# Reinforcement Learning: Assignment #3

Due on Sunday, Mai 07, 2023

**Group 4 - Abu El Komboz, Tareq 3405686** │ **Jain, Likhit 3678905** │ **Wurm, Marcel 3695946**

## Task 1

### Proofs (5 Points)

*(a) Show that the Bellman **optimality** operator $T$ is a $\gamma$-contraction. This is similar to but not the same as the Bellman **expectation backup** operator from lecture 3 slide 20. Be able to explain all the steps! (3P)*

$(Tv)(s) = \max_a \sum_{s',r} p(s',r|s,a)\left[r + \gamma v(s')\right]$

$$\left\|(Tv)(s) - (Tv')(s)\right\|_\infty$$

$$= \left\|\max_a \sum_{s',r} p(s',r|s,a)\left[r + \gamma v(s')\right] - \max_{a'} \sum_{s',r} p(s',r|s,a')\left[r + \gamma v'(s')\right]\right\|_\infty$$

$$\leq \left\|\max_a \left[\sum_{s',r} p(s',r|s,a)\left[r + \gamma v(s')\right] - \sum_{s',r} p(s',r|s,a)\left[r + \gamma v'(s')\right]\right]\right\|_\infty$$

$$= \left\|\max_a \sum_{s',r} \left[p(s',r|s,a)\left[r + \gamma v(s')\right] - p(s',r|s,a)\left[r + \gamma v'(s')\right]\right]\right\|_\infty$$

$$= \left\|\max_a \sum_{s',r} \left[p(s',r|s,a)r + p(s',r|s,a)\gamma v(s') - p(s',r|s,a)r - p(s',r|s,a)\gamma v'(s')\right]\right\|_\infty$$

$$= \left\|\max_a \sum_{s',r} \left[p(s',r|s,a)\gamma v(s') - p(s',r|s,a)\gamma v'(s')\right]\right\|_\infty$$

$$= \left\|\max_a \sum_{s',r} p(s',r|s,a)\gamma\left[v(s') - v'(s')\right]\right\|_\infty$$

$$= \gamma\left\|\max_a \sum_{s',r} p(s',r|s,a)\left[v(s') - v'(s')\right]\right\|_\infty$$

$$= \gamma \max_s \left|\max_a \sum_{s',r} p(s',r|s,a)\left[v(s') - v'(s')\right]\right|$$

$$\leq \gamma \max_s \left|\max_a \sum_{s',r} p(s',r|s,a) \max_s\left[v(s) - v'(s)\right]\right|$$

$$= \gamma \max_s \left|\max_a \max_s \left[v(s) - v'(s)\right]\right|$$

$$= \gamma \max_s \left|\max_s \left[v(s) - v'(s)\right]\right|$$

$$= \gamma \max_s \left|v(s) - v'(s)\right|$$

$$= \gamma\|v(s) - v'(s)\|_\infty$$

*(b) Assuming a general finite MDP $(S, A, R, p, \gamma)$ where rewards are bounded: $r \in [r_{min}, r_{max}]$ for all $r \in R$. Prove the following equations. (2P)*

1. $\frac{r_{\min}}{1-\gamma} \leq v(s) \leq \frac{r_{\max}}{1-\gamma}$

$$v(s) = \max_a \sum_{s',r} p(s', r|s, a)\Big[r + \gamma v(s')\Big]$$

$$= \max_a \sum_{s'} p(s'|s, a)\Big[r(s, a, s') + \gamma v(s')\Big]$$

$$= \max_a \sum_{s'} \Big[p(s'|s, a)r(s, a, s') + p(s'|s, a)\gamma v(s')\Big]$$

$$= \max_a \Big[\sum_{s'} p(s'|s, a)r(s, a, s') + \sum_{s'} p(s'|s, a)\gamma v(s')\Big]$$

$$= \max_a \Big[\sum_{s'} p(s'|s, a)r(s, a, s') + \sum_{s'} p(s'|s, a)\gamma v(s')\Big]$$

$$\leq \max_a \Big[\max_{s'} r(s, a, s') + \gamma v(s)\Big]$$

$$= \max_a \Big[r_{\max} + \gamma v(s)\Big]$$

$$= r_{\max} + \gamma v(s)$$

$$v(s) \leq r_{\max} + \gamma v(s) \qquad |-\gamma v(s)$$
$$v(s) - \gamma v(s) \leq r_{\max}$$
$$(1-\gamma)v(s) \leq r_{\max} \qquad |:(1-\gamma)$$
$$v(s) \leq \frac{r_{\max}}{1-\gamma}$$

