

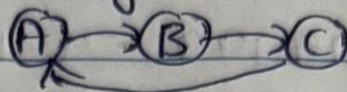
Homework - 6

Q1.1.

Policy 1



Policy 2



2. Policy 2

2.1

$$\gamma = 10^{-7}$$

Policy 1:

$$V^*(C) = 90 + \gamma \cdot 0 + \gamma^2 \cdot 90 + \gamma^3 \cdot 0 + \dots$$

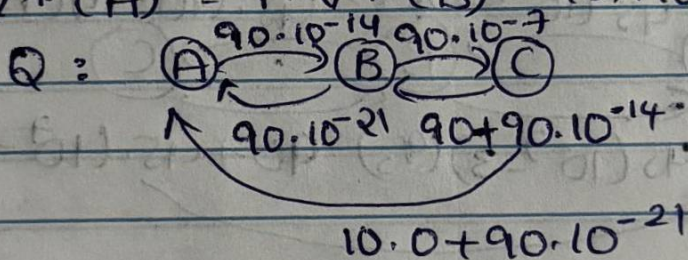
$$= 90 [1 + \gamma^2 + \gamma^4 + \dots]$$

$$= \frac{90}{1 - \gamma^2}$$

$$\approx 90$$

$$V^*(B) = \gamma \cdot V^*(C) = 90 \times 10^{-7}$$

$$V^*(A) = \gamma \cdot V^*(B) = 90 \times 10^{-14}$$



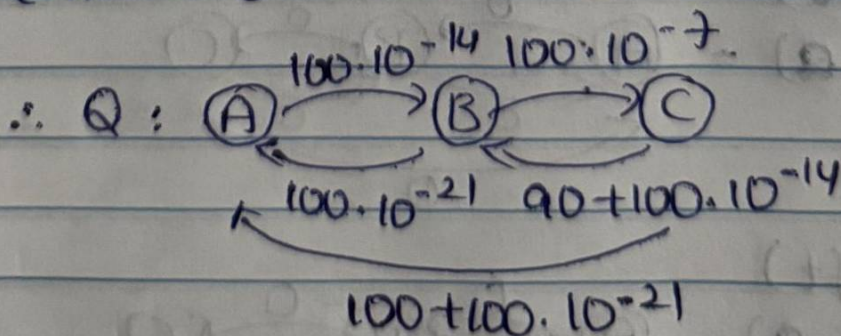
Policy 2:

$$V(C) = 100 + \gamma^3 100 + \gamma^6 100 + \dots$$

$$= \frac{100}{1 - \gamma^3} = \frac{100}{1 - 10^{-21}} \approx 100$$

$$V(B) = 100 \cdot 10^{-17}$$

$$V(A) = 100 \cdot 10^{-14}$$



\therefore Policy 2 is the optimal policy

3. Policy 1

3.1 $r = 0.9999999 = 1 - 10^{-7}$

$$V^*(C) = 90 + 90r^2 + 90r^4 + \dots$$

$$= \frac{90}{1-r^2} = 90$$

$$= \frac{90}{1-1-10^{-14}+2 \cdot 10^{-7}}$$

$$\approx \frac{90}{2} \cdot 10^7$$

$$\approx 45 \cdot 10^7$$

$$V^*(B) = r \cdot 45 \cdot 10^7 = (1 - 10^{-7}) 45 \cdot 10^7$$

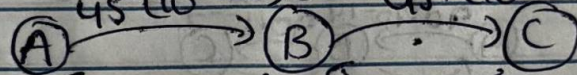
$$= 45(10^7 - 1)$$

$$V^*(A) = r \cdot 45(10^7 - 1)$$

$$= 45(10^7 - 1 - 1 + 10^{-7})$$

$$= 45(10^7 - 2)$$

$$45(10^7 - 2)(q_1) \quad 45(10^7 - 1)(q_2)$$



$$45(10^7 - 3)(q_3) \quad 90 + 45(10^7 - 2)(q_4)$$

$$100 + 45(10^7 - 3)(q_5)$$

$$q_4 = 90 + 45(10^7 - 2) = 45 \cdot 10^7$$

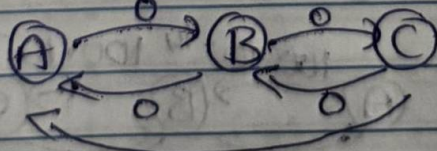
$$q_5 = 100 + 45(10^7 - 3) = 100 + 45 \times 10^7 - 150 - 45 \times 10^7 - 50$$

$$\therefore q_4 > q_5$$

\therefore Policy 1 is better.

