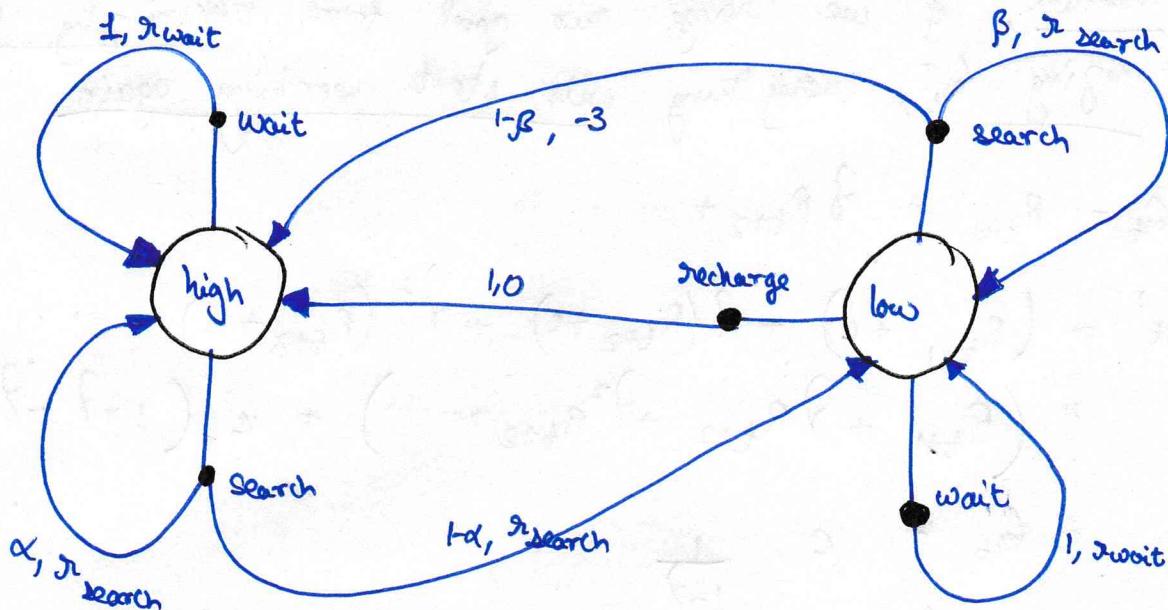


Seeta Gupta  
2018190

# Reinforcement Learning

Dynamic Programming  
ECE / CSE 564

Monsoon 20



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

This table formulates every thing that can happen in the environment.

Choose a  $(s, a)$  pair and then write down all possible  $(s', r)$  for it along with probabilities.

Q-3)

3.15:

Signs are important. Suppose the signs of every reward is flipped and nothing else is changed, then the agent will be motivated to hit the walls.

However, if we change our goal from maximizing  $G_t$  to minimizing  $G_t'$ , everything will start working again.

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

$$\begin{aligned} G_t' &= (R_{t+1} + c) + \gamma (R_{t+2} + c) + \gamma^2 (R_{t+3} + c) + \dots \\ &= (R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots) + c (1 + \gamma + \gamma^2 + \dots) \\ &= G_t + c \cdot \frac{1}{1-\gamma} \\ &= G_t + \text{c. } \frac{c}{1-\gamma} \rightarrow \text{a constant} \end{aligned}$$

$$V(\pi) = E_{\pi}[G_t | S_t = s]$$

$$V_{\pi'}(s) = E_{\pi}[G_t' | S_t = s]$$

$$= E_{\pi}[G_t + \frac{c}{1-\gamma} | S_t = s]$$

$$= E_{\pi}[G_t | S_t = s] + \frac{c}{1-\gamma} \quad (\text{as } \frac{c}{1-\gamma} \text{ is a constant})$$

and hence we can take it out of expectation)

$$\therefore V_c = \frac{c}{1-\gamma}$$

3.16

### Case 1: Undiscounted:

Original Task: Agent gets  $-1$  reward for every timestep that it stays in the maze.

It gets  $+1$  reward when it gets out of the maze.

The  $-1$  reward motivates agent to get out as quickly as possible

Modified task:  $c$  is added to all rewards. The agent gets  $c-1$  for every step it stays in the maze &  $c+1$  for exiting the maze.

If  $c-1 > 0$  i.e.  $c > 1$ , then the agent will be motivated to stay in the maze indefinitely and accumulates  $(c-1)$  rewards at each time step. The rewards for exiting doesn't change with time so agent has no motivation to exit early.  
Hence the agent never exits.

### Case 2: Discounted.

Original Task: Agent gets  $0$  reward for each time step it stays in the maze & gets reward  $+1$  discounted by  $\gamma$  at each time step upon exiting the maze

Modified Task: Agent gets  $c$  rewards for each time step it stays in the maze &  $(c+1)$  when it exits the maze discounted by  $\gamma$ .

In this case the agent has some motivation to exit the maze early because of  $\gamma$ .

However if  $c$  is high enough (example 1000), then it might be more beneficial to stay.

Example

Suppose the agent starts at position  $p$  which is at  $d$  distance away from exit. The agent can only move 1 unit distance every step.

If the agent exits as fast as possible, then

$$G_t = 1000 + \gamma 1000 + \gamma^2 1000 + \dots + \gamma^{d-1} 1000 + \gamma^d 1000$$

$$= 1000(1 + \gamma + \gamma^2 + \dots + \gamma^{d-1}) + \gamma^d 1000$$

$$= \frac{1000(1 - \gamma^{d+1})}{1 - \gamma} + \gamma^d 1000$$

If agent stays in the maze for infinite time and never exits, then

$$G_t = 1000 + \gamma 1000 + \gamma^2 1000 + \dots$$

$$= \frac{1000}{1 - \gamma}$$

Now if

$$\frac{1000(1 - \gamma^{d+1})}{1 - \gamma} + \gamma^d 1000 < \frac{1000}{1 - \gamma}$$

$$\begin{cases} \Rightarrow 1 - \gamma < \gamma \\ \Rightarrow 2\gamma > 1 \\ \Rightarrow \gamma > \frac{1}{2} \end{cases}$$

$$\Rightarrow \frac{\gamma^{d+1}}{1 - \gamma} + \gamma^d < 0$$

$$\Rightarrow \gamma^d < \frac{\gamma^{d+1}}{1 - \gamma}$$

$$\Rightarrow \cancel{\gamma^d} < \gamma^{d+1} \left(\frac{\gamma}{1 - \gamma}\right)$$

$$\Rightarrow 1 < \frac{\gamma}{1 - \gamma}$$

Then agent will be motivated to stay in the maze forever.

