

(i) Chearm

$$R_1=1, R_2=2, R_3=3, R_4=4, R_5=5, \text{ <Repeats>}$$

$$a) \quad Q_n = \frac{\sum_{i=1}^n R_i}{n}$$

$$b) \quad Q_1 = 1 ; \quad Q_{(n+1)} = Q_n + \frac{1}{n} [R_n - Q_n]$$

$$Q_2 = 1 + \frac{1}{2} [2-1] = 1.5$$

$$Q_3 = 1.5 + \frac{1}{3} [3-1.5] = 2$$

$$Q_4 = 2 + \frac{1}{4} [4-2] = 2.5$$

$$Q_5 = 2.5 + \frac{1}{5} [5-2.5] = 3$$

$$Q_6 = 3 + \frac{1}{6} [1-3] = 2.67$$

$$Q_7 = 2.67 + \frac{1}{7} [2-2.67] = 2.57$$

$$Q_8 = 2.57 + \frac{1}{8} [3-2.57] = 2.62$$

$$Q_9 = 2.62 + \frac{1}{9} [4-2.62] = 2.77$$

$$Q_{10} = 2.77 + \frac{1}{10} [5-2.77] = 2.993$$

$$Q_{11} = 2.993 + \frac{1}{11} [1 - 2.993] = 2.81$$

$$Q_{12} = 2.81 + \frac{1}{12} [2 - 2.81] = 2.74$$

$$Q_{13} = 2.74 + \frac{1}{13} [3 - 2.74] = 2.76$$

$$Q_{14} = 2.76 + \frac{1}{14} [4 - 2.76] = 2.84$$

$$Q_{15} = 2.84 + \frac{1}{15} [5 - 2.84] = 2.984$$

1c) Based on above 15 observations, we already find Q value to hover over $\approx 3 (\pm 0.5)$.

Because the sequence is stationary and (stepsize) \propto value is $\frac{1}{n}$, we can see the sequence

Converging to its expected value of $3 (\pm 0.5)$.

1d) From above observation, we see sequence converges to ≈ 3 .

This is because the average of (1, 2, 3, 4, 5) toward sequence is 3 and this sequence is repeated.

1e) Yes, we have used incremental Q_{n+1} in above calculations.

$$Q_{n+1} = \frac{\sum_{k=1}^{n+1} R_k}{(n+1)} = \frac{\sum_{k=1}^n R_k + R_{(n+1)}}{(n+1)}$$

$$Q_{n+1} = \frac{\sum_{k=1}^n R_k}{n} \left[\frac{n}{n+1} \right] + \frac{R_{(n+1)}}{(n+1)}$$

$$Q_{n+1} = Q_n \left[1 - \frac{1}{n+1} \right] + \frac{R_{(n+1)}}{n+1}$$

$$Q_{n+1} = Q_n + \frac{1}{n+1} [R_{(n+1)} - Q_n]$$

1f) $Q_1 = 1$; $Q_2 = \frac{1+2}{2} = 1.5$; $Q_3 = \frac{1+2+3}{3} = 2$;
 $Q_4 = \frac{1+2+3+4}{4} = 2.5$; $Q_5 = \frac{1+2+3+4+5}{5} = 3$

11y Using Incremental Update, refer values in (1b).

We see both values are equal.

1g) The new Update equation,

$$Q_{n+1} = Q_n + \frac{1}{(n+1)^{0.8}} [R_{n+1} - Q_n]$$

(ie) Step size $\alpha = \frac{1}{(n+1)^{0.8}}$

Trying observations,

$$Q_1 = 1 ; Q_2 = 1 + \frac{1}{2^{0.8}} [2 - 1] = 1 + \frac{1}{1.74} = 1.57$$

$$Q_3 = 1.57 + \frac{1}{3^{0.8}} [3 - 1.57] = 2.16$$

$$Q_4 = 2.16 + \frac{1}{4^{0.8}} [4 - 2.16] = 2.77$$

$$Q_5 = 2.77 + \frac{1}{5^{0.8}} [5 - 2.77] = 3.39$$

Thus, we observe as α (Stepsize) value is increased, the convergence rate is slow. But, still we can assume the series will converge to value ~ 3.39 .