$$r_{\min} + \gamma v(s) = \max_a \left[ r_{\min} + \gamma v(s) \right]$$

$$= \max_a \left[ \min_{s'} r(s, a, s') + \gamma v(s) \right]$$

$$\leq \max_a \left[ \sum_{s'} p(s'|s, a) r(s, a, s') + \sum_{s'} p(s'|s, a) \gamma v(s') \right]$$

$$= \max_a \left[ \sum_{s'} p(s'|s, a) r(s, a, s') + \sum_{s'} p(s'|s, a) \gamma v(s') \right]$$

$$= \max_a \sum_{s'} \left[ p(s'|s, a) r(s, a, s') + p(s'|s, a) \gamma v(s') \right]$$

$$= \max_a \sum_{s'} p(s'|s, a) \left[ r(s, a, s') + \gamma v(s') \right]$$

$$= \max_a \sum_{s',r} p(s', r|s, a) \left[ r + \gamma v(s') \right]$$

$$= v(s)$$

$$r_{\min} + \gamma v(s) \leq v(s) \qquad | - \gamma v(s)$$

$$r_{\min} \leq v(s) - \gamma v(s)$$

$$r_{\min} \leq (1 - \gamma) v(s) \qquad | : (1 - \gamma)$$

$$\frac{r_{\min}}{1 - \gamma} \leq v(s)$$

2. $|v(s) - v(s')| \leq \frac{r_{\max} - r_{\min}}{1 - \gamma}$

$$|v(s) - v(s')|$$

$$\leq \left| \frac{r_{\max}}{1 - \gamma} - \frac{r_{\min}}{1 - \gamma} \right|, \text{ because that is the maximal reachable distance}$$

$$= \left| \frac{r_{\max} - r_{\min}}{1 - \gamma} \right|$$

$$= \frac{r_{\max} - r_{\min}}{1 - \gamma}, \text{ because } r_{\max} > r_{\min}$$

## Task 2

### Value Iteration (5 points)

*(a) Implement the value iteration algorithm (see lecture 3 slide 28) in the function value iteration. Use the values for $\gamma$ and $\theta$ that are given in the code. Initialize the value function $V(s)$ to $0$ for all states. How many steps does it need to converge? What is the optimal value function? (3P)*

We worked with the standard 4x4 Frozen Lake environment. The Value Iteration algorithm needs 43 loops to converge close to the optimal value function with a maximum error of $\theta$.

The optimal value function is $v_* =$
$$
\begin{pmatrix}
v_*(s_0) \\
v_*(s_1) \\
v_*(s_2) \\
v_*(s_3) \\
v_*(s_4) \\
v_*(s_5) \\
v_*(s_6) \\
v_*(s_7) \\
v_*(s_8) \\
v_*(s_9) \\
v_*(s_{10}) \\
v_*(s_{11}) \\
v_*(s_{12}) \\
v_*(s_{13}) \\
v_*(s_{14}) \\
v_*(s_{15})
\end{pmatrix}
\approx
\begin{pmatrix}
0.015 \\
0.016 \\
0.027 \\
0.016 \\
0.027 \\
0 \\
0.060 \\
0 \\
0.058 \\
0.134 \\
0.197 \\
0 \\
0 \\
0.247 \\
0.544 \\
0
\end{pmatrix}
$$

*(b) Compute the optimal policy from the value function. (2P)*

The optimal policy is $\pi_* =$
$$
\begin{pmatrix}
\pi_*(s_0) \\
\pi_*(s_1) \\
\pi_*(s_2) \\
\pi_*(s_3) \\
\pi_*(s_4) \\
\pi_*(s_5) \\
\pi_*(s_6) \\
\pi_*(s_7) \\
\pi_*(s_8) \\
\pi_*(s_9) \\
\pi_*(s_{10}) \\
\pi_*(s_{11}) \\
\pi_*(s_{12}) \\
\pi_*(s_{13}) \\
\pi_*(s_{14}) \\
\pi_*(s_{15})
\end{pmatrix}
\approx
\begin{pmatrix}
1 \\
3 \\
2 \\
3 \\
0 \\
0 \\
0 \\
0 \\
3 \\
1 \\
0 \\
0 \\
0 \\
2 \\
1 \\
0
\end{pmatrix}
$$

| | | | |
|---|---|---|---|
| ↓ | ↑ | → | ↑ |
| ← | H | ← | H |
| ↑ | ↓ | ← | H |
| H | → | ↓ | G |