(1) (a)

States set = $\{S, 1, 3, 5, 6, 7, 8, W\}$

Transition matrix $\Rightarrow$

$T = $

|   | S | 1 | 3 | 5 | 6 | 7 | 8 | W |
|---|---|---|---|---|---|---|---|---|
| S | 0 | 0.25 | 0.25 | 0 | 0 | 0.25 | 0.25 | 0 |
| 1 | 0 | 0 | 0.25 | 0.25 | 0 | 0.25 | 0.25 | 0 |
| 3 | 0 | 0 | 0 | 0.25 | 0.25 | 0.25 | 0.25 | 0 |
| 5 | 0 | 0.25 | 0 | 0 | 0.25 | 0.25 | 0.25 | 0 |
| 6 | 0 | 0 | 0.25 | 0 | 0 | 0.25 | 0.25 | 0.25 |
| 7 | 0 | 0 | 0.25 | 0 | 0 | 0.25 | 0.25 | 0.25 |
| 8 | 0 | 0 | 0.25 | 0 | 0 | 0 | 0.5 | 0.25 |
| W | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

(b)

$$R(s) = \begin{cases} -1 & s \in \{S, 1, 3, 5, 6, 7, 8\} \\ 0 & s = W \end{cases}$$

Discount factor $r = 1$

we want to calculate the no. of throws required from states except $W$ to reach $W$.

Therefore taking states $k \in \{S, 1, 3, 5, 6, 7, 8\}$

$$R = [-1, -1, -1, -1, -1, -1]^T$$

$$\Rightarrow T = \begin{bmatrix} P & B \\ 0 & 1 \end{bmatrix}$$

$$V = (I - rP)^{-1} R = \begin{bmatrix} -7.118 \\ -7.050 \\ -6.711 \\ -6.77 \\ -5.35 \\ -5.35 \\ -5.35 \end{bmatrix}$$

③ ⓐ

$$P^{\pi_1} = \begin{array}{c} \\ A \\ B \\ C \\ D \end{array} \begin{array}{cccc} A & B & C & D \\ \left[\begin{array}{cccc} 0 & 0.9 & 0.1 & 0 \\ 0.1 & 0 & 0 & 0.9 \\ 0.9 & 0 & 0 & 0.1 \\ 0 & 0 & 0 & 1 \end{array}\right] \end{array}$$

$$P^{\pi_2} = \begin{array}{c} A \\ B \\ C \\ D \end{array} \left[\begin{array}{cccc} 0 & 0.1 & 0.9 & 0 \\ 0.9 & 0 & 0 & 0.1 \\ 0.1 & 0 & 0 & 0.9 \\ 0 & 0 & 0 & 1 \end{array}\right]$$

$$P^{\pi_3} = \begin{array}{c} A \\ B \\ C \\ D \end{array} \left[\begin{array}{cccc} 0 & 0.42 & 0.58 & 0 \\ 0.1 & 0 & 0 & 0.9 \\ 0.1 & 0 & 0 & 0.9 \\ 0 & 0 & 0 & 1 \end{array}\right]$$

$$R = \begin{array}{c} \\ A \\ B \\ C \\ D \end{array} \begin{array}{cc} a_1 & a_2 \\ \left[\begin{array}{cc} -10 & -10 \\ -10 & -10 \\ -10 & -10 \\ 100 & 100 \end{array}\right] \end{array} \qquad \gamma = 1$$

$$V^{\pi_1} = \left(I - \gamma P^{\pi_1}\right)^{-1} R = \begin{bmatrix} -24.39 & -24.39 \\ -12.43 & -12.43 \\ -31.95 & -31.95 \end{bmatrix}$$

$$V^{\pi_2} = \left(I - \gamma P^{\pi_2}\right)^{-1} R = \begin{bmatrix} -24.39 & -24.39 \\ & \\ -31.95 & -31.95 \\ -12.43 & -12.43 \end{bmatrix}$$

$$V^{\pi_3} = \begin{bmatrix} -22 \cdot 22 & -22 \cdot 22 \\ -12 \cdot 22 & \\ & -12 \cdot 22 \\ -12 \cdot 22 & -12 \cdot 22 \end{bmatrix}$$

(b) $\pi_3$ is the best policy because all the individual elements / ~~reward~~ values in value matrix with respect to $V^{\pi_1}$ and $V^{\pi_2}$ of $V^{\pi_3}$ is highest.

