

## Lecture 24

CS\_59000 - RL

Contextual Bandit & MDP

Agenda

- off policy policy evaluation
- off policy policy optimization

---

From last lecture:  $D := \{X_i, A_i, r_i\}_{i=1}^T$   
generated by  $\pi_b$

Estimate the performance of  $\pi$

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

$$\omega_{\max} := \text{esssup} \frac{\pi}{\pi_b} \leftarrow$$

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

with probability at least  $1 - \delta$ . (From Hoeffding)

what is the variance of  $\hat{\eta}(\pi)$ ?

let  $\sigma^2(x, a)$  denote  $\text{Var}[r(x, A) | x, A](x, a)$

$$\Rightarrow \text{Var}_{\mu, \pi_b}[\hat{\eta}(\pi)] = \frac{1}{T} \text{Var}_{\mu, \pi_b} \left[ \omega(x, A) r \right]$$

Importance weight  
function: i.e.  $\frac{\pi(x, A)}{\pi_b(x, A)}$

Using the law of total variance

$$\Rightarrow \text{Var}[X_1] = \underbrace{E[\text{Var}[X_1 | X_2]]}_{\text{a.s.}} + \text{Var}[E[X_1 | X_2]]$$

$$\begin{aligned} &\Rightarrow \frac{1}{2} \text{Var}_{\mu, \pi_b}[\omega(x, A) r] \\ &= \frac{1}{T} E_{\mu, \pi_b} \left[ \text{Var}[\underbrace{r}_{\omega(x, A)} | \underbrace{x, A}] \right] \\ &\quad + \frac{1}{T} \text{Var}_{\mu, \pi_b} \left[ E[\omega(x, A) r | x, A] \right] \\ &= \frac{1}{T} E_{\mu, \pi_b} \left[ \underbrace{\omega^2(x, A) \sigma^2(x, A)} \right] + \frac{1}{T} \text{Var}_{\mu, \pi_b} [\omega \bar{r}] \end{aligned}$$

⇒ One can use Bernstein or empirical Bernstein  
to come up with a better bound  
maybe a

Let's imagine some one gives us  $r'(x, a) \forall x, a \in \mathcal{X} \times \mathcal{A}$   
 ↪ a surrogate. Imagine  $r' \approx \bar{r}$   
 means, somewhat model based.

Doubly robust estimator

$$\hat{r}_b^{DR}(x_b, a) = \underbrace{r'(x_b, a)}_{\text{(finite action spaces)}} + \underbrace{\left( \bar{r}_+ - r'(x_b, a) \right)}_{\uparrow} \frac{1(A_i = a)}{\pi_b(A_i | x_b)}$$

HW: Similar to  $\hat{r}_t(x_t, a)$ , show that  $\hat{r}_b^{DR}(x_b, a)$

is also an unbiased estimator of  $\bar{r}(x_b, a)$

$$\begin{aligned} \Rightarrow \text{Var}^{DR}(\Pi) &= \frac{1}{T} E_{\mu, \Pi_b} \left[ \sigma^2(x, A) \omega^2(x, A) \right] \\ &+ \frac{1}{T} E_{\mu} \left[ \text{Var}_{\Pi_b} \left[ \omega(x, A) \left( \bar{r}(x, A) - r'(x, A) \right) \mid X \right] \right] \\ &+ \frac{1}{T} \text{Var}_{\mu} \left[ E_{\Pi_b} \left( \bar{r}(x, A) \omega(x, A) \mid X \right) \right] \end{aligned}$$

Someone gives us  $r'$   
we can learn it ourselves:

$$\text{E.g. given } D' \Rightarrow \inf_{f \in F} \frac{1}{T} \sum (f(x_t, A_t) - r(x_t, A_t))^2$$

$\hookrightarrow$  a class of functions.

$\Rightarrow$  Set  $r$  to be an  $f \in F$  with low loss

---

Sometime model based + IS  
might end up with better result.

---

policy optimization.

- Give a set of policies  $\Pi_N$  with  $N$  policies
- which one is the best?

we are looking for  $\max_{\pi \in \Pi_N} J(\pi) \rightarrow \pi^*$

E.g. w we IS.  $\Rightarrow \forall \pi \in \Pi_N \rightarrow \hat{J}(\pi)$

$\nearrow \max_{\pi \in \Pi_N} \hat{J}(\pi) \rightarrow \hat{\pi}^*$

How good is  $\hat{\Pi}^*$   $\Rightarrow \eta(\hat{\Pi}^*) \Leftrightarrow \eta(\Pi^*)$

Theorem

$$\eta(\hat{\Pi}^*) \geq \eta(\Pi^*) - \omega_{\max} \sqrt{\frac{2}{T} \log \frac{2N}{\delta}}$$

proof: using union bound and Hoeffding's

$$\rightarrow \max_{\Pi \in \Pi_N} |\hat{\eta}(\Pi) - \eta(\Pi)| \leq \omega_{\max} \sqrt{\frac{1}{2T} \log \frac{2N}{\delta}}$$

with prob at least  $1-\delta$ .

$$\text{where } \omega_{\max} = \max_{\Pi \in \Pi_N} \sup \frac{\Pi}{\Pi_b}$$

$$\begin{aligned} \rightarrow \eta(\hat{\Pi}^*) &\geq \hat{\eta}(\hat{\Pi}^*) - \omega_{\max} \sqrt{\frac{1}{2T} \log \left( \frac{2N}{\delta} \right)} \\ &\geq \hat{\eta}(\Pi^*) - \omega_{\max} \sqrt{\frac{1}{2T} \log \left( \frac{2N}{\delta} \right)} \\ &\geq \eta(\Pi^*) - \omega_{\max} \sqrt{\frac{2}{T} \log \left( \frac{2N}{\delta} \right)} \end{aligned}$$

Bingo

Now for MDPs.

