

# سوال 1

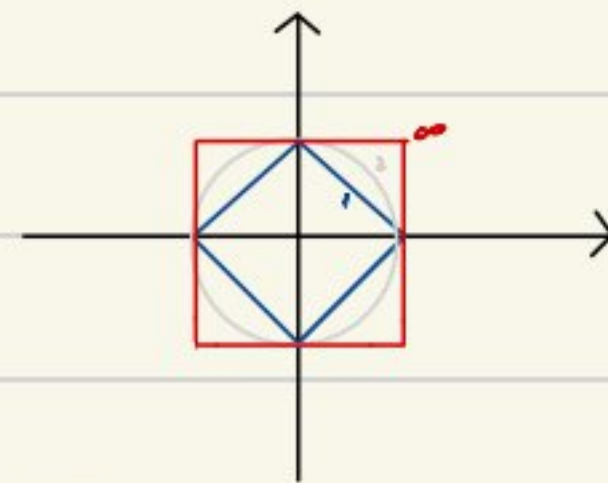
لف

To show that  $U(V)$  is a contraction mapping; we have to prove that  $\|U(V_1) - U(V_2)\| \leq \|V_1 - V_2\|$ ; i.e.,  $U(\cdot)$  makes two value function close to each other in any norm.

It's arbitrary to show that maximum difference between two vectors is their infinity norm distance; i.e.:

$$\|U(V_1) - U(V_2)\| \leq \|V_1 - V_2\|_\infty \Rightarrow \|U(V_1) - U(V_2)\| \leq \|V_1 - V_2\|$$

A short proof can be approved by visualization:



Thus The objection can be reduced to  $\infty$ -norm case.

$$\|U(V_1) - U(V_2)\|_\infty = \|R + \gamma P V_1 - R - \gamma P V_2\|_\infty = \gamma \|P(V_1 - V_2)\|_\infty \leq \|P(V_1 - V_2)\|_\infty \leq \|P\|_\infty \|V_1 - V_2\|_\infty \stackrel{*}{=} \|V_1 - V_2\|_\infty$$

where in  $\star$  I used the Cauchy-Schwartz theorem and in  $\star$  I used the fact that  $P$  is a transition matrix; thus its sum of rows is 1 and hence its  $\infty$ -norm is equal to one.

$$\begin{aligned} \lim_{n \rightarrow \infty} U^n(V) = V^* &\Leftrightarrow \lim_{n \rightarrow \infty} U^n(V) - V^* = 0 \Leftrightarrow \lim_{n \rightarrow \infty} U(U^{n-1}(V)) - V^* = 0 \Leftrightarrow \lim_{n \rightarrow \infty} U(U^{n-1}(V) - V^*) = 0 \Leftrightarrow \lim_{n \rightarrow \infty} \|U(U^{n-1}(V) - V^*)\|_\infty = 0 \\ &\Leftrightarrow \lim_{n \rightarrow \infty} \|U(U^{n-1}(V) - V^*)\|_\infty \leq \gamma \lim_{n \rightarrow \infty} \|U^{n-1}(V) - V^*\|_\infty \leq \dots \leq \lim_{n \rightarrow \infty} \gamma^n \|V - V^*\|_\infty \Leftrightarrow \lim_{n \rightarrow \infty} \gamma^n = 0 \end{aligned}$$

we know that  $\lim_{n \rightarrow \infty} \gamma^n$  converges to zero for any bound  $\gamma$  which is strictly less than one. Also in  $\star$  we used the fact that  $U(V^*) = V^*$  which is obvious; i.e.  $U(V^*) = V^*$  which means mapping has converged.

