

# ELEN E6885: Introduction to Reinforcement Learning

## Homework #3

Chenye Yang cy2540@columbia.edu

November 14, 2019

### P1

**Ans:**

The left-hand side of the equation can be written as

$$\begin{aligned}
 & \max_s \left| \mathbb{E}_\pi [G_{t:t+n} \mid S_t = s] - v_\pi(s) \right| \\
 &= \max_s \left| \mathbb{E}_\pi [R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n}) \mid S_t = s] - v_\pi(s) \right| \\
 &= \max_s \left| \mathbb{E}_\pi [R_{t:t+n} + \gamma^n V_{t+n-1}(S_{t+n}) \mid S_t = s] - v_\pi(s) \right|,
 \end{aligned}$$

where  $R_{t:t+n} \doteq R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n}$

Because

$$\begin{aligned}
 v_\pi(s) &\doteq \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right] \\
 &= \mathbb{E}_\pi \left[ R_{t:t+n} + \gamma^n \sum_{k=0}^{\infty} \gamma^k R_{t+n+k+1} \mid S_t = s \right] \\
 &= \mathbb{E}_\pi [R_{t:t+n} \mid S_t = s] + \gamma^n \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+n+k+1} \mid S_t = s \right]
 \end{aligned}$$

The left-hand side of the equation is

$$\begin{aligned}
 & \max_s \left| \mathbb{E}_\pi [R_{t:t+n} + \gamma^n V_{t+n-1}(S_{t+n}) \mid S_t = s] - v_\pi(s) \right| \\
 &= \max_s \left| \mathbb{E}_\pi [\gamma^n V_{t+n-1}(S_{t+n}) \mid S_t = s] - \gamma^n \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+n+k+1} \mid S_t = s \right] \right| \\
 &= \gamma^n \max_s \left| \mathbb{E}_\pi \left[ V_{t+n-1}(S_{t+n}) - \sum_{k=0}^{\infty} \gamma^k R_{t+n+k+1} \mid S_t = s \right] \right|
 \end{aligned}$$

Considering  $|\mathbb{E}[X]| \leq \mathbb{E}[|X|]$ , we get

$$\left| \mathbb{E}_\pi \left[ V_{t+n-1}(S_{t+n}) - \sum_{k=0}^{\infty} \gamma^k R_{t+n+k+1} \mid S_t = s \right] \right| \leq \mathbb{E}_\pi \left[ \left| V_{t+n-1}(S_{t+n}) - \sum_{k=0}^{\infty} \gamma^k R_{t+n+k+1} \right| \mid S_t = s \right]$$

Thus the left-hand side of the equation has

$$\max_s |\mathbb{E}_\pi [G_{t:t+n} | S_t = s] - v_\pi(s)| \leq \gamma^n \max_s \mathbb{E}_\pi \left[ \left| V_{t+n-1}(S_{t+n}) - \sum_{k=0}^{\infty} \gamma^k R_{t+n+k+1} \right| \mid S_t = s \right]$$

Given a completely free, unrestricted choice for  $S_t$ , the set of possible states  $S_{t+n}$  is a subset of the set of possible states  $S_t$ . Thus the following inequality is true.

$$\max_s \mathbb{E}_\pi \left[ \left| V_{t+n-1}(S_{t+n}) - \sum_{k=0}^{\infty} \gamma^k R_{t+n+k+1} \right| \mid S_t = s \right] \leq \max_s \mathbb{E}_\pi \left[ \left| V_{t+n-1}(S_t) - \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \right| \mid S_t = s \right]$$

Therefore, the left-hand side of the equation

$$\begin{aligned} \max_s |\mathbb{E}_\pi [G_{t:t+n} | S_t = s] - v_\pi(s)| &\leq \gamma^n \max_s \mathbb{E}_\pi \left[ \left| V_{t+n-1}(S_t) - \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \right| \mid S_t = s \right] \\ &= \gamma^n \max_s \mathbb{E}_\pi \left[ \left| V_{t+n-1}(S_t) - v_\pi(s) \right| \mid S_t = s \right] \\ &= \gamma^n \max_s \left| V_{t+n-1}(s) - v_\pi(s) \right| \end{aligned}$$

The error reduction property of  $n$ -step returns is proved.

**P2****1.****Ans:**

On-line updating methods update value function after each step in the episode.

In off-line updating methods, updates are accumulated within episode and applied in batch at the end of episode.