2) Given $\gamma \in (0,1)$, π , $V_\pi(s)$, $q_\pi(s,a)$, $p(s',r|s,a)$.

a)

$$q_\pi(s,a) = \sum_{s',r} p(s',r|s,a) [r + \gamma V_\pi(s')]$$

b)
$$V_\pi(s) = \sum_a \pi(a|s) q_\pi(s,a)$$

c)
$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V_\pi(s')]$$

d) Writing Bellman equation in Vector-matrix Notation

$$V_\pi = \gamma \pi + \gamma P^\pi V_\pi$$

$$\Rightarrow (I - \gamma P^\pi) V_\pi = \gamma \pi$$

$$V_\pi = (I - \gamma P^\pi)^{-1} \gamma \pi$$

$(I - \gamma P^\pi)$ is invertible, as Eigen values are $(0,1)$

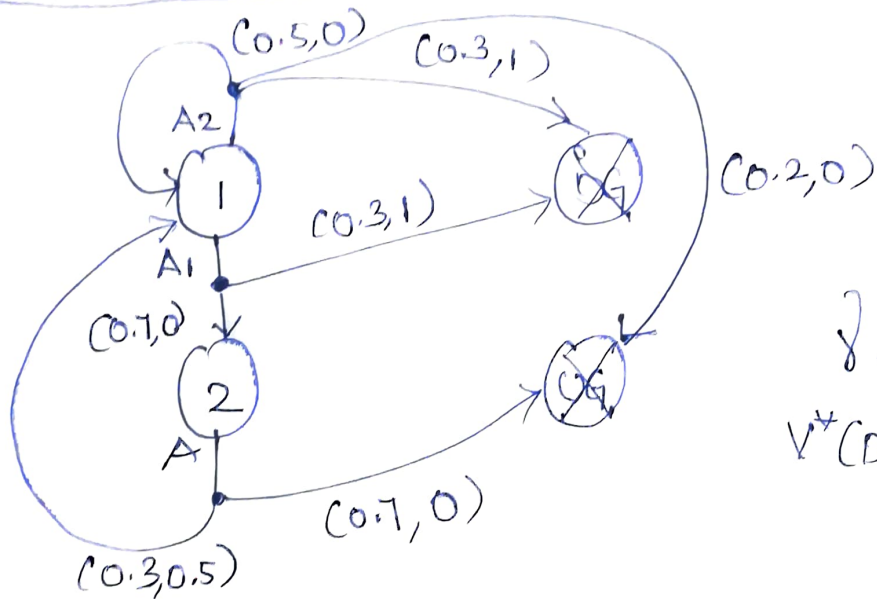
\Rightarrow Solution is Unique.

Thus Bellman equation has an Unique solution.

2e) Q-Bellman equation .

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) \left[r + \gamma \sum_b \pi(b | s') q_{\pi}(s', b) \right]$$

3 .



$$\gamma = 1$$

$$V^*(DG) = V^*(UG) = 0$$

$$8) \quad V(\text{Room 1}) = \frac{1}{2} [\text{Value from Action 1}] + \frac{1}{2} [\text{Value from Action 2}]$$

