

The MDP has 2 policies:

Policy 1:

$s_1$	a
$s_2$	c
$s_3$	c
$s_4$	d
$s_5$	d

Policy 2:

$s_1$	b
$s_2$	d
$s_3$	d
$s_4$	d
$s_5$	d

At infinite time, the value of  $V(S_i)$  will get constant.

So, we can write the following equations:  
for action a [when  $S_1$  is following policy 1]

$$V(S_1) = 0 + \gamma [(0.1)V(S_1) + (0.9)V(S_2)]$$

$$V(S_2) = 0 + \gamma [P_S \cdot V(S_2) + (1 - P_S) V(S_4)]$$

$$V(S_3) = 0 + \gamma [(0.1)V(S_3) + P_R \cdot V(S_5) + (0.9 - P_R) \cdot V(S_4)]$$

$$V(S_4) = 10 + \gamma [(0.1) \cdot V(S_4) + 0.9 \cdot V(S_1)]$$

$$V(S_5) = -10 + \gamma [(0.1) \cdot V(S_5) + 0.9 \cdot V(S_1)]$$

Solving them,

$$\rightarrow V(S_2) = \gamma P_S \cdot V(S_2) + \gamma (1 - P_S) V(S_4)$$

$$\therefore V(S_2) = \frac{\gamma (1 - P_S) V(S_4)}{(1 - \gamma P_S)} \quad \text{--- (1)}$$

$$\rightarrow V(S_4) = 10 + \gamma (0.1) V(S_4) + \gamma (0.9) V(S_1)$$

$$\therefore V(S_4) = \frac{[10 + (0.9)\gamma \cdot V(S_1)]}{1 - (0.1)\gamma} \quad \text{--- (2)}$$

Placing the value of  $V(S_4)$  from equation (2) in equation (1), we get

$$V(S_2) = \frac{(1 - p_s) \cdot \gamma \cdot (10 + 0.9\gamma \cdot V(S_1))}{(1 - \gamma p_s) \cdot (1 - 0.1\gamma)} \quad (3)$$

Placing this value in equation of  $V(S_1)$ , we get

$$\begin{aligned} V(S_1) &= 0 + \gamma [0.1 \cdot V(S_1) + 0.9 \cdot V(S_2)] \\ &= \gamma (0.1) \cdot V(S_1) + \frac{\gamma^2 (0.9) (1 - p_s) (10 + 0.9\gamma \cdot V(S_1))}{(1 - \gamma p_s) (1 - 0.1\gamma)} \\ \therefore (1 - (0.1)\gamma)^2 V(S_1) \cdot (1 - \gamma p_s) &= \gamma (0.9 - 0.9 p_s) (10 + 0.9\gamma \cdot V(S_1)) \end{aligned}$$

$$\therefore V(S_{10}) = \frac{(0.9 - 0.9 p_s) \gamma^2}{(1 - 0.1\gamma)^2 (1 - \gamma p_s) - \cancel{0.81} (\gamma^3 - p_s \gamma^3)} \quad (4)$$

We can similarly write the following equations for action b [i.e. when  $s_1$  is following policy 2]

$$V(s_1) = 0 + \gamma [(0.1)V(s_1) + (0.9)V(s_3)]$$

$$V(s_2) = 0 + \gamma [P_S V(s_2) + (1-P_S)V(s_4)]$$

$$V(s_3) = 0 + \gamma [(0.1)V(s_3) + P_R \cdot V(s_5) + (0.9 - P_R)V(s_4)]$$

$$V(s_4) = 10 + \gamma [(0.1)V(s_4) + (0.9)V(s_1)]$$

$$V(s_5) = -10 + \gamma [(0.1)V(s_5) + (0.9)V(s_1)]$$

Solving them,

$$\rightarrow V(s_5) = -10 + \gamma \cdot (0.1) \cdot V(s_5) + \gamma \cdot (0.9) \cdot V(s_1)$$

$$\therefore V(s_5) = \frac{-10 + \gamma \cdot (0.9) \cdot V(s_1)}{1 - (0.1)\gamma} \quad (5)$$

$$\rightarrow V(s_3) = \gamma \cdot (0.1) \cdot V(s_3) + \gamma P_R V(s_5) + \gamma (0.9 - P_R)V(s_4)$$

Putting the values of  $V(s_4)$  and  $V(s_5)$  from equations (2) and (5), we get

$$V(S_3) = \gamma \cdot (0.1) V(S_1) + \gamma P_R \cdot \left( \frac{-10 + 0.98 V(S_1)}{1 - 0.18} \right)$$

$$+ \gamma (0.9 - P_R) \left( \frac{10 + 0.98 \cdot V(S_1)}{1 - 0.18} \right)$$

$$\therefore V(S_3) = \frac{\gamma P_R (-10 + 0.98 V(S_1)) + \gamma (0.9 - P_R) (10 + 0.98 V(S_1))}{(1 - 0.18)^2}$$