**2.****Ans:**

$$\alpha = 0.1$$

On-line every-visit constant- $\alpha$  Monte Carlo method:

$$V_0(A) = 0$$

$$V_1(A) = V_0(A) + \alpha \times (1 + 2 + 1 - V_0(A)) = 0.4$$

$$V_2(A) = V_1(A) + \alpha \times (1 - V_1(A)) = 0.46$$

Off-line every-visit constant- $\alpha$  Monte Carlo method:

$$\alpha \times (1 + 2 + 1 - V_0(A)) = 0.4$$

$$\alpha \times (1 - V_0(A)) = 0.1$$

$$V(A) = V_0(A) + 0.4 + 0.1 = 0.5$$

**3.****Ans:**

On-line  $TD(0)$  method,  $\lambda = 0$ :

$$V_1(A) = V_0(A) + \alpha (1 - V_0(A)) = 0.1$$

$$V_2(A) = V_1(A) + \alpha (1 - V_1(A)) = 0.19$$

Off-line  $TD(0)$  method,  $\lambda = 0$ :

$$\alpha (1 - V_0(A)) = 0.1$$

$$\alpha (1 - V_0(A)) = 0.1$$

$$V(A) = V_0(A) + 0.1 + 0.1 = 0.2$$

**4.****Ans:**

$$\lambda = 0.5$$

On-line forward-view  $TD(\lambda)$  method:

$$\begin{aligned} V_1(A) &= V_0(A) + \alpha (G_t^\lambda - V_0(A)) \\ &= 0.225 \end{aligned}$$

$$\begin{aligned} V(A) &= V_1(A) + \alpha (G_t^\lambda - V_1(A)) \\ &= 0.225 + 0.1 \times (1 - 0.225) \\ &= 0.3025 \end{aligned}$$

Off-line forward-view  $TD(\lambda)$  method:

$$\alpha (G_{t_1}^\lambda - V_0(A)) = 0.225$$

$$\alpha (G_{t_2}^\lambda - V_0(A)) = 0.1$$

$$V(A) = V_0(A) + 0.225 + 0.1 = 0.325$$

**5.**

**Ans:**

On-line backward-view  $TD(\lambda)$  method,  $\gamma = 1, \lambda = 0.5$ :

Because

$$E_0(A) = E_0(B) = E_0(T) = 0$$

$$E_1(A) = \gamma\lambda E_0(A) + \mathbf{1}_{S_t=s} = 1$$

$$E_1(B) = \gamma\lambda E_0(B) + \mathbf{1}_{S_t=s} = 1$$

$$\delta_A = R + \gamma V_0(B) - V_0(A) = 1$$

Thus

$$V_1(A) = V_0(A) + 0.1 \times 1 \times 1 = 0.1$$

Then

$$\delta_B = R + \gamma V_1(A) - V_0(B) = 2 + 0.1 = 2.1$$

$$V_1(B) = V_0(B) + 0.1 \times 2.1 \times 1 = 0.21$$

Then

$$E_2(A) = (\gamma\lambda)^2 E_1(A) + \mathbf{1}_{S_t=s} = 1.25$$

$$\delta_A = R + \gamma V_0(T) - V_1(A)$$

$$= 1 + 0 - 0.1$$

$$= 0.9$$

Therefore

$$V(A) = V_1(A) + 0.1 \times 0.9 \times 1.25 = 0.2125$$

Off-line backward-view  $TD(\lambda)$  method:

$$\begin{aligned} \sum_{t=1}^T \alpha \delta_t E_{t(s)} &= \sum_{t=1}^T \alpha (G_t^\lambda - V(s_t)) \mathbf{1}_{S_t=s} \\ &= 0.1 \times (0.5 + 0.75 + 1) + 0.1 \times 1 \\ &= 0.1 \times 2.25 + 0.1 \\ &= 0.325 \end{aligned}$$

$$V(A) = V_0(A) + 0.325 = 0.325$$

## P3

### 1.a

**Ans:**

Left-hand side of the equation is backward view (Off-line  $TD(1)$ ) while right-hand side is  $\lambda = 1$  forward view (Off-line every-visit constant- $\alpha$  Monte Carlo method). According to the original statement, they are equivalent.