$$\|U^k(V) - U^{k-1}(V)\|_\infty < \epsilon \stackrel{?}{\Rightarrow} \|V^* - U^k(V)\|_\infty < \frac{\epsilon}{1-\gamma}$$

$$\begin{aligned} \|V^* - U^{k+1}(V)\|_\infty &\leq \|V^* - U^k(V)\|_\infty + \|U^k(V) - U^{k+1}(V)\|_\infty \leq \|V^* - U^k(V)\|_\infty + \gamma \|U^k(V) - U^{k-1}(V)\|_\infty \\ &\leq \|V^* - U^k(V)\|_\infty + \gamma \epsilon \leq \|V^* - U^{k+2}(V)\|_\infty + \|U^{k+2}(V) - U^{k+1}(V)\|_\infty + \gamma \epsilon \\ &\leq \|V^* - U^{k+2}(V)\|_\infty + \|U(U^k(V)) - U(U^{k-1}(V))\|_\infty + \gamma \epsilon \leq \\ &\leq \|V^* - U^{k+2}(V)\|_\infty + \gamma^2 \epsilon + \gamma \epsilon \\ &\vdots \\ &\leq \lim_{n \rightarrow \infty} \|V^* - U^n(V)\|_\infty + \sum_{i=1}^{\infty} \gamma^i \epsilon = \frac{\epsilon}{1-\gamma} \end{aligned}$$



$$\begin{aligned}
 G_{14} &= +10 & G_{12} &= 10\gamma & G_{13} &= 10\gamma^2 & G_{11} &= 10\gamma^3 & G_{10} &= 10\gamma^4 \cdot 10 & G_9 &= 10\gamma^5 \cdot 10\gamma & G_8 &= 10\gamma^6 \cdot 10\gamma^2 \cdot 10 \\
 G_7 &= 10\gamma^7 \cdot 10\gamma^3 \cdot 10\gamma & G_6 &= 10\gamma^8 \cdot 10\gamma^4 \cdot 10\gamma^2 & G_5 &= 10\gamma^9 \cdot 10\gamma^5 \cdot 10\gamma^3 & G_4 &= 10\gamma^{10} \cdot 10\gamma^6 \cdot 10\gamma^4 & G_3 &= 10\gamma^{11} \cdot 10\gamma^7 \cdot 10\gamma^5 \\
 G_2 &= 10\gamma^{12} \cdot 10\gamma^8 \cdot 10\gamma^6 & G_1 &= 10\gamma^{13} \cdot 10\gamma^9 \cdot 10\gamma^7
 \end{aligned}$$

الف

Every Visit Monte Carlo

$$\begin{aligned}
 V(1) &= G_1 & V(2) &= G_2 & V(3) &= G_3 & V(7) &= G_5 & V(8) &= G_4 & V(12) &= G_6 & V(16) &= G_7 & V(17) &= G_{13} \\
 V(20) &= G_8 & V(18) &= G_{14} & V(21) &= \frac{G_9 + G_{11}}{2} & V(22) &= \frac{G_{10} + G_{12}}{2}
 \end{aligned}$$

ب

First Visit Monte Carlo

$$\begin{aligned}
 V(1) &= G_1 & V(2) &= G_2 & V(3) &= G_3 & V(7) &= G_5 & V(8) &= G_4 & V(12) &= G_6 & V(16) &= G_7 & V(17) &= G_{13} \\
 V(20) &= G_8 & V(18) &= G_{14} & V(21) &= G_{11} & V(22) &= G_{12}
 \end{aligned}$$

## سوال ③

الف

$$\begin{aligned}
 V^{\pi}(s) &= E\left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0=s\right] = E\left[r_0 + \sum_{t=1}^{\infty} \gamma^t r_t \mid s_0=s\right] = E\left[r_0 + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_0=s\right] = E[r_0 \mid s_0=s] + \gamma E\left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_0=s\right] \\
 &= E\left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0=s\right] = E_{s' \sim p(s'|s, a)}\left[E\left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0=s' \mid s_0=s\right]\right] \xrightarrow{MPP} = E_{s' \sim p(s'|s, a)}\left[E\left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0=s'\right]\right] = E_{s' \sim p(s'|s, a)}[V^{\pi}(s')] \\
 &\Rightarrow V^{\pi}(s) = E\left[E_{s' \sim p(s'|s, a)}[r_0 + \gamma V^{\pi}(s')]\right] \xrightarrow{\text{independent of } s} \text{independent of } s. \quad * p(s_2|s_1, s_0) \xrightarrow{MPP} = p(s_2|s_1)
 \end{aligned}$$

