



CS234 - Winter 2023

Assignment 1

1) Reward Choices

$$a) \hat{r}(s, a) = r(s, a) + c$$

New Optimal Value function

$$\begin{aligned} V'(s) &= \max_a E_{\pi_\sigma} \left[\sum_{k=0}^{\infty} \gamma^k \hat{R}_{s+k+1} \mid s_t=s, A_t=a \right] \\ &= \max_a E_{\pi^*} \left[\sum_{k=0}^{\infty} \gamma^k r(R_{s+k+1} + c) \mid s_t=s, A_t=a \right] \\ &= V(s) + \frac{c}{1-\gamma} \quad \left(c \sum_{k=0}^{\infty} \gamma^k \right) \end{aligned}$$

$$\Rightarrow V'(s) - \gamma = V(s) \quad \text{or} \quad V'(s) = V(s) + \frac{\gamma c}{1-\gamma}$$

④ The optimal policy does not change.

Because fraction, the value of state has been increased
the same amount \Rightarrow Optimal point doesn't change!

$$\begin{aligned}\pi'(s) &= \arg\max_a \sum_{s'} p(s'|s,a) (\hat{r}(s,a) + \gamma V'(s')) \\ &= \arg\max_a \sum_{s'} p(s'|s,a) (r(s,a) + c + \gamma V(s') + \frac{\gamma c}{1-\gamma}) \\ &= \pi(s) \text{ (remove the constant)}\end{aligned}$$

$$b) \quad \hat{r}(s, a) = c \times r(s, a)$$

$$V'(s) = \max_a E_{\pi_a} \left[\sum_{k=0}^{\infty} \gamma^k \hat{R}_{t+k+1} \mid S_t = s, A_t = a \right]$$

$$= \max_a E_{\pi^*} \left[\sum_{k=0}^{\infty} c \times \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

$$= c \times V(s)$$

$$\textcircled{A} \quad c > 0 \Rightarrow \pi' = \pi^*$$

$$\textcircled{B} \quad c < 0 \Rightarrow \pi' = \pi_1^* \quad (V(s, \pi_1^*) \text{ has min value})$$

$$\textcircled{C} \quad c = 0 \Rightarrow V'(s) = 0 + a \Rightarrow \text{All policies are optimal.}$$

c) If MDP is finite, does it affect the answer from part a?

Because MDP is finite, we can ignore γ

$$V'(s) = \max_a E_{\pi^*} \left[\sum_{k=0}^{M_s} \hat{R}_{t+k+1} \mid S_t = s, A_t = a \right]$$

$$= \max_a E_{\pi^*} \left[\sum_{k=0}^{M_s} R_{t+k+1} + c \mid S_t = s, A_t = a \right]$$

$$= c \cdot M_s + V(s)$$

M_s : number of steps from state s to the end state

If take γ into account

$$\rightarrow V'(s) = V(s) + c \cdot \sum_{k=0}^{M_s} \gamma^k$$

Note: M_s is not fixed if the process is stochastic

So the optimal policy will change!

2) Bellman Residuals and performance bounds

B: Bellman backup with fixed point V^* (the optimal value function)

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

B^π : Bellman backup with fixed point V^π : (value function follows policy π)

$$(B^\pi V)(s) = \underset{a \sim \pi}{\mathbb{E}} [r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V(s')]$$

We did show that $\|BV - BV'\| \leq \gamma \|V - V'\|$

Similarly, we can show $\|B^\pi V - B^\pi V'\| \leq \gamma \|V - V'\|$

a) Prove fixed point of B^π is unique

Assume $m \neq n$ satisfies $B^\pi m = m, B^\pi n = n$

$$\Rightarrow \|B^\pi m - B^\pi n\| = \|m - n\|$$

But B^π is a contraction $\Rightarrow \|B^\pi m - B^\pi n\| \leq \gamma \|m - n\| < \|m - n\|$
 (contradiction)

-1 Done

b) $(BV)(s) \geq (B^\pi V)(s) \quad \forall s \in S, V \in \mathbb{R}^{|S|}$

c) When $V = V^*$ (optimal value function) then

$$\|BV - V\| = 0$$

d)

