

Q13) < b >

$P(Q|S)$  - probability of queuing in state state  
 $P(Q|F)$  " " " " fresh

we will find

$$V_{t+i}(\text{state} = S) = E[R_t + R_{t+i} \gamma^i]$$

$$V_{t+i}(\text{state} = S) = E[R_{t+i} + \gamma^2 R_{t+2+3}]$$

$$V_{t+i}(\text{states} = S) = E[R_{t+i} + \gamma R_{t+i+1} | S_t = S] \\ + P(S_{t+3} = F | S_t = S) [-10] + P(S_{t+3} = S | S_t = S) [-10]$$

$$E[R_{t+i} + \dots | S_t = S] = \sum_{R=0}^2 \gamma^R E[R_{t+2+R} | \text{state} = S]$$

I have calculated  $V_3(F)$ .

policy will be dynamic because  $V_i(F, S)$

changes with time-step (i).

Q13) <6>

$$V_3(F) = \text{average Reward after 3 time steps}$$

$$= E[R_{t+1}|F] + \gamma E[R_{t+2}|S_t=F] + \gamma^2 E[R_{t+3}|S_t=F]$$

$$= [1 - p(S_{t+3}=F|S_t=F)] [-10] +$$

$$p(S_{t+3}=F|S_t=F) [10]$$

Last two terms represent additional penalty or reward

$$E[R_{t+1}|F] = p(Q|F)(-4) + (1 - p(Q|F))(+4)$$

$$= 4 - 8p(Q|F) \quad \text{--- (i)}$$

$$E[R_{t+2}|F] = p(S_{t+1}=F|S_t=F) E[R_{t+2}|S_{t+1}=F]$$

$$+ [1 - p(S_{t+1}=F|S_t=F)] E[R_{t+2}|S_{t+2}=S]$$

$$= p(S_{t+1}=F|S_t=F) \{ E[R|F] + E[R|S] \} + E[R|S]$$

$$\rightarrow p(S_{t+1}=F|S_t=F) = p(Q|F)(0.9) + (1 - p(Q|F))(0.5)$$

$$p(S_{t+1}=F|S_t=F) = 0.5 - 0.4p(Q|F)$$

$$E[R|S] = p(Q|S)(-8) + [1 - p(Q|S)]4$$

$$= 4 - 12p(Q|S) \quad \text{--- (ii)}$$

$$E[R_{t+2}|S_t=F] = [0.5 - 0.4p(Q|F)] [1 - 8p(Q|F) - 4 + 12p(Q|S)]$$

$$+ 4 - 12p(Q|S)$$

$$= -0.15 + p(Q|F)[-40 - 0.4 + 3.2p(Q|F)] + 4$$

$$+ p(Q|S)[-6] - 3.6p(Q|S)p(Q|S)$$



$$E[R_{t+3} | S_t = F] =$$

$$p[S_{t+2} = F | S_t = F] E[R_{t+3} | S_{t+2} = F] + (1 - p[S_{t+2} = F | S_t = F]) E[R_{t+3} | S_{t+2} = S]$$

$$\Rightarrow p[S_{t+2} = F | S_t = F] = p(Q | S_t = F)(0.9) p(Q | S_t = F)(0.9) \\ + p(Q | F)(0.5) p(Q | S)(0.8) + \\ [1 - p(Q | F)](0.5) p(Q | F)(0.9) + \\ [1 - p(Q | F)](0.5) [1 - p(Q | F)](0.5)$$

$$\Rightarrow p(Q | F) [0.81 p(Q | F) - 0.45 p(Q | F)] + 0.45 \\ + p(Q | S) p(Q | F)(4) + \\ + p(Q | S) [4 - 1 - p(Q | F) 0.25 [p(Q | F)^2 - 2p(Q | F) + 1]]$$

$$\Rightarrow p(Q | F) [-1.19 p(Q | F)]$$

$$E[R_{t+3} | S_t = F] = p(Q | F)$$

$$= p(Q | F) [1.36 p(Q | F) - 2] +$$

$$p[S_{t+2} = F | S_t = F] = \\ = p(Q | F) [1.36 p(Q | F) - 1.55] + 0.25 + p(Q | S) p(Q | F)$$

put value of  $p[S_{t+2} = F | S_t = F]$  in eq (ii)

