

• CSE564 - Reinforcement Learning HW 2

Name: Ansh Arora

Roll No.: 2019022

- ① Given arms 1 to 10 with Gaussian distributions with mean 1 to 10 respectively and variance 1.

$$p[a] = \begin{cases} 2/15 & a=2, 5, 6, 8, 10 \\ 1/15 & a=1, 3, 4, 7, 9 \end{cases}$$

Expected Reward on each pick

$$= E[q_{\pi^*}(a)]$$

$$= \sum p(a) \times q_{\pi^*}(a) \quad (q_{\pi^*}(a) = \frac{E[R_t | A_t = a]}{1})$$

$$= \frac{2}{15} \times (2+5+6+8+10) + \frac{1}{15} \times (1+3+4+7+9)$$

$$= \frac{2}{15} \times (30) + \frac{1}{15} \times (25)$$

$$= 4 + \frac{5}{3} = \frac{17}{3}$$

Total reward after 10 picks = $10 \times E[q_{\pi^*}(a)]$

$$= 10 \times \frac{17}{3} = 56.66$$

② $q_{\pi^*}(a) = 0 \times 0.5 + 1 \times 0.5$
 $= 0.5 \quad \forall a \in \{1, 2, 4, 5, 7, 9, 10\}$

Similarly, $q_{\pi^*}(a) = 0 \times 0.3 + 0.2 \times 0.3 + 1 \times 0.4$
 $= 0.46 \quad \forall a \in \{3, 6, 8\}$

Thus, for stochastic optimal policies, we need to make use of arms with $q_{\pi}(a) = 0.5$. Any combination of these arms would be optimal. Six such-

Arm #

- 1) Select 9 and 10 with 0.5 each
- 2) Select 7, 8 and 10 with 0.33 each
- 3) Select 5, 7, 9 and 10 with 0.25 each
- 4) Select 7, 5, 7, 8 and 10 with 0.2 each
- 5) Select 1 with 0.3 and 2 with 0.7
- 6) Select 1 with 0.1, 2 with 0.1 and 4 with 0.8.

③ Let at time $t=1$,

$$Q_1(1) = 0.2, Q_1(2) = 0.3, Q_1(3) = 0.5$$

In the given ϵ -greedy -

1) We explore at $t=1, 3, 5$

2) We exploit at $t=2, 4, 6$

3) We explore over non-greedy actions only.

At $t=1$, explore among arms 1, 2 (non-greedy)

Suppose we select Arm #1.

Let it return value 1.

Then, after $t=1$, $A_1 = 1, R_1 = 1$

$$Q_2(1) = \frac{0.2+1}{2} = 0.6, Q_2(2) = 0.3, Q_2(3) = 0.5$$

At $t=2$, choose greedy -

We select arm #1.

Let it return value 0.

Then, after $t=2$, $A_2 = 1, R_2 = 0$

$$Q_3(1) = \frac{0.6+0}{3} = 0.2, Q_3(2) = 0.3, Q_3(3) = 0.5$$

Date: ___/___/___

Date: ___/___/___

At $t=3$, explore from arms 1, 2
We select arm #2.

Let it return value 0.

Then, after $t=3$, $A_3 = 2, R_3 = 0$

$$Q_4(2) = \frac{0.3+0}{2} = 0.15, Q_4(1) = 0.2, Q_4(3) = 0.5$$

At $t=4$, choose greedy -

We select arm #3.

Let it return value 1.

Then, after $t=4$, $A_4 = 3, R_4 = 1$

$$Q_5(3) = \frac{0.5+1}{2} = 0.75, Q_5(1) = 0.2, Q_5(2) = 0.15$$

At $t=5$, choose from arms (1, 2) (explore)

We select arm #1.

Let it return value 1.

After $t=5$, $A_5 = 1, R_5 = 1$

$$Q_6(1) = \frac{0.2+1}{2} = 0.3, Q_6(2) = 0.15, Q_6(3) = 0.75$$

At $t=6$, greedy.

We select arm #3

Let it return value 0.

After $t=6$, $A_6 = 3, R_6 = 0$.

#earnthesmarterway

#earnthesmarterway

s	a	s'	r	$p(s', r s, a)$
high	search	high	search	α
high	search	low	search	$1-\alpha$
high	wait	high	wait	1
low	search	high	-3	$1-\beta$
low	search	low	search	β
low	wait	low	wait	1
low	recharge	high	0	1

This table was derived from the MDP table and using the transition graph given for the example 3.3.