$$Q-5) \quad V_*(s) = \max_{a \in A(s)} q_*(s, a)$$

$$Q-8) \quad p(s', a | s, a) = \Pr [S_{t+1} = s', R_{t+1} = a | S_t = s, A_t = a]$$

$$\pi(a | s, a) = \Pr [R_{t+1} = a | S_t = s, A_t = a]$$

$$= \sum_{s' \in S} \Pr [S_{t+1} = s', R_{t+1} = a | S_t = s, A_t = a]$$

$$= \sum_{s' \in S} p(s', a | s, a) \quad \text{--- } ①$$

$\tilde{v}_\pi(s') | s, a$  is the probability of getting reward  $R_{t+2} = a'$

starting from  $S_t = s$  and taking action  $A_t = a$  under policy  $\pi$ .

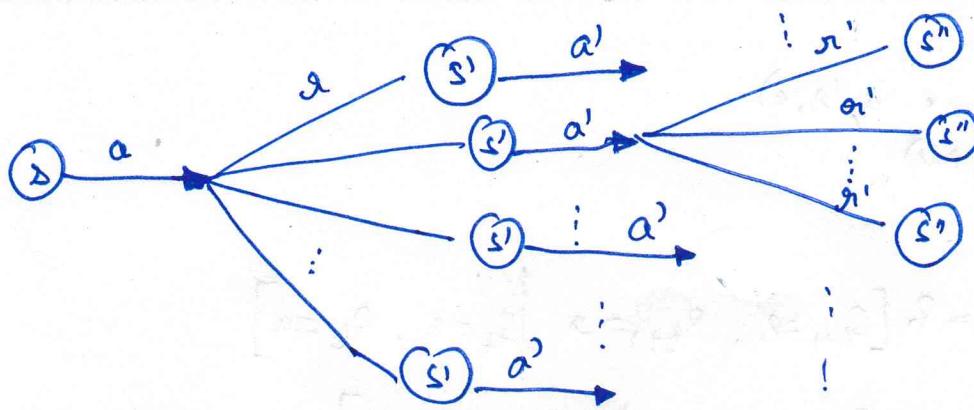
$$\tilde{v}_\pi(s') | s, a = \Pr [R_{t+2} = a' | S_t = s, A_t = a]$$

$$= \sum_{s' \in S} \Pr [S_{t+1} = s' | S_t = s, A_t = a] \times \sum_{a' \in A(s')} \Pr [A_{t+1} = a' | S_{t+1} = s', S_t = s, A_t = a]$$

$$\times \Pr [R_{t+2} = a' | S_{t+1} = s', A_{t+1} = a', S_t = s, A_t = a] \quad ②$$

Probability of getting reward at  $t+2$ ,  $R_{t+2} = a'$  is

probability of landing at  $s'$   $\times$  taking an action  $a'$  at  $s'$   $\times$  getting a reward  $a'$  due to  $a'$ .



$$\Pr [S_{t+1} = s' \mid S_t = s, A_t = a]$$

$$= \sum_{a \in R} \Pr [S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a]$$

$$= \sum_{a \in R} p(s', r \mid s, a)$$

— ③

$$\Pr [A_{t+1} = a' \mid S_{t+1} = s', S_t = s, A_t = a] = \Pr [A_{t+1} = a' \mid S_{t+1} = s'] \\ = \pi(a' \mid s) \quad — ④$$

(Because  $\pi$  only depends on current state)

$$\Pr [R_{t+2} = r' \mid S_{t+1} = s', A_{t+1} = a', S_t = s, A_t = a]$$

$$= \sum_{s'' \in S} \Pr [S_{t+2} = s'', R_{t+2} = r' \mid S_{t+1} = s', A_{t+1} = a', S_t = s, A_t = a]$$

$$= \sum_{s'' \in S} \Pr [S_{t+2} = s', R_{t+2} = r' \mid S_{t+1} = s', A_{t+1} = a'] \quad (\text{Markov property})$$

$$= \sum_{s'' \in S} p(s'', r' \mid s', a) \quad — ⑤$$

∴ By ② ③, ④ and ⑤

$$\tilde{\pi}_\pi(a'|s, a) = \sum_{\substack{s' \in S \\ a \in R}} \phi(s', a | s, a) \sum_{\substack{a' \in A(s) \\ s'' \in S}} \pi(a' | s) \cdot \phi(s'', a' | s', a')$$

∴  $\tilde{\pi}_\pi$  depends on  $s, a$  and  $\pi$ .

Q-9)  $E_\pi[R_{t+2} | S_t = s, A_t = a]$  (Expectation is dependent on policy  $\pi$ )

$$= \sum_{g \in R} g \cdot \tilde{\pi}_\pi(g | s, a)$$

$$= \sum_{\substack{s' \in R \\ g \in R}} g \cdot \sum_{\substack{s' \in S \\ a \in R}} \phi(s', g | s, a) \sum_{\substack{a' \in A(s) \\ s'' \in S}} \pi(a' | s) \cdot \phi(s'', a' | s', a)$$

Q-10)  $v_\pi(s) = E_\pi[G_t | S_t = s]$

$$= E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= E_\pi[R_{t+1} | S_t = s] + \gamma E_\pi[G_{t+1} | S_t = s]$$

$$= E_\pi[R_{t+1} | S_t = s] + \gamma E_\pi[E_\pi[G_{t+1} | S_{t+1} = s'] | S_t = s]$$

