

REINFORCEMENT LEARNING - CSE564

Homework 2

Suchet Aggarwal

2018105

** Note: All code questions are implemented in 2018105_code.ipynb, with comments specifying each question number.

Q1.

B1.

s	a	s'	τ	$p(s', \tau s, a)$
high	Search	high	r _{scan}	α
high	search	low	r _{scan}	$1 - \alpha$
low	Search	high	-3	$1 - \beta$
low	Search	low	r _{scan}	β
high	wait	high	r _{wait}	1
high	wait	low	-	0
low	wait	high	-	0
low	wait	low	r _{wait}	1
low	recharge	high	0	1
low	recharged	low	-	0

In general :

$$p(\tau | s', s, a) = \frac{p(\tau, s', s, a)}{p(s', s, a)}$$

$$= \frac{p(\tau, s', s, a)}{\frac{p(s, a)}{p(s', s, a)}}$$

$$= \frac{p(s', \tau | s, a)}{p(s' | s, a)}$$

$$\Rightarrow p(s', \tau | s, a) = p(s' | s, a) \cdot p(\tau | s', s, a)$$

In this example the reward is fixed for given transition $s \xrightarrow{a} s'$, $\therefore p(\tau | s', s, a) = 1$
 $\therefore p(s', \tau | s, a) = p(s' | s, a)$

Q2. Code Attached



Q3.

CLASSMATE
Date _____
Page _____

Q3. (a) The signs of the rewards aren't important rather the intervals between them are.

From Eq. (3.8)

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

adding C to each $R_i + 2$

$$G'_t = \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + C) \\ = G_t + C \cdot \sum_{k=0}^{\infty} \gamma^k \\ = G_t + \frac{C}{1-\gamma}$$

Now $v_{\pi}(s) = E(G_t | s_t = s)$

$$v_{\pi}'(s) = E(G'_t | s_t = s)$$

$$= E(G_t + \frac{C}{1-\gamma} | s_t = s) \\ = E(G_t | s_t = s) + \frac{C}{1-\gamma}$$

(Second term is constant)

$$v_{\pi}'(s) = v_{\pi}(s) + \frac{C}{1-\gamma}$$

The last step proves adding constant C to all rewards just adds a constant v_C to all $v_{\pi}(s)$, and doesn't affect relative values

Q3 (b)

In case of an Episodic task, the signs of reward are important, mainly because it conveys to the agent to complete the task within a time frame by giving negative rewards.

If such rewards are made positive the length of the episode may change that can alter the relative difference in rewards.

$$G_T = R_{T+1} + R_{T+2} + \dots + R_{T+T}$$

$$G'_T = (R_{T+1} + c) + \dots + (R_{T+K} + c)$$

~~where~~

~~here $K \neq T$ for all states~~
 here, K may not be same as T for $s \in S$
 thus the v_π values would get changed without maintaining proper relative order.

For Example : Consider a robot (agent), and the task is to move from position 0 to position 3, with allowed actions left and right.

(A)



with reward -1 at all steps.

Here obviously, optimal strategy is to move right always, with return = -3

However if all rewards are made +1, then practically maximum return becomes infinity with agent just moving between 0, 1, 2

Q4. Code Attached

Q5. The optimal value function is the max over all actions in the optimal action value function

$$v^*(s) = \max_{a \in A(s)} q^*(s, a) \quad \forall s \in S$$

Q6. Code Attached

Q7. Code Attached

Q8.

CLASSMATE

Date _____

Page _____

Q8. Yes, R_{t+2} is dependant on S_t, A_t

$$\text{consider } P(R_{t+2} | S_t, A_t)$$

$$= \sum_{S_{t+2}} P(S_{t+2}, R_{t+2} | S_t, A_t)$$

$$= \sum_{S_{t+2}} \sum_{S_{t+1}} \sum_{R_{t+1}} \sum_{A_{t+1}} P(S_{t+2}, R_{t+2} | S_{t+1}, R_{t+1}, A_{t+1}, S_t, A_t) \\ \cdot P(S_{t+1}, R_{t+1}, A_{t+1} | S_t, A_t)$$

$$= \sum_{S_{t+2}} \sum_{S_{t+1}} \sum_{R_{t+1}} \sum_{A_{t+1}} P(S_{t+2}, R_{t+2} | S_{t+1}, A_{t+1}),$$