Here, in the case where s is low and s' is high, we assume that the robot got discharged and to get recharged again took a -ve reward of -3.

- (5) Prove using (3.8) that adding a constant c to all rewards only adds a constant v_c to values of all states, and thus doesn't affect relative values of any states under any policies. Write v_c in terms of c and γ .

For a given state at time t ,

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots$$

$$\begin{aligned} G_{t, \text{new}} &= (c + R_{t+1}) + (\cancel{(c + R_{t+2})}) + (\cancel{\gamma(c + R_{t+3})}) + \dots \\ &= (R_{t+1} + \gamma R_{t+2} + \dots) + (c + \gamma c + \gamma^2 c + \dots) \end{aligned}$$

$$\Rightarrow G_{t, \text{new}} = G_t + \frac{c}{1-\gamma} \quad : \text{Relative value remain unchanged.}$$

#learnthesmarterway

Date ___/___/___

Date ___/___/___

Add a constant c to an episodic task. What effect does it have?

For an episodic task,

$$G_T = R_{T+1} + \gamma R_{T+2} + \dots + \overset{T \text{ (last)}}{\gamma^{T-1} R_T} \quad \text{where } T \text{ is final time step.}$$

Now, progressing like before

$$\begin{aligned} G_{T, \text{new}} &= (c + R_{T+1}) + \gamma(c + R_{T+2}) + \dots + \overset{T \text{ (last)}}{\gamma^{T-1}(c + R_T)} \\ &= (c + \gamma c + \dots + \overset{T \text{ (last)}}{\gamma^{T-1} c}) + G_T \\ &= c \frac{1 - \gamma^{T-1}}{1 - \gamma} + G_T \end{aligned}$$

v_c .

∴ If the value of T increases, the value of v_c would also increase. So, the added constant would promote the agent to take more steps, thus having an effect.

#learnthesmarterway

⑧ Equation for v_* in terms of q_{V*} .

$$v_*(s) = \max_{a \in A(s)} q_{V*}(s, a)$$

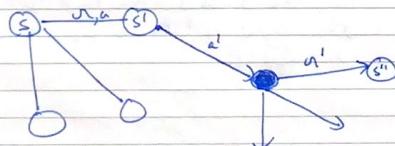
$$= \max_{a \in A(s)} \sum_{s' \in r} p(s', r | s, a) \left[r + \gamma \max_{a' \in A(s')} q_{V*}(s', a') \right]$$

⑨ We need to figure out if random variable R_{t+2} is dependent on s_t and A_t , and if yes then how.

Let $r = R_{t+2}$ and $v = R_{t+1}$.
Then,

$$P_r(r = R_{t+2} | s_t = S_t, A_t = a)$$

$$= \sum_{\substack{s' \in s \\ r \in R}} p(s', r | s, a) \sum_{a' \in A(s')} \pi(a' | s') \cdot p(s'', r' | s', a')$$



⑫ $E_{\pi}[R_{t+2} | s_t = S_t, A_t = a]$

$$= \sum_{r'} P_r(r = R_{t+2} | s_t = S_t, A_t = a)$$

$$= \sum_{r'} \sum_{s', a} p(s', r | s, a) \sum_{a', s''} \pi(a' | s') p(s'', r' | s', a')$$

$$\begin{aligned}
 \textcircled{13} \quad v_{\pi}(s) &= E[G_t | S_t = s] \\
 &= E_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\
 &= E_{\pi}[R_{t+1} | S_t = s] + \gamma E_{\pi}[G_{t+1} | S_t = s] \\
 &= E_{\pi}[R_{t+1} | S_t = s] + \gamma E_{\pi}[E_{\pi}[G_{t+1} | S_t = s'] | S_t = s] \\
 &= E_{\pi}[R_{t+1} | S_t = s] + \gamma E_{\pi}[v_{\pi}(s') | S_t = s] \\
 &= E_{\pi}[R_{t+1} + \gamma v_{\pi}(s') | S_t = s] \\
 &= \sum_{a \in A(s)} \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]
 \end{aligned}$$