ب

$$\begin{aligned}
 \text{define } \pi_1^* &= \arg \max_{\pi} E_{s \sim p(s|s_0, a)}[r_0 + \gamma V^{\pi}(s)] \xrightarrow{\times \gamma} \pi_2^* = \arg \max_{\pi} E_{s \sim p(s|s_0, a)}[\gamma r_0 + \gamma \gamma V^{\pi}(s)] \\
 &= \gamma E_{s \sim p(s|s_0, a)}[r_0 + \gamma V^{\pi}(s)]
 \end{aligned}$$

$\pi_1^*$  and  $\pi_2^*$  are the same because their objective function has been multiplied by a constant factor and hence can be factorized and ignored in an optimization seeking the optimal argument.

ج

$V^{\pi}(s) = E\left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0=s\right] \Rightarrow V_2^{\pi}(s) = E\left[\sum_{t=0}^{\infty} \gamma^t (r_t + c) \mid s_0=s\right] = E\left[\gamma^t r_t \mid s_0=s\right] + c E\left[\sum_{t=0}^{\infty} \gamma^t \mid s_0=s\right] = V_1^{\pi}(s) + c E\left[\frac{1-\gamma^{T+1}}{1-\gamma} \mid s_0=s\right]$   
 It's obvious that  $E\left[\frac{1-\gamma^{T+1}}{1-\gamma} \mid s_0=s\right]$  is not a constant term, it obviously depends not only on  $\gamma$ , but also on where we start which is not considered case in Bellman equation. Thus  $V_2^{\pi}(s) \neq V_1^{\pi}(s)$  i.e. we are dealing with a new value function in second definition and there no constraints that optimal policy for  $V_2^{\pi}(s)$  will converge  $\pi^*(s)$  as it's dealing with different objective function.

$$V_2^{\pi}(s) - V_1^{\pi}(s) = c E\left[\lim_{T \rightarrow \infty} \frac{1-\gamma^{T+1}}{1-\gamma} \mid s_0=s\right] = c E\left[\frac{1}{1-\gamma} \mid s_0=s\right] = \frac{c}{1-\gamma} \xrightarrow{\text{independent of } \gamma \text{ and } s} \text{like part b; } \pi_1^*(s) \text{ and } \pi_2^*(s) \text{ will be the same policies. proof}$$



also can be obtained by defining  $\pi_d^*(s) = \pi_2^*(s) - \pi_1^*(s) = \arg \max_a E_{s \sim p(s', s, a)} [r - r + \gamma V_1^{\pi}(s') + \frac{c}{1-\gamma} - \gamma V_1^{\pi}(s)] = \dots \Rightarrow \pi_d^*(s) = 0$

(5)

For first MDP; we have  $\pi_1^*(s) = \arg \max_a E_s [r + \gamma V_1^{\pi}(s)]$ . This implies that for any state  $s$ ;  $\pi_1^*(s)$  chooses the optimal action  $a^*$ . To prove that  $\pi_2^*(s) = \pi_1^*(s)$ , we can show that for any states;  $E_s [r + \gamma V_2^{\pi}(s) | a = a^*] > E_s [r + \gamma V_2^{\pi}(s) | a \neq a^*]$ . Thus we have:

$$E_s [r + \gamma V_2^{\pi}(s) | a = a^*] = E_s [r - c + \gamma V_2^{\pi}(s)] = E_s [r + \gamma V_1^{\pi}(s)] - c \xrightarrow{*} \pi_2^*(s) = \pi_1^*(s)$$

Also at  $*$  we assumed that value functions are the same if we choose the optimal action; which is shown above.

