## Question 3

(a) For each state we have 4 possible actions, so the number of policies is $4^6 = 4096$.

(b) State transition matrix $T$:

$$T = \begin{pmatrix} 0.2 & 0.8 & 0 & 0 & 0 & 0 \\ 0 & 0.2 & 0.8 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.2 & 0.8 & 0 \\ 0 & 0 & 0 & 0 & 0.2 & 0.8 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}$$

$$R = (0, 0, 20, 10, 0, -10)^T$$

$$\Rightarrow \quad V^{\pi_0} = (I - \gamma T)^{-1} R = \begin{pmatrix} 154.19 \\ 175.61 \\ 200 \\ -64.9 \\ -87.8 \\ -100 \end{pmatrix} \begin{matrix} S_1 \\ S_2 \\ S_3 \\ S_4 \\ S_5 \\ S_6 \end{matrix}$$

(C) For state $S_1$:

$R(S_1, up) = V^0(S_1) = 154.19$

$R(S_1, down) = 0.2V^0(S_1) + 0.8V^0(S_4) = -21.08$

$R(S_1, left) = V^0(S_1) = 154.19$

$R(S_1, right) = 0.2V^0(S_1) + 0.8V^0(S_2) = 171.326$

$\Rightarrow \pi'(S_1) = \underline{right}$

For state $S_2$:

$R(S_2, up) = V^0(S_2) = 175.61$

$R(S_2, down) = 0.2V^0(S_2) + 0.8V^0(S_5) = -35.12$

$R(S_2, left) = 0.2V^0(S_2) + 0.8V^0(S_1) = 158.474$

$R(S_2, right) = 0.2V^0(S_2) + 0.8V^0(S_3) = 195.122$

$\Rightarrow \pi'(S_1) = \underline{right}$

For state $S_3$:

$\pi'(S_3) = \max\{V^0(S_3), 0.2V(S_3) + 0.8V^0(S_2)\}$

$= \max\{V^0(S_3), v_tS_2)\} = V^0(S_3) = R(S_3, down)/\pi^0(S_3, right)$

$\Rightarrow \pi'(S_3) = \underline{down}.$

For state $S_4$:

$\pi(\max\{V^0(S_4), V^0(S_5)\} = V(S_4) = R(S_4, up/down/left)$

$\Rightarrow \pi'(S_4) = \underline{down}$

For S5:

$$\pi' \qquad \max \{ V^0(S_4), V^0(S_3), V^0(S_1), V^0(S_6) \} = V_0(S_2) = R(S_5, up).$$

$$\rightarrow \underline{\pi'(S_5) = up}$$

For S6:

$$\max_{S} \{ V^0(S_1), V^0(S_6), V^0(S_3) \} = V^0(S_3) = R(S_6, up)$$

$$\rightarrow \underline{\pi'(S_3) = up}$$

(d)  We can continue the iterative process in (C), over new transition matrix:
policy evaluation:

$$T' = \begin{pmatrix} 0.2 & 0.8 & 0 & 0 & 0 & 0 \\ 0 & 0.2 & 0.8 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0.8 & 0 & 0 & 0.2 & 0 \\ 0 & 0 & 0.8 & 0 & 0 & 0.2 \end{pmatrix}$$

$$\Rightarrow V^{\pi'} = (I - \gamma T')^{-1} R = \begin{pmatrix} 154.19 \\ 175.61 \\ 200 \\ 100 \\ 154.19 \\ 163.41 \end{pmatrix} \begin{matrix} S_1 \\ S_2 \\ S_3 \\ S_4 \\ S_5 \\ S_6 \end{matrix}$$

policy improvement:

① For $S_1$: $\max\{V'(S_1), V'(S_2), V'(S_4)\} = V'(S_2) = R(S_1, \text{right})$

$\longrightarrow \pi''(S_1) = \boxed{\text{right}}$

② For $S_2$: $\max\{V'(S_1), V'(S_2), V'(S_3), V'(S_5)\} = V'(S_3) = R(S_2, \text{right})$

$\longrightarrow \pi''(S_2) = \boxed{\text{right}}$.

③ For $S_3$: $\max\{V'(S_3), V'(S_2)\} = V'(S_3) = R(S_3, \text{up/down/right})$

$\longrightarrow \pi''(S_3) = \boxed{\text{down}}$

④ For $S_4$: $\max\{V'(S_4), V'(S_5)\} = V'(S_5) = R(S_4, \text{right})$

$\longrightarrow \pi''(S_4) = \boxed{\text{right}}$

⑤ For $S_5$: $\max\{V'(S_5), V'(S_4), V'(S_6), V'(S_2)\} = V'(S_2) = R(S_5, \text{up})$

$\longrightarrow \pi''(S_5) = \boxed{\text{up}}$

⑥ For $S_6$: $\max\{V'(S_3), V'(S_5), V'(S_6)\} = V'(S_3) = R(S_6, \text{up})$

$\longrightarrow \pi''(S_6) = \boxed{\text{up}}$.