[Assuming Equiprobability of actions]

$$V(\text{Room 1}) = \frac{1}{2} [(1+0) 0.3 + 0.7 (0 + V(\text{Room 2}))] + \frac{1}{2} [0.5 (0 + V(RI)) + 0.3 (1+0) + 0.2 (0 + 0)]$$

$$2 \quad V(RI) = 0.3 + 0.7 V(R2) + 0.5 V(RI) + 0.3 \cdot$$

$$2V(R_1) = 0.6 + 0.7V(R_2) + 0.5V(R_1)$$

$$\boxed{1.5V(R_1) - 0.7V(R_2) = 0.6} \quad (1)$$

$V(R_2) =$ ~~$\frac{1}{2}$~~ $\frac{1}{2}$

$$V(R_2) = (0.0 + \frac{1}{2}V(OG)) 0.7 + (0.5 + \frac{1}{2}V(R_1)) 0.3$$

$$V(R_2) = 0.15 + 0.3V(R_1)$$

$$\boxed{V(R_2) - 0.3V(R_1) = 0.15} \quad (2)$$

Solving (1) & (2)

$$V(R_1) = 0.546$$

$$V(R_2) = 0.314$$

b) Value Iteration Scheme

$$V(R_1) = \max \left(\begin{aligned} &(0.3(1+0) + 0.7(0 + V(R_2))) / 2, \\ &(0.5(0 + V(R_1)) + 0.3(1+0) + \\ &0.2(0+0)) / 2 \end{aligned} \right)$$

$$V(R_1) = \max \left(\begin{aligned} &(0.3 + 0.7V(R_2)) / 2, \\ &(0.5V(R_1) + 0.3) / 2 \end{aligned} \right)$$

~~$$V(R_2) = \max(0.6 + 0.7V(R_1))$$~~

$$V(R_2) = \max(0.7(0+0) + 0.3(0.5 + V(R_1)))$$

$$V(R_2) = 0.15 + 0.3V(R_1)$$

c)

$$V(R_1) = \max \left(\begin{aligned} &(0.15 + 0.35V(R_2)), \\ &(0.25V(R_1) + 0.15) \end{aligned} \right)$$

$$V(R_2) = 0.15 + 0.3V(R_1) \quad \# \text{No max - as only one action}$$

Step 0 , $V_0(R) = 0$; $V_0(R_2) = 0$

Step 1

$$V_1(R_1) = \max((0.15 + 0), (0 + 0.15))$$

$$V_1(R_1) = 0.15$$

$$V_1(R_2) = 0.15$$

Step 2

$$V_2(R_1) = \max((0.15 + (0.35 \times 0.15)), ((0.25 \times 0.15) + 0.15))$$

$$= \max(0.2025, 0.1875)$$

$$= 0.2025$$

$$V_2(R_2) = 0.15 + (0.3 \times 0.15)$$

$$= 0.195$$

Step 3

$$V_3(R_1) = \max((0.15 + (0.35 \times 0.195)), ((0.25 \times 0.2025) + 0.15))$$

$$= \max(0.218, 0.201) = 0.218$$

$$V_3(R_2) = 0.15 + 0.3(0.2025)$$

$$V_3(R_2) = 0.21075$$

Step 4

$$V_4(R_1) = \max \left((0.15 + (0.35 \times 0.201)), \right. \\ \left. ((0.25 \times 0.218) + 0.15) \right)$$

$$= \max(0.224, 0.205) = 0.224$$

$$V_4(R_2) = 0.15 + (0.3 \times 0.218) = 0.2154$$

$$\neq 0.154$$

Step 5

$$V_5(R_1) = \max \left((0.15 + (0.35 \times 0.2154)), \right. \\ \left. ((0.25 \times 0.224) + 0.15) \right)$$

$$= \max(0.2254, 0.206)$$

$$= 0.2254$$

$$V_5(R_2) = 0.15 + (0.3 \times 0.224) = 0.2172$$

3d)

$$V_{\pi}(R1) = 0.3(1+0) + 0.7(0 + \gamma V_{\pi}(R2))$$

$$\boxed{V_{\pi}(R1) = 0.3 + 0.7V_{\pi}(R2)} \rightarrow (1)$$

$$V_{\pi}(R2) = 0.7(0+0) + 0.3(0.5 + \gamma V_{\pi}(R1))$$

$$\boxed{V_{\pi}(R2) = 0.15 + 0.3V_{\pi}(R1)} \rightarrow (2)$$

Solving (1) & (2),

$$V_{\pi}(R1) = 0.513 \quad ; \quad V_{\pi}(R2) = 0.304$$

3e)