$$P(A_{t+1} | S_{t+1}) \cdot P(S_{t+1}, R_{t+1} | S_t, A_t)$$

(by markov assumption :

$$P(S_{t+2}, R_{t+2} | S_{t+1}, R_{t+1}, A_{t+1}, S_t, A_t) =$$

$$P(S_{t+2}, R_{t+2} | S_{t+1}, A_{t+1})$$

probability is independent of previous state reward)

$$\text{and } (P(S_{t+1}, R_{t+1}, A_{t+1} | S_t, A_t) =$$

$$P(A_{t+1} | S_{t+1}) \cdot P(S_{t+1}, R_{t+1} | S_t, A_t)$$

as action at $t+1$ depends only on S_{t+1})

$$P(R_{t+2} | S_t, A_t) = \sum_{S_{t+2}} \sum_{S_{t+1}} \sum_{R_{t+1}} \sum_{A_{t+1}} P(S_{t+2}, R_{t+2} | S_{t+1}, A_{t+1}) \cdot \\ \pi(A_{t+1} | S_{t+1}) \cdot \\ P(S_{t+1}, R_{t+1} | S_t, A_t)$$

Q9.

classmate

Date _____

Page _____

Q9.

$$\Rightarrow E(R_{t+2} | S_t, A_t) = \sum \cancel{s} \cancel{r} \cancel{\gamma}$$

$$= \sum_r r \cdot p(R_{t+2} = r | S_t, A_t)$$

From Question 8.

$$= \sum_r r \cdot \sum_{S_{t+2}} \sum_{R_{t+1}} \sum_{A_{t+1}} p(S_{t+2}, R_{t+2} = r | S_{t+1}, A_{t+1}) \\ \cancel{\sum_r} \quad S_{t+2} \quad R_{t+1} \quad A_{t+1} \\ \cdot P(A_{t+1} | S_{t+1}) \\ \cdot P(S_{t+1}, R_{t+1} | S_t, A_t)$$

(For ease of read; read S_{t+2} as $S_{t+2} = s''$

R_{t+2} as $R_{t+2} = r$

R_{t+1} as $R_{t+1} = r'$

A_{t+1} as $A_{t+1} = a'$

S_{t+1} as $S_{t+1} = s'$

S_t as $S_t = s$

A_t as $A_t = a$)

Q10.

$$\begin{aligned}v_{\pi}(s) &= E_{\pi}(G_t \mid S_t = s) \\&= E_{\pi}(R_t + \gamma G_{t+1} \mid S_t = s) \\&= \sum_a \pi(a|s) \sum_{s',r} p(s',r \mid s,a) [r + \gamma E_{\pi}(R_t + \gamma G_{t+1} \mid S_{t+1} = s')] \\&= \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_{\pi}(s')]\end{aligned}$$

Q11.

classmate

Date _____

Page _____

Q11. Given $R_1 = 2, R_2 = -1, R_3 = 10, R_4 = -3$

Also

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$$G_4 = 0$$

$$\Rightarrow G_3 = -3$$

$$G_2 = 10 + 0.5(-3) = 8.5$$

$$G_1 = -1 + 0.5(8.5) = 3.25$$

$$G_0 = 2 + 0.5(3.25) = 3.625$$

(Considering the episode length to be 4)

Now in general:

