \frac{D\sqrt{L|\mathcal{X}|}}{2} \sum_{x_{1:a} \in \mathcal{X}_{1:A}} \sum_{k=1}^K \frac{T_k(x_{1:a})}{\sqrt{1 + T_{k-1}(x_{1:a})}}$$

Using Azuma inequality

we have:  $\mathbb{P}\left(\mathcal{F} \text{ and } \left| \sum_{t=1}^n E[V_{k_b}(x_{t+1})] - V(x_{t+1}) \right| \leq \frac{D\sqrt{T \log \frac{2}{\delta}}}{2} \right) \leq 1 - \frac{\delta}{2}$

Monday, November 9, 2020

$$\text{Now: } \sum_{n \in X} \sum_{a \in A} \frac{T_n(n, a)}{\sqrt{1 + \sum_{\tau_{k-1}} T_{\tau_{k-1}}(n, a)}}$$

Lemma 19 (UCRL2)

For any sequence of numbers  $z_i \in \mathbb{N} \cup \{0\}$

$$\text{with } 0 \leq z_k \leq \bar{z}_{k-1} = 1 + \sum_{i=1}^{k-1} z_i$$

$$\sum_{k=1}^K \frac{z_k}{\sqrt{\bar{z}_{k-1}}} \leq (\sqrt{2} + 1) \sqrt{\bar{z}_K}$$

Proof. HW.

Using this lemma:

$$\sum_{k=1}^K \frac{T_{\tau_k}(n, a)}{\sqrt{1 + \sum_{\tau_{k-1}} T_{\tau_{k-1}}(n, a)}} \leq (\sqrt{2} + 1) \sqrt{T_{\tau}(n, a)}$$

$$\text{Since } \sum_{(n, a) \in X \times A} T_{\tau}(n, a) = T$$

$$\Rightarrow \sum_{n \in X} \sum_{a \in A} \sqrt{T_{\tau}(n, a)} \leq \sqrt{|X| |A| T}$$

Monday, November 9, 2020

$$\Rightarrow \hat{R}_T \leq KD + \frac{D\sqrt{2|\mathcal{X}|}}{2} \sqrt{|\mathcal{X}| |A| T} + \frac{D\sqrt{T \log \frac{2}{\delta}}}{2}$$

with probability at least  $1 - \frac{\delta}{2}$

what is  $K$ ?

How many times for  $(n, a)$ , we can double the epoch?

$$\underbrace{1 + \log_2 \frac{T}{T}(n, a)}$$

$$\Rightarrow K \leq \sum_{(n, a) \in \mathcal{X} \times A} \underbrace{\left(1 + \log_2 \left(\frac{T}{T}(n, a)\right)\right)}$$

since  $\sum_{(n, a) \in \mathcal{X} \times A} T(n, a) = T \Rightarrow$  the max for the right hand side:

$$K \leq |\mathcal{X}| |A| \left(1 + \log_2 \left(\frac{T}{|\mathcal{X}| |A|}\right)\right)$$

we just proved UCR2 achieves regret of

$$\rightarrow \frac{D|\mathcal{X}| |A| \left(1 + \log_2 \frac{T}{|\mathcal{X}| |A|}\right)}{2} + \frac{D\sqrt{2|\mathcal{X}|}}{2} \sqrt{|\mathcal{X}| |A| T} + \frac{D\sqrt{T \log \frac{2}{\delta}}}{2}$$

i.e.  $\tilde{O}(D|\mathcal{X}| \sqrt{|A| T})$



Monday, November 9, 2020

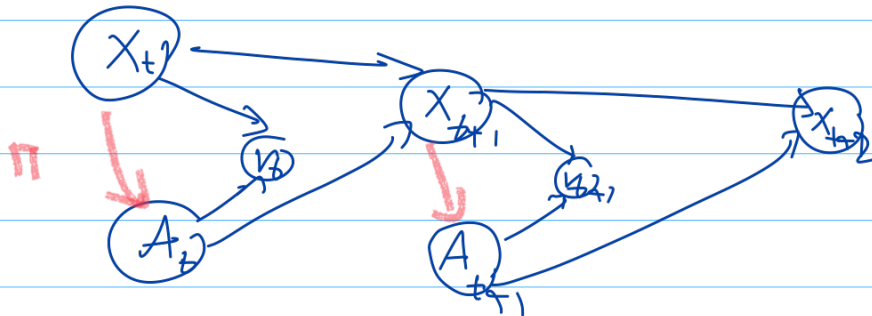
Note: the lower bound is  $\sqrt{D(X) |A| T}$ ,  $\Rightarrow$   
we are off with  $\sqrt{D(X)}$