$G_{t+1}$  only depends on  $S_{t+1}$  &  $\pi$  not  $S_t$

$$\therefore = E_\pi[R_{t+1} | S_t = s] + \gamma E_\pi[E_\pi[G_{t+1} | S_{t+1} = s'] | S_t = s]$$

$$= E_\pi[R_{t+1} | S_t = s] + \gamma E_\pi[v_\pi(S_{t+1}) | S_t = s]$$

$$= E_\pi[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]$$

$$= \sum_{\substack{a \in R \\ s' \in S \\ a \in \pi(s')}} p(a, s' | s, a) \cdot p_\pi(a | s) \cdot [a + \gamma v_\pi(s')]$$

Q-11.)

$$G_4 = 0$$

Part 1

$$\begin{aligned} G_3 &= R_4 + \gamma G_4 \\ &= -3 + (0.5) 0 \\ &= -3 \end{aligned}$$

$$\begin{aligned} G_2 &= R_3 + \gamma G_3 \\ &= 10 + (0.5)(-3) \\ &= 10 - 1.5 \\ &= 8.5 \end{aligned}$$

$$\begin{aligned} G_1 &= R_2 + \gamma G_2 \\ &= -1 + 0.5(8.5) \\ &= -1 + 4.25 \\ &= 3.25 \end{aligned}$$

$$\begin{aligned} G_0 &= R_1 + \gamma G_1 \\ &= 2 + 0.5(3.25) \\ &= 2 + 1.625 \\ &= 3.625 \end{aligned}$$

Part 2:

$$R_t = c + r_t$$

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= c + \gamma c + \gamma^2 c + \gamma^3 c + \dots \\ &= c (1 + \gamma + \gamma^2 + \gamma^3 + \dots) \\ &= c \cdot \frac{1}{1-\gamma} \\ &= \frac{c}{1-\gamma} \end{aligned}$$

$$Q-12) \quad \pi_*(s) = \operatorname{argmax}_{a \in A(s)} \left\{ p(s'|s, a) [r(s', a, s) + \gamma r_*(s')] \right\}$$

where :

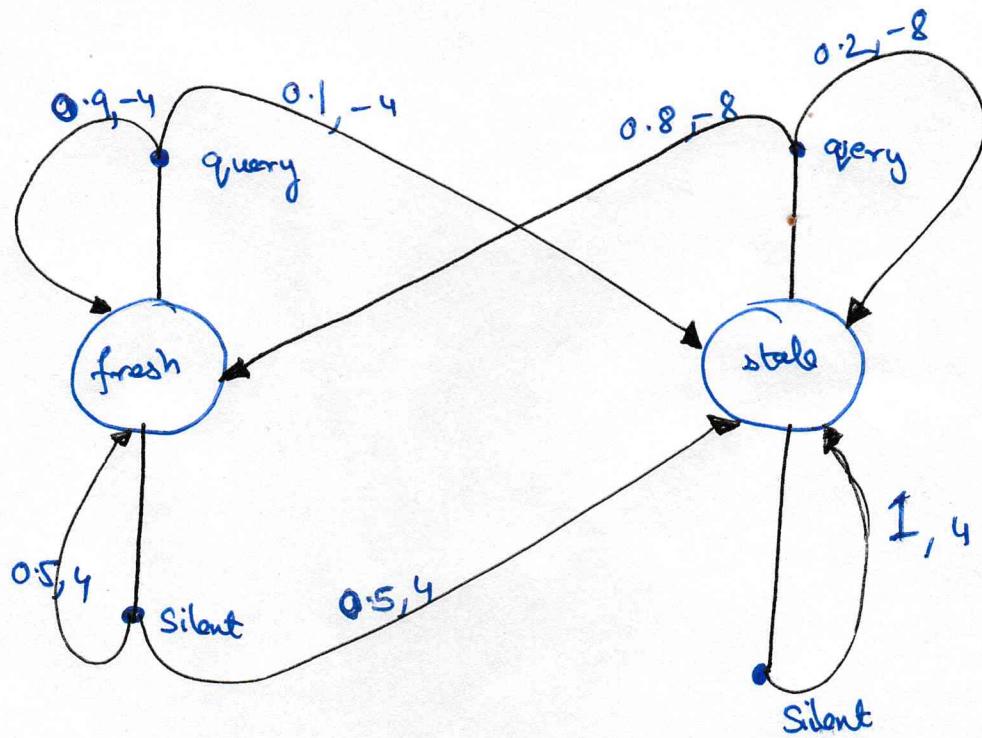
$p(s'|s, a)$  = probability of landing at  $s'$  if I take action  $a$  at  $s$ .

$r(s', a, s)$  = expected reward while going to  $s'$  from  $s$  via  $a$ .

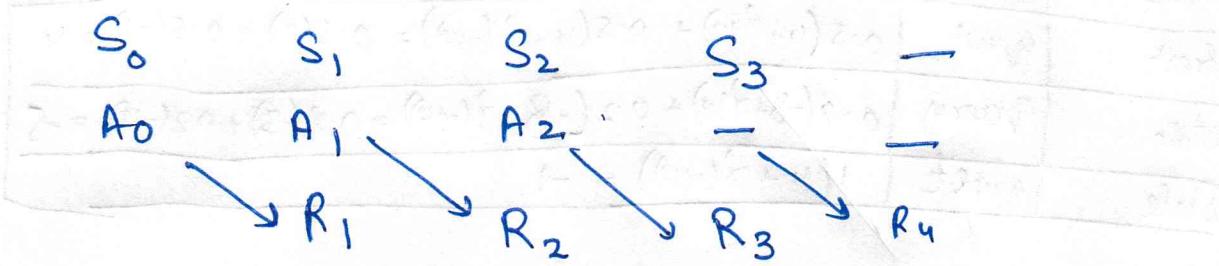
$$p(s'|s, a) = \sum_{a' \in R} p(s'|s, a)$$

$$r(s, a, s) = \sum_{a' \in R} r_{a'} \cdot \frac{p(s'|s, a)}{p(s'|s, a)}$$

Q-13.) a)