④ Proof of $\|V - V^\pi\| \leq \frac{\|V - B^\pi V\|}{1-\gamma} + V$

We know that V^π is the fixed point of the contraction B^π

$$\Rightarrow \lim_{n \rightarrow \infty} (B^\pi)^n V = V^\pi$$

$$\underline{\text{Lemma}} : \|(B^\pi)^n v - (B^\pi)^{n-1} v\| \leq \gamma^{n-1} \|B^\pi v - v\|$$

Because B^π is a contraction

$$\begin{aligned} \Rightarrow \|(B^\pi)^n v - (B^\pi)^{n-1} v\| &\leq \gamma \|(B^\pi)^{n-1} v - (B^\pi)^{n-2} v\| \\ &\quad \dots \\ &\leq \gamma^{n-1} \|B^\pi v - v\| \end{aligned}$$

Back to the main :

$$\begin{aligned} \|v - v^\pi\| &= \lim_{n \rightarrow \infty} \|(B^\pi)^n v - v\| \\ &\leq \lim_{n \rightarrow \infty} \|(B^\pi)^n v - (B^\pi)^{n-1} v\| + \|(B^\pi)^{n-1} v - (B^\pi)^{n-2} v\| + \dots \\ &\quad + \|B^\pi v - v\| \\ &\leq \lim_{n \rightarrow \infty} (\gamma^{n-1} \|B^\pi v - v\| + \gamma^{n-2} \|B^\pi v - v\| + \dots + \|B^\pi v - v\|) \\ &< \text{Apply the lemma} \end{aligned}$$

$$= \|B^\pi v - v\| \cdot \lim_{n \rightarrow \infty} \sum_{k=0}^{n-1} \gamma^k$$

$$= \|v - B^\pi v\| \cdot \frac{1}{1-\gamma} \quad (\text{Q.E.D})$$

Similarly, because v^* is the fixed point of contraction B

$$\Rightarrow \|v - Bv^*\| \leq \frac{\|v - Bv\|}{1-\gamma}$$

e) Proof: $v^\pi(s) \geq v^*(s) - \frac{\alpha \varepsilon}{1-\gamma}$ with $\varepsilon = \|Bv - v\|$

$$(=) \frac{\alpha \varepsilon}{1-\gamma} \geq v^*(s) - v^\pi(s)$$

$$\begin{aligned}
 V^*(s) - V^\pi(s) &= (V^*(s) - V(s)) + (V(s) - V^\pi(s)) \\
 &\leq \|V^*(s) - V(s)\| + \|V(s) - V^\pi(s)\| \\
 &\leq \frac{\|BV - V\|}{1-\gamma} + \frac{\|B^\pi V - V\|}{1-\gamma} \\
 &\leq \frac{\|BV - V\|}{1-\gamma} + \frac{\|BV - U\|}{1-\gamma}
 \end{aligned}$$

(Because π is greedy policy)
 $B^\pi V$ converges faster than BV

Note : If T is the fixed point of contraction B

\Rightarrow If V then $\{V, BV, B^2V, \dots\}$ converges to T

- π^* is the optimal policy with Bellman operator B
 π is an arbitrary policy with Bellman operator B^π

\Rightarrow Convergent Speed of $\{V, BV, B^2V, \dots\} \geq \dots$ of $\{V, B^\pi V, B^{\pi^2}V, \dots\}$

f) Prove $V^\pi(s) \geq v^*(s) - \frac{\epsilon}{1-\gamma}$ if $V \leq v^*$

we have $v^*(s) - V^\pi(s) \leq V(s) - V^\pi(s)$
 $\leq \|V - V^\pi\| \leq \frac{\epsilon}{1-\gamma}$ (Q.E.D)

g) Show that if $BV \leq V$ then $v^* \leq V$
 (Từ đây suy ra π^* là BV là một tách hàn)

vì nếu không thì V^* (optimal) $\leq V$
 $(V \delta V)$

We prove $B^k V \leq V + k \in \mathbb{N}^\infty$

$$B^k(V)(s) = \max_a [r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) (B^{k-1}V)(s')]$$

$$\leq \max_a [r(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V(s')] \\ (\text{because } B^{k-1}V \leq V)$$

$$\Rightarrow B^k V \leq V$$

$$\Rightarrow V \geq \lim_{k \rightarrow \infty} B^k V = V^*$$

Because V^* is the fixed point of the contraction B .

$$\Rightarrow V \geq V^* \text{ (Q.E.D.)}$$