Bellman eqⁿ for $v_{\pi}(s)$

\textcircled{14} We start from the back. $\gamma = 0.5$

$$\begin{aligned}
 G_3 &= R_3 + \gamma G_4 \\
 &= -3 + 0.5 \times 0 \\
 &= -3
 \end{aligned}$$

$$\begin{aligned}
 G_2 &= R_2 + \gamma G_3 \\
 &= 10 + 0.5 \times (-3) \\
 &= 8.5
 \end{aligned}$$

$$\begin{aligned}
 G_1 &= R_1 + \gamma G_2 \\
 &= -1 + 0.5 \times (8.5) \\
 &= 8.25
 \end{aligned}$$

$$\begin{aligned}
 G_0 &= R_0 + \gamma G_1 \\
 &= 2 + 0.5 \times 8.25 \\
 &= 2 + 4.125 \\
 &= 6.125
 \end{aligned}$$

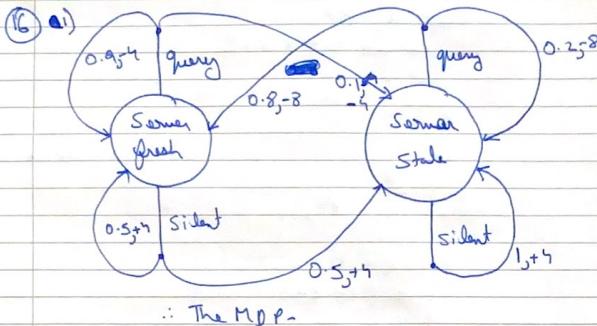
If an agent receives a constant reward c ,

$$\begin{aligned}
 G_t &= R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots \\
 &= c + \gamma c + \gamma^2 c + \dots \\
 &= c(1 + \gamma + \gamma^2 + \dots) \\
 &= c \left(\frac{1}{1-\gamma} \right) \\
 &= \frac{c}{1-\gamma}
 \end{aligned}$$

$$\textcircled{15} \quad \pi^*(s) = \arg \max_{a \in A(s)} Q_{\pi}(s, a)$$

$$\Rightarrow \pi^*(s) = \arg \max_{a \in A(s)} \sum_{s'} \sum_{r} p(s', r | s, a) [r + \gamma \cdot v_{\pi}(s')]$$

Thus, what we do is simply go over all actions that can be taken at a given state s and find the best action that returns the highest value.



s	a	s'	r	$p(s', r s, a)$
fresh	query	fresh	-4	0.9
fresh	query	stale	-4	0.1
fresh	silent	fresh	+4	0.5
fresh	silent	stale	+4	0.5
stale	query	fresh	-8	0.8
stale	query	stale	-8	0.2
stale	silent	stale	+7	1

3) Value iteration-

$$\gamma = 0.5$$

$$\text{Initialization } v(\text{fresh}) = 0, v(\text{stale}) = 0$$

$$v_{k+1}(s) = \max_a \sum_{s', r} p(s', r | s, a) (r + \gamma v_k(s'))$$

1st iteration-

$$v(\text{fresh}) = \max \left\{ 0.9 \times (-4 + 0.5 \times 0) + 0.1 \times (-4 + 0.5 \times 0), 0.5 \times (4 + 1/2 \times 0) + 0.5 \times (4 + 1/2 \times 0) \right\}$$

$$= \max \{-4, 4\} = 4$$

#learnthesmarterway

$$v(\text{stale}) = \max \left\{ 0.8 \times (-8 + 0.5 \times 0) + 0.2 \times (-8 + 0.5 \times 0), 1 \times (4 + 0.5 \times 0) \right\}$$

$$= \max \{-8, 4\} = 4$$

$$\therefore v(\text{fresh}) = 4, v(\text{stale}) = 4.$$

2nd iteration-

$$v(\text{fresh}) = \max \left\{ 0.9 \times (-4 + 0.5 \times 4) + 0.1 \times (-4 + 0.5 \times 4), 0.5 \times (4 + 0.5 \times 4) + 0.5 \times (4 + 0.5 \times 4) \right\}$$

$$= \max \{-2, 6\} = 6$$

$$v(\text{stale}) = \max \left\{ 0.8 \times (-8 + 0.5 \times 4) + 0.2 \times (-8 + 0.5 \times 4), 1 \times (4 + 0.5 \times 4) \right\}$$