$\Rightarrow \pi'' = (\text{right, right, down, right, up, up}) \neq \pi'$,

so we keep doing iteration.

policy evaluation: new transition matrix:

$$T'' = \begin{pmatrix} 0.2 & 0.8 & 0 & 0 & 0 & 0 \\ 0 & 0.2 & 0.8 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0.2 & 0.8 & 0 \\ 0 & 0.8 & 0 & 0 & 0.2 & 0 \\ 0 & 0 & 0.8 & 0 & 0 & 0.2 \end{pmatrix}$$

$$\Rightarrow V^{\pi''} = (I - \gamma T'')^{-1} R = \begin{pmatrix} 154.19 \\ 175.61 \\ 200 \\ 147.58 \\ 154.19 \\ 162.41 \end{pmatrix}$$

policy improvement:

① For $S_1$: $\max\{V^3(S_1), V^3(S_2), V^3(S_4)\} = V^3(S_2) = R(S_1, right)$

$\rightarrow \pi^3(S_1) = \boxed{right}$

② For $S_2$: $\max\{V^3(S_1), V^3(S_2), V^3(S_3), V^3(S_5)\} = V^3(S_3) = R(S_2, right)$

$\rightarrow \pi^3(S_2) = \boxed{right}$

③ For $S_3$: $\max\{V^3(S_3), V^3(S_2)\} = V^3(S_3) = R(S_3, up/down/right)$

$\rightarrow \pi^3(S_3) = \boxed{down}.$

④ For S4: $\max\{V^3(S_4), V^3(S_5)\} = V^3(S_5) = R(S_4, right)$

$\longrightarrow \pi^3(S_4) = \boxed{right}$

⑤ For S5: $\max\{V^3(S_5), V^3(S_4), V^3(S_6), V^3(S_2)\} = V^3(S_2) = R(S_5, up)$

$\longrightarrow \pi^3(S_5) = \boxed{up}$

⑥ For S6: $\max\{V^3(S_3), V^3(S_5), V^3(S_6)\} = V^3(S_3) = R(S_6, up)$

$\longrightarrow \pi^3(S_6) = \boxed{up}$

$\Longrightarrow \pi^3 = ( right, right, down, right, up, up) = \pi''$

$\Longrightarrow$ Now policy iteration algorithm converges

$\Longrightarrow V^* = \begin{pmatrix} 154.19 \\ 175.61 \\ 200 \\ 147.58 \\ 154.19 \\ 163.41 \end{pmatrix}$

(e) The optimal value function is unique since $V^*$ is defined as the best

value that can be achieved by any state:

$$V^*(s) = \max_\pi V^\pi(s).$$

(f) By (d) we know the optimal policy is:

$$\pi^* = (\text{right, right, down, right, up, up})$$
$$\quad\quad S_1 \quad S_2 \quad S_3 \quad S_4 \quad S_5 \quad S_6$$

(g) No. Since there could be other policies $\tilde{\pi}$ that also returns the same optimal value function. For instance, $\tilde{\pi} = (\text{right, right, up, right, up, up})$
$$\quad\quad\quad\quad\quad\quad\quad\quad S_1 \quad S_2 \quad S_3 \quad S_4 \quad S_5 \quad S_6$$

( obtained by changing $\pi^*(S_3)$ into "up", which yields the same reward)
is also an optimal policy.

(h). If we scale all rewards for each state by a constant factor k, then all deductions would remain unchanged, which means we could yield the same optimal policy, with the corresponding value function scaled by k.

Proof of this proposition:

Let $R' = kR$ (where k is a real constant), then by the new value function under policy $\pi$ is:

$$V'_\pi = (I - \gamma T^\pi)^{-1} \cdot kR = k \cdot (I - \gamma T^\pi)^{-1} R = k \cdot V_\pi$$

When we do policy improvement for a particular state S, we compute:

$$\max_{S_i} \left\{ P V_\pi(S_i) + (1-P) V_\pi(S) \right\} = \max_{S_i} \left\{ V_\pi(S_i) \right\}$$

If we scale $V_\pi$ by k, $\max_{S_i} \left\{ V_\pi(S_i) \right\}$ wouldn't change, implying the new transition matrix $T^{\pi'}$ is also unchanged. Thus, by induction, the entire policy iteration process would yield the same optimal policy, with the corresponding $V^*$ scaled by a factor k.