

Name: Kartik Bhargadwaj

Roll No.: CS205020

"I Pledge that I haven't copied or given any unauthorized assistance on this exam".

---

Q1) c

Q3) Let the discount factor be  $\alpha$  for a disc-MDP.

We know the cum. disc. reward  $R \in \mathbb{R}$  is:

$$R = \sum_{t=0}^{\infty} \alpha^t r(x_t).$$

(a)

$$J(x) = E[R | x_0 = x]$$

$$J(x) = E\left[\sum_{t=0}^{\infty} \alpha^t r(x_t) | x_0 = x\right]$$

$$J(x) = r(x) + E\left[\sum_{t=1}^{\infty} \alpha^t r(x_t) | x_0 = x\right]$$

$$J(x) = r(x) + \alpha E\left[E\left[\sum_{t=1}^{\infty} \alpha^{t-1} r(x_t) | x_0 = x, x_1 = x'\right]\right]$$

$$J(x) = r(x) + \alpha \sum_{x' \in \mathcal{X}} P(x' | x) \cdot J(x').$$

~~16~~

$$M(x) = E[R^2 | x_0 = x]$$

$$\Rightarrow M(x) = E\left[\left(\sum_{t=0}^{\infty} \alpha^t r(x_t)\right)^2 \mid x_0 = x\right]$$

$$\Rightarrow M(x) = E\left[\left(r(x) + \sum_{t=1}^{\infty} \alpha^t r(x_t)\right)^2 \mid x_0 = x\right]$$

$$\Rightarrow M(x) = E\left[r(x)^2 + 2r(x) \left[\sum_{t=1}^{\infty} \alpha^t r(x_t)\right] + \left[\sum_{t=1}^{\infty} \alpha^t r(x_t)\right]^2 \mid x_0 = x\right]$$

$$\Rightarrow M(x) = r(x)^2 + 2r(x) E\left[\sum_{t=1}^{\infty} \alpha^t r(x_t) \mid x_0 = x\right]$$

$$+ E\left[\left(\sum_{t=1}^{\infty} \alpha^t r(x_t)\right)^2 \mid x_0 = x\right]$$

$$\Rightarrow M(x) = r(x)^2 + 2\alpha r(x) \sum_{x' \in X} p(x'|x) J(x')$$

$$+ \alpha^2 \sum_{x'} p(x'|x) M(x') \rightarrow \text{Using } J(x) \text{ from (a).}$$

Part(b)

We know variance =  $E(X^2) - E(X)^2$

$$\Rightarrow V(x) = M(x) - J(x)^2$$

$$\Rightarrow V(x) = r(x)^2 + 2\alpha r(x) \sum_{x'} p(x'|x) J(x') + \alpha^2 \sum_{x'} p(x'|x) M(x') - \underbrace{\left[r(x) + \alpha \sum_{x'} p(x'|x) J(x')\right]^2}_{J(x)^2}$$

$$\Rightarrow V(x) = r(x)^2 + 2\alpha r(x) \sum_{x'} P(x'|x) J(x')$$

$$+ \alpha^2 \left[ \sum_{x'} P(x'|x) V(x') + \sum_{x'} P(x'|x) J(x')^2 \right] - J(x)^2$$

$$\Rightarrow V(x) = r(x)^2 + 2\alpha r(x) \sum_{x'} P(x'|x) J(x') + \alpha^2 \sum_{x'} P(x'|x) J(x')^2 - J(x)^2 + \alpha^2 \sum_{x'} P(x'|x) V(x')$$

