

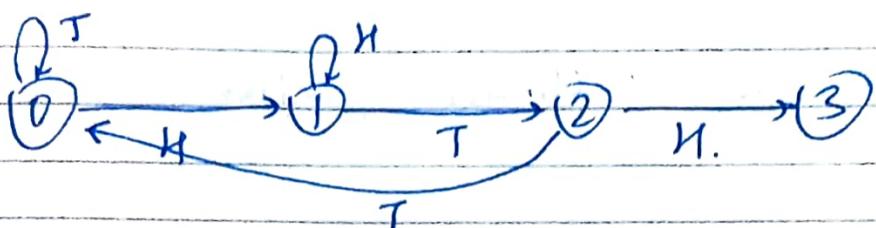
Question 1

State 0: → No part of pattern has been matched.

State 1: - The first 'H' has been matched

State 2: - The first 'H' followed by 'T' has been matched.

State 3: - The pattern 'HTH' has been matched



Transition
matrix

$$P = \begin{bmatrix} 0 & 1/2 & 1/2 & 0 & 0 \\ 1 & 0 & 1/2 & 1/2 & 0 \\ 2 & 1/2 & 0 & 0 & 1/2 \\ 3 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Absorbing state is State 3

States we can write the states of MRP as
< S, P, R, V > as follows

$$S = \{\emptyset, 'H', 'HT', 'HTH'\}$$

$$P = \begin{bmatrix} \emptyset & H & HT & HTH \\ \emptyset & 1/2 & 1/2 & 0 & 0 \\ H & 0 & 1/2 & 1/2 & 0 \\ HT & 1/2 & 0 & 0 & 1/2 \\ HTH & 0 & 0 & 0 & 1 \end{bmatrix}$$

Ques I Set a reward for all states = -1 except absorbing state 'NTH'.

i.e

$$R(s) = -1 \text{ for } s \in \{0, 1N, 'NT'\} \text{ for every extra tors.}$$

$$R(s) = 0 \text{ for } s = 'NTH'.$$

The discount factor I take is $\gamma = 1$.

Now as per the Bellman equation

$$V(s) = E(\text{rew}) | s_t = s + \gamma \sum_{s' \in S} P_{ss'} V(s').$$

applying it to all states.

$$V(0) = -1 + \frac{1}{2} V(1) + \frac{1}{2} V(0)$$

$$V(1) = -1 + \frac{1}{2} V(1) + \frac{1}{2} V(0)$$

$$V(2) = -1 + \frac{1}{2} V(0) + \frac{1}{2} V(3)$$

$$V(3) = 0. \quad \text{as the 'NTH' is the absorbing state.}$$

Solving the above system of equations.

$$V(2) = -1 + \frac{1}{2} V(0) \rightarrow (1)$$

$$V(1) = -1 + \frac{1}{2} V(1) + \frac{1}{2} V(2)$$

$$\Rightarrow V(1) = -2 + V(2) \rightarrow (2)$$

$$v(0) = -2 + v(1)$$

$$= -2 + (-2 + v(2))$$

$$\Rightarrow v(0) = -4 + v(2)$$

using ③ equation ①

$$x(2) = v(2) = -4 + (-1 + \gamma_2 v(0))$$

$$\therefore v(0) = -5 + \gamma_2 v(0)$$

$$\therefore v(0) = -10$$

Now as we had taken Reward $R(S) = -1$ for every tors with -ve value being acting as a penalization factor.

Also $v(0)$ signifies the expected number of tors required to reach goal state i.e. value from state 0.

Hence No of Tors required to reach goal State = 10. {as -1 is awarded for every extra tors}.

Therefore dividing -10 by -1 i.e $R(S)$.

Question 2.

2(a)

(i) State Space :- defines the number of machines that are currently working.

$$S: \{0, 1, 2, 3, \dots - N\}$$

where N = total machines present in the production facility.

Any state s has m machines working where $0 \leq m \leq N$. as per the problem statement.

(ii) Action Space :- Actions available at each state m are

→ Repair :- call the Repairman to repair all machines at the cost of $-\frac{N}{2}$ \$

→ No Repair :- Do nothing & continue with current number of working machines.