$$G_t = \sum_{k=0}^{\infty} \gamma^k \cdot R_{t+1+k}$$

$$\text{Here } R_{t+1+k} \forall k = c$$

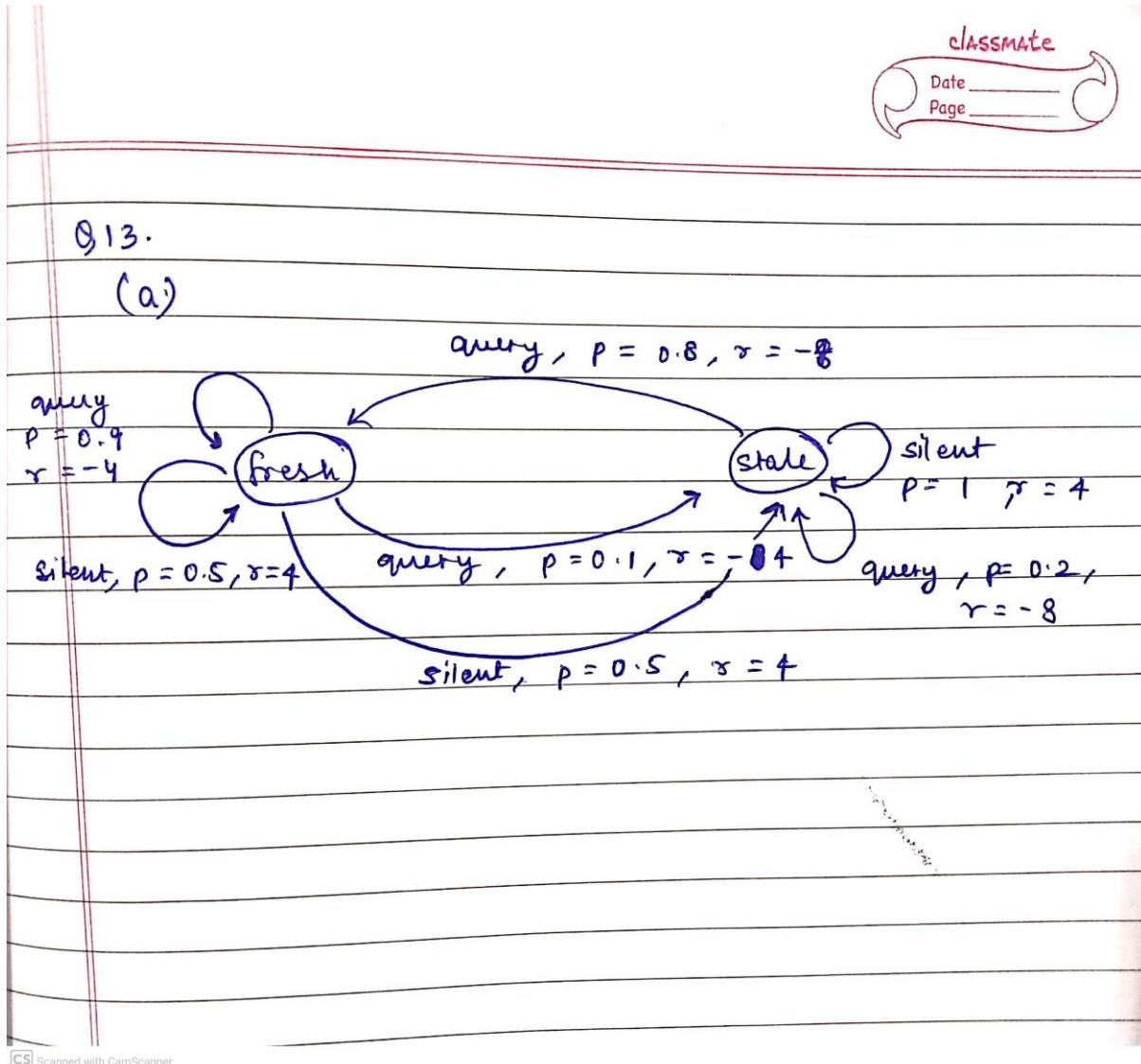
$$\therefore G_t = c \cdot \sum_{k=0}^{\infty} \gamma^k$$

$$G_t = \frac{c}{1-\gamma}$$

Q12. The optimal policy can be obtained by greedily picking all actions from a state a that maximise the overall return using the optimal value function itself as follows:

$$\pi^*(s) = \operatorname{argmax}_{a \in A(s)} (E(R_t + \gamma v^*(s') | S_t = s, A_t = a)) \quad \forall s \in S$$

Q13.



(b) Given a finite horizon problem (i.e $T=3$)

the rewards are R_1, R_2, R_3 and finally
depending on the final state we get R .

(assuming a discounting factor of 0.5)

let $dp(s, t)$ denote the maximum expected return
in state "s" at time t when beginning to sit at
time t .