Subsequent work show  $\sqrt{D(X) |A| T} + \text{poly}(\log(X, |A|))$   
 $\hookrightarrow \tilde{O}(\sqrt{D(X) |A| T})$

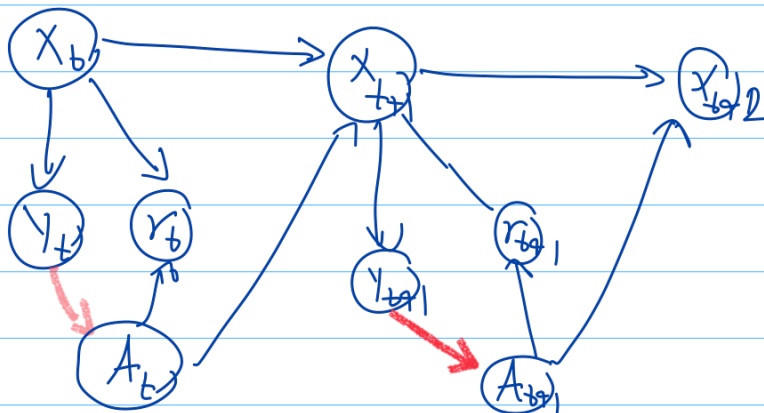
---

## Rich observable MDP (ROMDP)

→ MDP



→ ROMDP



There is a onto mapping from  $Y \rightarrow X$   
i.e. given  $Y$  we know the underlying  $X$

Monday, November 9, 2020

in other words, there exist a function  $f: Y \rightarrow X$

$$\text{RCMDP: } M := (X, Y, A, P, P, R, O)$$

$$\left. \begin{array}{l} \text{Infinite case - } X_1 \sim P \\ X_{t+1} \sim P(\cdot | X_t, A_t) \\ X_t \sim O(\cdot | X_t) \\ r_t \sim R(X_t, A_t) \end{array} \right\} \begin{array}{l} O \text{ is a stochastic} \\ \text{mapping from } X \text{ to } Y \\ \text{(surjective map)} \\ f \text{ is its inverse map.} \\ \text{which onto.} \end{array}$$

RCMDP  $\rightarrow$  results in MDP on  $Y$  space  
with  $|Y|$  number of states  
and  $D_y$  as its diameter.  
 $\rightarrow Y$

$\Rightarrow$  UCB style algorithms give regret of  $\tilde{O}(|X| D \sqrt{KT})$   
However UCRL2 guarantees  $\tilde{O}(|Y| D_y \sqrt{KT})^x$   
where  $|Y| \gg |X|$   $D_y \gg D_x$

Azizzadenesheli 2017

Monday, November 9, 2020

## Partially Observable MDPs (POMDP)

simple case

$$\rightarrow M := (X, A, Y, P, P_1, O, R, \gamma)$$

where

$X$  := State space

$A$  := Action space

$Y$  := Observation space

$P$  := Transition kernel

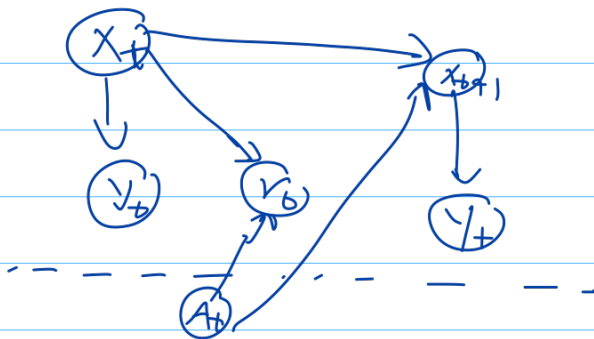
$P_1$  := Initial state measure

$R$  := Reward kernel

$\gamma$  := Discount factor.

The process:

$$- X_1 \sim P_1$$



$$Y_t \sim O(X_t)$$

$$r_t \sim R(X_t, A_t)$$

$$X_{t+1} \sim P(X_t, A_t)$$

How  $A_t$  is generated?

Monday, November 9, 2020

Policy:

- History dependent policy  $A_t \sim \Pi(\mathcal{H}_t)$

- Memory less policy  $A_t \sim \Pi(y_t)$  history ( $\dots, A_{t-1}, y_{t-1}$ )

- Markovian policy  $A_t \sim \Pi(y_{t-1}, t)$

→ Note MDP  $\subset$  RMDP  $\subset$  PGMDP

---