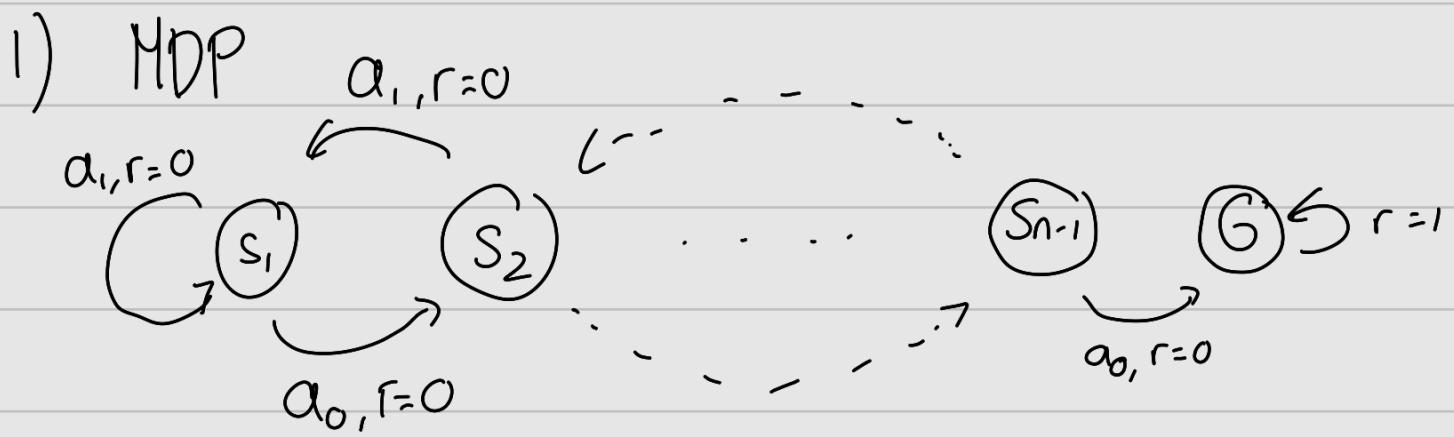


Reinforcement Learning

Aditya Kelvianto Siahaan

aks 2266

Homework 1



- Starting from s_1 , move left (a_1) or right (a_0)
- action deterministic $p(s_2 | s_1, a_0) = 1$
- reward only if reach G , $r < 1$

a) optimal policy at any state $s_i \neq G$

$$\pi(s_i) = a_0 \quad \forall i = 1 \dots n-1$$

optimal value function

$$V^*(G) = 1 + r + r^2 + r^3 + r^4 + \dots = \frac{1}{1-r}$$

$$V^*(S_{n-1}) = 0 + r \frac{1}{1-r} = \frac{r}{1-r}$$

$$V^*(S_{n-2}) = 0 + r \times 0 + r^2 \times \frac{1}{1-r} = \frac{r^2}{1-r}$$

$$V^*(S_{n-k}) = \frac{r^k}{1-r} \quad \forall k = 1 \dots n-1$$

$$V^*(S_1) = \frac{r^{n-1}}{1-r}$$

f) Assuming that the discount factor is less than 1, then the optimal policy will not be changed, because the value function will not change, i.e. $Q(s_k, a_0) > Q(s_k, s_1)$
 $\forall k = 1 \dots n-1$ because the longer it reaches goal state, the more discounted the reward will be.

$$V^*(s_1) < V^*(s_2) < \dots < V^*(G)$$

$\forall r < 1$

Nevertheless, if the discount factor is 1, then $Q(s_k, a_0) = Q(s_k, a_1) \quad \forall k = 1 \dots n-1$, because it no longer matters when do we actually reach G, as $V^*(s_k) = Q(s_k, a_j) = \infty \quad \forall k = 1 \dots n-1, j = 0, 1$.
 Thus, the optimal policy will for all s_k states will be both a_0 and a_1 .

$$V^*(s_1) = V^*(s_2) = \dots = V^*(G)$$

if $r = 1$

c) After adding constant c to all rewards,
because $c < 1+c \neq c$



The Optimal Value for each of the state will be following

$$V^*(G) = 1+c + \gamma(1+c) + \gamma^2(1+c) + \dots \\ = \frac{1+c}{1-\gamma}$$

$$V^*(S_{n-1}) = c + \gamma \frac{1+c}{1-\gamma}$$

$$V^*(S_{n-2}) = c + \gamma \times c + \gamma^2 \times \frac{1+c}{1-\gamma}$$

$$V^*(S_{n-k}) = \left(\sum_{i=0}^{k-1} \gamma^i c \right) + \gamma^k \frac{1+c}{1-\gamma}$$

$$V^*(S_1) = \left(\sum_{i=0}^{n-2} \gamma^i c \right) + \gamma^{n-1} \left(\frac{1+c}{1-\gamma} \right)$$