(c) No, because if many policy $\pi_1$ and $\pi_2$ ~~its~~ ~~any two values~~

$$V^{\pi_1} = \begin{bmatrix} 0.9 \\ 0.1 \end{bmatrix} \qquad V^{\pi_2} = \begin{bmatrix} 0.8 \\ 0.7 \end{bmatrix}$$

then $\qquad V^{\pi_1}[1] > V^{\pi_2}[1]$

$$V^{\pi_2}[2] < V^{\pi_2}[2]$$

then thy two policy are not comparable

(a) We will take values fruitions value corresponding to each policy $\pi_1$ and $\pi_2$ then we will pick maximum value actions corresponding to each state from policy $\pi_1$ and $\pi_2$ and create a new policy $\pi$ with this actions (maximum value) for each state (can have mixture of actions aswell)
( This will be similar for the question ~~this~~ MDP $\pi_3$ as we have different actions for different states)

(4) (a) for policy $\pi_1$, the agent is preferring close exit and risk the cliff. Which means the agent is short sighted that means $r$ will be low (0.1) and not giving importance to future (distant). We will be putting noise $\eta$ to be zero in the environment so that there is no danger of tripping of the cliff.

for $\pi_3$, the agent is preferring close exit and not risk the cliff. Which means the agent is short sighted taking $r$ to be low (0.1) and the noise $\eta$ to be more (0.5) in the environment so that there is danger of tripping in the cliff.

for $\pi_2$, the agent is preferring distant exit and risk the cliff which means the agent is far sighted taking $r$ to be high (0.9) and the noise $\eta$ to be low (0) in the environment so that the danger of the tripping is none.

for $\pi_4$, the agent is preferring distant exit and not risk the cliff. which means the agent is far sighted taking $r$ to be

high loss and the noise to be more loss in the environment so that the danger of trapping is exists.

$$Q_1^\pi(s,a) = \mathbb{E}_\pi \left[ \dot{e}_{t+1} + \gamma \dot{e}_{t+2} \cdots \mid s_t = s, a_t = a \right]$$

$$Q_2^\pi(s,a) = \mathbb{E}_\pi \left[ \dot{e}_{t+1}^2 + \gamma \dot{e}_{t+2}^2 \cdots \mid s_t = s, a_t = a \right]$$

$$\Rightarrow Q_1^\pi(s,a) + Q_2^\pi(s,a)$$

$$\Rightarrow \mathbb{E}_\pi \left[ \dot{e}_{t+1}' + \gamma \dot{e}_{t+2}' \cdots \mid s_t = s, a_t = a \right]$$

$$+ \mathbb{E}_\pi \left[ \dot{e}_{t+1}^2 + \gamma \dot{e}_{t+2}^2 \cdots \mid s_t = s, a_t = a \right]$$

$$\Rightarrow \mathbb{E}_\pi \left[ (\dot{e}_{t+1}' + \dot{e}_{t+2}') + \gamma(\dot{e}_{t+1}' + \dot{e}_{t+2}^2) \cdots \mid s_t = s, a_t = a \right]$$

$$\Rightarrow Q_3^\pi(s,a)$$

$$Q_3^\pi(s,a) = Q_2^\pi(s,a) + Q_1^\pi(s,a)$$

from ⓐ

$$V_3^\pi(s,a) = V_1^\pi(s,a) + V_2^\pi(s,a)$$

when $\pi^*$ is the optimal policy for
$Q_1(s,a)$ and $a$   MDP $M_1$ and $M_2$
then we can say that $M_3$ will
have $\pi^*$ as the optimal policy