$$= \{-6, 6\} = 6$$

3rd iteration-

$$v(\text{fresh}) = \max \left\{ 0.9 \times (-4 + 0.5 \times 6) + 0.1 \times (-4 + 0.5 \times 6), 0.5 \times (4 + 0.5 \times 6) + 0.5 \times (4 + 0.5 \times 6) \right\}$$

$$= \max \{-1, 7\} = 7$$

$$v(\text{stale}) = \max \left\{ 0.8 \times (-8 + 0.5 \times 6) + 0.2 \times (-8 + 0.5 \times 6), 1 \times (4 + 0.5 \times 6) \right\}$$

$$= \max \{-5, 7\} = 7$$

4th iteration-

$$v(\text{fresh}) = \max \left\{ 0.9 \times (-4 + 0.5 \times 7) + 0.1 \times (-4 + 0.5 \times 7), 0.5 \times (4 + 0.5 \times 7) + 0.5 \times (4 + 0.5 \times 7) \right\}$$

$$= \max \{-0.5, 7.5\} = 7.5$$

#learnthesmarterway

$$v(\text{stale}) = \max \left\{ -8 + \frac{1}{2} (0.8 \times 7 + 0.2 \times 7), 7 + \frac{1}{2} \times 1 \times 7 \right\}$$

$$= \max \{-4.5, 7.5\} = 7.5$$

∴ Final after 4 iterations,

$$v(\text{fresh}) = 7.5 \quad v(\text{stale}) = 7.5$$

Policy Iteration -

~~$v_{k+1}(s) = \sum \pi_k(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')]$~~

$$v_{k+1}(s) = \sum_{a \in \mathcal{A}(s)} \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')]$$

and then,

$$\pi_{k+1}(s) = \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma v_{k+1}(s')]$$

$$v(\text{fresh}) = \pi(\text{query | fresh}) \times [0.9(-4 + 0.5v_k(\text{fresh})) + 0.1 \times (-4 + 0.5v_k(\text{stale}))] \\ + \pi(\text{quiet | fresh}) \times [0.5 \times (-4 + 0.5v_k(\text{fresh})) + 0.5 \times (-4 + 0.5v_k(\text{stale}))]$$

↓↓↓↓↓

$$= \pi(\text{query | fresh}) [-4 + 0.5 \left(\frac{0.9v_k(\text{fresh}) + 0.1v_k(\text{stale})}{2} \right)] \\ + \pi(\text{quiet | fresh}) [-4 + 0.5 \left(\frac{0.9v_k(\text{fresh}) + 0.1v_k(\text{stale})}{2} \right)]$$

$$v(\text{stale}) = \pi(\text{query | stale}) [0.8 \times (-8 + 0.5v_k(\text{fresh})) \\ + 0.2 \times (-8 + 0.5v_k(\text{stale}))] \\ + \pi(\text{quiet | stale}) [-4 + 0.5v_k(\text{stale})]$$

$$= \pi(\text{query | stale}) [-8 + 0.5 \left(\frac{0.8v_k(\text{fresh}) + 0.2v_k(\text{stale})}{2} \right)] \\ + \pi(\text{quiet | stale}) [-4 + 0.5v_k(\text{stale})]$$

#learnthesmarterway

1st Iteration

Policy evaluation -

Let π_0 be -

$$\pi_0(\text{query | stale}) = 0$$

$$\pi_0(\text{query | fresh}) = 0$$

$$\pi_0(\text{quiet | stale}) = 1$$

$$\pi_0(\text{quiet | fresh}) = 1$$

$$v_{\pi_0}(\text{stale}) = 0 \rightarrow y.$$

$$v_{\pi_0}(\text{fresh}) = 0 \rightarrow z.$$

$$v_{\pi_0}(\text{stale}) = x = 4 + 0.25x + 0.25y$$

↓↓↓↓↓

$$y = 4 + 0.5z.$$

$$\Rightarrow y = 8, z = 8$$

$$v_{\pi_0}(\text{stale}) = v_{\pi_0}(\text{fresh})$$

Policy Improvement -

$$\pi_1(\text{fresh}) = \arg \max_{\text{query, quiet}} \left\{ -4 + 0.5 (0.9 \times 8 + 0.1 \times 8), 4 + 0.5 (8 + 8) \right\}$$