By adding constant c to all of the states' rewards, we actually do not change the relative reward between any of the states, i.e

$$V^*(s_1) < V^*(s_2) < \dots < V^*(G)$$
$$\forall r < 1$$

Thus, the optimal policy will still be

$$\pi(s_i) = a_0 \quad \forall i = 1 \dots n-1$$

In other words,

$$Q(s_i, a_i) < Q(s_i, a_0)$$

$$\forall i = 1 \dots n-1$$

d) After scaling the reward with constant α ,



when $\alpha > 0$, $\alpha c < \alpha + \alpha c$

$\alpha = 0$, $\alpha c = \alpha c$

$\alpha < 0$, $\alpha c > \alpha + \alpha c$

Intuitively speaking, when $\alpha = 0$ or $\alpha < 0$, rushing to the goal state is no longer the optimal policy!

Lets explore each of the possible values of α .

① When $\alpha \neq 0$:

$$V^*(G) = 0$$

$$V^*(S_i) = 0 \quad \forall i = 1 \dots n-1$$

Therefore

$$\pi^*(S_i) : \operatorname{argmax}_a Q^*(S_i, a) = a_0 \text{ or } a, \quad \forall i = 1 \dots n-1,$$

in other words, any policy chosen by the system is an

optimal policy, as it maximizes the reward of the system which is zero throughout.

② When $\alpha > 0$

$$V^*(G) = \alpha(1+c) + \gamma(1+c) + \gamma^2(1+c) + \dots \\ = \frac{\alpha(1+c)}{1-\gamma}$$

$$V^*(S_{n-1}) = ac + \gamma \frac{\alpha(1+c)}{1-\gamma}$$

$$V^*(S_{n-2}) = ac + \gamma ac + \gamma^2 \frac{\alpha(1+c)}{1-\gamma}$$

$$V^*(S_{n-k}) = \left(\sum_{i=0}^{k-1} \gamma^i ac \right) + \gamma^k \frac{\alpha(1+c)}{1-\gamma}$$

$$V^*(S_1) = \left(\sum_{i=0}^{n-2} \gamma^i ac \right) + \gamma^{n-1} \frac{\alpha(1+c)}{1-\gamma}$$

In this setting,

$$Q(s_i, a_1) < Q(s_i, a_0)$$

$$\forall i = 1 \dots n-1$$

It also can be seen from the equation above that

$$V^*(s_1) < V^*(s_2) < \dots < V^*(G)$$

$\text{if } r < 1$

Thus, the optimal policy will still be

$$\pi(s_i) = a_0 \quad \text{if } i = 1 \dots n-1$$

when $a < 0$, the optimal value function is

$$\begin{aligned} V^*(G) &= a(1+c) + ar(1+c) + ar^2(1+c) + \dots \\ &= \frac{a(1+c)}{1-r} \end{aligned}$$

$$\begin{aligned} V^*(s_1) &= V^*(s_2) = V^*(s_{n-1}) \\ &= ac + r ac + r^2 ac + \dots \\ &= \frac{ac}{1-r} \end{aligned}$$

Intuitively speaking, the negative reward obtained by staying in the state $s_1 \dots s_{n-1}$ is much smaller than

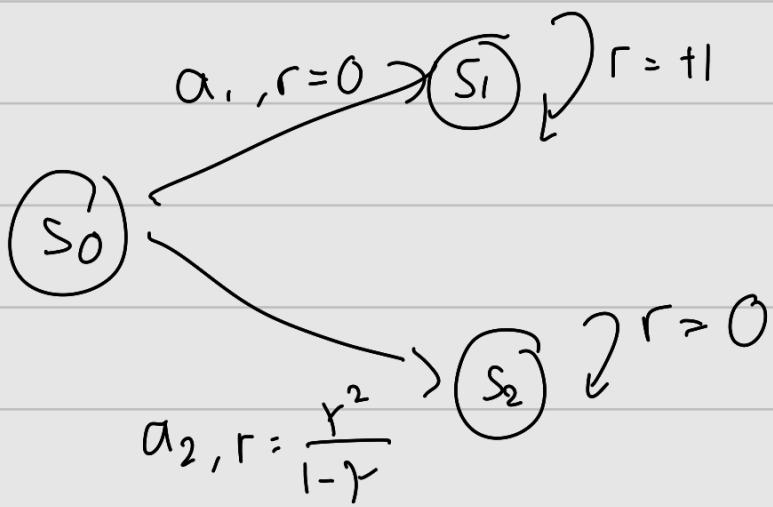
negative reward obtained by reaching the goal state G . Therefore,
the optimal solution is as follows:

$$\pi^*(s_k) = a_0 \text{ or } a_1 \quad \forall k = 1 \dots n-2$$

$$\pi^*(s_{n-1}) = a_1$$

2

- Infinite MDP, discount factor $\gamma \leq 1$
- Rewards are given upon taking action from the state



a) discounted reward of taking action a_1 :

$$\begin{aligned} Q(S_0, a_1) &= 0 + \gamma + \gamma^2 + \gamma^3 + \dots \\ &= \frac{\gamma}{1-\gamma} \end{aligned}$$

b) discounted reward of taking action a_2

$$\begin{aligned} Q(S_0, a_2) &= \frac{\gamma^2}{1-\gamma} + \gamma \cdot 0 + \gamma^2 \cdot 0 + \dots \\ &= \frac{\gamma^2}{1-\gamma} \end{aligned}$$

The optimal action on time step 0 is
to take action a_1 ,

c) Assume that $\alpha \in V_{=0}, V_{n_0}(s) = 0$
 Show that performing Value iteration
 will always choose suboptimal action n^*

$$V_{n_0}(s_0) = 0, V_{n_0}(s_1) = 0, V_{n_0}(s_2) = 0$$

We will always choose a_2 in the start of the iteration
 because $V(s_1)$ is underrepresented. Therefore, for
 the optimal action a_1 to be chosen, at iteration n^* ,

$$\pi^*(s_0) = \underset{a}{\operatorname{argmax}} Q^*(s_0, a) = a_1$$

in other words,

$$Q_{n^*}(s_0, a_1) > Q_{n^*}(s_0, a_2)$$

$$r(s_0, a_1) + \gamma V_{n^*}(s_1) > r(s_0, a_2) + \gamma V_{n^*}(s_2)$$

When $n=1$

$$Q_{n_1}(s_0, a_1) = r(s_0, a_1) + \gamma V_{n_0}(s_1) = 0$$

$$Q_{n_1}(s_0, a_2) = r(s_0, a_2) + \gamma V_{n_0}(s_2)$$

$$= \frac{\gamma^2}{1-\gamma}$$

$$\begin{aligned} V_{n_1}(s_1) &= r(s_1, a) + \gamma V_{n_1}(s_1) \\ &= 1 + 0 = 1 \end{aligned}$$

$$V_{n_1}(s_2) = 0$$

When $n=2$

$$Q_{n_2}(s_0, a_1) = 0 + \gamma \cdot 1 = \gamma$$

$$Q_{n_2}(s_0, a_2) = \frac{\gamma^2}{1-\gamma}$$

$$\begin{aligned} V_{n_2}(s_1) &= r(s_1, a) + \gamma V_{n_2}(s_1) \\ &= 1 + \gamma \end{aligned}$$

$$V_{n_2}(s_2) = 0$$

when $n=3$

$$Q_{n_3}(S_0, a_1) = 0 + \gamma(1+\gamma) = \gamma + \gamma^2$$

$$Q_{n_3}(S_0, a_2) = \frac{\gamma^2}{1-\gamma}$$

$$V_{n_3}(S_1) = 1 + \gamma(1+\gamma) = 1 + \gamma + \gamma^2$$

$$V_{n_3}(S_2) = 0$$

when $n=4$

$$Q_{n_4}(S_0, a_1) = 0 + \gamma(1+\gamma+\gamma^2) = \gamma + \gamma^2 + \gamma^3$$