(a)

$$V_1^{\pi}(s) = E_{\pi}\left( \sum_{k=0}^{\infty} \gamma^k r'_{t+k+1} \right)$$

$$= E_{\pi}\left[ \sum_{k=0}^{\infty} \gamma^k \left( r^2_{t+k+1} + \varepsilon \right) \right]$$

$$= E_{\pi}\left[ \sum_{k=0}^{\infty} \gamma^k r^2_{t+k+1} \right] + E_{\pi}\left[ \sum_{k=0}^{\infty} \gamma^k \varepsilon \right]$$

$$V_1^{\pi}(s) = V_2^{\pi}(s) + \frac{\varepsilon}{1-\gamma}$$

(b) It is not possible to do so by combining $\pi_i^*$ and $\pi_2^*$ because the rewards for the optimal policies may be totally different and the optimal policy to be obtained ~~with~~ for $\pi_2^*$ will have totally different reward sums.

(2)(a)  State space = $\{0, 1, \ldots \ldots N\}$

$\qquad\qquad$ = No. of machines working

Action = [ Repair , No repair ]

$\qquad$ repairing $\qquad\qquad\qquad$ Repairing
$\qquad$ will be $\qquad\qquad\qquad\qquad$ will not be
$\qquad\qquad$ done $\qquad\qquad\qquad\qquad\qquad$ done

Rewards = $\begin{cases} K - \frac{N}{2} & \text{( Repairing was done} \\ & \text{with K working} \\ K & \text{machines)} \end{cases}$

$\qquad\qquad\qquad\qquad\qquad$ ( working
$\qquad\qquad\qquad\qquad\qquad\qquad$ machines count)

Transition
Probability

## Transition probability matrix with no repair action

$$
\Rightarrow \quad
\begin{array}{c|ccccccc}
 & 0 & 1 & 2 & \cdots & N-1 & N \\
\hline
0 & 1 & 0 & 0 & \cdots & 0 & 0 \\
1 & 1/2 & 1/2 & 0 & \cdots & 0 & 0 \\
2 & 1/3 & 1/3 & 1/3 & \cdots & 0 & 0 \\
\vdots & & & & & & \\
N-1 & 1/N & 1/N & 1/N & \cdots & 1/N & 0 \\
N & 1/N+1 & 1/N+1 & 1/N+1 & \cdots & 1/N+1 & 1/N+1
\end{array}
$$

## Transition matrix with repair action

$$
\Rightarrow \quad
\begin{array}{c|cccccc}
 & & & & & & \\
0 & 0 & - & - & - & - & \\
1 & 0 & - & - & & - & \\
2 & 0 & - & - & & & \\
\vdots & & & & & & \\
N-1 & - & - & - & & - & \\
N & - & - & - & & - &
\end{array}
$$

(b) We are going to use discounted setting with $r < 1$ as we don't have any absorbing states (can be thought as indefini horizon problem)

$$
\begin{array}{c} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{array}
\begin{bmatrix}
\phantom{x}0 & \phantom{x}1 & \phantom{x}2 & \phantom{x}3 & \phantom{x}4 & \phantom{x}5 \\
1 & 0 & 0 & 0 & 0 & 0 \\
\tfrac{1}{2} & \tfrac{1}{2} & 0 & 0 & 0 & 0 \\
\tfrac{1}{3} & \tfrac{1}{3} & \tfrac{1}{3} & 0 & 0 & 0 \\
\tfrac{1}{4} & \tfrac{1}{4} & \tfrac{1}{4} & \tfrac{1}{4} & 0 & 0 \\
\tfrac{1}{5} & \tfrac{1}{5} & \tfrac{1}{5} & \tfrac{1}{5} & \tfrac{1}{5} & 0 \\
\tfrac{1}{6} & \tfrac{1}{6} & \tfrac{1}{6} & \tfrac{1}{6} & \tfrac{1}{6} & \tfrac{1}{6}
\end{bmatrix}
$$

$$V = (I - \gamma P)^{-1} R$$

$$R = \begin{bmatrix} 0 & 1 & 2 & 3 & 4 & 5 \end{bmatrix}^T$$

$$\pi = \text{no repair policy is choosen}$$

$$V^{\pi} = \begin{bmatrix} 0 & 1 \cdot 904 & 3 \cdot 783 & 5 \cdot 706 & 7 \cdot 61 & 9 \cdot 302 \end{bmatrix}^T$$