Therefore $A = \{\text{Repair}, \text{No Repair}\}$.

(iii) Transition Probabilities. $P(n|m, a)$

n = no of working machines of next day

m = no of working machines of current day

a = under the action a .
i.e Repair / No repair

$$P(A = \text{No Repair}) =$$

$$P(n|m, a = \text{"No repair"}) =$$

	$n=0$	$n=1$	$n=2$	\dots	$n=N-1$	$n=N$
$m=0$	$\frac{1}{2}$	$\frac{1}{2}$	0	0	\dots	
$m=1$		$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$		
$m=2$			$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	0
$m=3$				$\frac{1}{4}$	0	0
$m=N$		$\frac{1}{N+1}$	$\frac{1}{N+1}$	$\frac{1}{N+1}$	\dots	$\frac{1}{N+1}$

where m = currently correctly working machines.
 n = no of machines that would start up correctly next day.

After repairing state always transitions to $m=N$
 (all machines working)

$$P(n|m, \text{"Repair"}) =$$

$$P(n=N|m, \text{Repair}) = 1$$

i.e

	$n=0$	$n=1$	$n=2$	$n=3$	\dots	$n=N$
$m=0$	0	0	0	0	0	1
$m=1$	0	0	0	0	\dots	$\frac{1}{2} - \frac{1}{2}$
$m=2$	0	0	0	0	\dots	$\frac{1}{3} - \frac{1}{3}$
$m=3$						
$m=N$		$\frac{1}{N+1}$	$\frac{1}{N+1}$	$\frac{1}{N+1}$	\dots	$\frac{1}{N+1}$

(iv) Rewards.

$R(m, a = "No Repair") = m$ { since each working machine generates a revenue of 1 \$. hence Reward for state when m machines are working is m .

$R(m, a = "Repair") = -\frac{N}{2}$ { cost - lumpsum to repair all incorrectly working machines.

Note:- it is independent of m .

Q(b) I would use undiscouted setting, since the problem statement focusses on maximizing total long term profit, without any indication that future rewards should be less prioritized.

Hence Total Rewards

$$\sum_{t=0}^{\infty} r^t R_t$$

$$= \sum_{t=0}^{\infty} R_t \text{ for } r=1 \text{ i.e undiscouted setting.}$$

(c) Given MDP $\langle S, A, P, R, \gamma \rangle$.

$$N=5$$

S : State space S i.e. $m \in \{0, 1, 2, 3, 4, 5\}$

	0	1	2	3	4	5
$P_{m=0}$	0	0	0	0	0	0
$m=1$	$\frac{1}{2}$	$\frac{1}{2}$	0	0	0	0
2	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	0	0	0
3	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	0	0
4	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	0
5	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

$A = \text{"No Repair"} = \text{Action.}$

$\gamma = 1$ \downarrow as undiscounted setting assumed.

