## Lecture 25: Reinforcement Learning

*Lecturer: Bhuvesh Kumar*      *Scribes: Alen Polakof, Yinglun Xu*

**Disclaimer**: *These notes have not been subjected to the usual scrutiny reserved for formal publications.*

## 25.1 Reinforcement Learning Introduction

**Problem Setting:** The problem is formulated as an agent, which takes actions in an environment changing its state so as to maximize a cumulative reward. Same as the Stochastic Bandit Setting presented in class before, we have the exploration vs. exploitation trade-off.

**Problem Formulation:** MDP (Markov Decision Process)

- A model of the environment and the interactions of the environment.

- A set $S$ of states.

- A start state $s_0$.

- A set $A$ of actions.

- Transition Probability: $Pr(s' \mid s, a) \implies s' = (s, a)$. Probability of taking action $a$ in state $s$ and arriving to the new state $s'$.

- Reward Probability: $Pr(r' \mid s, a) \implies r' = r(s, a)$. Probability of taking action $a$ in state $s$ and obtaining a reward $r'$.

- Policy: A policy $\pi : s \to A$ or $\pi(s)$ is a mapping from $S$ to probability distributions over the action space $A$.

We use MDPs because the problem is only dependent on the current state and action taken, not on its history.

**Problem Objective:** Maximize the reward, where the reward depends on the policy taken. Hence, the objective is to find the policy that produces the highest cumulative reward.

**Examples**

1. Agent Mario interacts with environment and gets rewards as coins or superpowers. Actions can be moving, or jumping, which change the environment.

2. Alpha Go. Computer program developed by DeepMind that plays the game of Go.

**Type of horizons:** In Reinforcement Learning there are two different games:

- Finite Horizon Game:

$$\text{Reward} = \sum_{t=0}^{T} r(s_t, \pi(s_t))$$

- Infinite Horizon Game: Let $\gamma \in (0, 1)$ be the discount factor. In other words, $\gamma$ is a weighting term that represents the importance of future rewards, and decreases the further we look into the future.

$$\text{Reward} = \sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t))$$

**Policy Value:** The value $V_\pi(s)$, which is the expected reward if $s_0 = s$ and you follow the policy $\pi$. As a result, is an expectation of the reward that will be accumulated starting at state $s$ and following the policy $\pi$.

- For finite horizon case: $V_\pi(s) = \mathbb{E}[\sum_{t=0}^{T} r(s_t, \pi(s_t)) | s_0 = s]$

- For infinite horizon case: $V_\pi(s) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t)) | s_0 = s]$

(From now no, only focus on infinite horizon cases)

## 25.2 Bellman Equation:

**Theorem 25.1** *For all $s$ in $S$, where $S$ is the set of all possible states, and $s'$ is the states that can be reached from $s$. We define the Bellman Equation as:*

$$V_\pi(s) = \mathbb{E}[r(s, \pi(s))] + \gamma \sum_{s' \in S} Pr(s' \mid s, \pi(s)) V_\pi(s')$$

**Proof:**

$$V_\pi(s) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, \pi(s_t) | s_0 = s)]$$

$$= \mathbb{E}[r(s, \pi(s))] + \gamma \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_{t+1}, \pi(s_{t+1})) | s_0 = s]$$

$$= \mathbb{E}[r(s, \pi(s))] + \gamma \mathbb{E}[\sum_{s' \in S} \sum_{t=0}^{\infty} \gamma^t r(s_{t+1}, \pi(s_{t+1})) | s_1 = s', s_0 = s] Pr(s_1 = s' | s_0 = s)$$

$$= \mathbb{E}[r(s, \pi(s))] + \gamma \sum_{s' \in S} \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_{t+1}, \pi(s_{t+1})) | s_1 = s'] Pr(s' | s, \pi(s))$$

$$= \mathbb{E}[r(s, \pi(s))] + \gamma \sum_{s' \in S} V_\pi(s') Pr(s' | s, \pi(s))$$

■

As a result, this demonstrates that $V = R + \gamma P \cdot V$, where:

$$P \in \mathbb{R}^{|s| x |s|} \qquad\qquad P_{ss'} = Pr(s' | s, \pi(s))$$
$$V \in \mathbb{R}^{|s|} \qquad\qquad V_s = V_\pi(s)$$
$$R \in \mathbb{R}^{|s|} \qquad\qquad R_s = \mathbb{E}[r(s, \pi(s))]$$

**Theorem 25.2** *For finite MDP, the Bellman Equation has a unique solution. Therefore, $V = (I - \gamma P)^{-1} \cdot R$.*
**Proof:** *Note that:*
$$\|P\|_\infty = max_s \sum_{s' \in S} |P_{s,s'}| = max_s 1 = 1 \ and \ (I - \gamma P) \cdot V = R$$

$$\|\gamma P\|_\infty = \gamma$$
$$\implies \textit{All eigenvalues of } \gamma P \leq \gamma$$
$$\implies \textit{Eigenvalues of } (I - \gamma P) \geq 1 - \gamma > 0 \qquad\qquad \gamma \in (0, 1)$$

*As a result, since all eigenvalues are greater than 0, the matrix $I - \gamma P$ is invertible, and $V = (I - \gamma P)^{-1} R$* ∎

## 25.3   Optimal Policy and Q-Function

A policy $\pi^*$ is optimal if for all $s \in S$ and for all possible policies $\pi$,

$$V_{\pi^*}(s) \geq V_\pi(s)$$

As a result, the value of the optimal policy must be always greater all equal to the value of all other possible policies.

