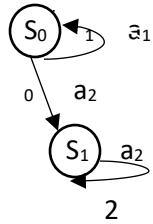


Assignment 3

Bellman optimality equation for a finite-horizon MDP :



$T = 1$; Reward = 1, Optimum Value [1 2]

$T = 2$; Optimum Value [2 2]

Bellman Optimality Equation for a finite horizon MDP,

$$V_t(s) = \max_a (r(s,a) + \sum P(s, a, s') V_{t+1}(s')) \quad , \text{ where } \max_a = \text{best action}$$

Problem in state, $S \rightarrow$ action, $a \rightarrow$ goto s' . $P(s,a,s')$

Backward Function :

$t \rightarrow t+1$

Objective :

R_t – Random variable of reward at time ‘t’

$$E [\sum_{t=1}^T R_t]$$

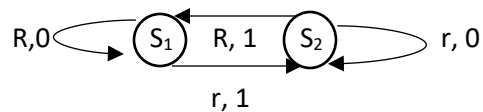
Part b:

Different Policies : $\pi : 1)$ History Dependent

2) Deterministic Markov policy

$$\pi: \mathcal{S} \rightarrow a \quad (\text{Randomized Policy})$$

\rightarrow To check which policy we are in,



Probability : r (right) : 0 1

l (left) : 0 1

$\{ \frac{1}{2} r \text{ and } \frac{1}{2} l \}$ Half the time right and half left, reward greater than ‘1’.

Randomized policy :

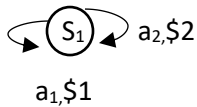
$$\Pi : \mathcal{F} \rightarrow \Delta^a$$

$$\mathcal{F} * a [0,1]$$

$$V_t(s) = r(s, \Pi(s)) + \sum_{s'} P(s, \Pi(s), s') V_{t+1}(s')$$

If horizon applied to Infinite loops, MDPs applied in discounted state,

Discounted MDPs :



For example: Day1 Day2 Day3 a_2 is not best since Infinite

$$\begin{array}{lll} \$1 & \$1.25 & \$1.25^2 + 1.25 \end{array}$$

$$\begin{array}{lll} \$2 & \$2.50 & \$2.50^2 + 2.50 \end{array}$$

Predicting in backward direction (right to left)

$$\$1 \leftarrow \$5$$

$$\$0.2 \leftarrow \$1$$

$$E \left[\sum_{t=1}^{\infty} r^t R_t \right]$$

$$a_1: \sum_{t=0}^{\infty} r^t = 1 / (1-r) \quad r = 0.9, 1 \rightarrow 0.9$$

$$a_2: 2 \sum_{t=0}^{\infty} r^t = 2 / (1-r) \quad 2 \rightarrow 1.8$$

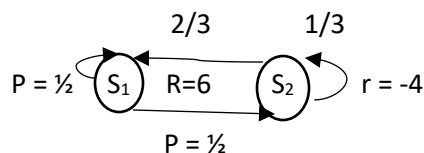
Predicting 10%

$$\text{Geometric Series, } ([1+r+r^2+r^3+r^4] (1-r)) / (1-r)$$

Bellman's

$$V(s) = \max_a (v(s,a) + \gamma \sum_{s'} p(s,a,s') v(s')), \text{ where } \gamma = \text{discount}$$

$$\Pi(V_{\Pi}(S)) = r(s, \Pi(s)) + \gamma \sum_{s'} P(S, \Pi(s), s') v_{\Pi}(s')$$



$$V(S_1) = 6 + \frac{1}{2} V(S_1) + \frac{1}{2} V(S_2)$$

$$V(S_2) = -4 + \frac{2}{3} (V(S_1) + \frac{1}{3} V(S_2))$$

2 equations, 2 values, 1 V value

Writing in vector-matrix :

From equations, $V = V + \gamma P_v$, where v = vector

$$P = \begin{bmatrix} 1/2 & 1/2 \\ 2/3 & 1/3 \end{bmatrix} \quad r = \begin{bmatrix} 6 \\ -4 \end{bmatrix}$$

$$P \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Matrix, multiply by 1, should get 1

Backward multiplication,

$$\alpha = \frac{1}{0} \quad \alpha^T P = [1/2 \quad 1/2]$$

Probability of matrices,

$$V - \gamma P_v = r$$

$$(I - \gamma P) v = r$$

$$V = (I - \gamma P)^{-1} r \quad I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Methods include using LV, QR, SVD, Gauss sieidel except Inverse which should not be used.

$$\Rightarrow \text{Eigenvalue} \neq 0$$

$$\Rightarrow (P \in (-\gamma, \gamma))$$

$$\Rightarrow (I - \gamma P) \in (1 - \gamma, 1 + \gamma)$$

$$\Rightarrow \gamma = \frac{21}{12}$$

$$P_0 = \frac{1}{0} \quad P = 21$$

$$P_0 = \frac{0}{1} \quad p = 12$$

$$P_0 = \frac{2/3}{1/3} \quad P = p_0^T v \quad \sum_{t=0}^{\infty} E[r^t R_t], \text{ where } p_0^T v = \text{reward}$$

$$P_1 = p_0^T P = \frac{1/2}{1/2}, \text{ where } P_1^T r = \text{reward}$$

$$P_3 = p_2^T P = P_1 P^3, \text{ where } P_1 P^3 = \text{reward}$$

$$p_0^T \sum_{t=0}^{\infty} r^t p^t r$$

$$P_0^T r + \gamma P_0^T P r + \gamma^2 P_0^T P^2 r + \dots$$

$$\sim 1 + \gamma + \gamma^2 + \gamma^3 + \dots = (1 - \gamma)^{-1}$$

Neumann series, $= P_0^T(I - \gamma P)^{-1} v$, where $(I - \gamma P)^{-1} v = P_0^T r$

$$3) V - \gamma P_v = r$$

Difference between a stationary and Markovian policy

a) Stationary policies are not optimal always. Used if we want to continue taking actions indefinitely. A stationary policy is a mapping from states to actions.

$\pi: S \rightarrow A$, $\pi(s)$ is action to do at state s (regardless of time)

b) Markovian policy is not dependent on history

$\Pr(S_{t+1} | A_t, S_t, A_{t-1}, S_{t-1}, \dots, S_0) = \Pr(S_{t+1} | A_t, S_t)$, Next state only depends on current state and current action.

Is the optimal policy in finite-horizon MDP guaranteed to be Markovian or stationary (or both)?

a) To be optimal under finite horizon, a non-stationary policy should be used, hence it is guaranteed to be Markovian.

What happens to the value function when you add a constant to all the rewards?

Only the Relative size of the rewards are vital, for $K = 0$, we have value of state : $V = \sum_{t=0}^{\infty} \gamma^{(t-t_0)} r_t$. A value for K gives,

$$\sum_{t=0}^{\infty} \gamma^{(t-t_0)} (r_t + K) = V + (K/(1 - \gamma)).$$

In a sequential task, the sum does not move to infinity, the result is less simple in particular if the length is not fixed.

What happens to the policy. How do the answers differ for finite and infinite-horizon problem?

For the policy, the constant does not change the relative values, only the relative size is vital for all the policies. For example, adding a constant 'c' to all the rewards adds a constant, V_c to all the values of all states and this does not affect the relative values of any states under policies. In an episodal task, The result will be simple in particular if the length of the episodes is not fixed.