

Ravina Lad.

EX 3 → Dynamic Programming.

Q.1

$$a) V^*(s) = \sum_a \pi^*(a|s) q^*(s, a)$$

$$b) q^*(s, a) = \sum_{s', r} p(s'|s, a) (r + \gamma V^*(s'))$$

$$c) \pi^*(s) = \arg\max_a q^*(s, a)$$

$$d) \pi^*(s) = \arg\max_a \sum_{s', r} p(s'|s, a) (r + \gamma V^*(s'))$$

e) Bellman eqⁿ (V^* , π^* , q^*) in terms of 3 args & p
& 2 args r .

$$V^*(s) = \max_a [r(s, a) + \gamma \sum_{s'} p(s'|s, a) V^*(s')]$$

$$\pi^*(s) = \sum_a \pi(a|s) [r(s, a) + \gamma \sum_{s'} p(s'|s, a) \cdot V^*(s')]$$

$$q^*(s, a) = r(s, a) + \sum_{s'} p(s'|s, a) \left(\gamma \sum_a \pi(a|s') q^*(s') \right)$$

$$q^*(s, a) = r(s, a) + \sum_{s'} p(s'|s, a) \gamma \max_a q^*(s', a')$$

Q.2 Fixing Policy Iteration -

a) We can fix the bug in policy iteration by every time selecting optimal actions with equal probability.

e.g. if policy has 4 action $a = \{ \text{Left, Right, Up, Down} \}$ at each iteration will select optimal policy with equal probability. If Left & Up are optimal then \rightarrow

$$\pi^*(a) = \{ 0.5, 0, 0.5, 0 \}$$

Here's modified pseudo code \rightarrow

1. Initialization - $V(s) \in \mathbb{R}$ $\forall s \in S$

2. Policy evaluation.

$$\text{loop } \Delta \leftarrow 0$$

loop for each $s \in S$

$$V(s) \leftarrow V(s)$$

$$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} P(s',r|s,a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |V - V(s)|)$$

until $\Delta < \epsilon$ (ϵ a small +ve number)

3. Policy improvement

policy-stable \leftarrow True

for each $s \in S$:

$$\text{old-action} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \frac{1}{\text{No. of optimal actions}} \cdot \underset{s',r}{\operatorname{argmax}} \sum_a P(s',r|s,a) [r + \gamma V(s')]$$

if $\text{old-action} \neq \pi(s)$:

policy stable \leftarrow False.

if policy-stable:

stop & return $V \approx V^*$ & $\pi \approx \pi^*$
(Break)

else:

go to 2.

b). I don't think so that value iteration has any bug similar to policy iteration.

- I think optimal value f^* should be unique.
- It should also unique regardless of how many different optimal policies we have.

Q.3) Policy Iterations for action values:

a) 1. Initialization:

$$q(s, a) \in R, \pi(s) \in A(s) \text{ arbitrarily } \forall s \in S.$$

2. Policy Evaluation:

Loop: $\Delta \leftarrow 0$

Loop for each $s \in S$:

Loop for each $a \in A$:

$$q(s, a) \leftarrow \sum_{s', r} p(s', r | s, a) [r + \gamma \sum_{a'} \pi(a' | s') q(s', a')]$$

$$\Delta \leftarrow \max(\Delta; \text{abs}(q(s, a) - q(s, a)))$$

until $\Delta < \epsilon$
(small positive number)

3. Policy Improvement:

policy-stable \leftarrow True

For each $s \in S$: for each $a \in A$:

$$\begin{aligned} \text{old-action} &\leftarrow \pi \\ \pi(a) &\leftarrow \underset{a}{\operatorname{argmax}} q(s, a) \end{aligned}$$