### 1.b

**Ans:**

Eligibility trace

$$\begin{aligned} E_0(s) &= 0 \\ E_t(s) &= \gamma\lambda E_{t-1}(s) + \mathbf{1}(S_t = s) \end{aligned}$$

Assume all the  $k$ -s, which let  $S_k = s$ , are in order  $\{k_1, k_2, \dots, k_n\}$ , then we have

$$E_{k_1}(s) = (\gamma\lambda)^i E_0(s) + \mathbf{1}(S_t = s) = 1$$

Thus

$$\begin{aligned} E_{k_n}(s) &= (\gamma\lambda)^{k_n - k_{n-1}} E_{k_{n-1}}(s) + 1 \\ &= (\gamma\lambda)^{k_n - k_{n-1}} [(\gamma\lambda)^{k_{n-1} - k_{n-2}} E_{k_{n-2}}(s) + 1] + 1 \\ &= (\gamma\lambda)^{k_n - k_{n-1}} [(\gamma\lambda)^{k_{n-1} - k_{n-2}} ((\gamma\lambda)^{k_{n-2} - k_{n-3}} E_{k_{n-3}}(s) + 1) + 1] + 1 \\ &= \dots \\ &= (\gamma\lambda)^{k_n - k_1} E_{k_1}(s) + (\gamma\lambda)^{k_n - k_2} + \dots + (\gamma\lambda)^{k_n - k_{n-1}} + 1 \\ &= (\gamma\lambda)^{k_n - k_1} + (\gamma\lambda)^{k_n - k_2} + \dots + (\gamma\lambda)^{k_n - k_{n-1}} + (\gamma\lambda)^{k_n - k_n} \end{aligned}$$

Therefore,

$$E_t(s) = \sum_{k=0}^t \gamma^{t-k} \cdot \mathbf{1}(S_k = s)$$

### 1.c

**Ans:**

Accumulating eligibility trace can be written as

$$E_t(s) = \sum_{k=0}^t (\gamma\lambda)^{t-k} \mathbf{1}(S_k = s)$$

Thus, the left-hand side of the equation

$$\begin{aligned}
 \sum_{t=0}^{T-1} \alpha \delta_t E_t(s) &= \sum_{t=0}^{T-1} \alpha \delta_t \sum_{k=0}^t (\gamma \lambda)^{t-k} \mathbf{1}(S_k = s) \\
 &= \sum_{k=0}^{T-1} \alpha \delta_k \sum_{t=0}^k (\gamma \lambda)^{k-t} \mathbf{1}(S_k = s) \\
 &= \sum_{k=t}^{T-1} \alpha \delta_k \sum_{t=0}^{T-1} (\gamma \lambda)^{k-t} \mathbf{1}(S_k = s) \\
 &= \sum_{t=0}^{T-1} \alpha \mathbf{1}(S_t = s) \sum_{k=t}^{T-1} (\gamma \lambda)^{k-t} \delta_k
 \end{aligned}$$

Consider an individual update of the  $\lambda$ -return algorithm

$$\begin{aligned}
 G_t - V_t(S_t) &= -V_t(S_t) + (1 - \lambda) \lambda^0 [R_{t+1} + \gamma V_t(S_{t+1})] \\
 &\quad + (1 - \lambda) \lambda^1 [R_{t+1} + \gamma R_{t+2} + \gamma^2 V_t(S_{t+2})] \\
 &\quad + (1 - \lambda) \lambda^2 [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 V_t(S_{t+3})] \\
 &\quad \dots \\
 &= -V_t(S_t) \\
 &\quad + (\gamma \lambda)^0 [R_{t+1} + \gamma V_t(S_{t+1}) - \gamma \lambda V_t(S_{t+1})] \\
 &\quad + (\gamma \lambda)^1 [R_{t+2} + \gamma V_t(S_{t+2}) - \gamma \lambda V_t(S_{t+2})] \\
 &\quad + (\gamma \lambda)^2 [R_{t+3} + \gamma V_t(S_{t+3}) - \gamma \lambda V_t(S_{t+3})] \\
 &\quad \dots \\
 &= (\gamma \lambda)^0 [R_{t+1} + \gamma V_t(S_{t+1}) - V_t(S_t)] \\
 &\quad + (\gamma \lambda)^1 [R_{t+2} + \gamma V_t(S_{t+2}) - V_t(S_{t+1})] \\
 &\quad + (\gamma \lambda)^2 [R_{t+3} + \gamma V_t(S_{t+3}) - V_t(S_{t+2})] \\
 &\quad \dots \\
 &\approx \sum_{k=t}^{\infty} (\gamma \lambda)^{k-t} \delta_k \\
 &\approx \sum_{k=t}^{T-1} (\gamma \lambda)^{k-t} \delta_k
 \end{aligned}$$