$$\text{Also given } dp(\text{stale}, 4) = -10 \\ dp(\text{fresh}, 4) = 10$$

$$dp(\text{stale}, 3) = \max_{t \in \{0, 1\}} \left(0.8 \left(-8 + \frac{1}{2}(10) \right) + 0.2 \left(-8 + \frac{1}{2}(-10) \right) \right. \\ \left. + + \frac{1}{2}(-10) \right) : :$$

$$\pi(\text{stale}, 3) = \text{silent} \quad (\text{argmax of above})$$

$$dp(\text{fresh}, 3) = \max_{t \in \{0, 1\}} \left(0.1 \left(-4 + \frac{-10}{2} \right) + 0.9 \left(-4 + \frac{1}{2}(10) \right) \right. \\ \left. + 0.5 \left(4 + \frac{1}{2}(10) \right) + 0.5 \left(4 + \frac{1}{2}(-10) \right) \right) \\ = 4$$

$$\pi(\text{fresh}, 3) = \text{silent} \quad (\text{argmax of above})$$

$$dp(\text{stale}, 2) = \max_{\{0, 5\}} (0.8)(-8 + \frac{1}{2}(dp(\text{fresh}, 3))$$

$$+ 0.2(-8 + \frac{1}{2}(dp(\text{stale}, 3)))$$

$$1(4 + \frac{1}{2}dp(\text{stale}, 3))$$

$$= \max_{\{0, 5\}} (0.8(-8 + 2) + 0.2(-8 + -0.5))$$

$$9, 4 + \frac{1}{2}(-1)$$

$$= \max_{\{-2.5, 3.5\}} -6.5$$

$$dp(\text{stale}, 2) = 3.5$$

$$\pi(\text{stale}, 2) = \text{Silent}$$

$$dp(\text{fresh}, 2) = \max_{\{0, 5\}} (0.9(-4 + \frac{1}{2}dp(\text{fresh}, 3)) +$$

$$0.1(-4 + \frac{1}{2}dp(\text{stale}, 3))$$

$$0.5(4 + \frac{1}{2}dp(\text{stale}, 3))$$

$$= \max_{\{-4 + 2, -4 + 0.5\}} (0.9(-4 + 2) + 0.1(-4 + 0.5))$$

$$0.5(4 + 2) + 0.5(4 + -0.5)$$

$$dp(\text{fresh}, 2) = \max(-2.25, 4.75) = 4.75$$

$$\pi(\text{fresh}, 2) = \text{Silent}$$

$$dp(\text{state}, 1) = \max_{\text{S1}, \text{S2}} \left(0.8 \left(-8 + \frac{1}{2} dp(\text{fresh}, 2) \right) + 0.2 \left(-8 + \frac{1}{2} dp(\text{state}, 2) \right), 1 \left(4 + \frac{1}{2} dp(\text{state}, 2) \right) \right)$$

$$= \max \left(0.8(-8 + 2.375) + 0.2(-8 + 1.75), 4 + 1.75 \right)$$

$$dp(\text{state}, 1) = 5.75$$

$$\pi(\text{state}, 1) = \text{silent}$$

$$dp(\text{fresh}, 1) = \max_{\text{S1}, \text{S2}} \left(0.9 \left(-4 + \frac{1}{2} dp(\text{fresh}, 2) \right) + 0.1 \left(-4 + \frac{1}{2} dp(\text{state}, 2) \right) \right)$$

$$0.5 \left(4 + \frac{1}{2} dp(\text{fresh}, 2) \right) + 0.5 \left(4 + \frac{1}{2} dp(\text{state}, 2) \right)$$

$$= \max \left(0.9(-4 + 2.375) + 0.1(-4 + 1.75), 0.5(4 + 2.375) + 0.5(4 + 1.75) \right)$$

$$dp(\text{fresh}, 1) = 6.0625$$

$$\pi(\text{fresh}, 1) = \text{silent}.$$

Thus the value function for each start state :-

$$v(\text{fresh}) = dp(\text{fresh}, \perp) = 6.0625$$

$$v(\text{stale}) = dp(\text{stale}, \perp) = 5.75$$

π :-

at each time step , and at each state always choose "silent".

i.e.

$$\pi(\text{fresh}) = \text{silent}$$

$$\pi(\text{stale}) = \text{silent} .$$

(C) When the server optimises costs over an infinite interval :- the task becomes a continuing task

$$g_t = \sum_{n=0}^{\infty} R_{t+n} \pi_n \cdot \gamma^n \quad (\gamma = 1/2)$$

Assuming an initial policy π_0 :-

$$\pi_0(\text{silent | fresh}) = \pi_0(\text{query | fresh}) = 0.5$$

$$\pi_0(\text{silent | stale}) = \pi_0(\text{query | stale}) = 0.5$$

~~Policy iteration~~ Policy iteration:

