

n-state MDP

(a) $V^*(n) = ?$

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \cdot V^*(s')]$$

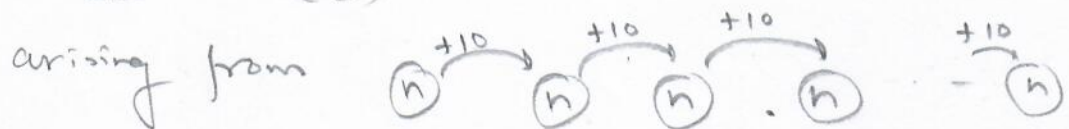
At state n , there is only one action (Go Right).

$$V^*(n) = 1.0 \times [+10 + 0.5 \times V^*(n)]$$

$$\Rightarrow V^*(n) = \boxed{20}$$

Another approach to arrive at this conclusion is to view $V^*(n)$ as the maximum expected discounted reward that the agent can gather from state n .

$$\Rightarrow 10 + 10 \times \left(\frac{1}{2}\right) + 10 \times \left(\frac{1}{2}\right)^2 \dots \text{discounting}$$



$$\Rightarrow 10 \times \left(1 + \frac{1}{2} + \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^3 \dots \right)$$

$$\Rightarrow 10 \times \left[\frac{1}{1 - (1/2)} \right] = \boxed{20}$$

(b) Optimal value function $V^*(k)$ for all $k \in \{1, 2, \dots, (n-1)\}$
 Given $V^*(k) < V^*(k+1)$ for this MDP

Examine,

$$V^*(n-1) = \max_{\{GoRight, Reset\}} \begin{cases} \text{For GoRight: } 1 + \frac{1}{2} \times V^*(n) \\ \text{For Reset: } 0 + \frac{1}{2} \times V^*(1) \end{cases}$$

Since, $V^*(1) < V^*(n)$

$$\text{Hence, } \max \left[\left(1 + \frac{1}{2} \times V^*(n) \right), \left(\frac{1}{2} \times V^*(1) \right) \right]$$

$$\text{will return } \left[1 + \frac{1}{2} \times V^*(n) \right]$$

Similarly,

$$V^*(n-2) = \max_{\{GoRight, Reset\}} \begin{cases} \text{For GoRight: } 1 + \frac{1}{2} \times V^*(n-1) \\ \text{For Reset: } 0 + \frac{1}{2} \times V^*(1) \end{cases}$$

Since, $V^*(1) < V^*(n-1)$

$$\text{Hence, } \max \left[\left(1 + \frac{1}{2} \times V^*(n-1) \right), \left(\frac{1}{2} \times V^*(1) \right) \right]$$

$$\text{will return } \left[1 + \frac{1}{2} \times V^*(n-1) \right]$$

$$\begin{aligned} V^*(n-2) &= 1 + \frac{1}{2} \times V^*(n-1) = 1 + \frac{1}{2} \left(1 + \frac{1}{2} \times V^*(n) \right) \\ &= 1 + \frac{1}{2} + \left(\frac{1}{2} \right)^2 V^*(n) \end{aligned}$$

$$V^*(n-k) = 1 + \left(\frac{1}{2} \right) + \left(\frac{1}{2} \right)^2 + \dots + \left(\frac{1}{2} \right)^{k-1} + \left(\frac{1}{2} \right)^k V^*(n)$$

$$= \left[\frac{1 - \left(\frac{1}{2} \right)^k}{1 - \left(\frac{1}{2} \right)} \right] + \left(\frac{1}{2} \right)^k V^*(n)$$

(c) Value Iteration initialised to zero for all states.