b) Let's try to back calculate. The sequence is

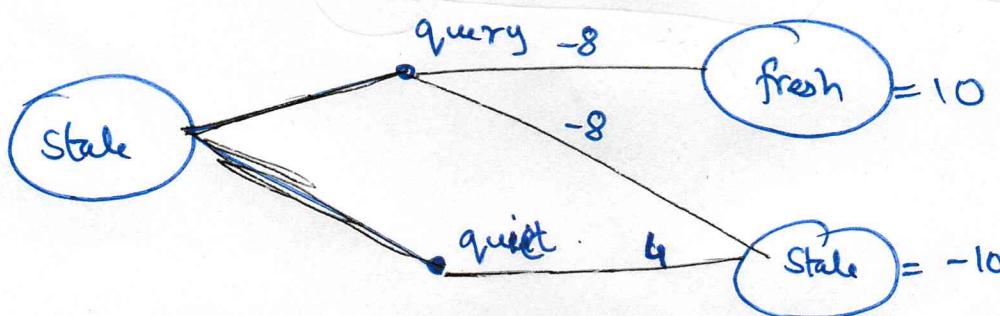
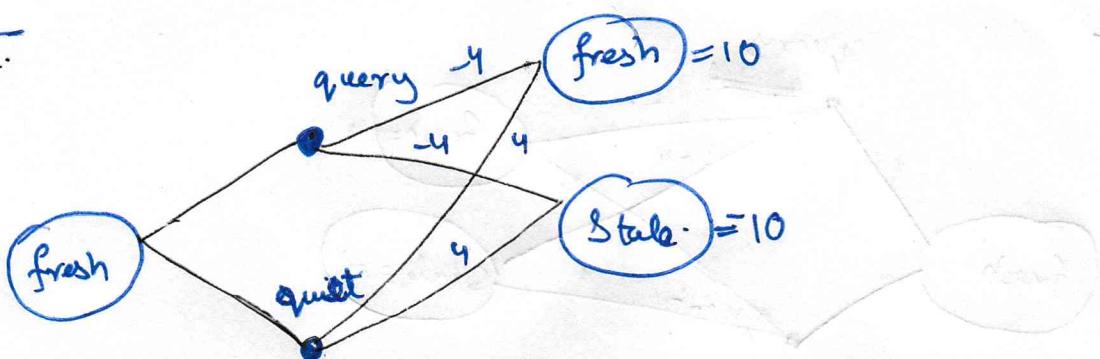


$S_3$ :

$$\text{If } S_3 = \text{fresh}, \quad v^*(S_3 = \text{fresh}) = 10$$

$$S_3 = \text{stale}, \quad v^*(S_3 = \text{stale}) = -10$$

$S_2$ :



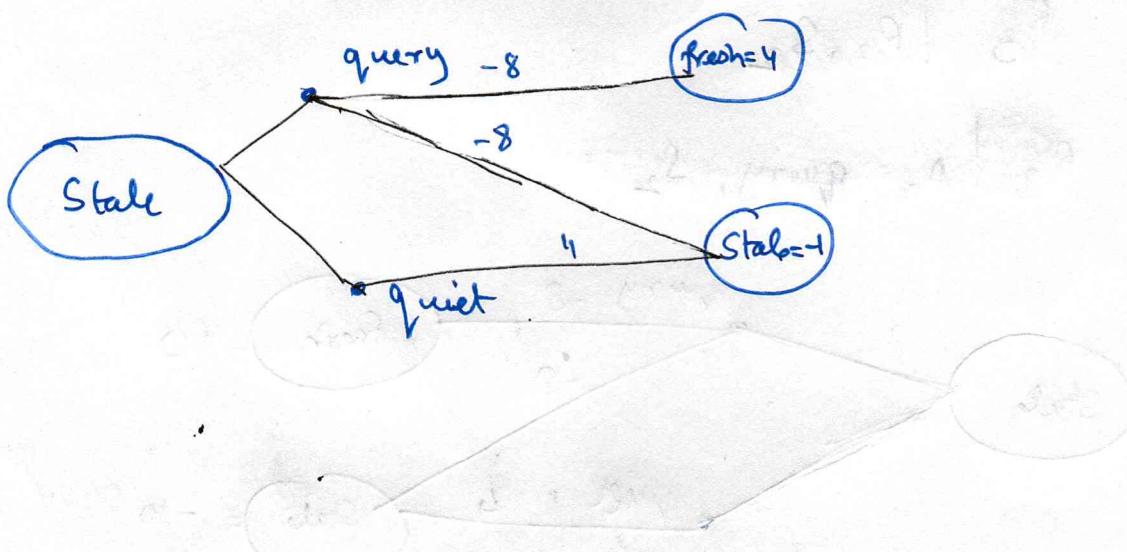
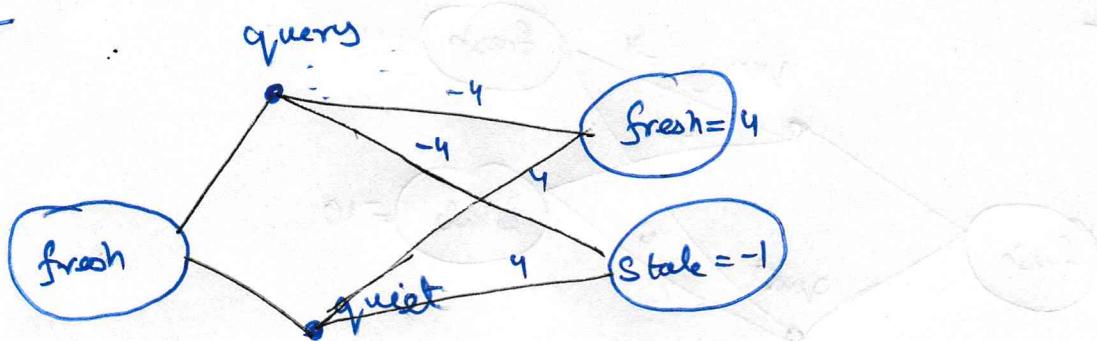
$S_2$	$A_2$	$B_2$
fresh	query	$0.9(-4+7(10)) + 0.1(-4+7(-10)) = 0.9(1) + 0.1(-9) = 0$
fresh	quiet	$0.5(4+7(10)) + 0.5(4+7(-10)) = 0.5(9) + 0.5(-1) = 4$
stale	query	$0.8(-8+7(10)) + 0.2(-8+7(-10)) = 0.8(-3) + 0.2(-13) = -5$
stale	quiet	$1(4+7(-10)) = -1$