(itr 1) Policy Evaluation :

$$v_{\pi_0}(\text{stale}) = 0.5 \left(4 + \frac{1}{2} v_{\pi_0}(\text{stale}) \right) + 0$$

$$0.5 \left(0.2 (-8 + \frac{1}{2} v_{\pi_0}(\text{stale})) + 0.8 (-8 + \frac{1}{2} v_{\pi_0}(\text{fresh})) \right)$$

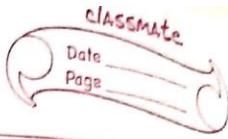
$$v_{\pi_0}(\text{fresh}) = 0.5 \left(0.5 (4 + \frac{1}{2} v_{\pi_0}(\text{fresh})) + 0.5 (4 + \frac{1}{2} v_{\pi_0}(\text{stale})) \right) +$$

$$0.5 \left(0.9 (-4 + \frac{1}{2} v_{\pi_0}(\text{fresh})) + 0.1 (-4 + \frac{1}{2} v_{\pi_0}(\text{stale})) \right)$$

on simplifying :-

$$0.49 v_{\pi_0}(\text{stale}) - 0.4 v_{\pi_0}(\text{fresh}) = \underline{3.2} \rightarrow ①$$

$$0.05 v_{\pi_0}(\text{stale}) - 0.55 v_{\pi_0}(\text{fresh}) = \underline{2.4} \rightarrow ②$$



on solving :

$$v_{\pi_0}(\text{stale}) = 0.64$$

$$v_{\pi_0}(\text{fresh}) = -7.21$$

Policy update :

$$\pi_1(\text{fresh}) = \underset{\{Q, S\}}{\operatorname{argmax}} \left(\begin{array}{l} \cancel{0.5} \\ 0.9(-4 + \cancel{0.5}) + \\ 0.1(-4 + \cancel{0.64}) \end{array} \right),$$

greater → $0.5(4 - \cancel{7.21}) +$
 $0.5(4 + \cancel{0.64})$

$$\pi_1(\text{fresh}) = \text{silent}$$

$$\pi_1(\text{stale}) = \underset{\{Q, S\}}{\operatorname{argmax}} \left(\begin{array}{l} 0.2(-8 + \overset{0.5}{\cancel{(0.64)})} + \\ 0.8(-8 - \cancel{7.21}) \end{array} \right)$$

$$\text{greater } \rightarrow (4 + (0.64)(0.5))$$

$$\pi_1(\text{stale}) = \text{silent}$$

Iteration 2: Policy evaluation:

$$v_{\pi_1}(\text{stale}) = 4 + 0.5(v_{\pi_1}(\text{stale}))$$

$$v_{\pi_1}(\text{fresh}) = 0.5(4 + 0.5(v_{\pi_1}(\text{fresh}))) \\ + 0.5(4 + 0.5 v_{\pi_1}(\text{stale}))$$

$$v_{\pi_1}(\text{stale}) = 8$$

$$v_{\pi_1}(\text{fresh}) = 8$$

Policy update:

$$\pi_2(\text{stale}) = \underset{\{Q, S\}}{\operatorname{argmax}} \left(0.2(-8 + \frac{1}{2}(8)) + 0.8(-8 + \frac{1}{2}(8)) \right)$$

$$\text{greater } \xrightarrow{1} 4 + \frac{1}{2}(8) \xrightarrow{2}$$

$$\pi_2(\text{stale}) = \text{silent}$$

$$\pi_2(\text{fresh}) = \underset{\{Q, S\}}{\operatorname{argmax}} \left(0.9(-4 + \frac{1}{2}(8)) + 0.1(-4 + \frac{1}{2}(8)) \right)$$

$$\text{greater } \xrightarrow{1} 0.5(4 + \frac{1}{2}(8)) \xrightarrow{2} 0.5(4 + \frac{1}{2}(8))$$

$$\therefore \pi_2(\text{fresh}) = \text{silent}$$

$$\pi_1 = \pi_2 \therefore \text{Terminate.}$$

$$\text{Optimal policy } \rightarrow \pi^*(\text{fresh}) = \pi^*(\text{stale}) = \text{silent}$$

$$\text{Value function } \rightarrow V^*(\text{fresh}) = V^*(\text{stale}) = 8$$

Value iteration?