(6)

$$\rightarrow V(S_1) = \gamma (0.1) V(S_1) + \gamma (0.9) \cdot V(S_3)$$

Placing the value of  $V(S_3)$  from equation (6), we get

$$(1 - 0.1\gamma) S_1 = \gamma \cdot 0.9 \left[ \gamma P_R (-10 + 0.98 V(S_1)) + \gamma (0.9 - P_R) (10 + 0.98 V(S_1)) \right]$$

$$\frac{(1 - 0.1\gamma)^2}{(1 - 0.1\gamma)^2}$$

$$\therefore V(S_1) = \frac{\gamma (0.9) [-20\gamma P_R + 9\gamma + 0.81\gamma^2 V(S_1)]}{(1 - 0.1\gamma)^3}$$

$$v(S_1) \cdot (1 - 0.18)^3 - \gamma(0.9)(0.81\gamma^2 v(S_1)) \\ = \gamma(0.9)(-20\gamma p_R + 9\gamma)$$

$$\therefore v(S_{1B}) = \frac{0.9\gamma^2(9 - 20p_R)}{(1 - 0.18)^3 - (0.9)^3\gamma^3} \quad (7)$$

(1) (a)  $p_S = 0.2, p_R = 0.01$

From eq.(4), we get  $v(S_1)$  when following the first policy.

$$v(S_{1B}) = \frac{(9 - 9p_S)\gamma^2}{(1 - 0.18)^2(1 - p_S\gamma) - 0.81(\gamma^3 - p_S\gamma^3)} \\ = \frac{6.498}{0.10783} = 60.2608242$$

Similarly, if  $S_1$  is following policy second policy, from eq.(7)

$$v(S_{1B}) = \frac{0.9\gamma^2(9 - 20p_R)}{(1 - 0.18)^3 - (0.9)^3\gamma^3} \\ = \frac{7.1478}{0.11619125} = 61.5175411$$

As  $v(s_{1a}) < v(s_{1b})$ , we select  
the second policy.

$$\therefore v(s_1) = v(s_{1b}) = 61.5175411$$

~~∴~~

From eq (2)

$$v(s_4) = \frac{10 + 0.9\gamma \cdot v(s_1)}{1 - (0.1)\gamma}$$

$$= 69.1685056801$$

From eq (5)

$$v(s_5) = \frac{-10 + \gamma(0.9) \cdot v(s_1)}{1 - (0.1)\gamma}$$

$$= 47.0690581663$$

From eq.(1)

$$v(s_2) = \frac{\gamma(1-p_s) \cdot v(s_4)}{(1-\gamma p_s)}$$

$$= 64.8988448356$$

$$v(s_3) = \gamma[(0.1)v(s_3) + p_r \cdot v(s_5) + (0.9 - p_r) \cdot v(s_4)]$$

$$= 65.1150581272$$

$\therefore$  For  $P_S = 0.2, P_R = 0.01,$

The optimal policy is

$s_1 \quad b$

$s_2 \quad d$

$s_3 \quad d$

$s_4 \quad d$

$s_5 \quad d.$

and value function is

$$v(s_1) = 61.5175411$$

$$v(s_2) = 64.8988448356$$

$$v(s_3) = 65.1150581272$$

$$v(s_4) = 69.1685056801$$

$$v(s_5) = 47.069058166$$

(b)  $P_S = 0.2 \quad P_R = 0.03.$

$$v(s_{1a}) = 60.2608242 \quad (\text{from eq.(4)})$$

from eq.(7),

$$v(s_{1b}) = \frac{0.98^2 (9 - 20(0.03))}{(1-0.18)^3 - (0.9)^3 8^3}$$

$$= \frac{6.8229}{0.11619125} = 58.7212893$$

As  $V(S_{1a}) > V(S_{1b})$  we select  
First policy.

$$\therefore V(S_1) = V(S_{1a}) = 60.2608242.$$

From eq.(2),

$$V(S_4) = 67.98122.$$

From eq.(5),

$$V(S_5) = 45.88177314.$$

From eq.(1)

$$V(S_2) = 63.78484901$$

$$V(S_3) = 63.52940266.$$

$\therefore$  For  $P_S = 0.2$  &  $P_R = 0.03$ ,

Optimal policy

$S_1 \quad a$

$S_2 \quad d$

$S_3 \quad d$

$S_4 \quad d$

$S_5 \quad d$

Value function is:

$$V(S_1) = 60.2608242.$$

$$V(S_2) = 63.78484901.$$

$$V(S_3) = 63.52940266.$$

$$V(S_4) = 67.98122.$$

$$V(S_5) = 45.88177314.$$

(C) If the agent is indifferent between taking action a and b in state  $S_1$ , the value function of  $S_1$  should remain the same.