$R(m, a = \text{"No Repair"}) = m.$

$$\text{i.e } R(0) = 0$$

$$R(1) = 1$$

$$R(2) = 2$$

$$R(3) = 3$$

$$R(4) = 4$$

$$R(5) = 5.$$

using the Bellman ~~equation~~ evaluation equation

$$v^n = r^n + \gamma p^n V^n.$$

$$V(0) = 0 \rightarrow ① \quad \{ \text{as no machines are working} \}$$

$$V(1) = 1 + \frac{1}{2} V(0) + \frac{1}{2} V(1) \rightarrow ②$$

$$V(2) = 2 + \frac{1}{3} \{ V(0) + V(1) + V(2) \} \rightarrow ③$$

$$V(3) = 3 + \frac{1}{4} \{ V(0) + V(1) + V(2) + V(3) \} \rightarrow ④$$

$$V(4) = 4 + \frac{1}{5} \{ V(0) + V(1) + V(2) + V(3) + V(4) \}$$

$$V(5) = 5 + \frac{1}{6} \{ V(0) + V(1) + V(2) + V(3) + V(4) + V(5) \}$$

Solving ① & ② we get

$$V(1) = 1 + \frac{1}{2} V(0) + \gamma_2 V(1)$$

$$\Rightarrow V(1) = 2 \times 1 = 2 \rightarrow ⑦$$

Solving ③ using values of $V(0)$ & $V(1)$

$$\Rightarrow V(2) = 2 + \frac{1}{3} \{ 0 + 2 + V(2) \}$$

$$\Rightarrow \frac{2}{3} V(2) = \frac{6+2}{3}$$

$$\Rightarrow V(2) = 8/2 = 4.$$

Solving (4) using values of $v(0), v(1), v(2)$

$$v(3) = 3 + \frac{1}{4} \{ 0 + 2 + 4 + v(3) \}$$

$$\Rightarrow \frac{3}{4} v(3) = \frac{12 + 2 + 4}{4}$$

$$\Rightarrow v(3) = 18/3 = 6.$$

Solving (4) using values of $v(0), v(1), v(2), v(3)$

$$v(4) = 4 + \frac{1}{3} \{ 0 + 2 + 4 + 6 + v(4) \}$$

$$\Rightarrow \frac{4}{3} v(4) = \frac{20 + 2 + 4 + 6}{3}$$

$$\Rightarrow v(4) = 32/4 = 8.$$

Solving (5) using values of $v(0), v(1), v(2), v(3)$
and $v(4)$

$$v(5) = 6 + \frac{1}{6} \{ 0 + 2 + 4 + 6 + 8 + v(5) \}$$

$$\Rightarrow \frac{5}{6} v(5) = \frac{30 + 2 + 4 + 6 + 8}{6}$$

$$\Rightarrow v(5) = \frac{50}{6} = 10$$

Hence

$$\left\{ \begin{array}{l} v(0) = 50 \\ v(1) = 2 \\ v(2) = 4 \\ v(3) = 6 \\ v(4) = 8 \\ v(5) = 10. \end{array} \right.$$

Hence value of the policy for number of working machines
 $m=5$ $v(5)=10.$

2(a) Step 1 :- Initial Policy & its Evaluation

Let's assume the initial policy π_1 to be a "No repair" policy for all the states

i.e. q for states

$m=1, m=2, m=3, m=4, m=5$
all has No repair, No repair as the action

i.e. $\{0, 1, 2, 3, 4, 5\}$ = State Space

Initial	$m=0$	$m=1$	$m=2$	$m=3$
	{ No repair, No repair, No repair, No repair, No repair, No repair }			

and as from previous answer

$$V^0(0) = 0$$

$$V^0(1) = 2$$

$$V^0(2) = 4$$

$$V^0(3) = 6$$

$$V^0(4) = 8$$

$$V^0(5) = 10. \quad \text{for the No repair policy}$$

The above values represent the expected long term reward when starting with m working machines under "No repair" policy.

Step 2 Under the "Repair" action. And using the Bellman evaluation equation

$$V^n = \mathbb{E} R^n + \gamma P^n V^n.$$

transition matrix

under

Action "Repair"

	$m=0$	$n=1$	$m=2$	$n=3$	$n=4$	$n=5$
$m=0$	0	0	0	0	0	1
$m=1$	0	0	0	0	$\frac{1}{2}$	$\frac{1}{2}$
$m=2$	0	0	0	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
$m=3$	0	0	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
$m=4$	0	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{1}{5}$
$m=5$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

Initial policy

π (No Repair, No Repair, No Repair, No Repair,
No Repair, No Repair)

Step 2 Improve the policy for one iteration.

~~$$V(n+1) = V^n + R^n + \gamma P^n V^n.$$~~

for Action $a = \text{"Repair"}$.

$R(m, \text{"Repair"}) = m - N/2$ where $m >$
no of currently
correctly working machines

Discount factor $r = 1$.