initialise $v_0(\text{fresh}) = v(\text{stale}) = 0$

Iteration 1:

$$v_1(\text{fresh}) = \max_{\{-8, 5, 3\}} \left(0.9 \left(-4 + \frac{1}{2}(0) \right) + 0.1 \left(-4 + \frac{1}{2}(0) \right) \right)$$

$$= 0.5 \left(4 + \frac{1}{2}(0) \right) + 0.5 \left(4 + \frac{1}{2}(0) \right)$$

$$= \max_{\{-8, 5, 3\}} (-4, 4) = 4$$

$$v_1(\text{stale}) = \max_{\{-8, 5, 3\}} \left(0.2 \left(-8 + \frac{1}{2}(0) \right) + 0.8 \left(-8 + \frac{1}{2}(6) \right) \right)$$

$$= \max_{\{-8, 5, 3\}} (-8, 4) = 4$$

$$v_1(\text{fresh}) = v_1(\text{stale}) = 4$$

Iteration 2:

$$v_2(f) = \max_{\{-8, 5, 3\}} \left(0.9 \left(-4 + \frac{1}{2}(4) \right) + 0.1 \left(-4 + \frac{1}{2}(4) \right) \right)$$

$$= \max_{\{-8, 5, 3\}} \left(0.5 \left(4 + \frac{1}{2}(4) \right) + 0.5 \left(4 + \frac{1}{2}(4) \right) \right)$$

$$= \max_{\{-8, 5, 3\}} (-2, 6) = 6$$

v_1

v_1

$$v_2(\text{stale}) = \max_{\in Q, S_2} \left(0.2(-8 + \frac{1}{2}(4)) + 0.8(-8 + \frac{1}{2}(4)), 4 + \frac{1}{2}(4) \right),$$

$$= \max_{\in Q, S_2} (-6, 6) = 6$$

$$\therefore v_2(\text{fresh}) = v_2(\text{stale}) = 6$$

Iteration 3:

v_2

$$v_3(\text{fresh}) = \max_{\in Q, S_3} \left(0.9(-4 + \frac{1}{2}(6)) + 0.1(-4 + \frac{1}{2}(6)), 0.5(4 + \frac{1}{2}(6)) + 0.5(4 + \frac{1}{2}(6)) \right)$$

$$v_3(\text{stale}) = \max_{\in Q, S_3} \left(0.2(-8 + \frac{1}{2}(6)) + 0.8(-8 + \frac{1}{2}(6)), 4 + \frac{1}{2}(6) \right),$$

$$= 7$$

$$v_3(\text{fresh}) = v_3(\text{stale}) = 7.$$

Iteration 4:

v_3

$$v_4(\text{fresh}) = \max_{\in Q, S_4} \left(0.9(-4 + \frac{1}{2}(7)) + 0.1(-4 + \frac{1}{2}(7)), 0.5(4 + \frac{1}{2}(7)) + 0.5(4 + \frac{1}{2}(7)) \right),$$

$$= 7.5$$

$$v_4(\text{stale}) = \max_{\in Q, S_4} \left(0.2(-8 + \frac{1}{2}(7)) + 0.8(-8 + \frac{1}{2}(7)), 4 + \frac{1}{2}(7) \right),$$

$$= 7.5$$

$$v_4(\text{fresh}) = v_4(\text{stale}) = 7.5$$

Q14.

Q14. To prove:

1) Policy improvement step either improves the policy or the current policy is optimal.

Equivalently we prove:

- ① Policy improvement either improves the policy or leaves it unchanged.
- ② If the policy is unchanged, then current policy is optimal.

Proof for ①:

From policy evaluation

$$v_{\pi_k}(s) = T_{\pi_k}(v_{\pi_k}(s))$$

After policy improvement. $T_{\pi_{k+1}} v_{\pi_k}(s) = T v_{\pi_k}$

$$T_{\pi_{k+1}}(v_{\pi_k}(s)) \geq v_{\pi_k}(s)$$

$$(i.e. q_{\pi_k}(s, \pi_{k+1}(s)) \geq q_{\pi_k}(s, \pi_k(s)) = v_{\pi_k}(s))$$

since $T_{\pi_{k+1}}$ is monotonically increasing

$$T_{\pi_{k+1}}^2(v_{\pi_k}(s)) \geq T_{\pi_{k+1}}(v_{\pi_k}(s)) \geq v_{\pi_k}(s)$$

⋮

$$T_{\pi_{k+1}}^n(v_{\pi_k}(s)) \geq v_{\pi_k}(s)$$

$$\text{as } N \rightarrow \infty \quad T_{\pi_{k+1}}^N(v_{\pi_k}(s)) = v_{\pi_{k+1}}$$

(i.e. $\because v_{\pi_{k+1}}$ is a stationary pt. of $T_{\pi_{k+1}}$)

$$\Rightarrow \lim_{N \rightarrow \infty} T_{\pi_{k+1}}^N(v_{\pi_k}(s)) \geq v_{\pi_k}(s)$$

$$\Rightarrow v_{\pi_{k+1}}(s) \geq v_{\pi_k}(s) \quad \forall s$$

$\therefore \pi_{k+1}$ (new policy) improves the value function
 or leaves it unchanged.