$$R1 \xrightarrow{0} R2 \xrightarrow{0.5} R1 \xrightarrow{0} R2 \xrightarrow{0.5} R1 \xrightarrow{0} R2 \xrightarrow{0.5} R1 \xrightarrow{0} R2$$

$$\xrightarrow{0.5} R1 \xrightarrow{1} DG$$

$$\gamma = 1$$

i) First visit method.

$$V_{\pi}(R1) = G_1 = \frac{1 + 0.5 + 0.5 + 0.5 + 0.5}{1}$$

$$V_{\pi}(R1) = 3$$

$$V_{\pi}(R2) = G_1 = \frac{1 + 0.5 + 0.5 + 0.5 + 0.5}{1}$$

$$= 3$$

ii) Every Visit Method

$$V_{\pi}(S) = \frac{\sum_{i=1}^k G_i}{k}$$

k : Number of occurrences

G_i : Return at i th Occurrence

~~W(R1)~~

For R1

$$k=5$$

$$G_1 = ((0.5) \times 4) + 1 = 3$$

$$G_2 = ((0.5) \times 3) + 1 = 2.5$$

$$G_3 = ((0.3) \times 2) + 1 = 2$$

$$G_4 = 1.5$$

$$G_5 = 1$$

$$\therefore V_{\pi}(R1) = \frac{1+1.5+2+2.5+3}{5} = 2$$

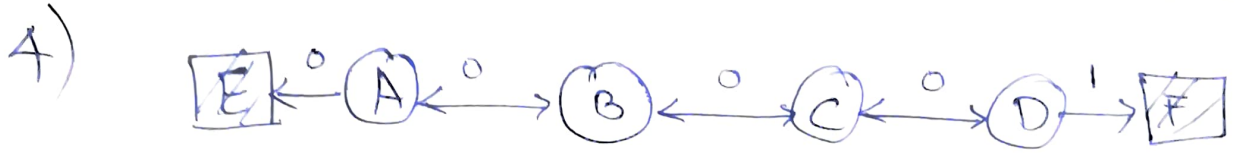
For R2

$$k=4$$

$$G_1 = ((0.5) \times 4) + 1 = 3$$

$$G_2 = 2.5; G_3 = 2; G_4 = 1.5$$

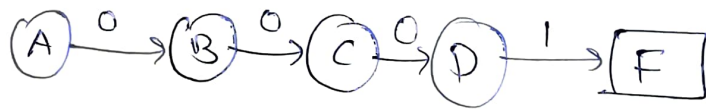
$$V_{\pi}(R2) = \frac{1.5 + 2 + 2.5 + 3}{4} = \frac{9}{4} = 2.25$$



Action Probability $\gamma = 1$

$\leftarrow \frac{1}{2} \quad \frac{1}{2} \rightarrow$
L $\frac{1}{2}$ $\frac{1}{2}$ R