$$Q_{n_4}(S_0, a_2) = \frac{\gamma^2}{1-\gamma}$$

$$V_{n_3}(S_1) = 1 + \gamma + \gamma^2 + \gamma^3$$

$$V_{n_3}(S_2) = 0$$

Thus, when $n=k$

$$Q_{n_k}(S_0, a_1) = \sum_{i=1}^{i=k-1} \gamma^i$$

$$Q_{n_k}(S_0, a_2) = \frac{\gamma^2}{1-\gamma}$$

$$V_{n_k}(S_1) = 1 + \sum_{i=1}^{i=k-1} \gamma^i$$

$$V_{n_k}(S_2) = 0$$

$$\sum_{i=1}^{k-1} r^i \geq \frac{r^2}{1-r}$$

$$\sum_{i=0}^{k-1} r^i - 1 \geq \frac{r^2}{1-r}$$

$$\frac{r^{k-1}}{1-r} \geq \frac{r^2}{1-r}$$

$$r^k \geq \frac{r^2}{1-r} (r-1) + r$$

$$r^k \geq 1-r$$

$$\log(r^n) \geq 1-r$$

$$k \log r \geq \log(1-r)$$

$$n^* \text{ or } k \geq \frac{\log(1-r)}{\log r} \quad \checkmark$$

3a) The value algorithm that we use indicates that we stop updating the value function once we reach the following threshold

$$\|v_{k+1} - v_k\| < \epsilon \cdot \frac{1-\gamma}{2\gamma}$$

Theorem A) If we do this, then it is guaranteed that $\|v_{k+1} - v^*\| < \epsilon$

Proof:

$$\begin{aligned} \|v^{\pi_k} - v^*\| &\leq \|v^{\pi_k} - v_{k+1}\| + \\ &\quad \|v_{k+1} - v^*\| \end{aligned}$$

(1)

where $k+1$ is a point of convergence.

Expressing it as a vector notation,

$$L_\pi v_{k+1} = L v_{k+1}$$

Therefore,

$$\begin{aligned}\|V^{\pi_k} - V_{k+1}\| &\leq \|L_\pi V^{\pi_k} - L_{V_{k+1}}\| \\ &\quad + \|L_{V_{k+1}} - L_{V_k}\| \\ &\leq \gamma \|V^{\pi_k} - V_{k+1}\| \\ &\quad + \gamma \|V_{k+1} - V_k\|\end{aligned}$$

$$\begin{aligned}\|V^{\pi_k} - V_{k+1}\| - \gamma \|V^{\pi_k} - V_{k+1}\| &\leq \gamma \|V_{k+1} - V_k\|\end{aligned}$$

$$\|V^{\pi_k} - V_{k+1}\| \leq \frac{\gamma}{1-\gamma} \|V_{k+1} - V_k\| \quad (2)$$

From (1), it also follows that

$$\|V_{k+1} - V^*\| \leq \frac{\gamma}{1-\gamma} \|V_{k+1} - V_k\|$$

Adding the two terms up,

$$\|V^{\pi_k} - V^*\| \leq \frac{2\gamma}{1-\gamma} \|V_{k+1} - V_k\| < \varepsilon$$

Using Theorem ④,

if $\|v_{k+1} - v_k\| < \varepsilon \cdot \frac{1-\gamma}{\gamma}$,
then $\|v_{k+1} - v^*\| \leq \varepsilon$

letting $\varepsilon = \frac{\gamma^n}{1-\gamma} \|v_1 - v_0\|$

$$\|v_{k+1} - v_k\| \leq \gamma^{k-1} \quad (3)$$
$$\|v_{k+1} - v^*\| \leq \frac{\gamma^k}{1-\gamma} \|v_1 - v_0\| \quad *$$

$$3b) \|v^{\pi_k} - v^*\| \leq \|v^{\pi_k} - v^{k+1}\| + \|v^{k+1} - v^*\| \quad \text{from (1)}$$

from (2)

$$\|v^{\pi_k} - v_{k+1}\| \leq \frac{\gamma}{1-\gamma} \|v_{k+1} - v_k\|$$

from (3)

$$\|v_{k+1} - v^*\| \leq \frac{\gamma}{1-\gamma} \|v_1 - v_0\|$$

Therefore

$$\| V^{\pi_k} - V^* \| \leq \frac{\gamma}{1-\gamma} \| V_{k+1} - V_k \|$$

$$+ \frac{\gamma^k}{1-\gamma} \| V_1 - V_0 \|$$

$$\leq \frac{\gamma^k}{1-\gamma} \| V_1 - V_0 \|$$

$$+ \frac{\gamma^k}{1-\gamma} \| V_1 - V_0 \|$$

$$\| V^{\pi_k} - V^* \| \leq \frac{2\gamma^k}{1-\gamma} \| V_1 - V_0 \|$$