$\therefore$  At  $S_2 = \text{fresh}$  choose  $A_2 = \text{quiet}$

$S_2 = \text{stale}$  choose  $A_2 = \text{quiet}$

$$v_* (S_2 = \text{fresh}) = 4 \quad v_* (S_2 = \text{stale}) = -1$$

$S_1$ :



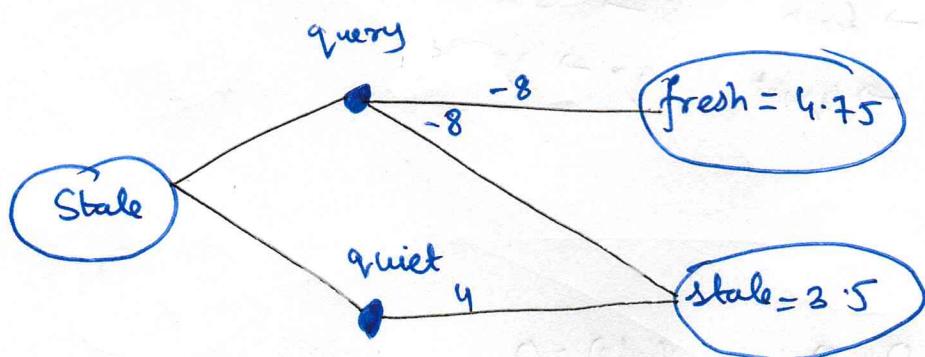
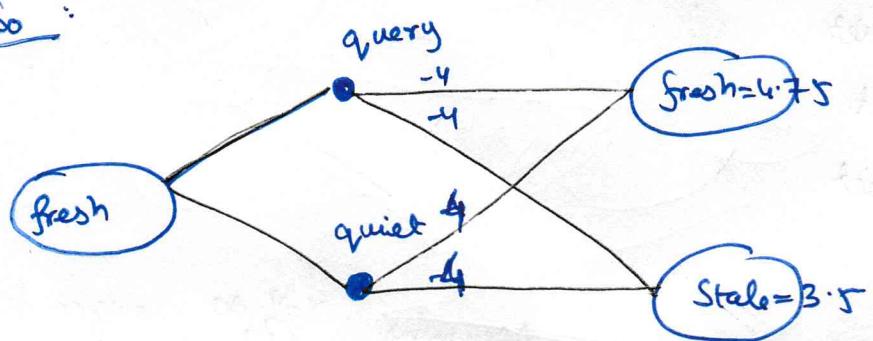
$S_1$	$A_1$	$G_1$
fresh	query	$0 \cdot 9(-4 + 7(4)) + 0 \cdot 1(-4 + 7(-1)) = 0 \cdot 9(-2) + 0 \cdot 1(-4 \cdot 5) = -2.25$
fresh	quiet	$0 \cdot 5(4 + 7(4)) + 0 \cdot 5(4 + 7(-1)) = 0 \cdot 5(6) + 0 \cdot 5(3 \cdot 5) = 4.75$
stale	query	$0 \cdot 8(-8 + 7(4)) + 0 \cdot 2(-8 + 7(-1)) = 0 \cdot 8(-6) + 0 \cdot 2(-8 \cdot 5) = -6.5$
stale	quiet	$1(4 + 7(-1)) = 3.5$

$\therefore$  At  $S_1 = \text{fresh}$  choose  $A_1 = \text{quiet}$

$S_2 = \text{stale}$  choose  $A_1 = \text{quiet}$

$$v_* (S_1 = \text{fresh}) = 4.75 \quad v_* (S_2 = \text{stale}) = 3.5$$

So:



$S_0$	$A_0$	$G_0$
fresh	query	$0.9(-4 + \gamma(4.75)) + 0.1(-4 + \gamma(3.5)) = -1.6875$
fresh	quiet	$0.5(-4 + \gamma(4.75)) + 0.5(-4 + \gamma(3.5)) = 6.0625$
stale	query	$0.8(-8 + \gamma(4.75)) + 0.2(-8 + \gamma(3.5)) = -5.75$
stale	quiet	$1(-4 + \gamma(3.5)) = 5.75$

$\therefore$  Take quiet action for  $A_0$

$$\therefore V_*(S_0 = \text{fresh}) = 6.0625$$

$$V_*(S_0 = \text{stale}) = 5.75$$

Optimal policy:

$\pi_0 : A_0 = \text{quiet}$

$A_1 : \text{quiet}$

$A_2 : \text{quiet}$

c)

Value iterations:

Let  $f \rightarrow \text{fresh}$      $s \rightarrow \text{stale}$   
 $q_t \rightarrow \text{quiet}$      $q_r \rightarrow \text{query}$ .

Assuming  $\gamma = 0.9$ , horizon problem.

Initialization:

$$V(f) = 0, \quad V(s) = 0$$

$$V_{R+1}(s) = \max_a \sum_{s', r} p(s', r | s, a) (r + \gamma V_R(s'))$$

Iter ①:

for  $f$ :

$$v(f) = \max \left\{ 0.9(-4 + v_2(0)) + 0.1(-4 + v_2(0)), \right. \\ \left. 0.5(4 + v_2(0)) + 0.5(4 + v_2(0)) \right\} \\ = \max \{-4, 4\} = 4$$

$$v(s) = \max \left\{ -8 + 0.8(0) + 0.2(0), \right. \\ \left. 4 + 1(0) \right\} \\ = \max \{-8, 4\} = 4$$

Iter ②

$$v(f) = \max \left\{ -4 + (0.9)(v_2(6)) + (0.1)(v_2(6)), \right. \\ \left. 4 + (0.5)(v_2(6)) + (0.5)(v_2(6)) \right\} \\ = \max \{-4+2, 4+2\} = 6$$

$$v(s) = \max \left\{ -8 + (v_2)(0.8(6) + 0.2(6)), \right. \\ \left. 4 + (v_2)(1)(6) \right\} \\ = \max \{-8+2, 4+2\} = 6$$