After 1st iteration:

$$\bullet V_1(n) = 10 + \frac{1}{2} \times V_0(n) = 10 + \frac{1}{2} \times 0 = \boxed{10}$$

$$\bullet V_1(n-1) = \max_{\{GoRight, Reset\}} \left\{ \begin{array}{l} \text{For GoRight: } 1.0 + \frac{1}{2} \times V_0(n) \\ \text{For Reset: } 0.0 + \frac{1}{2} \times V_0(1) \end{array} \right.$$
$$= \max \left\{ \begin{array}{l} 1.0 \\ 0.0 \end{array} \right\}$$

$$\Rightarrow \boxed{V_1(n-1) = 1}$$

Similarly, for all states $k \in \{1, \dots, n-1\}$,

$$\boxed{V_1(n-k) = 1}$$

After 2nd iteration:

$$\bullet V_2(n) = 10 + \frac{1}{2} \times 10 = \boxed{15}$$

$$\bullet V_2(n-1) = \max_{\{GoRight, Reset\}} \left\{ \begin{array}{l} \text{For GoRight: } 1 + \frac{1}{2} \times V_1(n) \\ \text{For Reset: } 0 + \frac{1}{2} \times V_1(1) \end{array} \right.$$

$$= \max_{\{GoRight, Reset\}} \left\{ \begin{array}{l} \text{For GoRight: } 1 + \frac{1}{2} \times 10 \\ \text{For Reset: } 0 + \frac{1}{2} \times 1 \end{array} \right.$$

$$= \max \left\{ \begin{array}{l} 6 \\ 1 \end{array} \right\}$$

$$\Rightarrow \boxed{V_2(n-1) = 6}$$

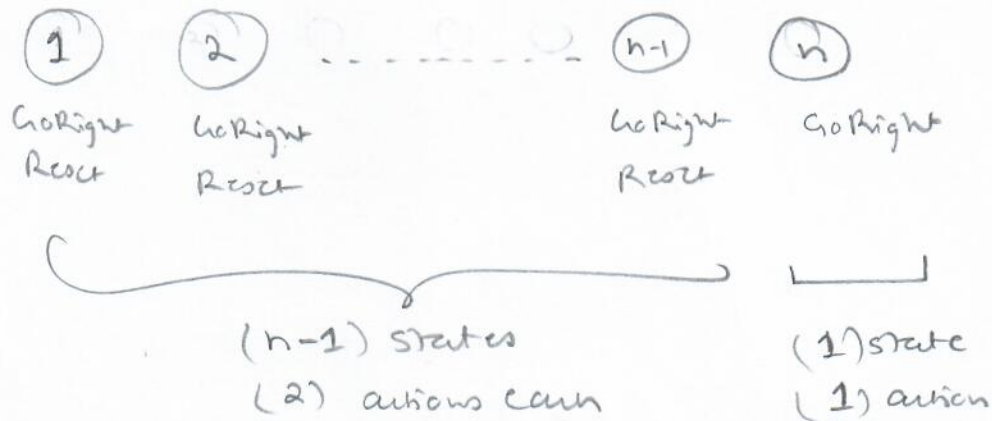
$$\bullet V_2(n-2) = \max_{\{GoRight, Reset\}} \left\{ \begin{array}{l} \text{For GoRight: } 1 + \frac{1}{2} \times V_1(n-1) \\ \text{For Reset: } 0 + \frac{1}{2} \times V_1(1) \end{array} \right.$$

$$= \max_{\{GoRight, Reset\}} \left\{ \begin{array}{l} \text{For GoRight: } 1 + \frac{1}{2} \times 1 \\ \text{For Reset: } 0 + \frac{1}{2} \times 1 \end{array} \right.$$

$$V_2(n-2) = \max\{1.5, 0.5\} = \boxed{1.5}$$

The same for $V_2(n-k) = 1.5 \quad \forall k \in \{1, \dots, n-1\}$

(d) How many policies are there for this MDP?



of policies $2^{n-1} \times 1 = \boxed{2^{n-1}}$

An optimal policy:

$\pi(k) = \text{Go Right} \quad \text{for } k \in \{1, 2, \dots, n\}.$