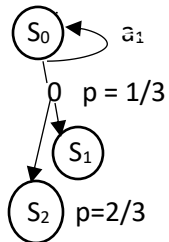


Assignment 5

MDP : $r(s,a)$

$$r_2(s,a,s') = r_1(s,a)$$



$$r_1(s,a,s_1) = 5$$

$$r_2(s,a,s') = 10$$

$$r(s,a_1) = 5 \cdot (1/3) + 10 \cdot (2/3)$$

$$= 25/3$$

$$V(s) = \max_a r(s,a) + \gamma \sum p(s,a,s') \cdot r(s')$$

$$V(s) = \max_a \sum_{s'} p(s,a,s') \cdot (r(s,a,s') + \gamma V(s'))$$

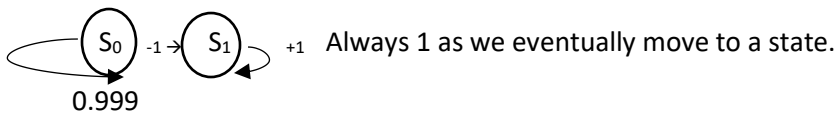
Average Reward MDPs:

$$\text{Discounted } (\sum_{t=0}^{\infty} \gamma^t E[R_t])$$

Average Reward MDPs

$$\lim_{T \rightarrow \infty} (1/T) (\sum_{t=1}^T E[R_t]) \quad (1/1) / (2/2) / (3/3) / \dots$$

$$S_0 \rightarrow S_1 \quad (-1000/1) / (-991/2) / (-998/3) / \dots$$



Gain ~ Value function

Gain(s) = average reward starting from state s.

$$G(s) = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T E[R_t]$$

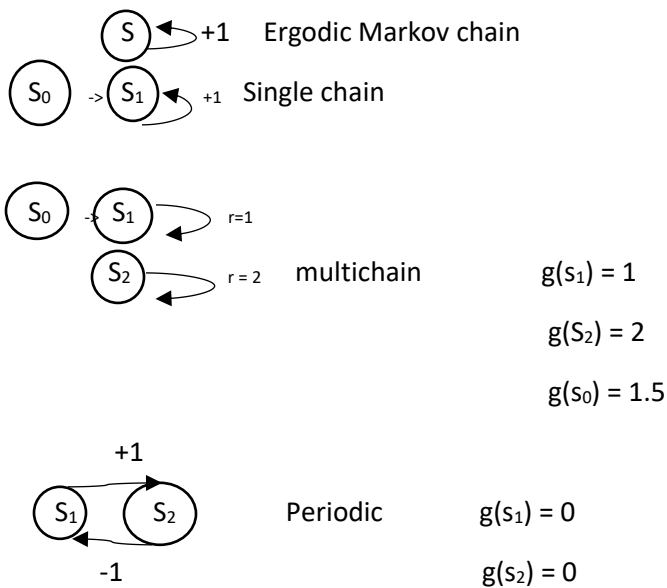
$$G(s_1) = 1$$

$$G(s_0) = 1$$

1. Base stock is maximum mandatory value that one has to hold in inventory management. For example, if the base stock is 35, and end of day count is 25, one has to order 10 more to satisfy base stock policy.
2. When we increase the discount factor, the order amount decreases, very marginally.
3. Computation time decreases when we increase the discount factor

Advantage of TD(0) over Monte-carlo algorithm:

- TD methods do not require a model of the environment, only experience
- TD, but not MC, methods can be fully incremental.

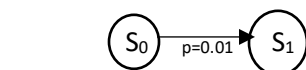


$$P = p_0^T p^\infty$$

$$p^1 = p$$

$$P^2 = p.p$$

$$P^3 = p.p.p$$



$$P = \begin{pmatrix} 0.99 & 0.01 \\ 0 & 1 \end{pmatrix}$$

$$\sum_{t=0}^{\infty} \gamma^t p^t = (I - \gamma P)^{-1}$$

$$P^\infty = \lim_{T \rightarrow \infty} (1/T) \sum_{t=0}^T p^t = (1/T) (I - P)^{-1} \text{ (Inverse doesn't exist)}$$

Stochastic matrix has an eigen value 1.

$$P_s^T = P_s^T P$$

$$P_s = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$