Action = Repair.

$$V(t_0) = \infty$$

using:

$$V^n(s) : R(s, "Repair") + r [P_{S0}V(S_0) + P_{S1}V(S_1) + P_{S2}V(S_2) + P_{S3}V(S_3) \\ + P_{S4}V(S_4) + P_{S5}V(S_5)]$$

$$V(5) = 5 - \frac{5}{2} + 1 * \left[\frac{1}{6} * 0 + \frac{1}{6} * 2 + \frac{1}{6} * 4 + \frac{1}{6} * 6 + \frac{1}{6} * 8 + \frac{1}{6} * 10 \right]$$

$$\Rightarrow V(5) = 2.5 + \frac{30}{6}$$

$$= 7.5.$$

Similarly for

$$V(4) = 4 - \frac{5}{2} + 1 * \left[\frac{1}{5} * 0 + \frac{1}{5} * 20 + \frac{1}{5} * 4 + \frac{1}{5} * 6 + \frac{1}{5} * 8 + \frac{1}{5} * 10 \right]$$

$$= 4 - 2.5 + \frac{30}{5} = 7.5.$$

for calculating for $V(3)$

$$V(3) = 3 - \frac{5}{2} + 1 * \left[\frac{1}{4} * 0 + \frac{1}{4} * 0 + \frac{1}{4} * 0 + \frac{1}{4} * 6 + \frac{1}{4} * 8 + \frac{1}{4} * 10 \right]$$

$$= 3 - 2.5 + 7 = 7.5.$$

$$V(2) = \frac{2-5}{2} + 1 * \left[\frac{1}{3} \times 0 + \frac{1}{3} \times 0 + \frac{1}{3} \times 0 + \frac{1}{3} \times 6 + \frac{1}{3} \times 8 \right. \\ \left. + \frac{1}{3} \times 10 \right] \\ = 2 - 2.5 + 8 = 7.5.$$

$$V(1) = \frac{1-5}{2} + 1 * \left[\frac{1}{2} \times 0 + \frac{1}{2} \times 0 + \frac{1}{2} \times 0 + \frac{1}{2} \times 0 + \frac{1}{2} \times 8 \right. \\ \left. + \frac{1}{2} \times 10 \right] \\ = 1 - 2.5 + 9 = 7.5.$$

$$V(0) = 0 - \frac{5}{2} + 1 * \left[0 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 + 0 \times 0 \right. \\ \left. + 1 \times 10 \right] \\ = 10 - 5/2 = 7.5.$$

Once

	Repair	No Repair	Better Action
	$V(S, "Repair")$	$V(S, "No Repair")$	
S0	7.5	0	Repair
S1	7.5	2	Repair
S2	7.5	4	Repair
S3	7.5	6	Repair.
S4	7.5	8	No Repair
S5	7.5	10.	No Repair.

So the initial policy

$\pi \{ "No\ Repair", "No\ Repair", "No\ Repair", "No\ Repair",$
 $"No\ Repair", "No\ Repair" \}$

get's changed to

$\pi \{ "Repair", "Repair", "Repair", "Repair", "Repair",$
 $"No\ Repair" \}$

after one iteration

Section 3(a)

(a). At the terminal time step N player must take the reward based on the value of the dice as no further rolling is possible.

Therefore at $n=N$, the value function $V^N(s)$ for each state s is just the reward corresponding to the dice value s .

$$V^N(s) = R(s) = 3s^2 + 5.$$

$$V^N(s=1) = 3(1)^2 + 5 = 8$$

$$V^N(s=2) = 3(2)^2 + 5 = 17.$$

$$V^N(s=3) = 3(3^2) + 5 = 32$$

$$V^N(s=4) = 3(4^2) + 5 = 53.$$

(b) To compute $Q^{N-1}(s, a)$ for $a = \text{"Quit"}$ or $a = \text{"continue"}$,

we need to evaluate the expected value of taking each action (Quit/continue) at the second to last time step $N-1$ given state s .

$$Q^{N-1}(s, \text{"Quit"}) = R(s) = 3s^2 + 5.$$

$$Q^{N-1}(s, \text{"continue"}) =$$

$$E[V_N(s')] = \frac{1}{4} \sum_{s'=1}^4 V^N(s')$$

$$\Rightarrow E[V^N(s')] = \frac{1}{4} (V^N(1) + V^N(2) + V^N(3) + V^N(4)) \\ = \frac{1}{4} (8 + 17 + 32 + 53) = \frac{1}{4} \times 110 = 27.5.$$