Proof for ②

for policy improvement

$$T_{\pi_{k+1}}(v_{\pi_k}(s)) = T(v_{\pi_k}(s))$$

$$(i.e. \pi_{k+1}(s) = \underset{a \in A(s)}{\operatorname{argmax}} E[R_{t+1} + \gamma v_{\pi}(s') | s_t = s, a_t = a])$$

~~Suppose~~ policy remains unchanged

$$\therefore i.e. v_{\pi_k}(s) = v_{\pi_{k+1}}(s)$$

$$\Rightarrow T_{\pi_{k+1}}(v_{\pi_k}(s)) = T_{\pi_{k+1}}(v_{\pi_{k+1}}(s))$$

(stationary pt.)

$$= v_{\pi_{k+1}}(s)$$

$$= v_{\pi_k}(s)$$

Also

$$T_{\pi_{K+1}}(v_{\pi_K}(s)) = T(v_{\pi_K}(s))$$

$$v_{\pi_K}(s) = T(v_{\pi_K}(s))$$

i.e. v_{π_K} is a stationary pt for T

i.e.

$$v_{\pi_K}(s) = \max_{a \in A(s)} \sum_{s', r} p(s', r | s, a) (r + \gamma v_{\pi_K}(s'))$$

which is the bellman optimality condition.

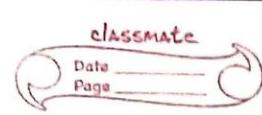
$\therefore \pi_K$ is the optimal policy

($\because v_{\pi_K}$ is the optimal value function)

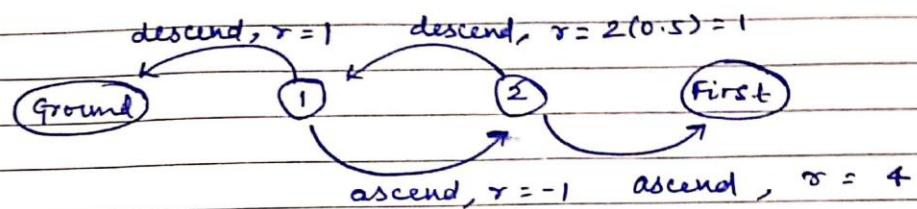
$$\therefore \pi^* = \pi_K$$

$$v^* = v_{\pi_K}$$

Q15.



Q15. for $n = 2$
the MDP is given by:



Also the robot ascends, descends with equal prob = 0.5

$$\therefore \pi_0^*(\text{ascend} \mid S=1) = 0.5$$

$$\pi_0^*(\text{descend} \mid S=1) = 0.5$$

$$\pi_0^*(\text{ascend} \mid S=2) = 0.5$$

$$\pi_0^*(\text{descend} \mid S=2) = 0.5$$

Policy iteration:

Iteration 1: Policy evaluation:

$$v_{\pi_0}(\text{Ground}) = v_{\pi_0}(\text{First}) = 0 \quad (\text{terminal states})$$

$$v_{\pi_0}(1) = \pi(\text{descend} \mid 1) \cdot (1 + v_{\pi_0}(\text{Ground})) \\ + \pi(\text{ascend} \mid 1) \cdot (-1 + v_{\pi_0}(2))$$

$$v_{\pi_0}(2) = \pi(\text{descend} \mid 2) \cdot (1 + v_{\pi_0}(1)) \\ + \pi(\text{ascend} \mid 2) \cdot (4 + v_{\pi_0}(\text{First}))$$

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

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

Solving $v_{\pi_0}(1)$ and $v_{\pi_0}(2)$

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

$$v_{\pi_0}(2) = 2 - 5 + 0.5 v_{\pi_0}(1)$$

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

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

Policy improvement :

$$\pi_1(1) = \underset{\text{ascend, descend 3-}}{\operatorname{argmax}} \left(1, \underbrace{-1 + v_{\pi_0}(2)}_{\frac{7}{3}} \right)$$

= ascend

$$\pi_1(2) = \underset{\text{descend, ascend 3-}}{\operatorname{argmax}} \left(1 + v_{\pi_0}(1), \underbrace{4}_{\frac{8}{3}} \right)$$

= ascend

iteration 2 : Policy eval.