$$\therefore V(S_{1a}) = V(S_{1b}).$$

$$\therefore (9 - 9P_S)\gamma^2$$

$$(1 - 0.18)^2 (1 - \gamma P_S) - 0.81(\gamma^3 - P_S \gamma^3)$$

$$= \frac{0.9\gamma^2(9 - 20P_R)}{(1 - 0.18)^3 - (0.9)^3\gamma^3}$$

$$\therefore 60.2608242 = \frac{(9 - 20P_R) \cdot 0.9\gamma^2}{0.11619125}$$

$$\therefore (9 - 20P_R) = 8.62022837$$

$$\therefore P_R = 0.0189885815$$

$$(2)(a) P_S = 0.6 \quad P_R = 0.1$$

from eq.(4), we get  $V(S_1)$  when following the first policy as follows

$$\begin{aligned} V(S_{1a}) &= \frac{(q - qP_S)\gamma^2}{(1 - 0.1\gamma)^2 (1 - P_S\gamma)} - 0.81(\gamma^3 - P_S\gamma^3) \\ &= \frac{3.249}{0.35218075 - 0.2777895} \\ &= 43.6744913. \end{aligned}$$

Similarly, if  $S_1$  is following second policy, we get the following from eq.(7)

$$\begin{aligned} V(S_{1b}) &= \frac{0.9\gamma^2(q - 20P_R)}{(1 - 0.1\gamma)^3 - (0.9)^3\gamma^3} \\ &= 48.9344077. \end{aligned}$$

$$\therefore V(S_{1b}) > V(S_{1a})$$

We select the second policy

$$V(S_1) = V(S_{1b}) = 48.9344077.$$

$$V(S_4) = \frac{10 + 0.9\gamma \times V(S_1)}{1 - (0.1)\gamma}$$

$$= 57.28057303.$$

$$V(S_5) = \frac{-10 + \gamma(0.9)V(S_1)}{1 - (0.1)\gamma}$$

$$= 35.18112552.$$

$$V(S_2) = \frac{\gamma(1-p_s)V(S_4)}{(1-\gamma p_s)}$$

$$= 50.62004128$$

$$V(S_3) = \gamma[(0.1)V(S_3) + p_r V(S_5) + (0.9 - p_r)V(S_4)]$$

$$= 51.79606898.$$

$\therefore$  for  $p_s = 0.6$  and  $p_r = 0.1$ ,

the optimal policy is

$S_1$	b
$S_2$	d
$S_3$	d
$S_4$	d
$S_5$	d.

And the value function is

$$V(S_1) = 48.9344077$$

$$V(S_2) = 50.62004128$$

$$V(S_3) = 51.79606898$$

$$V(S_4) = 57.28057303$$

$$V(S_5) = 35.18112552$$

(b)  $P_S = 0.6 \quad P_\gamma = 0.2$

From eq. (4)

$$V(S_{1a}) = 43.6744913$$

From eq. (7),

$$V(S_{1b}) = \frac{0.9\gamma^2(9 - 20(0.2)^2)}{(1 - 0.1\gamma)^3 - (0.9)^3\gamma^3}$$

$$(1 - 0.1\gamma)^3 - (0.9)^3\gamma^3$$

$$= \underline{4.06125}$$

$$0.741217625 - 0.625026375$$

$$= 34.9531484$$

As  $V(S_{1a}) > V(S_{1b})$ , we select the first policy.

$$\therefore V(S_1) - V(S_{1a}) = 43.6744913$$

From eq(2),  $\gamma = 0.9$ ,  $\alpha = 0.1$

$$V(S_4) = \frac{10 + 0.9\gamma \cdot V(S_1)}{1 - 0.1\gamma}$$

$$= 52.31125972$$

From eq. (5),

$$V(S_5) = \frac{-10 + \gamma(0.9) \cdot V(S_1)}{1 - 0.1\gamma}$$

$$= 30.21181221$$

From eq. (1)

$$V(S_2) = \frac{\gamma(1 - p_S) V(S_4)}{(1 - \gamma p_S)}$$

$$= 46.22855511$$

$$V(S_3) = \gamma[(0.1)V(S_3) + p_R V(S_5) + (0.9 - p_R)V(S_4)]$$

$$= 44.78147186$$

$\therefore f$

$\therefore$  For  $P_S = 0.6$ ,  $P_R = 0.2$ ,

the optimum policy is

$S_1$	a
$S_2$	d
$S_3$	d
$S_4$	d
$S_5$	d

And the value function is:

$$V(S_1) = 43.67449129.$$

$$V(S_2) = 46.22855511$$

$$V(S_3) = 42.78147186.$$

$$V(S_4) = 52.31125972$$

$$V(S_5) = 30.21181221.$$

(C) If the agent is indifferent, then

$$V(S_{1a}) = V(S_{1b})$$

$$\therefore \underbrace{(9 - 9P_S)\gamma^2}_{(1 - 0.1\gamma)^2(1 - \gamma P_S) - 0.81(\gamma^3 - P_S\gamma^3)}$$

$$= \frac{0.9\gamma^2(9 - 20P_R)}{(1 - 0.1\gamma)^3 - (0.9)^3\gamma^3}$$

$$\therefore 43.6744913 = \underbrace{(0.9)(0.95)^2(9 - 20P_R)}_{0.11619125}.$$

$$\boxed{\therefore P_R = 0.137621192.}$$