Since, We have calculated

$$E[R_{t+i} | S_t = F] \quad 1 \leq i \leq 3$$

$$P(S_{t+2} = F | S_t = F) \text{ in term's of}$$

$$P(0|S) \text{ and } P(0|F)$$

Now calculate

$$P[S_{t+3} = F | S_t = F]$$

$$\begin{aligned} &= P[S_{t+3} = F | S_{t+2} = F] P[S_{t+2} = F | S_t = F] \\ &\Rightarrow + P[S_{t+3} = F | S_{t+2} = S] (1 - P(S_{t+2} = F | S_t = F)) \end{aligned}$$

#  ~~$P(0|F)(0.9) + P$~~   $P(S_{t+2} = F | S_t)$  can be taken from eq (IV)

$$P[S_{t+3} = F | S_{t+2} = F] = P(0|F)(0.9) + (1 - P(0|F))(0.5)$$

$$P[S_t = F | S_{t+2} = S] = P(0|S)(0.8) + (1 - P(0|S))(0.5)$$

Now we can put all value's in

equation of  $V_3(F)$  to get state value,

same approach for  $V(S)$



follow same approach to calculate  $V_2(F)$ ,  $V_1(F)$   
 $V_i(\text{state} = \text{State}) \quad 1 \leq i \leq 3$

$$V_i(\text{state} = s) = V_i(s_{t+3} = s_{t+i} = s) \quad s \in \{F, S\} \quad 1 \leq i \leq 3$$

$V_i(\cdot)$  represent average ~~over~~ reward starting from state  $t+i$  to  $t+3$  from starting time-step state 's'.

optimal policy calculation :-

at time step  $t+i$  and in state  $= s$

$$\arg \max_a \text{Reward}(a|s) + \gamma$$

Take Action 'a' such that

$$\arg \max_{a \in \{Q, NQ\}} P \left( \sum_{r \in R} \sum_{s' \in S} p(r, s' | s, a) [r + \gamma V_{3-i}(s')] \right)$$

~~Q~~  $Q \rightarrow$  Query

$NQ \rightarrow$  not Query

$R \rightarrow$  set of reward's

$S \rightarrow$  set of states

$p(r, s' | s, a)$  is Given in question

$V_{3-i}(s')$  can be calculated by

above method

Q13) &lt;c&gt;

6

$$V_F(F) = p(Q|F)(0.9)[-4 + \gamma V_F] + p(Q|F)(0.1)[-4 + \gamma V_F] \\ + (1 - p(Q|F))(0.5)[4 + \gamma V_F] + (1 - p(Q|F))(0.5)[4 + \gamma V_F]$$

$$V_F = V_F[\gamma 0.9 p(Q|F) + 0.5\gamma - 0.5 p(Q|F)] \\ + V_S[0.1\gamma p(Q|F) + 0.5\gamma - 0.5 p(Q|F)] \\ + p(Q|F)[-0.45 - 0.4 - 2 - 2] \\ + 2 + 2$$

$$V_S = p(Q|S)[0.8][-8 + \gamma V_F] + p(Q|S)(0.2)[-8 + \gamma V_S] \\ + [1 - p(Q|S)](0)[4 + \gamma V_F] + [1 - p(Q|S)](1.0)[4 + \gamma V_S]$$

$$V_S = V_F[-0.64\gamma p(Q|S)] + V_S[0.2\gamma p(Q|S) + \gamma - p(Q|S)] \\ + [-0.64 + (-1.6) + 4] p(Q|S) \\ + 4$$

Solve by Value Policy iteration

$p(Q|S)$  → probability of Query in state state  
 $p(Q|F)$  → " " " " " fresh "

$\gamma$  → discounting factor

Ans 14) In policy improvement :-  
for a given state  $s \in S$ .

Calculate action state-action value for all actions expect deterministic action defined by current policy.

either  $a \in A(s)$  such that  $a \neq \pi(s)$

let this set be of new state-action be SA.

if any action in SA give's more reward than by current policy.

As either  $q(s, a) > \pi(a|s) \sum p(s',$   
 $q(s, a) > q(s, a')$

where  $a' = \pi(a|s)$   
then change  $\pi(a|s) = a$

otherwise do nothing.

Since, updating  $\pi(s)$  [ $s \in S$  and  $s \neq s$ ]  
does not change other policies so,  
it only update current ~~state value~~  
policy.

so, it leads to policy improvement