

Chapter 2 State Value and Bellman Equation

State Value, Bellman Equation, Action Value

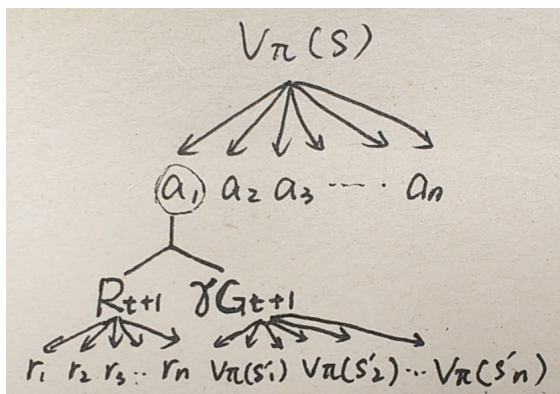
1. Why did we introduce State Value? Isn't return enough to value a policy?
 2. Explain the understanding of Random Variables in RL.
 3. Explain how state values depend on each other(Bootstrapping)?
 4. What is G_t ? What is the difference between G_t and v_t ?
 5. What does Bellman Equation tell us?
 6. What is policy evaluation?
 7. Derive the Bellman Equation. How to understand the ultimate format of Bellman Equation?
 8. Did find any principle when dealing with multiple Σ ?
 9. In the Bellman Equation, what are unknown, what represent the model, is π given?
 10. Derive matrix-vector form of Bellman Equation.
 11. Why need matrix-vector form beyond elementwise form?
 12. How to solve Bellman Equation? Give 2.
 13. Prove that the iterative solution is converged.
 14. What is action value? Is action value based on state or has no relationship with state?
 15. What is the relationship between state value and action value?
 16. The elementwise form of action value.
 17. Derive the matrix-vector form of action value.
 18. In the matrix-vector form of state value and action value, we have vector $v = [v_1, v_2, v_3, \dots, v_n]^T$ and $q = [q_1, q_2, q_3, \dots, q_n]^T$. Explain what do these 2 sequences represent for? Does sequence v represent for the state values of the whole state space? Does sequence q represent for the action values of the action space of a state?
-
1. In stochastic system, the model (state transition and reward) is stochastic and can be described using conditional probabilities. Return can only describe a deterministic trajectory, but state value can describe the expected return of a stochastic trajectory.
 2. Random Variables do not have a deterministic value, but have a row of values with corresponded probabilities(Probability Distribution). Such as action, state transition, return, they are stochastic and are described with probability distribution.
 3. $v_t = r_{t+1} + \gamma \cdot r_{t+2} + \gamma^2 \cdot r_{t+3} + \dots = r_{t+1} + \gamma \cdot v_{t+1}$
 4. G_t is a stochastic variable describes the discounted return of state s_t following a policy. v_t is the expectation of G_t and it is a deterministic value.

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) \\ &= R_{t+1} + \gamma G_{t+1}, \end{aligned}$$

$$\begin{aligned}
 v_{\pi}(s) &= \mathbb{E}[G_t | S_t = s] \\
 &= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\
 &= \mathbb{E}[R_{t+1} | S_t = s] + \gamma \mathbb{E}[G_{t+1} | S_t = s].
 \end{aligned}$$

5. Bellman Equation solve state values based on the model and a given policy. It explains the relationship between $v_{\pi}(s)$ and $v_{\pi}(s')$ based on the model. Also, Bellman Equation is fundamentally a formula calculating the AVERAGE DISCOUNTED RETURN on state s , it shows the general situation.
6. The process of solving Bellman Equation to get state values is policy evaluation.
- 7.

$$\begin{aligned}
 \textcircled{1} \quad V_{\pi}(s) &= \mathbb{E}[G_t | S_t = s] \\
 &= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\
 &= \mathbb{E}[R_{t+1} | S_t = s] + \gamma \mathbb{E}[G_{t+1} | S_t = s] \\
 \textcircled{2} \quad \mathbb{E}[R_{t+1} | S_t = s] &= \sum_{a \in A(s)} \pi(a|s) \cdot \sum_{r \in R(s,a)} r \cdot p(r|s,a) \\
 \textcircled{3} \quad \mathbb{E}[G_{t+1} | S_t = s] &= \sum_{s' \in S} p(s'|s) \cdot \mathbb{E}[G_{t+1} | S_t = s, S_{t+1} = s'] \\
 &= \sum_{s' \in S} p(s'|s) \cdot \mathbb{E}[G_{t+1} | S_{t+1} = s'] \Rightarrow \text{Markov Property} \\
 &= \sum_{s' \in S} p(s'|s) \cdot V_{\pi}(s') \\
 &= \sum_{s' \in S} \sum_{a \in A(s)} \pi(a|s) \cdot p(s'|s,a) \cdot V_{\pi}(s') \\
 \textcircled{4} \quad V_{\pi}(s) &= \sum_{a \in A} \pi(a|s) \cdot \sum_{r \in R} r \cdot p(r|s,a) + \gamma \sum_{s' \in S} p(s'|s,a) \cdot V_{\pi}(s') \\
 &= \sum_{a \in A} \pi(a|s) \cdot \left[\sum_{r \in R} p(r|s,a) \cdot r + \gamma \sum_{s' \in S} p(s'|s,a) \cdot V_{\pi}(s') \right]
 \end{aligned}$$



$$\sum_{a \in A} \pi(a|s) \left[\sum_{r \in R} p(r|s,a) r + \gamma \sum_{s' \in S} p(s'|s,a) v_{\pi}(s') \right]$$