Iter ③

$$v(f) = \max \left\{ -4 + (v_2)(0.9(6) + 0.1(6)), \right. \\ \left. 4 + (v_2)(0.5(6) + 0.5(6)) \right\} \\ = \max \{-4+3, 4+3\} = 7$$

$$v(s) = \max \left\{ -8 + (v_2)(0.8(6) + 0.2(6)), \right. \\ \left. 4 + (v_2)(1)(6) \right\} = \{-8+3, 4+3\} = 7$$

Iter ④

$$\begin{aligned} v(f) &= \max \left\{ -4 + 1/2 \left[ (0.9)(7) + (0.1)(7) \right], \right. \\ &\quad \left. 4 + 1/2 \left[ (0.5)(7) + (0.5)(7) \right] \right\} \\ &= \max \{-4 + 3.5, 4 + 3.5\} = 7.5 \end{aligned}$$

$$\begin{aligned} v(s) &= \max \left\{ -8 + 1/2 \left( 0.8(7) + 0.2(7) \right), \right. \\ &\quad \left. 4 + 1/2 (1)(7) \right\} \\ &= \max \{-8 + 3.5, 4 + 3.5\} = 7.5 \end{aligned}$$

Policy Iteration:

$$\begin{aligned} v_{\pi}(s) &= \sum_{a \in A(s)} \sum_{s', r} \pi(a|s) \cdot p(s', r | s, a) \cdot (r + \gamma v_{\pi}(s')) \\ v_{\pi}(f) &= \pi(q_r|f) \left[ (0.9)(-4 + \gamma v_{\pi}(f)) + (0.1)(4 + \gamma v_{\pi}(f)) \right] + \\ &\quad \pi(q_t|f) \left[ (0.5)(4 + \gamma v_{\pi}(f)) + (0.5)(-4 + \gamma v_{\pi}(f)) \right] \\ &= \pi(q_r|f) (-4 + 0.45v_{\pi}(f) + 0.05v_{\pi}(s)) \\ &\quad + \pi(q_t|f) (4 + 0.25v_{\pi}(f) + 0.25v_{\pi}(s)) \quad \text{--- (A)} \end{aligned}$$

$$\begin{aligned} v_{\pi}(s) &= \pi(q_r|s) \left[ -8 + 0.4 v_{\pi}(f) + 0.1 v_{\pi}(s) \right] \\ &\quad + \pi(q_t|s) \left[ 4 + 0.5 v_{\pi}(s) \right] \quad \text{--- (B)} \end{aligned}$$

## Initialization

$$P_{q_t}(q_t | s) = 1 ; P_{q_t}(q_t | f) = 1 ; v_{\pi_0}(s) = 0$$

$$P_{q_t}(q_f | s) = 0 ; P_{q_t}(q_f | f) = 0 ; v_{\pi_0}(f) = 0$$

## Iter ① of policy iteration

$$v_{\pi_0}(s) = 4 + 0.25v_{\pi_0}(f) + 0.25v_{\pi_0}(s) \quad \text{--- (i)}$$

$$v_{\pi_0}(s) = 4 + 0.5v_{\pi_0}(s) \quad \text{--- (ii)}$$

Solving we get

$$v_{\pi_0}(f) = 8 \quad v_{\pi_0}(s) = 8$$

## Policy improvement :

$$\pi_1(f) = \operatorname{argmax} \{-4 + 0.5(0.9 \times 8 + 0.1 \times 8)\},$$

$$= \operatorname{argmax} \{0, 8\}$$

$$\therefore \pi_1(f) = q_t$$

$$\pi_1(s) = \operatorname{argmax} \{-8 + 0.5(0.8 \times 8 + 0.2 \times 8)\},$$

$$= \operatorname{argmax} \{4, 8\} = q_t.$$

Policy stable

Done!

Q-14) At policy improvement:

$$T_{\pi_{\tau_k}} v_{\pi_{\tau_k}} = T_{\pi_{\tau_k+1}} v_{\pi_{\tau_k}} \quad \text{--- (i)}$$

Also

$$T_{\pi_{\tau_k}} > v_{\pi_{\tau_k}} \quad (\text{As we take argmax}) \quad \text{--- (ii)}$$

∴ By (i) & (ii)

$$T_{\pi_{\tau_k+1}} v_{\pi_{\tau_k}} > v_{\pi_{\tau_k}} \quad \text{--- (iii)}$$

Note if

$$f(s) \geq g(s) + \delta \\ \Rightarrow T_{\pi_{\tau_k+1}} f(s) \geq T_{\pi_{\tau_k+1}} g(s) \quad \text{--- (iv)}$$

As

$$T_{\pi_{\tau_k+1}} f(s) = \sum_{a \in A(s)} \pi_{\tau_k+1}(a|s) p(s', a | s, a)(r + \gamma f(s')) \\ \geq \sum_{a \in A(s)} \pi_{\tau_k+1}(a|s) p(s', a | s, a)(r + \gamma g(s')) \quad \text{for } \gamma > 0$$

$$\therefore T_{\pi_{\tau_k+1}} g(s)$$

By (iii) and (iv)

$$T_{\pi_{\tau_k+1}} T_{\pi_{\tau_k+1}} v_{\pi_{\tau_k}} > T_{\pi_{\tau_k+1}} v_{\pi_{\tau_k}} \quad \text{--- (v)}$$

By (v) and (ii)

$$T^2 \pi_{R+1} v_{\pi_R} \geq v_{\pi_R}$$

Repeating the same procedure:

$$\lim_{M \rightarrow \infty} T^M \pi_{R+1} v_{\pi_R} \geq v_{\pi_R} \quad \text{--- (vi)}$$

But  $T_{\pi_{R+1}}$  has a fix pt.  $= v_{\pi_{R+1}}$

$$\therefore v_{\pi_{R+1}} \geq v_{\pi_R} \quad \text{--- (vii)}$$

$v_{\pi_{R+1}}$  either improves  $v_{\pi_R}$  i.e.  $\pi_{R+1}$  is an improvement over  $\pi_R$  or it doesn't change  $v_{\pi_R}$ .