سؤال 4

الف

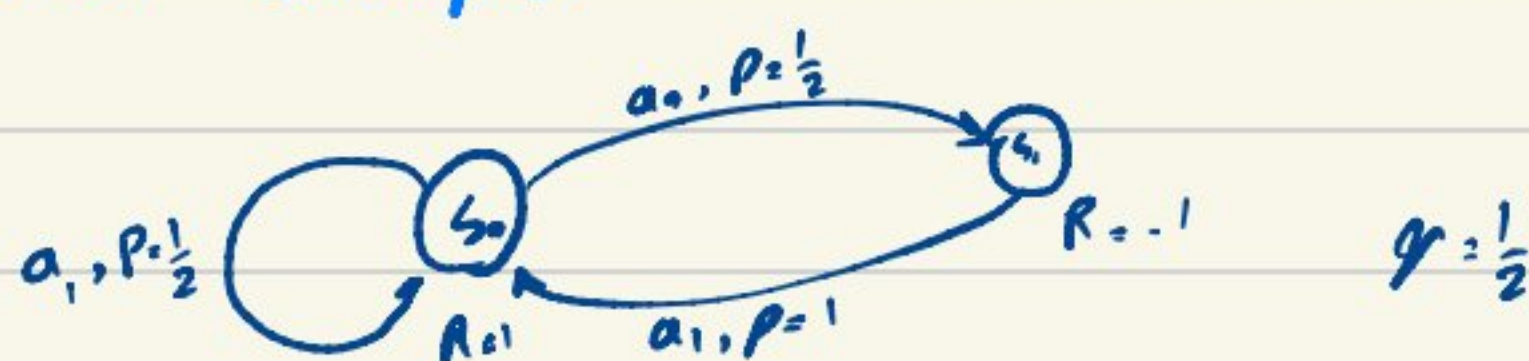
$$V^{\pi}(s') = \sum_s \sum_a p(s, a, s') \pi(s, a) \left[ \frac{V^{\pi}(s) - R(s, a)}{\gamma} \right]$$

$$\text{Bellman equation: } V^{\pi}(s) = \sum_{s'} \sum_a p(s', s, a) \pi(a|s) (R(s, a) + \gamma V^{\pi}(s'))$$

$p(s', a, s)$  describes path probability and  $\pi(a|s)$  explores joint probability of states and actions.

Consider the case where  $s'$  and  $s$  both have action  $a$ . then  $p(s, s', a)$  doesn't have specific meaning because it can't obtain 2 different values.

Counter example:



$$\text{Bellman equation: } \begin{cases} V^{\pi}(s_0) = \frac{1}{2} (1 + \frac{1}{2} V^{\pi}(s_0)) + \frac{1}{2} (-1 + \frac{1}{2} V^{\pi}(s_1)) \Rightarrow \frac{3}{4} V^{\pi}(s_0) + \frac{1}{4} V^{\pi}(s_1) \Rightarrow V^{\pi}(s_0) = \frac{1}{3} V^{\pi}(s_1) \\ V^{\pi}(s_1) = \frac{1}{2} V^{\pi}(s_0) \Rightarrow \frac{1}{2} V^{\pi}(s_0) = \frac{1}{3} V^{\pi}(s_1) \Rightarrow V^{\pi}(s_0) = \frac{2}{3} V^{\pi}(s_1) \end{cases}$$

$$\text{Suggested equation: } \begin{cases} V^{\pi}(s_0) = \frac{2}{3} (V^{\pi}(s_1) + 1) + \frac{1}{3} (V^{\pi}(s_0) - 1) \Rightarrow V^{\pi}(s_0) = 0 \\ V^{\pi}(s_1) = 2 (V^{\pi}(s_0) - 1) \Rightarrow V^{\pi}(s_1) = -2 \end{cases}$$