Therefore for "continue" action

$$Q^{N-1}(s, \text{"continue"}) = 27.5$$

for $s=1$

$$Q^{N-1}(1, \text{"quit"}) = 3(1^2) + 5 = 8$$

$$Q^{N-1}(1, \text{"continue"}) = 27.5$$

for $s=2$

$$Q^{N-1}(2, \text{"quit"}) = 3(2^2) + 5 = 17$$

$$Q^{N-1}(2, \text{"continue"}) = 27.5$$

for $s=3$

$$Q^{N-1}(3, \text{"quit"}) = 3(3^2) + 5 = 32$$

$$Q^{N-1}(3, \text{"continue"}) = 27.5$$

for $s=4$

$$Q^{N-1}(4, \text{"quit"}) = 3(4^2) + 5 = 53$$

$$Q^{N-1}(4, \text{"continue"}) = 27.5$$

s	$Q^{N-1}(s, \text{"quit"})$	$Q^{N-1}(s, \text{"continue"})$
1	8	27.5
2	17	27.5
3	32	27.5
4	53	27.5

(c) as $V^{N+1}(s) = \max (\alpha^{N+1}(s, "quit"), \alpha^{N+1}(s, "continue"))$

therefore

$$\text{for } s=1 \quad V^{N+1}(1) = \max(8, 27.5) = 27.5$$

$$V^{N+1}(2) = \max(17, 27.5) = 27.5$$

$$V^{N+1}(3) = \max(32, 27.5) = 32$$

$$V^{N+1}(4) = \max(53, 27.5) = 53$$

(d) Applying Bellman equation for $V^{N+1}(s)$

$$V^{N+1}(s) = \max (\alpha^{N+1}(s, "quit"), \alpha^{N+1}(s, "continue"))$$

$$= \max (\cancel{\alpha^{N+1}})$$

$$\alpha^{N+1}(s, "quit") = R(s) = 3s^2 + 5$$

$$\alpha^{N+1}(s, "continue") = E[V_{s'}^N(s')] = \frac{1}{4} \sum_{s'=1}^4 V^N(s')$$

Hence

$$V^{N+1}(s) = \max \left(3s^2 + 5, \frac{1}{4} (V^N(1) + V^N(2) + V^N(3) + V^N(4)) \right)$$

(e)

When the player chooses to "Continue" at time $n-1$ they roll the dice again.

Since the dice is fair the future state s' will be uniformly distributed over the possible outcomes.

Therefore the expected value of continuing at time $n-1$ is the average of the action value of continuing at time n .

$$\mathbb{Q}_s^{n-1}(\delta, \text{"Continue"}) = \frac{1}{4} (\mathbb{Q}^n(1, \text{"Continue"}) + \mathbb{Q}^n(2, \text{"Continue"}) \\ + \mathbb{Q}^n(3, \text{"Continue"}) + \mathbb{Q}^n(4, \text{"Continue"})$$

Thus $\mathbb{Q}^{n-1}(\delta, \text{"Continue"})$ depends on average value of $\mathbb{Q}^n(\delta', \text{"Continue"})$.

(f)

$$\mathbb{Q}^n(\delta, \text{"Quit"}) = 3\delta^2 + 5$$

$$\mathbb{Q}^n(\delta, \text{"Continue"}) = \frac{1}{4} \sum_{\delta'=1}^4 V^{n+1}(\delta')$$

optimal policy

$$\pi^*(\delta, n) = \begin{cases} \text{"Quit"} & \text{if } \mathbb{Q}^n(\delta, \text{"Quit"}) \geq \mathbb{Q}^n(\delta, \text{"Continue"}) \\ \text{"Continue"} & \text{if } \mathbb{Q}^n(\delta, \text{"Quit"}) < \mathbb{Q}^n(\delta, \text{"Continue"}) \end{cases}$$

$$= \begin{cases} \text{"Quit"} & \text{if } 3\delta^2 + 5 \geq \frac{1}{4} \sum_{\delta'=1}^4 V^{n+1}(\delta') \\ \text{"Continue"} & \text{if } \frac{1}{4} \sum_{\delta'=1}^4 V^{n+1}(\delta') > 3\delta^2 + 5 \end{cases}$$

(g) The optimal policy in the dice game is non stationary as the decision to "stop" or "continue" depends not only on the current state s , but also on the time step n .

As the game progresses & the number of remaining rolls decreases the player's optimal choice can change making the policy time dependent.