In the case where

$$v_{\pi_{R+1}} = v_{\pi_R} \quad \text{--- (viii)}$$

By (i) & (viii)

$$T v_{\pi_{R+1}} = T \pi_{R+1} v_{\pi_{R+1}} = v_{\pi_{R+1}} \quad \text{--- (ix)}$$

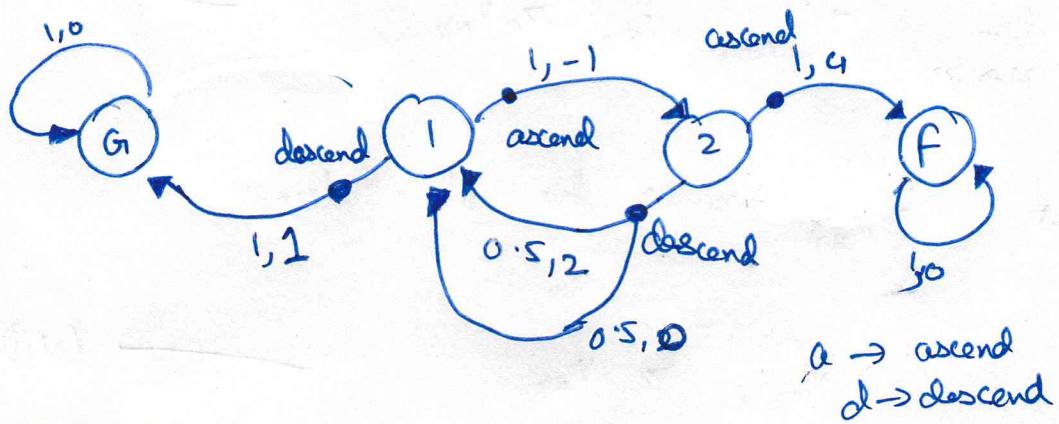
(As  $v_{\pi_{R+1}}$  is fix pt. of  $\pi_{R+1}$ )

$$\Rightarrow v_{\pi_{k+1}} = T v_{\pi_{k+1}} \Rightarrow v_{\pi_{k+1}} \text{ is the fix pt. of } T$$

But  $T$  has fix pt. of  $v_*$

$$\therefore v_* = v_{\pi_{k+1}} \therefore \text{optimal}$$

(8-15.)



Initialization:

$$\pi_0(a|s) \quad \forall s \in \{G, 1, 2, F\} = 0.5$$

$$\pi_0(d|s) \quad \forall s \in \{G, 1, 2\} = 0.5$$

$$v_{\pi_0}(s) \quad \forall s \in \{G, 1, 2, F\} = 0$$

Note:  $v_{\pi_0}(G) = 0$   $\left. \right\} \text{As terminal}$   
 $v_{\pi_0}(F) = 0$

$$\begin{aligned} v_{\pi}(1) &= \pi(d|1) \cdot 1 [1 + v_{\pi}(0)] + \\ &\quad \pi(a|1) \cdot 1 [-1 + 1 + v_{\pi}(2)] \\ &= \pi(d|1) + \pi(a|1) (-1 + v_{\pi}(2)) - \textcircled{A} \end{aligned}$$

$$\begin{aligned} v_{\pi}(2) &= \pi(d|1) [0.5(2 + 1 \cdot v_{\pi}(1)) + 0.5[0 + v_{\pi}(1)]] + \\ &\quad \pi(a|1)[1][4 + 1 \cdot (0)] \\ &= \pi(d|1) (1 + v_{\pi}(1)) + 4\pi(a|1) \end{aligned}$$

— \textcircled{B}

## Iter ①

$$v_{\pi_0}(1) = 0 \cdot 5 + 6 \cdot 5 (-1 + v_{\pi_0}(2)) = v_{\pi_0}(2)/2$$

$$v_{\pi_0}(2) = 0 \cdot 5 (1 + v_{\pi_0}(1)) + 4 (0 \cdot 5) = 2 \cdot 5 + v_{\pi_0}(1)/2$$

$$\therefore v_{\pi_0}(1) = 5/3, \quad v_{\pi_0}(2) = 10/3$$

### Policy improvement

$$\pi_1(1) = \arg \max \{ 1[1+0], 1[-1+10/3] \}$$

$$= \arg \max \{ 1, 7/3 \} = \text{ascend.}$$

$$\pi_1(2) = \arg \max \{ 0 \cdot 5 [0 + \frac{5}{3}] + 0 \cdot 5 [2 + \frac{5}{3}] + 4 \}$$

$$= \arg \max \{ \frac{8}{3}, 4 \} = \text{ascend.}$$

## Iter ②

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

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

### Policy improvement

$$\pi_2(1) = \arg \max \{ 1, -1+4 \} = \text{ascend.}$$

$$\pi_2(2) = \arg \max \{ 1+3, 4 \} = \text{ascend} \Delta \text{descend.}$$

### Policy is stable!

$\pi^*(1) = \text{ascend}$

$\pi^*(2) = \text{ascend or descend. We can choose } \Pr(a|z) \text{ & } P(a|z)$   
arbitrarily

The optimal policy in this environment doesn't give any preference to ascending or descending from any state  $i > 1$ .

The reason is that  $\gamma = 1$ . I have no motivation to reach  $F$  faster.

If at state  $i$ , I traverse to  $i-1$  and then come back to  $i$  and then go to  $F$ , I get the same reward if I had directly went from  $i$  to  $F$ .

