## Using value iteration

$$V_{opt}^{(t)}(s) \leftarrow \max_{a \in Actions(s)} \sum_{s'} T(s,a,s')[Reward(s,a,s') + \gamma V_{opt}^{(t-1)}(s')]$$

## prob 1a)

① $V_{opt}^{(1)}(S_A) = \max \begin{cases} a = + \rightarrow 0 + 0.00|x0 \\ a = - \rightarrow 5 + 0.00|x0 \end{cases} = 5 \quad (-)$

$V_{opt}^{(1)}(S_B) = \max \begin{cases} a = + \rightarrow 0 + 0.00|x0 \\ a = - \rightarrow 0 + 0.00|x0 \end{cases} = 0$

$V_{opt}^{(1)}(S_c) = \max \begin{cases} a = + \rightarrow 16 + 0.00|v0 \\ a = - \rightarrow 0 + 0.00|x0 \end{cases} = 16 \quad (+)$

$V_{opt}^{(1)}(S_d) = \max \begin{cases} 0 \\ 0 \end{cases} = 0$

② $V_{opt}^{(2)}(S_A) = \max(0, 5 + 0.005) = 5.005 \quad (-)$

$V_{opt}^{(2)}(S_B) = \max(0 + 0.016, 0 + 0.005) = 0.016 \quad (+)$

$V_{opt}^{(2)}(S_c) = \max(16, 0) = 16 \quad (+)$

$V_{opt}^{(2)}(S_D) = \max(0, 0) = 0$

③ $V_{opt}^{(3)}(S_A) = \max(0 + 0.016 \times 0.0|, 5 + 0.01 \times 5.005) \approx 5.005 \quad (-) \text{ almost same}$

$V_{opt}^{(3)}(S_B) = \max(0.00| \times 16, 0.01 \times 5.005) = 0.016 \quad (+) \text{ same}$

$V_{opt}^{(3)}(S_c) = \max(16, 0.01 \times 5.005) \approx 16 \quad (+) \text{ same}$

$V_{opt}^{(3)}(S_D) = 0 \quad same$

$\Rightarrow$ ∴ optimal policy for $S_A : (-)$

prob 1 b )

① $V_{opt}^{(1)}(S_A) = \max \begin{cases} a=+ \to 0+0.999 \times 0 \\ a=- \to 5+0.999 \times 0 \end{cases} = 5$  (-)

$V_{opt}^{(1)}(S_B) = \max \begin{cases} a=+ \to 0+0.999 \times 0 \\ a=- \to 0+0.999 \times 0 \end{cases} = 0$

$V_{opt}^{(1)}(S_c) = \max \begin{cases} a=+ \to 16+0.999 \times 0 \\ a=- \to 0+0.999 \times 0 \end{cases} = 16$  (+)

$V_{opt}^{(1)}(S_d) = \max \begin{cases} 0 \\ 0 \end{cases} = 0$

② $V_{opt}^{(2)}(S_A) = \max(0, 5+0.999 \times 5) = 9.995$  (-)

$V_{opt}^{(2)}(S_B) = \max(0+0.999 \times 16, 0+0.999 \times 5) = 15.984$  (+)

$V_{opt}^{(2)}(S_c) = \max(16, 0) = 16$  (+)

$V_{opt}^{(2)}(S_D) = \max(0, 0) = 0$

③ $V_{opt}^{(3)}(S_A) = \max(0+0.999 \times 15.984, 0.999 \times 9.995) = 15.968016$ (+)

$V_{opt}^{(3)}(S_B) = \max(0+0.999 \times 16 +0.0999 \times 9.995) = 15.984$ (+) same

$V_{opt}^{(3)}(S_c) = \max(16, 0+0.999 \times 15.984) = 16$ (+) same

$V_{opt}^{(3)}(S_D) = 0$  same

④ $V_{opt}^{(4)}(S_A) = \max(0+0.999 \times 15.984, 0.999 \times 15.968016) = 15.968016$ (+) same.

∴ optimal policy for $S_A$ : (+)

## prob 1 c )

$$V_{opt}^{(t)}(S_A) = \max\left(0 + \gamma V_{opt}^{(t-1)}(S_B), \; 5 + \gamma V_{opt}^{(t-1)}(S_A)\right)$$

$$V_{opt}^{(t)}(S_B) = \max\left(0 + \gamma V_{opt}^{(t-1)}(S_C), \; 0 + \gamma V_{opt}^{(t-1)}(S_A)\right)$$

$$V_{opt}^{(t)}(S_C) = \max\left(16 + \gamma V_{opt}^{(t-1)}(S_D) + 0 + \gamma V_{opt}^{(t-1)}(S_B)\right) = \max\left(16, \; \gamma V_{opt}^{(t-1)}(S_B)\right)$$

$$V_{opt}^{(t)}(S_D) = 0$$

$\Rightarrow$ ① $V_{opt}^{(1)}(S_A) = 5$ 　② $V_{opt}^{(2)}(S_A) = (1+\gamma)5$

$\quad\quad V_{opt}^{(2)}(S_B) = 0$ 　　　$V_{opt}^{(2)}(S_B) = 16\gamma$

$\quad\quad V_{opt}^{(3)}(S_C) = 16$ 　　　$V_{opt}^{(3)}(S_C) = 16$

$\quad\quad V_{opt}^{(4)}(S_D) = 0$ 　　　$V_{opt}^{(4)}(S_D) = 0$

③ $V_{opt}^{(3)}(S_A) = \max\left(16\gamma^2, \; 5(1 + \gamma(1+\gamma))\right)$

$\quad V_{opt}^{(3)}(S_B) = \max\left(16\gamma, \; \gamma(1+\gamma)5\right)$ 　　　$\text{---}$

$\quad V_{opt}^{(3)}(S_C) = \max\left(16, \; 16\gamma^2\right)$

$\quad V_{opt}^{(4)}(S_D) = 0$

because $16\gamma$ is always bigger than $\gamma V_{opt}^{(t-1)}(S_A)$,

<u>optimal policy of $s_B$ is (+)</u>