**Definition 25.3 (Q Function)**

$$Q_\pi(s, a) = \mathbb{E}[r(s, a)] + \gamma \sum_{s'} P[s' \mid s, \pi(s)] V_\pi(s')$$

$$Q_{\pi^*}(s, a) = \mathbb{E}[r(s, a)] + \gamma \sum_{s'} P[s' \mid s, \pi^*(s)] V_{\pi^*}(s')$$

*Bellman optimality condition:*

$$V_{\pi^*}(s) = \max_a Q_{\pi^*}(s, a)$$

Two Settings:

- Planning: $P[s' \mid s, \pi(s)]$ and $r(s, a)$ are known

- Learning: Don't know anything

---

**Algorithm 1** Value Iteration

---
**while** $\|V - \phi(V)\|_\infty \geq \frac{(1-\gamma)\epsilon}{\gamma}$ **do**
  $V = \phi(V)$
   where $\phi(V)(s) = \max_a \{\mathbb{E}[r(s, a)] + \gamma \sum_{s'} P[s' \mid s, a] V(s')\}$
**end while**
**return**  $\pi(s) = \text{argmax}_a \{\mathbb{E}[r(s, a)] + \gamma \sum_{s'} P[s' \mid s, a] V(s')\}$

---

**Lemma 25.4** $\phi$ *is* $\gamma - Lipschitz$ *in* $\|\cdot\|_\infty$, *i.e.,*

$$\|\phi(V) - \phi(U)\|_\infty \leq \gamma \|V - U\|_\infty.$$

**Proof:** $\forall s \in S$, Let $a^*(s) = \text{argmax}_a \{\mathbb{E}[r(s, a)] + \gamma \sum_{s'} P[s' \mid s, a] V(s')\}$. Then

$$\phi(V)(s) - \phi(U)(s) \leq \phi(V)(s) - [\mathbb{E}[r(s, a^*)] + \gamma \sum_{s'} P[s' \mid s, a^*] V(s')]$$

$$= \gamma \sum_{s'} P[s' \mid s, a^*][V(s') - U(s')]$$

$$\leq \gamma \sum_{s'} P[s' \mid s, a^*] \|V - U\|_\infty$$

$$= \gamma \|V - U\|_\infty$$

Similarly, Let $a^{*2}(s) = \text{argmax}_a \{\mathbb{E}[r(s,a)] + \gamma \sum_{s'} P[s' \mid s, a]U(s')\}$, it can be proved that

$$\phi(U)(s) - \phi(V)(s) \leq \gamma \|V - U\|_\infty.$$

Combine both then

$$\|\phi(V) - \phi(U)\|_\infty \leq \gamma \|V - U\|_\infty.$$

∎

**Lemma 25.5** $V^* = \phi(V^*)$ *where* $V^* = V_{\pi^*}$

**Proof: (Exercise)** Hint: Bellman Optimality ∎

**Theorem 25.6** *For any* $V_0$, *the sequence* $V_{n+1} = \phi(V_n) = \phi(V_n)$ *converge to* $V^*$.

**Proof:**
$$\|V^* - v_{n+1}\|_\infty = \|\phi(V_* - \phi(V_n))\|_\infty \leq \gamma \|V_* - V_n\|_\infty \leq \gamma^n \|V_* - V_1\|_\infty$$

So $V_n$ converges to $V^*$ ∎

**Theorem 25.7** *Value Iteration halts in* $O(\log\frac{1}{\epsilon})$ *steps and return* $\pi$ *such that* $\|V_n - V^*\|_\infty \leq \epsilon$

**Proof:** $\forall n \in N$,

$$\begin{aligned}
\|V^* - V_{n+1}\|_\infty &\leq \|V^* - \phi(V_{n+1})\|_\infty + \|\phi(V_{n+1}) - V_{n+1}\|_\infty \\
&= \|\phi(V^*) - \phi(V_{n+1})\|_\infty + \|\phi(V_{n+1}) - \phi(V_n)\|_\infty \\
&\leq \gamma \|V^* - V_{n+1}\|_\infty + \gamma \|\phi(V_n) - V_n\|_\infty
\end{aligned}$$

∎

Then

$$\|V^* - V_{n+1}\|_\infty \leq \frac{\gamma}{1-\gamma} \|\phi(V_n) - V_n\|_\infty = \epsilon$$

Let n be the largest index such that $\|\phi(V_n) - v_n\|_\infty \leq \frac{1-\gamma}{\gamma}\epsilon$. Since $\|\phi(V_n) - V_n\|_\infty \leq \gamma^n \|V_1 - V_0\|_\infty$, then

$$\frac{1-\gamma}{\gamma}\epsilon \leq \gamma^n \|V_1 - V_0\|_\infty$$

$$n \leq O(\log\frac{1}{\epsilon})$$