Contradiction

ب

$$\begin{aligned} V^{\pi}(s) &= E[G_t | s_t = s, \pi] \stackrel{\text{As shown in part 3.a}}{=} E^{\pi} [r_t + E_{s' \sim p(s', s, \pi(s))} [V^{\pi}(s')] | s_t = s] \\ &= E^{\pi} [r_t | s_t = s] + E^{\pi} [E_{s' \sim p(s', s, \pi(s))} [V^{\pi}(s')] | s_t = s] \\ &\stackrel{\text{from Markov property}}{=} E^{\pi} [R(s, \pi(s))] + E^{\pi} [E_{s' \sim p(s', s, \pi(s))} [V^{\pi}(s')]] \\ &= E^{\pi} [r_0 | s_0 = s] + E^{\pi} [E_{s' \sim p(s', s, \pi(s))} [V^{\pi}(s') | s_0 = s]] \\ &= E^{\pi} [r_0 + E_{s' \sim p(s', s, \pi(s))} [V^{\pi}(s') | s_0 = s] | s_0 = s] = E[G_0 | s_0 = s, \pi] = V^{\pi}(s) \end{aligned}$$



(ج)

$$\begin{aligned}
 V(s_t) &= E^{\pi} [r_t + E_{s_{t+1} \sim p(s_{t+1}|s_t, \pi(s_t))} [V(s_{t+1})]] \\
 &= E^{\pi} [r_t + V(s_{t+1})] \\
 &= E^{\pi} [r_t] + V(s_{t+1}) \\
 &\leftarrow V(s_{t+1})
 \end{aligned}$$

transition is deterministic  
actions don't effect transition  
rewards are definitely negative.

سوال 5

الف

n-step TD learning;  $V^{k+1}(s_t) \leftarrow V^k(s_t) + \alpha (\sum_{k=1}^n \gamma^{k-1} R_k + \gamma^n V^k(s_{t+n}) - V^k(s_t))$

$n=1 \rightarrow V(E)$  gets updated

$n=2 \rightarrow V(D), V(E)$  gets updated

$n=3 \rightarrow$  Value function for all steps gets updated

(ب)

As we increase  $\alpha$ ; at first TD learns more and converges better for limited iteration. After wards; by increasing  $\alpha$ ; temporal error makes bigger impact in algorithm and hence for big values of  $\alpha$ ; value functions get unstable and don't converge to optima.

For large  $n$ ; we have more precise temporal difference and better converges. If we increase  $n$  however makes algorithm needy for more iterations and thus converges slowly.

So proper  $n, \alpha$  for TD algorithm is nothing huge and not neglectable amount.

(ج)

A) Increasing states causes more leafs for state-action tree; i.e. we have bigger environment to explore in contrast to before. Thus more probably the error rate for same episode and  $n$  will increase.

B) From law of large numbers we know that more episodes mean more samples and thus less error probably. In fact; for good amount of  $n$ ; we converge to  $V^*(s)$  if we take  $\infty$  samples.

C) For a proper amount of  $\alpha$  (recall part ب); more repetition means Convergence! i.e. if  $\alpha$  is small; running algorithm for  $V(s)$  makes maximum usage of the explored data. However; if  $\alpha$  is large; more repeat has no good case it may vary.

$$E_t(s) = \gamma \wedge E_{t-1}(s) + 1\{s_t = s_{t-1}\}$$

(د)

$E_t(s)$  reaches it's maximum increase " $E_{t+1}(s)$ " when the state  $s$  is selected for accumulating trans.

$$\text{thus: } E_t(s) = \frac{1}{4} + \frac{1}{5} E_{t+1}(s) + 1 = \frac{1}{5} E_{t+1}(s) + 1 \rightarrow E_{t+1}(s) = \frac{5}{4}$$