$$= \arg \max \{ 0, 8 \}.$$

$$\therefore \pi_1(\text{fresh}) = \text{quiet} \Rightarrow \pi_1(\text{quiet | fresh}) = 1$$

$$\pi_1(\text{stale}) = \arg \max_{\text{query, quiet}} \left\{ -8 + 0.5 (0.8 \times 8 + 0.2 \times 8), 4 + 0.5 (8) \right\}$$

$$= \arg \max \{ -4, 8 \} = \text{quiet}$$

$$\therefore \pi_1(\text{stale}) = \text{quiet} \Rightarrow \pi_1(\text{quiet | stale}) = 1$$

∴ $\pi_0 = \pi_1$, by policy improvement theorem
we have achieved the optimal policy.

(17) Proving the correctness of policy improvement step-

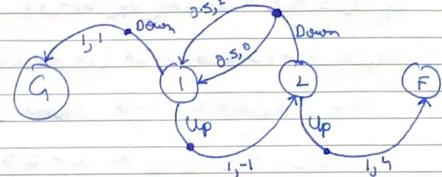
$$\begin{aligned}
 v_{\pi}(s) &\leq q_{\pi}(s, \pi'(s)) \quad \text{where } \pi'(s) \\
 &= E_{\pi'}[R_{t+1} + \gamma v_{\pi}(s_{t+1}) | S_t = s, A_t = \pi'(s)] \\
 &= E_{\pi'}[R_{t+1} + \gamma v_{\pi}(s_{t+1}) | S_t = s] \\
 &\leq E_{\pi'}[R_{t+1} + \gamma q_{\pi}(s_{t+1}, \pi'(s_{t+1})) | S_t = s] \\
 &= E_{\pi'}[R_{t+1} + \gamma E_{\pi}[R_{t+2} + \gamma v_{\pi}(s_{t+2}) | S_{t+1}]] \\
 &\quad A_{t+1} = \pi'(S_{t+1}) \\
 &= E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_{\pi}(s_{t+2}) | S_t = s]
 \end{aligned}$$

On further expanding:

$$\begin{aligned}
 &\leq E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} \dots | S_t = s] \\
 &= E_{\pi'}[G_+ | S_t = s] \\
 &= v_{\pi'}(s)
 \end{aligned}$$

$$\therefore v_{\pi}(s) \leq v_{\pi'}(s) \quad \text{if } \pi(s) \leq \pi'(s)$$

(18) Modelling taking a stairwell of 2 steps b/w G and F.



Initialization -

$$\begin{aligned}
 \pi(U_p|1) &= 0.5 & \pi(U_p|2) &= 0.5 \\
 \pi(Down|1) &= 0.5 & \pi(Down|2) &= 0.5
 \end{aligned}$$

$$v_{\pi_0}(s) = 0 \quad \forall s \in \{G, I, L, F\}$$

$$Up = \uparrow \quad Down = \downarrow$$

General form -

$$v_{\pi}(s) = \sum_{a \in \mathcal{A}(s)} \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$

$$\begin{aligned}
 v_{\pi}(1) &= \pi(\uparrow|1)(1 \times (-1 + v_{\pi}(2))) \\
 &\quad + \pi(\downarrow|1)(1 \times (1 + v_{\pi}(2)))
 \end{aligned}$$

$$= \pi(\downarrow|1) + \pi(\uparrow|1) (v_{\pi}(2) - 1)$$

$$\begin{aligned}
 v_{\pi}(2) &= \pi(\uparrow|2)(1 \times (0 + v_{\pi}(F))) \\
 &\quad + \pi(\downarrow|2)(0.5 \times (2 + v_{\pi}(1)) + 0.5 \times (0 + v_{\pi}(1)))
 \end{aligned}$$