Episode



$$V(A) = V(B) = V(C) = V(D) = 0.6 \quad \alpha = 0.1$$

2) n-step TD with $n=1$

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

3-iterations of Regular TD

e.g. $V(S_1) \leftarrow V(S_1) + \alpha [R_2 + \gamma V(S_2) - V(S_1)]$

Iteration-1

$$V(A) = V(A) + \alpha [R_{A \rightarrow B} + \gamma V(B) - V(A)]$$

$$V(A) = 0.6 + 0.1 [0 + 0.6 - 0.6]$$

$$\boxed{V(A) = 0.6}$$

$$V(B) = 0.6 + 0.1 [0 + 0.6 - 0.6]$$

$$\boxed{V(B) = 0.6}$$

$$V(C) = 0.6 + 0.1 [0 + 0.6 - 0.6]$$

$$\boxed{V(C) = 0.6}$$

$$V(D) = 0.6 + 0.1 [1 + 0 - 0.6]$$

$$V(D) = 0.6 + 0.1 [0.4]$$

$$\boxed{V(D) = 0.64}$$

Iteration-2

$$V(A) = 0.6 + 0.1 [0 + 0.6 - 0.6] = 0.6$$

$$V(B) = 0.6 + 0.1 [0 + 0.6 - 0.6] = 0.6$$

$$V(C) = 0.6 + 0.1 [0 + 0.64 - 0.6] = 0.604$$

$$V(D) = 0.64 + 0.1 [1 + 0 - 0.64] = 0.676$$

Iteration-3

$$V(A) = 0.6 + 0.1[0 + 0.6 - 0.6] = 0.6$$

$$V(B) = 0.6 + 0.1[0 + 0.604 - 0.6] = 0.6004$$

$$V(C) = 0.604 + 0.1[0 + 0.676 - 0.604] = 0.6112$$

$$V(D) = 0.676 + 0.1[1 + 0 - 0.676] = 0.7084$$

4b) 2-step TD

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+3}) - V(S_t)]$$

Iteration-1

$$V(A) \leftarrow V(A) + \alpha [R_{A \rightarrow B} + \gamma R_{B \rightarrow C} + \gamma^2 V(C) - V(A)]$$

$$V(A) = 0.6 + 0.1[0 + 0 + 0.6 - 0.6] = 0.6$$

$$V(B) = 0.6 + 0.1[0 + 0 + 0.6 - 0.6] = 0.6$$

$$V(C) = 0.6 + 0.1[0 + 1 + 0 - 0.6] = 0.64$$

$$V(D) = 0.6 + 0.1[1 + 0 + 0 - 0.6] = 0.64$$

Iteration-2

$$V(A) = 0.6 + 0.1 \left[\overset{A \rightarrow B}{0} + \overset{B \rightarrow C}{0} + \overset{C \rightarrow D}{0.64} - \overset{V(A)}{0.6} \right] = 0.604$$

$$V(B) = 0.6 + 0.1 \left[\overset{B \rightarrow C}{0} + \overset{C \rightarrow D}{0} + \overset{D \rightarrow F}{0.64} - \overset{V(B)}{0.6} \right] = 0.604$$

$$V(C) = 0.64 + 0.1 \left[\overset{C \rightarrow D}{0} + \overset{D \rightarrow F}{1} + \overset{F \rightarrow H}{0} - \overset{V(C)}{0.64} \right] = 0.676$$

$$V(D) = 0.64 + 0.1 \left[\overset{D \rightarrow F}{1} + \overset{F \rightarrow H}{0} + \overset{H \rightarrow I}{0} - \overset{V(D)}{0.64} \right] = 0.676$$

Iteration-3

$$V(A) = 0.604 + 0.1 \left[\overset{A \rightarrow B}{0} + \overset{B \rightarrow C}{0} + \overset{C \rightarrow D}{0.676} - \overset{V(A)}{0.604} \right] = 0.6112$$

$$V(B) = 0.604 + 0.1 \left[\overset{B \rightarrow C}{0} + \overset{C \rightarrow D}{0} + \overset{D \rightarrow F}{0.676} - \overset{V(B)}{0.604} \right] = 0.6112$$

$$V(C) = 0.676 + 0.1 \left[\overset{C \rightarrow D}{0} + \overset{D \rightarrow F}{1} + \overset{F \rightarrow H}{0} - \overset{V(C)}{0.676} \right] = 0.7084$$

$$V(D) = 0.676 + 0.1 \left[\overset{D \rightarrow F}{1} + \overset{F \rightarrow H}{0} + \overset{H \rightarrow I}{0} - \overset{V(D)}{0.676} \right] = 0.7084$$