If old.action $\neq \pi(a)$:

policy-stable \leftarrow False

If policy-stable, stop & return $q \approx q^*$ & $\pi \approx \pi^*$

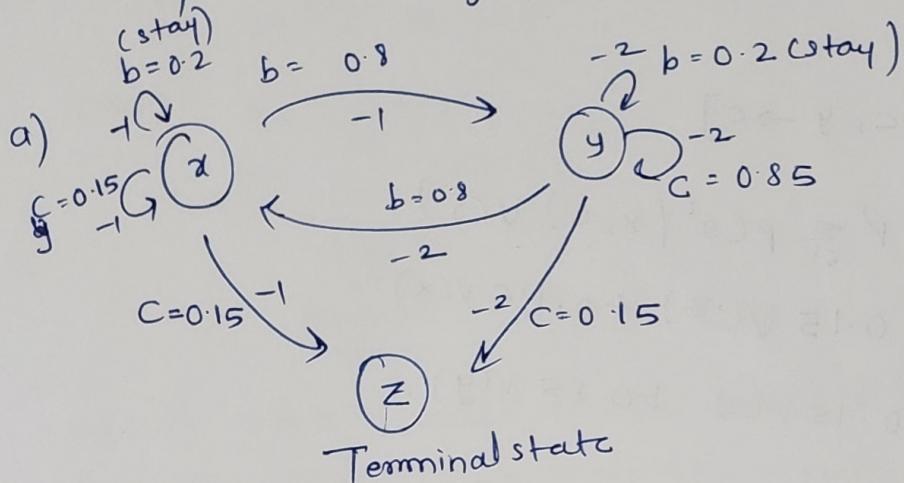
else go to 2.

$$\rightarrow \text{If } \left(\sum_{s', r} p(s', r | s, \pi(a)) [r + \gamma \sum_{a'} \pi(a' | s') q(a', s')] \neq \sum_{s', r} p(s', r | s, \text{old-action}) [r + \gamma \sum_{a'} \pi(a' | s') q(a', s')] \right)$$

b)

$$q_{t+1}(s, a) \leftarrow \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_t(s', a')]$$

Q4 Policy - iteration by hand,



- agent's goal is to reach $\rightarrow Z \rightarrow$ Terminal state
- from any state $X/Y \rightarrow$ prob \rightarrow Terminal state $= C = 0.15$
- considering reward values,
 - State X \rightarrow receive -1
 - State Y \rightarrow receive -2
- To minimize -ve rewards \rightarrow agent should stay at X.
- optimal policy would be-
 - starting in Y \rightarrow move to X \rightarrow take C from X \hookrightarrow Goal
OR
 - starting in X \rightarrow move to C \rightarrow reach goal

(B) iteration 1

$$V(x) = -1 + 0.15 V(z) + 0.85 V(x)$$

$$V(y) = -2 + 0.2 V(y) + 0.8 V(x).$$

$$V(z) = 0$$

$$V(x) = \frac{-1}{0.15} = \frac{-20}{3} = -6.67$$

$$0.8 V(y) = -2 + 0.8 \left(\frac{20}{3} \right)$$

$$V(y) = \frac{-7.33}{0.8} = -9.1625$$

Policy improvement

$$\begin{aligned} \pi(x) &= \arg \max_a \begin{cases} -1 + 0.8 V(y) + 0.2 V(x) \\ -1 + (0.15 \times 0) + 0.85 V(x) \end{cases} \\ &= \arg \max_a \begin{cases} -1 + 0.8(-9.16) + 0.2(-6.67) \\ -1 + 0.85(-6.67) \end{cases} \end{aligned}$$

$$\pi(x) = c$$

$$\begin{aligned} \pi(y) &= \arg \max_a \begin{cases} -2 + 0.8 (V(x) + 0.2 V(y)) \\ -2 + 0.85 V(y) \end{cases} \\ &, \arg \max_a \begin{cases} -2 + 0.8 (-6.67) + 0.2 (-9.16) \\ -2 + 0.85 (-9.16) \end{cases} \end{aligned}$$

$$\pi = [x \rightarrow c, y \rightarrow b] \rightarrow \text{same as prev. policy is converged.}$$

c) Initial policy $\pi = [x \rightarrow b, y \rightarrow b]$

$$V(x) = -1 + 0.8 V(y) + 0.2 V(x)$$

$$V(y) = -2 + 0.8 V(x) + 0.2 V(y)$$

$$V(z) = 0$$

$$0.8 V(x) = -1 + 0.8 V(y)$$

$$0.8 V(y) = -2 + 0.8 V(x)$$

$$0.8 V(x) = -1 + [-2 + 0.8 V(x)]$$

$0 = -3$ which is inconsistent.