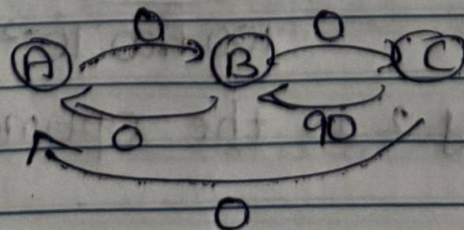
4.1 (B, left)

$$\hat{Q}(s, a) =$$

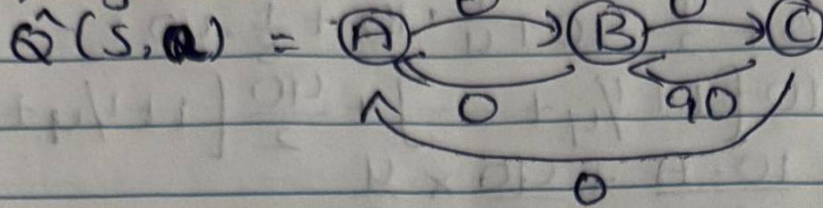


2. (C, left)

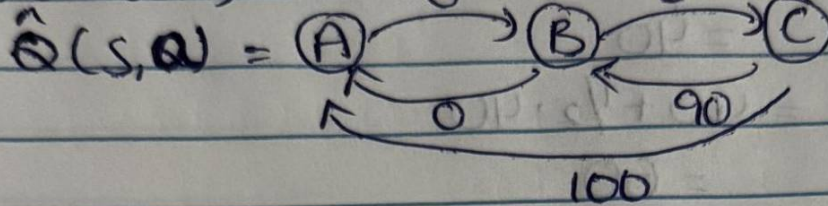
$$\hat{Q}(s, a) =$$



3. (A, right)

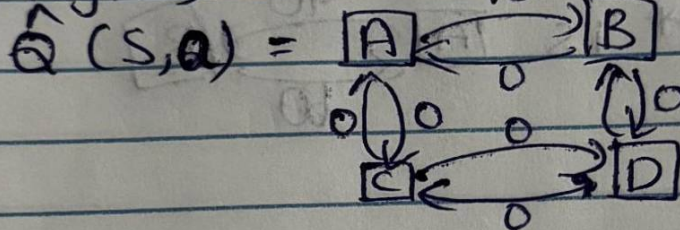


4. (C, down)

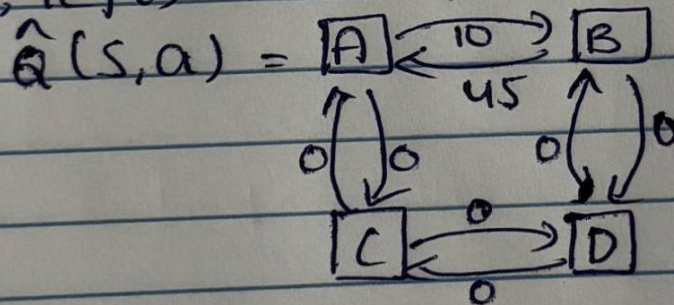


Q2. PART-A

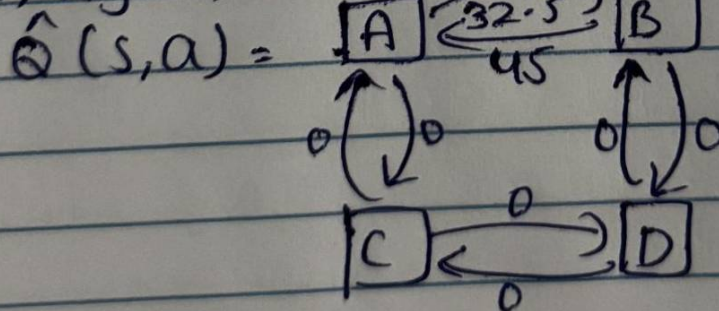
1. (A, right)



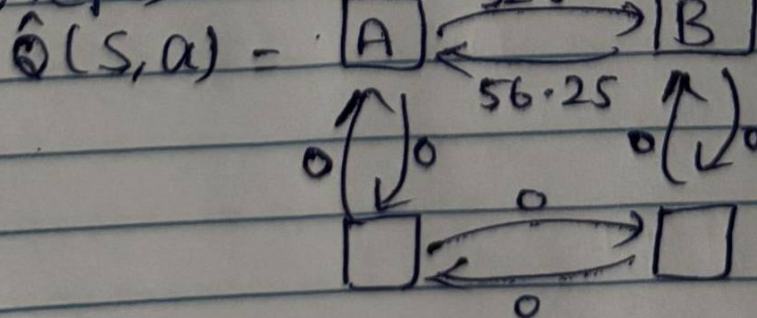
2. (B, left)



3. (A, right)



4. (B, left)



PART-B

1. ~~V_A~~ $V_A^* = 10 + \frac{1}{2} \cdot 40 + \frac{1}{4} \cdot 10 + \frac{1}{8} \cdot 40 + \dots$
 $= 10 \left[1 + \frac{1}{4} + \dots \right] + \frac{40}{2} \left[1 + \frac{1}{4} + \dots \right]$
 $= \frac{10 \cdot A}{3} + \frac{40}{2} \times \frac{4}{3}$
 $= 40$

$$V_B^* = 40 + \frac{1}{2} \cdot 40$$
$$= 60$$

\therefore The optimal policy