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

$$v_{\pi_1}(2) = 4$$

$$\Rightarrow v_{\pi_1}(1) = 3$$

Policy update :

$$\pi_2(1) = \underset{\in \{\text{ascend, descend}\}}{\operatorname{argmax}} (1, -1 + 4)$$

= ascend

$$\pi_2(2) = \underset{\in \{\text{descend, ascend}\}}{\operatorname{argmax}} (v_1(1) + 4)$$

= ~~descend~~ ∈ ascend, descend 3

we pick ascend

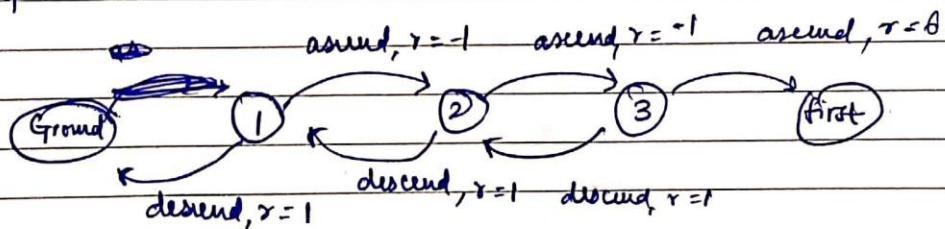
* Note if we pick descend, then next time

We'll again have a choice, this way $\pi_{n+1} \neq \pi_n$
and policy iteration would never end.

The no. of iterations would always be = 2,

for any n. at ~~for the last step~~

Proof : Consider n = 3



assuming $\pi(\text{asc} | s) = 0.5$
 $\pi(\text{desc} | s) = 0.5 \quad \text{for } s \in \{1, 2, 3\}$

iteration 1 : Policy eval

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

$$v_{\pi_0}(2) = 0.5(1 + v_{\pi_0}(1)) + 0.5(-1 + v_{\pi_0}(3))$$

$$v_{\pi_0}(3) = 0.5(1 + v_{\pi_0}(2)) + 0.5(6)$$

on solving:

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

$$v_{\pi_0}(2) = 0.5 v_{\pi_0}(1) + 0.5 v_{\pi_0}(3)$$

$$v_{\pi_0}(3) = 3.5 + 0.5 v_{\pi_0}(2)$$

$$v_{\pi_0}(1) = 7/4$$

$$v_{\pi_0}(2) = 7/2$$

$$v_{\pi_0}(3) = 21/4$$

Policy update

$$\pi_1(1) = \underset{\in \{d, a_3\}}{\operatorname{argmax}} (1, -1 + 7/2)$$

= ascend

$$\pi_1(2) = \underset{\in \{d, a_3\}}{\operatorname{argmax}} (1 + \frac{7}{4}, -1 + \frac{7}{2})$$

= Ascend

$$\pi_1(3) = \underset{\in \{d, a_3\}}{\operatorname{argmax}} (1 + \frac{7}{2}, 6)$$

= descend

Iteration 2: Policy eval

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

$$v_{\pi_1}(2) = -1 + v_{\pi_1}(3)$$

$$v_{\pi_1}(3) = 6$$

\Rightarrow

$$v_{\pi_2}(1) = 4$$

$$v_{\pi_2}(2) = 5$$

$$v_{\pi_2}(3) = 6$$

Policy update

$$\pi_1(1) = \underset{\{d, a\}}{\operatorname{argmax}} (1, -1 + 5)$$

= ascend

$$\pi_2(2) = \underset{\{d, a\}}{\operatorname{argmax}} (1 + 4, -1 + 6)$$

* Can pick both, same as
 $n = 2$; we pick ascend for termination

$$\pi_2(3) = \underset{\{d, a\}}{\operatorname{argmax}} (1 + 5, 6)$$

~~↑~~
 Same for $s = 3$ = ascend

Thus it ends after two iterations.

This can be generalised for any n where always going up is optimal.