Our initial policy always takes action b which causes it to loop between states x & y .

Hence we are taking undiscounted MDP,
will get sum of ∞ -ve no.

• Does discounting help?

→ discounting factor can help in this case.

As values in future will not be significant.

→ it will be difficult to achieve optimal solution as problem changes. Also, MDP optimal policy is independent of discounting factor, policy might be able to converge.

Q.7) Proving convergence of value iteration.

for any $f^n \rightarrow f \wedge g$.

$$|\max_a f(a) - \max_a g(a)| \leq \max_a |f(a) - g(a)|$$

Hint: consider the cases $\max_a f(a) - \max_a g(a) \geq 0$.
 $\& \max_a f(a) - \max_a g(a) < 0$.

→ ① Case 1 →

$$\max_a f(a) - \max_a g(a) \geq 0. \quad \text{---(1)}$$

Let, $\max_a f(a) = a^* f$, $\max_a g(a) = a^* g$

$$f(a^*_f) - g(a^*_g) \geq 0.$$

$$\therefore |f(a^*_f) - g(a^*_g)| \leq \max_a |f(a) - g(a)| \quad \text{---(2)}$$

$g(a^*_f) \leq g(a^*_g) \rightarrow a^*_g \rightarrow \text{optimal action for } g.$

from eqn (1),

$$f(a^*_f) - g(a^*_f) \geq f(a^*_f) - g(a^*_g) \geq 0.$$

$$f(a^*_f) - g(a^*_g) = |f(a^*_f) - g(a^*_g)|$$

$$f(a^*_f) - g(a^*_g) \leq f(a^*_f) - g(a^*_f)$$

$$= |f(a^*_f) - g(a^*_g)| \leq \max_a |f(a) - g(a)|$$

2nd case \rightarrow

$$\max_a f(a) - \max_a g(a) = f(a_f^*) - g(a_g^*) < 0$$

$$f(a_g^*) < f(a_f^*)$$

$$\rightarrow g(a_g^*) - f(a_g^*) \geq g(a_g^*) - f(a_f^*) > 0$$

$$g(a_g^*) - f(a_g^*) = |g(a_g^*) - f(a_g^*)|$$

$$|g(a_g^*) - f(a_f^*)| \leq |g(a_g^*) - f(a_g^*)|$$

$$= |g(a_g^*) - f(a_g^*)| \leq \max_a |g(a) - f(a)|$$

$$= \max_a |f(a) - g(a)|$$

b) $|\max_a f(a) - \max_a g(a)| \leq \max_a |f(a) - g(a)|$

$V_k \leftarrow$ 12th value iteration

$B \leftarrow$ Bellman Backup operator

$V_{k+1} \leftarrow B V_k$

$$V_{k+1}(s) \leftarrow \max_a p(s', r | s, a) [r + \gamma V_k(s')]$$

$\|V\|_\infty = \max$ absolute value of it's element.

→ from given definition ~~in~~ in the problem,

$$\begin{aligned}\|\beta v_i - \beta v'_i\|_\infty &= \|R(s) + \max_a \gamma \sum p(s'|s,a) v(s') \\ &\quad - R(s) - \max_a \gamma \sum p(s'|s,a) v'(s')\|_\infty \\ &= \left\| \max_a \gamma \sum p(s'|s,a) \cdot v(s') - \max_a \gamma \sum p(s'|s,a) v'(s') \right\|\end{aligned}$$

$$\begin{aligned}\rightarrow &\leq \max_a \left| \gamma \sum p(s'|s,a) v(s') - v \sum p(s'|s,a) v'(s') \right| \\ &\leq \max_a \gamma \sum p(s'|s,a) \|v - v'\| \\ &= \max_a \gamma \|v - v'\| \\ &\leq \sum_{s'} p(s'|s,a) = 1 \\ \therefore \max_a \gamma \|v - v'\| &= \gamma \|v - v'\| \\ \therefore \|\beta v_i - \beta v'_i\|_\infty &\leq \gamma \|v_i - v'_i\|_\infty\end{aligned}$$