$\therefore$  As long as I don't descend from  $I$ , every policy is optimal.  $\Rightarrow V_{\pi}(I) = \text{descend}$  is the only criteria

To promote  $\pi_1$  to be optimal, I must set up  $V_{\pi_0}$  such that

$$\arg \max \{1, -1 + r_{\pi_1}(2)\} = \text{ascend}$$

In other words,

$$V_{\pi_0}(2) - 1 > 1$$

$$\Rightarrow V_{\pi_0}(2) > 2$$

Let's calculate  $V_{\pi_0}(i)$

By observation we find that  $\sum_{i=1}^n v_{\pi_0}(i) = \frac{5n}{2}$  — (1)

We also observe

$$v_{\pi_0}(i) = \Delta i \quad \text{where} \quad \Delta = \frac{5}{n+1} \quad — (2)$$

Let's verify this claim: (Under  $\pi_0$  being equiprobable)

$$v_{\pi_0}(1) = 0.5 + 0.5[-1 + v_{\pi_0}(2)] = v_{\pi_0}\left(\frac{2}{2}\right) \quad — (3)$$

$$\begin{aligned} v_{\pi_0}(i) &= 0.5[-1 + v_{\pi_0}(i+1)] + 0.5[1 + v_{\pi_0}(i-1)] \\ &= \underbrace{v_{\pi_0}(i+1) + v_{\pi_0}(i-1)}_2 \end{aligned} \quad — (4)$$

$$\begin{aligned} v_{\pi_0}(n) &= 0.5[1 + v_{\pi_0}(n-1)] + 4(0.5) \\ &= \underbrace{5 + v_{\pi_0}(n-1)}_2 \end{aligned} \quad — (5)$$

Let's put (2) in (3), (4), (5)

$$(3) : \underbrace{v_{\pi_0}(2)}_2 = \frac{2\Delta}{2} = \Delta = v_{\pi_0}(1) \quad \text{Holds!}$$

$$(4) : \underbrace{v_{\pi_0}(i+1) + v_{\pi_0}(i-1)}_2 = \frac{\Delta(i+1) + \Delta(i-1)}{2} = \cancel{\Delta} = i\Delta \quad \text{Holds!}$$

$$(5) : \underbrace{5 + v_{\pi_0}(n-1)}_2 = \frac{5 + (n-1)\Delta}{2} = 2.5 + \frac{5(n-1)}{(n+1)2}$$

$$= \frac{5n + \cancel{5} + \cancel{5n - \cancel{5}}}{2(n+1)}$$

$$= \frac{10n}{2(n+1)} = \frac{5n}{n+1} = n\left(\frac{5}{n+1}\right) = n\Delta$$

Holds!

Now for

$$v_{\pi_0}(2) > 2$$

$$2\Delta > 2$$

$$\Delta > 1$$

$$\Rightarrow \frac{5}{n+1} > 1$$

$$\Rightarrow 5 > n+1$$

$$\Rightarrow n < 4$$

$\therefore$  for  $n < 4$   $\pi_g$  will be optimal & we will only need 1 iteration!

for  $n \geq 4$ : if  $n > 4$ ,  $\Delta < 1$

$$\pi_1(1) = \arg \max \{1, -1 + v_{\pi_0}(2)\}$$

$$= \arg \max \{1, -1 + 2\Delta\}$$

= descend.

$$\pi_1(i) = \arg \max \{-1 + (i+1)\Delta, 1 + (i-1)\Delta\} \quad 1 < i < n$$

$$= \arg \max \{i\Delta + (\Delta - 1), i\Delta - (\Delta - 1)\}$$

= descend.

$$\pi_1(n) = \arg \max \{2n, 1 + \Delta(n-1)\}$$

$$= \arg \max \left\{2n, 1 + \frac{5n-5}{n+1}\right\}$$

$$= \arg \max \left\{2n, \frac{6n-4}{n+1}\right\}$$

$$2n > \frac{6n-4}{n+1}$$

$$\Rightarrow 2n^2 + 2n > 6n - 4$$

$$\Rightarrow 2n^2 > 4n - 4$$

$$\Rightarrow n^2 > 2n - 1$$

$$\Rightarrow n^2 - 2n + 1 > 0$$

$$n^2 - 2n + 1 = 0$$

$$\Rightarrow n = 2 \pm \frac{\sqrt{4-4(1)(1)}}{2} = 1$$

$$\therefore n^2 - 2n + 1 > 0 \text{ for } n > 1, \text{ holds.}$$

$$\therefore \bar{n}_1(n) = \text{asconds.}$$

$\therefore$  Only  $\bar{n}_1(1)$  is wrongly chosen.

Let's calculate  $v_{\bar{n}_1}(i)$

$$v_{\bar{n}_1}(1) = 1 \times 1 = 1$$

$$v_{\bar{n}_1}(i) = 1[-1 + v_{\bar{n}_1}(i+1)] = v_{\bar{n}_1}(i+1) - 1 \quad | < i < n$$

$$v_{\bar{n}_1}(n) = 1[2n-1] = 2n-1$$

$$\begin{aligned} \therefore v_{\bar{n}_1}(i) &= 2n-1 - (n-i) = 2n-n-1+i \\ &= n+i-1 \end{aligned}$$

Let's calculate  $\pi_2$

$$\begin{aligned}\pi_2(i) &= \arg\max \{1, -1 + v_{\pi_1}(2)\} \\ &= \arg\max \{1, -1 + (n+2-i)\} \\ &= \arg\max \{1, n\}\end{aligned}$$

Since  $n > 4$ , we choose ascend  $\therefore \pi_2$  is optimal.

∴ 2 iterations:

for  $n=4$   $\Rightarrow \Delta = 1$

$$\begin{aligned}\pi_1(1) &= \arg\max \{1, -1 + v_{\pi_0}(2)\} \\ &= \arg\max \{1, 1\} = \text{ascend or descend.}\end{aligned}$$

$$\begin{aligned}\pi_1(i) &= \arg\max \{1 + (i+1)\Delta, 1 + (i-1)\Delta\} \\ &= \arg\max \{1, 1\} = \text{ascend or descend.}\end{aligned}$$

$$\begin{aligned}\pi_2(n) &= \arg\max \{2n, 1 + \Delta(n-1)\} \\ &= \arg\max \{2n, 1 + n\Delta\} = \text{ascend.}\end{aligned}$$

In this case we may need more iterations!



$\alpha$   $\alpha$   $\alpha$   $\alpha$