4c) 3-step TD

$$V(S_t) \leftarrow V(S_t) + \alpha \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 V(S_{t+4}) - V(S_t) \right]$$

Iteration-1

$$V(A) \leftarrow V(A) + \alpha \left[R_{A \rightarrow B} + \gamma R_{B \rightarrow C} + \gamma^2 R_{C \rightarrow D} + \gamma^3 V(D) - V(A) \right]$$

9

$$V(A) = 0.6 + 0.1 [0 + 0 + 0 + 0.6 - 0.6] = 0.6$$

$$V(B) = 0.6 + 0.1 \overset{A \rightarrow B \quad B \rightarrow C \quad C \rightarrow D \quad V(B)}{[0 + 0 + 1 + 0 - 0.6]} = 0.64$$

$$V(C) = 0.6 + 0.1 \overset{A \rightarrow B \quad B \rightarrow C \quad C \rightarrow D \quad V(C)}{[0 + 1 + 0 + 0 - 0.6]} = 0.64$$

$$V(D) = 0.6 + 0.1 [1 + 0 + 0 + 0 - 0.6] = 0.64$$

Iteration-2

$$V(A) = 0.6 + 0.1 \overset{A \rightarrow B \quad B \rightarrow C \quad C \rightarrow D \quad V(C)}{[0 + 0 + 0 + 0.64 - 0.6]} = 0.604$$

$$V(B) = 0.64 + 0.1 [0 + 0 + 1 + 0 - 0.64] = 0.676$$

$$V(C) = 0.64 + 0.1 [0 + 1 + 0 + 0 - 0.64] = 0.676$$

$$V(D) = 0.64 + 0.1 [1 + 0 + 0 + 0 - 0.64] = 0.676$$

Iteration-3

$$V(A) = 0.604 + 0.1 [0 + 0 + 0 + 0.676 - 0.604] = 0.6112$$

$$V(B) = 0.676 + 0.1 [0 + 0 + 1 + 0 - 0.676] = 0.7084$$

$$V(C) = 0.676 + 0.1 [0 + 1 + 0 + 0 - 0.676] = 0.7084$$

$$V(D) = 0.676 + 0.1 [1 + 0 + 0 + 0 - 0.676] = 0.7084$$

4d) From above 3 steps, these are the observations

1-step TD

	A	B	C	D
Iteration-1	0.6	0.6	0.6	0.64
Iteration-2	0.6	0.6	0.604	0.676
Iteration-3	0.6	0.6004	0.6112	0.7084

2-step TD

Iteration	A	B	C	D
1	0.6	0.6	0.64	0.64
2	0.604	0.604	0.676	0.676
3	0.6112	0.6112	0.7084	0.7084

3-step TD

Iteration	A	B	C	D
1	0.6	0.64	0.64	0.64
2	0.604	0.676	0.676	0.676
3	0.6112	0.7084	0.7084	0.7084

We see 3-step TD is converging / learning at a faster rate [As 3 steps are updated in single iteration]

Sharmila

5a) Semi Gradient TD(0) Update

$$\omega_{t+1} = \omega_t + \alpha [R_{t+1} + \gamma \hat{V}(S_{t+1}, \omega_t) - \hat{V}(S_t, \omega_t)] * \nabla \hat{V}(S_t, \omega_t)$$

Here, instead of,

$V(S_{t+1})$, we use function approximation

$\hat{V}(S_{t+1}, \omega_t)$. And hence, it is called

"Semi-gradient".

