

# Lecture 23

Wednesday, November 11, 2020

CS59000 - RL

PGMDP

Agenda

- Protocol
- Belief
- Hardness

---

Let's consider finite state, action, observation spaces. Then we can write  $P, O, P_1$  in the following sense.

$$P(x'|x, a)$$

$$O(y|x)$$

$$P_1(x)$$

---

Protocol:

At time step 1, -  $X_1 \sim P_1$

-  $Y_1 \sim O(X_1)$

- choose action  $A_1$

- succeed to  $X_2 \sim P(X_1, A_1)$

⋮

---

At time  $t=1$ , without any observation

$P_1$  : prior

After observing  $Y_1$ , then what we think

about  $X_1$ ?  $b_1(x) = P_r(X_1=x | Y_1=y_1)$

↙ belief at time  $t=1$  ↗

How to compute  $b_1$ ?

$$b_1(x) = \frac{P_r(X_1=x, Y_1=y_1)}{P_r(Y_1=y_1)} = \frac{O(y_1|x) P_1(x)}{\sum_{x_1} O(y_1|x_1) P_1(x_1)}$$

note:  $b_1(x) = P_r(X_1=x | Y_1=y_1, A_1=a)$

if  $A_1$  is chosen irrespective to  $X_1$ .

---

$$b_t(x) = P_r(X_t=x | \underbrace{y_1, a_1, \dots, y_t}_{h_t})$$

$$\begin{aligned} &\rightarrow \frac{P_r(Y_t=y_t | X_t=x, h_{t-1}, a_{t-1}) P_r(X_t=x | h_{t-1}, a_{t-1})}{P_r(Y_t | h_{t-1}, a_{t-1})} \end{aligned}$$

$$b_t(n) = \frac{o(y_t | n) \sum_{n^i} P(x_t = n | x_{t-1} = n^i, a_{t-1}) P_r(x_{t-1} = n^i | h_{t-1}, a_{t-1})}{\sum_{n^i} \sum_{n^j} o(y_t | n^j) P(n^j | n^i, a_{t-1}) P_r(n^i | h_{t-1}, a_{t-1})}$$

$b_{t-1}(n^i)$

$$\Rightarrow P_r(x_{t-1} = n^i | h_{t-1}, a_{t-1}) = P_r(x_{t-1} = n^i | h_{t-1}) = b_{t-1}(n^i)$$


---

$b \in \Delta^X \Rightarrow$  we can have a policy as a function of  $b_b$  i.e.  $A \sim \pi(b_b) \Rightarrow$  MDP on belief space

---

Optimal policy:

- Infinite horizon PG MDP

$\hookrightarrow$  Undecidable Madani et al. 1999

- Fixed horizon:

given  $y_t \rightarrow b_1(n) ; |y|$  vectors.

$b_2(n) ; |y|^2 |A|$

$b_3(n) ; |y|^2 |A|^2$

Wednesday, November 11, 2020

$$L_H(n); |Y|^H |A|^{H-1}$$

$\Rightarrow$  PSPACE-Complete

Papadimitriou and Tsitsiklis 1987

---

Memory less policy.

NP-hard

Vlassis et al. 2002

The optimal policy is stochastic Singh et al 1994

---

Memory-less  $\Rightarrow$  Much better

Optimism in PG MDP

$$\tilde{O}(D|X|^{3/2} \sqrt{|A||Y|T})$$

$\rightarrow$  Azizzadenesheli 2016

---

So many open Problem.

## off-policy learning.

$\pi_b$  : behavioural policy ;

- using data generated by  $\pi_b$
- can we say how well a new  $\pi$  does
- can we find an optimal policy.

---

Remember the Contextual bandit, i.e. a Fixed horizon MDP with horizon 1.