$$= (0.5 \pi(\uparrow|2) + (1 + v_{\pi}(1))) \times \cancel{\pi(\downarrow|2)}$$

$$= 0.5 \pi(\uparrow|2) + \pi(\downarrow|2) (1 + v_{\pi}(1))$$

#learnthesmarterway

1st iteration

$$v_{\pi_0}(1) = 0.5 + 0.5 \times (v_{\pi_0}(2) - 1)$$

$$v_{\pi_0}(2) = 4 \times 0.5 + 0.5 \times (1 + v_{\pi_0}(1))$$

Taking $v_{\pi_0}(1) = x$ and $v_{\pi_0}(2) = y$

$$x = 0.5 + 0.5(y-1) \Rightarrow x = 0.5y$$

$$y = 2 + 0.5(1+x) \Rightarrow y = 2.5 + 0.5x$$

$$y = 2.5 + 0.25y$$

$$0.75y = 2.5$$

$$y = 2.5 = \frac{10}{3}$$

$$x = 0.5y = \frac{5}{3}$$

$$\therefore v_{\pi_0}(1) = \frac{5}{3}$$

$$v_{\pi_0}(2) = \frac{10}{3}$$

Policy improvement-

$$\pi_1(1) = \operatorname{argmax}_{\uparrow, \downarrow} \{ 1 \times (\frac{10}{3} - 1), 1 \}$$

$$= \operatorname{argmax} \{ \frac{7}{3}, 1 \} = \uparrow$$

$$\pi_1(2) = \operatorname{argmax}_{\uparrow, \downarrow} \{ 4, 1 + \frac{5}{3} \}$$

$$= \operatorname{argmax} \{ 4, \frac{8}{3} \} = \uparrow$$

$$\therefore \pi_1(\uparrow 1) = 1 \quad \pi_1(\downarrow 1) = 0$$

$$\pi_1(\uparrow 2) = 1 \quad \pi_1(\downarrow 2) = 0$$

Date ___/___/___

Date ___/___/___

2nd iteration

$$v_{\pi_1}(1) = 0 + 1 \times (v_{\pi_1}(2) - 1)$$

$$v_{\pi_1}(2) = 1 \times 4 + 0 = 4$$

$$\therefore v_{\pi_1}(1) = 3, v_{\pi_1}(2) = 4$$

Policy improvement-

$$\pi_2(1) = \operatorname{argmax}_{\uparrow, \downarrow} \{ 1 \times (4-1), 1 \}$$

$$= \operatorname{argmax} \{ 3, 1 \} = \uparrow$$

$$\pi_2(2) = \operatorname{argmax}_{\uparrow, \downarrow} \{ 4, 1+3 \} = \uparrow, \downarrow$$

$$\therefore \pi_2(\uparrow 1) = 1 \quad \pi_2(\downarrow 1) = 0$$

$$\pi_2(\uparrow 2) = 0.5 \quad \pi_2(\downarrow 2) = 0.5$$

3rd iteration

$$v_{\pi_2}(1) = 0 + 1 \times (v_{\pi_2}(2) - 1)$$

$$v_{\pi_2}(2) = 4 \times 0.5 + 0.5 \times (1 + v_{\pi_1}(1))$$

$$\Rightarrow x = y-1 \quad \text{and} \quad y = 2.5 + \frac{x}{2}$$

$$x+1 = 2.5 + \frac{x}{2}$$

$$\frac{2x}{2} = 1.5 \Rightarrow x = 3 \\ y = 4.$$

$$v_{\pi_2}(1) = 3 \quad v_{\pi_2}(2) = 4.$$

#learnthesmarterway

#learnthesmarterway

Policy improvement -

$$\pi_1(1) = \underset{\uparrow, \downarrow}{\operatorname{argmax}} \{ 1 \times (-1), 1 \} = \uparrow$$

$$\pi_3(2) = \underset{\uparrow, \downarrow}{\operatorname{argmax}} \{ -1, 1 + 3 \} = \uparrow, \downarrow$$

$$\therefore \pi_3 = \pi_2$$

~~∴ By policy improvement theorem,~~
 $\pi_2 = \pi_3 = \pi^*$

∴ Optimal policy is -

$$\begin{aligned}\pi^*(\uparrow|1) &= 1 & \pi^*(\uparrow|2) &= 0.5 \\ \pi^*(\downarrow|1) &= 0 & \pi^*(\downarrow|2) &= 0.5\end{aligned}$$

The states above 1 show an equal policy for both \uparrow and \downarrow . This is because $\gamma=1$, which means that the number of steps taken to reach the top don't play a role. Every policy where we don't go down from 1 plays out and returns the exact same value.

Date ___/___/___

Date ___/___/___

(a) Let $H = \text{Healthy}$, $S = \text{Sick}$.