Thus, the right-hand side of the equation

$$\sum_{t=0}^{T-1} \alpha (G_t - V(S_t)) \mathbf{1}(S_t = s) = \sum_{t=0}^{T-1} \alpha \mathbf{1}(S_t = s) \sum_{k=t}^{T-1} (\gamma \lambda)^{k-t} \delta_k$$

Therefore,

$$\sum_{t=0}^{T-1} \alpha \delta_t E_t(s) = \sum_{t=0}^{T-1} \alpha (G_t - V(S_t)) \mathbf{1}(S_t = s)$$

**2.**

**Ans:**

No. On-line methods will update value function once reaching state  $s$ . Backward and forward views use different information. They can't match exactly in every update.

**P4****1.****Ans:**

$$\hat{Q}(s, a, w) = 1$$

$$q_t^n = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n q(S_{t+n}, \Lambda_{t+n})$$

Thus

$$q_1^1 = -1 + 1 = 0, \quad q_1^2 = -1 - 1 + 1 = -1, \quad q_1^3 = -1 - 1 - 1 + 1 = -2, \quad q_1^4 = -1 - 1 - 1 - 1 + 0 = -4$$

$$q_2^1 = -1 + 1 = 0, \quad q_2^2 = -1 - 1 + 1 = -1, \quad q_2^3 = -1 - 1 - 1 + 0 = -3$$

$$q_3^1 = -1 + 1 = 0, \quad q_3^2 = -1 - 1 + 0 = -2$$

$$q_4^1 = -1 + 0 = -1$$

$$q_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} q_t^n$$

Thus, the  $\lambda$ -return  $q_t^\lambda$  corresponding to this episode is

$$q_1^\lambda = (1 - 0.5) \times (1 \times 0 + 0.5 \times (-1) + 0.25 \times (-2)) + 0.125 \times (-4) = -1$$

$$q_2^\lambda = (1 - 0.5) \times (1 \times 0 + 0.5 \times (-1)) + 0.25 \times (-3) = -1$$

$$q_3^\lambda = (1 - 0.5) \times (1 \times 0) + 0.5 \times (-2) = -1$$

$$q_4^\lambda = -1$$

**2.****Ans:**For forward-view  $TD(\lambda)$ , target is the action-value  $\lambda$ -return.

$$\Delta w = \alpha (q_t^\lambda - \hat{q}(S_t, A_t, w)) \nabla_w \hat{q}(S_t, A_t, w)$$

Thus, the sequence of updates to weight  $w_1$  is

$$\Delta w_1^1 = 0.5 \times (-1 - 1) \times 1 = -1$$

$$\Delta w_1^2 = 0.5 \times (-1 - 1) \times 1 = -1$$

$$\Delta w_1^3 = 0.5 \times (-1 - 1) \times 1 = -1$$

$$\Delta w_1^4 = 0.5 \times (-1 - 1) \times 0 = 0$$

The total update to weight  $w_1$  is  $-3$ **3.****Ans:**The  $TD(\lambda)$  accumulating eligibility trace  $\mathbf{e}_t$  when using linear value function approximation is

$$\mathbf{e}_t = \gamma \lambda \mathbf{e}_{t-1} + x(s, a)$$

The sequence of eligibility traces corresponding to *right* action are

$$1, \frac{3}{2}, \frac{7}{4}, \frac{7}{8}$$



**4.****Ans:**For back-view  $TD(\lambda)$  we have

$$\begin{aligned}\delta_t &= R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, w) - \hat{q}(S_t, \Lambda_t, w) \\ \Delta w &= \alpha \delta_t E_t\end{aligned}$$

Thus, the sequence of updates to weight  $w_1$  is

$$\begin{aligned}\Delta w_1^1 &= 0.5 \times (-1 + 1 - 1) \times 1 = -\frac{1}{2} \\ \Delta w_1^2 &= 0.5 \times (-1 + 1 - 1) \times \frac{3}{2} = -\frac{3}{4} \\ \Delta w_1^3 &= 0.5 \times (-1 + 1 - 1) \times \frac{7}{4} = -\frac{7}{8} \\ \Delta w_1^4 &= 0.5 \times (-1 + 0 - 1) \times \frac{7}{8} = -\frac{7}{8}\end{aligned}$$

The total update to weight  $w_1$  is  $-3$ **5.****Ans:**When using off-line updates and linear function approximation, forward-view and backward-view  $TD(\lambda)$  are equivalent to each other.