5b) Using Linear function approximation, we will have

$$\hat{V}(S_t, \omega_t) = \omega_t^T X(S_t)$$

$$\hat{V}(S_{t+1}, \omega_t) = \omega_t^T X(S_{t+1})$$

$$\nabla \hat{V}(S_t, \omega_t) = X(S_t)$$

Hence, Update rule of TD(0) under a linear approximation architecture:

$$\begin{aligned}
 W_{t+1} &= W_t + \alpha [R_{t+1} + \gamma \omega_t^T X(s_{t+1}) - \omega_t^T X(s_t)] X(s_t) \\
 &= W_t + \alpha [R_{t+1} + \gamma X(s_{t+1})^T W_t - X(s_t)^T W_t] X(s_t) \\
 &= W_t + \alpha [R_{t+1} X(s_t) + \gamma X(s_{t+1})^T W_t X(s_t) - X(s_t)^T W_t X(s_t)]
 \end{aligned}$$

$\therefore X(s_{t+1})^T W_t$ is a scalar,

we can write

$$X(s_{t+1})^T W_t X(s_t) = X(s_t) X(s_{t+1})^T W_t$$

$$\therefore W_{t+1} = W_t + \alpha [R_{t+1} X(s_t) - X(s_t) [\hat{X}(s_t)^T - \gamma X(s_{t+1})^T] W_t]$$

5c) For 1-step SARSA

$$\omega_{t+1} = \omega_t + \alpha [R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \omega_t) - \hat{q}(S_t, A_t, \omega_t)] * \nabla \hat{q}(S_t, A_t, \omega_t)$$

(Or)

Using Linear function Approximation,

$$W_{t+1} = W_t + \alpha [R_{t+1} + \gamma X(S_{t+1}, A_{t+1}) - X(S_t, A_t)]$$

$$[X(S_t, A_t)^T - \gamma X(S_{t+1}, A_{t+1})^T] W_t$$

5d) Q-learning Algorithm with linear function approximation

Input: a differential function

$$\hat{q} : \mathcal{S}^+, \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R} \text{ such that}$$

$$\hat{q}(\text{terminal}, \cdot, \cdot) = 0$$

Algorithm parameters: step size $\alpha > 0$,
small $\epsilon > 0$

Initialize value function weights $w \in \mathbb{R}^d$ arbitrarily
(eg., $w = 0$)

Loop for each episode:

Initialize s

Loop for each step of episode:

choose A from s using policy derived from \hat{q}
(eg. ϵ -greedy)

Take action A , observe R, s'

$$w \leftarrow w + \alpha (R + \gamma \max_a \hat{q}(s', a, w)$$

$$- \hat{q}(s, A, w)) \nabla \hat{q}(s, A, w)$$

$$s \leftarrow s'$$

Until s is terminal

5e) Input: 2 differential functions

$$\hat{q}_1: s^+, A \times \mathbb{R}^d \rightarrow \mathbb{R} \text{ such that}$$

$$\hat{q}_1(\text{terminal}, \cdot, \cdot) = 0$$

$$\hat{q}_2: s^+, A \times \mathbb{R}^d \rightarrow \mathbb{R} \text{ such that}$$

$$\hat{q}_2(\text{terminal}, \cdot, \cdot) = 0$$

Algorithm parameter : step size $\alpha \in (0, 1]$,
Small $\epsilon > 0$

Initialize Value-function weights $w_1 \in \mathbb{R}^d$ arbitrarily
(eg, $w_1 = 0$)
~~and~~ $w_2 \in \mathbb{R}^d$

Loop for each episode:

Initialize S

Loop for each step of episode:

Choose A from S using policy ϵ -greedy in $\hat{q}_1 + \hat{q}_2$

Take action A , observe R, s'

With 0.5 probability

$$w_1 \leftarrow w_1 + \alpha (R + \frac{1}{2} \hat{q}_2(s', \arg \max_a \hat{q}_1(s', a, w_1))$$

$$- \hat{q}_1(s, A, w_1)) \nabla \hat{q}_1(s, A, w_1)$$

else:

$$\omega_2 \leftarrow \omega_2 + \alpha (R + \gamma \hat{q}_1(s', \arg \max_a \hat{q}_2(s', a, \omega_2)) - \hat{q}_2(s, A, \omega_2)) \nabla \hat{q}_2(s, A, \omega_2)$$

$$S \leftarrow s'$$

Until S is terminal.