8. We can freely arrange the items and Σ . Just make sure that when the item and the Σ have the same variable notation, the items should follow behind the Σ but not before.
9. $v_{\pi}(s)$ and $v_{\pi}(s')$ is unknown, $p(r|s,a)$ and $p(s'|s,a)$ represent the model, π is a given policy and the process is policy evaluation because state values can evaluate the performance of the policy.
- 10.

$$\begin{aligned}
 v_{\pi}(s) &= r(s) + \gamma v_{\pi}(s') \\
 &= r(s) + \gamma \sum_{s' \in S} p(s'|s) \cdot v_{\pi}(s') \\
 \Rightarrow \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ \vdots \\ v_n \end{bmatrix} &= \begin{bmatrix} r_1 \\ r_2 \\ r_3 \\ \vdots \\ r_n \end{bmatrix} + \gamma \cdot \begin{bmatrix} p(s_1|s_1) & p(s_2|s_1) & p(s_3|s_1) & \dots & p(s_n|s_1) \\ p(s_1|s_2) & p(s_2|s_2) & p(s_3|s_2) & \dots & p(s_n|s_2) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ p(s_1|s_n) & p(s_2|s_n) & p(s_3|s_n) & \dots & p(s_n|s_n) \end{bmatrix} \cdot \begin{bmatrix} v_1 \\ v_2 \\ v_3 \\ \vdots \\ v_n \end{bmatrix} \\
 \vec{v} &= \vec{r} + \gamma \vec{P} \cdot \vec{v}
 \end{aligned}$$

11. Bellman Equation is a row of formula of every state. To solve it, we must put it together. Matrix-vector form can deal with this.

12. Closed-form solution and iterative solution.

13.

$$\begin{aligned}
 \text{define } \vec{\delta}_k &= \vec{v}_k - \vec{v}_{\pi} \\
 \vec{v}_{k+1} &= \vec{r}_{\pi} + \gamma \vec{P}_{\pi} \vec{v}_k \\
 \Rightarrow \vec{\delta}_{k+1} + \vec{v}_{\pi} &= \vec{r}_{\pi} + \gamma \vec{P}_{\pi} (\vec{\delta}_k + \vec{v}_{\pi}) \\
 \vec{\delta}_{k+1} &= \gamma \vec{P}_{\pi} \cdot \vec{\delta}_k \\
 \text{Thus, } \vec{\delta}_{k+1} &= \gamma \vec{P}_{\pi} \cdot \gamma \vec{P}_{\pi} \vec{\delta}_{k-1} \\
 &= \dots \\
 &= \gamma^{k+1} \vec{P}_{\pi}^{k+1} \cdot \vec{\delta}_0 \rightarrow 0 \\
 0 < \gamma < 1, \text{ then } \gamma^{k+1} &\rightarrow 0. \quad \vec{P}_{\pi} \text{ 特征值} \leq 1, \text{ then } \vec{P}_{\pi}^{k+1} \text{ is finite.}
 \end{aligned}$$

14. Action value is the discounted return of taking a specific action at a specific state. Same action on different state may have different action value.

15. State value is the expectation value of every action value on this state.

$$v_{\pi}(s) = \sum_{a \in A} \pi(a|s) q_{\pi}(s, a)$$

16.

$$q_{\pi}(s, a) = \sum_{r \in \mathcal{R}} p(r|s, a) r + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v_{\pi}(s')$$

17. Derive by myself. No record in the book. So maybe some mistakes.

$$\begin{aligned}
q_{\pi}(s, a) &= r(s, a) + \gamma \sum_{s' \in S} P(s' | s, a) \sum_{a' \in A(s')} \pi(a' | s') \cdot q_{\pi}(s', a') \\
\begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_n \end{bmatrix} &= \begin{bmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{bmatrix} + \gamma \cdot \begin{bmatrix} P(s'_1 | s, a_1) & P(s'_1 | s, a_2) & \dots & P(s'_1 | s, a_n) \\ P(s'_2 | s, a_1) & P(s'_2 | s, a_2) & \dots & P(s'_2 | s, a_n) \\ \vdots & \vdots & \ddots & \vdots \\ P(s'_n | s, a_1) & P(s'_n | s, a_2) & \dots & P(s'_n | s, a_n) \end{bmatrix}_{n \times m} \cdot \\
&\quad \begin{bmatrix} \pi(a_1 | s'_1) & \pi(a_1 | s'_2) & \dots & \pi(a_1 | s'_m) \\ \pi(a_2 | s'_1) & \pi(a_2 | s'_2) & \dots & \pi(a_2 | s'_m) \\ \vdots & \vdots & \ddots & \vdots \\ \pi(a_n | s'_1) & \pi(a_n | s'_2) & \dots & \pi(a_n | s'_m) \end{bmatrix}_{m \times n} \cdot \\
&\quad \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_n \end{bmatrix}_{n \times 1} \\
\vec{q}_{\pi} &= \vec{r} + \gamma \vec{P} \vec{\pi} \vec{q}_{\pi}
\end{aligned}$$

18. They don't represent the state values of the whole state space or the action values of the action space of a state. They come from bootstrapping. They are sequential records of the state values and action values obtained during a trajectory following the policy.