

Supplementary Material #1

1. **Bellman Equation:** Let's derive the Bellman Equation for the state-value function step by step.

$$\begin{aligned}
\mathbf{v}^\pi(s) &\stackrel{\text{def}}{=} \mathbb{E}[G_t | S_t = s] \\
&= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \\
&= \mathbb{E}[R_{t+1} | S_t = s] + \gamma \mathbb{E}[R_{t+2} | S_t = s] + \gamma^2 \mathbb{E}[R_{t+3} | S_t = s] + \dots
\end{aligned}$$

Bellman Equation is expressed by a recursion, that is, we aim to produce $\mathbf{v}^\pi(s') = \mathbb{E}[G_{t+1} | S_{t+1} = s']$ on the right hand side. To this end, we rewrite each term on the right hand side as follows:

$$\begin{aligned}
\mathbb{E}[R_{t+1} = r | S_t = s] &= \sum_{s'} \sum_a p(s', a | s) r \\
&= \sum_{s'} \sum_a \frac{p(s', a, s)}{p(s)} r \\
&= \sum_{s'} \sum_a \frac{p(s', a, s)}{p(s)} \frac{p(a, s)}{p(a, s)} r \\
&= \sum_{s'} \sum_a \frac{p(s', a, s)}{p(a, s)} \frac{p(a, s)}{p(s)} r \\
&= \sum_{s'} \sum_a p(s' | s, a) p(a | s) r \\
&= \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) r \\
\mathbb{E}[R_{t+1} | S_t = s] &= \sum_a \pi(a | s) \sum_{s'} \sum_r p(s', r | s, a) r \\
\mathbb{E}[R_{t+2} = r' | S_t = s] &= \sum_{s''} \sum_{a'} \sum_{s'} \sum_a p(s'', a' | s') p(s', a | s) r' \\
&= \sum_{s'} \sum_a p(s', a | s) \sum_{s''} \sum_{a'} p(s'', a' | s') r' \\
&= \sum_{s'} \sum_a p(s | s', a) p(a | s) \sum_{s''} \sum_{a'} p(s'', a' | s') r' \\
&= \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) \mathbb{E}[R_{t+2} = r' | S_{t+1} = s'] \\
\mathbb{E}[R_{t+2} | S_t = s] &= \sum_a \pi(a | s) \sum_{s'} \sum_r p(s', r | s, a) \mathbb{E}[R_{t+2} | S_{t+1} = s']
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}[R_{t+3} = r'' | S_t = s] &= \sum_{s'''} \sum_{a''} \sum_{s''} \sum_{a'} \sum_{s'} \sum_a p(s''', a'' | s'') p(s'', a' | s') p(s', a | s) r'' \\
&= \sum_{s'} \sum_a p(s', a | s) \sum_{s'''} \sum_{a''} \sum_{s''} \sum_{a'} p(s''', a'' | s'') p(s'', a' | s') r'' \\
&= \sum_{s'} \sum_a p(s' | s, a) p(a | s) \sum_{s'''} \sum_{a''} \sum_{s''} \sum_{a'} p(s''', a'' | s'') p(s'', a' | s') r'' \\
&= \sum_a \pi(a | s) \sum_{s'} p(s' | s, a) \mathbb{E}[R_{t+3} = r'' | S_{t+1} = s'] \\
\mathbb{E}[R_{t+3} | S_t = s] &= \sum_a \pi(a | s) \sum_{s'} \sum_r p(s', r | s, a) \mathbb{E}[R_{t+3} | S_{t+1} = s']
\end{aligned}$$

Repeat the same procedure for $t + 4, t + 5, \dots$, and plug them back in.

$$\begin{aligned}
\mathbf{v}^\pi(s) &= \mathbb{E}[R_{t+1} | S_t = s] + \gamma \mathbb{E}[R_{t+2} | S_t = s] + \gamma^2 \mathbb{E}[R_{t+3} | S_t = s] + \dots \\
&= \sum_a \pi(a | s) \sum_{s'} \sum_r p(s', r | s, a) \{r + \gamma \mathbb{E}[R_{t+2} + \gamma R_{t+3} + \dots | S_{t+1} = s']\} \\
&= \sum_a \pi(a | s) \sum_{s'} \sum_r p(s', r | s, a) \{r + \gamma \mathbf{v}^\pi(s')\}
\end{aligned}$$

Recall the notation we introduced in the lecture:

$$\begin{aligned}
\mathcal{P}_{ss'}^a &\stackrel{\text{def}}{=} P[S_{t+1} = s' | S_t = s, A_t = a] = \sum_r p(s', r | s, a) \\
\mathcal{R}_s^a &\stackrel{\text{def}}{=} \mathbb{E}[R_{t+1} = r | S_t = s, A_t = a] = \sum_{s'} \sum_r p(s', r | s, a) r
\end{aligned}$$

Thus, we get the desired equation:

$$\mathbf{v}^\pi(s) = \sum_a \pi(a | s) \left(\mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a \mathbf{v}^\pi(s') \right)$$

2. (exercise) Derive the Bellman Equation for the action-value function, analogously.

$$\begin{aligned}
\mathbf{q}^\pi(s, a) &= \mathbb{E}[G_t | S_t = s, A_t = a] \\
&= \vdots \\
&= \mathcal{R}_s^a + \gamma \sum_{s'} \mathcal{P}_{ss'}^a \sum_{a'} \pi(a' | s') \mathbf{q}^\pi(s', a')
\end{aligned}$$

Special Thanks to Michael. Indeed, the Markov property need to be taken into account. That is, the probability to reach a state s'' in time step $t + 2$ starting from a state s in t should be $p(s'', a' | s') p(s', a | s)$.