We now have  $V(x)$  in the following form,

$$V(x) = \psi(x) + \alpha^2 \sum_{x' \in X} P(x'|x) V(x')$$

$$\text{where } \psi(x) = r(x)^2 + 2\alpha r(x) \sum_{x'} P(x'|x) J(x') + \alpha^2 \sum_{x'} P(x'|x) J(x')^2 - J(x)^2$$

$$\Rightarrow \psi(x) = [r(x) + J(x)] \left[ r(x) - J(x) \right] + 2\alpha r(x) \sum_{x'} P(x'|x) J(x') + \alpha^2 \sum_{x'} P(x'|x) J(x')^2$$

$$\Rightarrow \psi(x) = [r(x) + J(x)] \left[ -\alpha \sum_{x'} P(x'|x) \cdot J(x') \right] \quad \text{From A3(a) can get } J(x)$$

$$+ 2\alpha r(x) \sum_{x'} P(x'|x) J(x')$$

$$+ \alpha^2 \sum_{x'} P(x'|x) J(x')^2$$

$$\Rightarrow \psi(x) = [r(x) - J(x)] \left[ \alpha \sum_{x'} P(x'|x) J(x') + \alpha^2 \sum_{x'} P(x'|x) J(x')^2 \right]$$

$$\Rightarrow \psi(x) = -\alpha^2 \left[ \sum_{x'} P(x'|x) J(x') \right]^2 + \alpha^2 \sum_{x'} P(x'|x) J(x')^2$$

$$\Rightarrow \boxed{\psi(x) = \alpha^2 \left[ \sum_{x'} P(x'|x) \psi(x')^2 - \left( \sum_{x'} P(x'|x) \psi(x') \right)^2 \right]}.$$

Remade, proved.

84)

States =  $\{A, B\}$ .

$$\left. \begin{aligned} P(B|A) &= p_1 \\ P(A|B) &= p_2 \\ P(B|B) &= q_2 = 1 - p_2 \\ P(A|A) &= 1 - p_1 = q_1 \end{aligned} \right\} \text{Transition probabilities}$$

~~g~~ # single stage cost

$$g_A = -A ; g_B = B.$$

$$J^*(A) = -A + \gamma [P(B|A) J^*(B) + P(A|A) J^*(A)]$$

$$J^*(B) = B + \gamma [P(A|B) J^*(A) + P(B|B) J^*(B)]$$

$$\Rightarrow J^*(A) = -A + \gamma p_1 J^*(B) + \gamma q_1 J^*(A) \quad - (1)$$

$$J^*(B) = B + \gamma p_2 J^*(A) + \gamma q_2 J^*(B) \quad - (2)$$

from (2), we have:

$$J^*(B) = \frac{B + \gamma p_2 J^*(A)}{1 - \gamma q_2} \quad - (3)$$



Q4) continued

$$\Rightarrow J^*(B) = \frac{B + \gamma p_2 (\gamma B p_1 - A(1 - \gamma q_2))}{1 - \gamma q_1 - \gamma q_2}$$


---


$$1 - \gamma q_2$$

$$\Rightarrow J^*(B) = \frac{B - \gamma B q_1 - \gamma B q_2 + \gamma^2 B p_1 p_2 - \gamma A p_2 + \gamma^2 A p_2 q_2}{(1 - \gamma q_2)(1 - \gamma q_1 - \gamma q_2)}$$

↳ (5)

Q4

(b)  $\min \{M + J^*(a), J^*(a)\}$

$$\Rightarrow \min \{M, 1\} + J^*(a)$$

$\Rightarrow$  whenever cost  $M < 1$ , it is optimal to buy a new machine.

Q2) True since  $\pi_3(i)$  is taking the optimal action w.r.t  $\pi_1(i)$  &  $\pi_2(i)$ . ~~However, it is not necessary that~~

However, it is not necessary that  $J_{\pi_3}(i) = \min \{J_{\pi_1}, J_{\pi_2}\}$

Since  $\pi_3$  is a better policy than  $\pi_1$  &  $\pi_2$ ,

$$J_{\pi_3} \leq \min \{J_{\pi_1}, J_{\pi_2}\}.$$

References for all my solution!

Q1) Internet (Google)

Q3) TD Methods for variance of reward to go (The authors prove for ~~the~~ an SSP. But I've converted it into a disc. MDP)