Consider a fixed horizon MDP

$$M(X, A, P, R, P_0, \gamma, H); \gamma \in [0, 1]$$

$$\text{Let } \tau_h = (x_1, A_1, r_1, \dots, x_h, A_h, r_h)$$

in an episode.

stage h

$$\rightarrow J(\pi) = E_{P_\pi} \left[ \sum_{h=1}^H \gamma^{h-1} r_h \right]$$

we are given data from  $\pi_b \Rightarrow J(\pi)$

$$\text{If } P_\pi < P_{\pi_b} \Rightarrow J(\pi) = E_{P_{\pi_b}} \left[ \frac{dP_\pi}{dP_{\pi_b}} \left( \sum_{h=1}^H \gamma^{h-1} r_h \right) \right]$$

E.g. For special case of Euclidean space:

$$dP_\pi = P_1(x_1) \pi(A_1; x_1) P(x_2; x_1, A_1) \dots dx_1 dA_1 \dots$$

$$\hookrightarrow \frac{dP_\pi}{dP_{\pi_b}} = \prod_{h=1}^H \frac{\pi(A_h; x_h)}{\pi_b(A_h; x_h)} \quad ; \quad \rho_h = \frac{\pi(A_h; x_h)}{\pi_b(A_h; x_h)}$$

Simply for step-IS we have ↓

$$v(\pi) = E_{\pi_b} \left[ \sum_{h=1}^H \gamma^{h-1} \left( \sum_{h'=1}^h \rho_{h'} \right) r_h \right]$$


---

$$V_H = \bar{r} \rightarrow V_{H-1} = r + P V_H$$

$$V_{H-2} = r + P V_{H-1}$$


---

We are given data  $\pi_b : D := \left\{ (x_h^t, A_h^t, r_h^t) \right\}_{h=1, b=1}^{H, T}$

$$\hat{V}_h(n) = \frac{1}{T_h(n)} \sum_{t=1}^T \rho_h^t (r_h^t + \gamma \hat{V}_{h+1}) \mathbb{I}(x_h^t = n)$$

$$\hookrightarrow \hat{V}_H(n) = \frac{1}{T_H(n)} \sum_{t=1}^T \underbrace{\rho_H^t(x_H^t, A_H^t)}_{\rho_H^t} \underbrace{r_H^t}_{r_H^t} \mathbb{I}(x_H^t = n)$$

$$\hat{V}_{H-1}(n) = \frac{1}{T_{H-1}(n)} \sum_{t=1}^T \rho_{H-1}^t(x_{H-1}^t, A_{H-1}^t) (r_{H-1}^t + \gamma \hat{V}_H(x_H^t)) \mathbb{I}(x_{H-1}^t = n)$$

DR:

$$\hat{V}_h^{DR}(x) = \underbrace{V'_h(x)} + \frac{1}{T_h(x)} \sum_{t=1}^T \underbrace{\rho_h^t(r_h^t + \gamma \hat{V}_{h+1}^{DR} - Q'(x_h, A_h))}_{I(x_h^t = x)}$$

---

Both of these estimators are unbiased.

---

What are the variances:

$$\text{Var}[\hat{V}_h^{DR}] = \text{Var}[V'_h(x_h)]$$

$$\begin{aligned} &+ E\left[\text{Var}\left[\rho_h\left(Q'(x_h, A_h) - Q(x_h, A_h)\right) \mid x_h\right]\right] \\ &+ E\left[\rho_h^2 \text{Var}[r_h \mid x_h, A_h]\right] \\ &+ E\left[\gamma^2 \rho_h^2 \text{Var}[V_{h+1}^{DR} \mid x_h, A_h]\right] \end{aligned}$$

---

Proof:

Final exam



Model based (UCRL2)

Model free  $\downarrow$   $\rightarrow$  Value based (Q-learning, SARSA)  
 $\downarrow$   $\rightarrow$  Policy based

model based

$\hookrightarrow M \rightarrow Q \rightarrow \pi \Rightarrow \text{good } \eta(\pi)$

model free - value based

$\rightarrow Q \rightarrow \pi \Rightarrow \text{good } \eta(\pi)$

model free - policy based

$\pi \Rightarrow \text{good } \eta(\pi)$

$\max_{\pi \in \Pi} \eta(\pi)$

$\Rightarrow$  Policy gradient

$$\nabla_{\pi} \eta(\pi) \Big|_{\pi=\pi_t} \rightarrow \pi_{t+1} \leftarrow \text{Proj} \left( \pi_t + \alpha \nabla_{\pi} \eta(\pi) \Big|_{\pi=\pi_t} \right)$$