---

Consider a stationary Contextual bandit, with context set  $\mathcal{X}$ , action set  $\mathcal{A}$ , and a measure  $\mu$  on  $\mathcal{X}$ . with  $R$ , the reward kernel and mean  $\bar{r}$  ;  $0 \leq \bar{r} \leq 1$   
 $0 \leq r \leq 1$

Protocol:

- Draw  $X \sim \mu$
- Draw  $A \sim \pi(X_t)$
- receive  $r \sim R(X_t, A_t)$

$$\begin{aligned} \text{For } \pi ; \quad \eta(\pi) &= E_{\mu} [E_{\pi} [r | X]] \\ &= E_{\mu} [E_{\pi} [\bar{r}(X, A) | X]] \end{aligned}$$

After following  $\pi_b$  for  $T$  time step,

we have  $D_T = \left\{ x_t, A_t, r_t \right\}_{t=1}^T$

what is  $\eta(\pi_b)$ ?

$$\hat{\eta}(\pi_b) = \frac{1}{T} \sum r_t \quad \left| \hat{\eta}(\pi_b) - \eta(\pi_b) \right| \leq \sqrt{\frac{1}{2T} \log \frac{2}{\delta}}$$

with probability at least  $1 - \delta$

How about  $\eta(\pi)$ ?

For simplicity, let  $\mathcal{A}$  to be a finite set.  
space and  $\pi(dA; x) = \pi(A; x) dA$

Define  $\hat{r}_t(x_t, a) = r_t \frac{1(A_t = a)}{\pi_b(A_t; x_t)}$

↳ the inverse propensity score (IPS)

$$\begin{aligned} E_{\pi_b} [\hat{r}_t(x_t, a) | x_t] &= \int \frac{\bar{r}(x_t, A_t) 1(A_t = a)}{\pi_b(A_t; x_t)} \pi_b(A_t; x_t) dA_t \\ &= \bar{r}(x_t, a) \end{aligned}$$

now we can use  $\hat{v}$  to estimate  $\eta(\pi)$

$$\begin{aligned} \eta_{IPS}(\pi) &= \frac{1}{T} \sum_{t=1}^T E_n[\hat{v} | x_t] \\ &= \frac{1}{T} \sum_{t=1}^T \int \hat{v}(x_t, A) \pi(A; x_t) dA \\ &= \frac{1}{T} \sum_{t=1}^T \frac{\pi(A_t; x_t)}{\pi_b(A_t; x_t)} \end{aligned}$$

why it is good?

$$\eta(\pi) = E_\mu[E_n[v | x]] = E_\mu E_{\pi_b} \left[ \underbrace{\frac{\pi}{\pi_b}}_{\text{Importance weight}} v | x \right]$$

justification,

$$\begin{aligned} \eta(\pi) &= E_\mu[E_n[\bar{v}(x, A) | x]] \\ &= E_\mu \left[ \int \bar{v}(x, A) \pi(A; x) dA \right] \\ &= E_\mu \left[ \int \underbrace{\bar{v}(x, A) \frac{\pi(A; x)}{\pi_b(A; x)}}_{\text{Importance weight}} \pi_b(A; x) dA \right] \\ &\rightarrow E_\mu \left[ E_{\pi_b} \left[ \bar{v}(x, A) \frac{\pi(A; x)}{\pi_b(A; x)} | x \right] \right] \end{aligned}$$

$\Rightarrow$  therefore, the empirical estimate

$$\hat{\eta}(\pi) = \sum_{t=1}^T V(X_t, A_t) \frac{\pi(A_t; X_t)}{\pi_b(A_t; X_t)}$$

i.e. an Monte Carlo estimate of the integral

$$\Rightarrow \hat{\eta}_{\text{IPS}}(\pi) = \hat{\eta}(\pi)$$

How good is this estimation?

Assume  $\omega_{\max} := \text{ess sup } \frac{\pi}{\pi_b} < \infty$

$$|\hat{\eta}(\pi) - \eta(\pi)| \leq \omega_{\max} \sqrt{\frac{1}{2T} \log \frac{2}{\delta}}$$

w.p. at least  $1-\delta$ .