

Lecture 3: Bellman consistency equation; MRPs; Optimal value function

Lecturer: Chicheng Zhang

Scribe: Yinan Li

1 Questions from last class

Recall that last time, we defined the value function of π , $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$, which is the expected discounted return of a policy, conditioned on the starting state:

$$V^\pi(s) = \mathbb{E}[G_0 \mid s_0 = s, \pi].$$

We also defined the action-value function of π , $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, which is conditioned additionally on the starting action:

$$Q^\pi(s, a) = \mathbb{E}[G_0 \mid s_0 = s, a_0 = a, \pi].$$

Using *Bellman consistency equation*, and eliminating variables $V^\pi(s)$, we have:

$$\forall s \in \mathcal{S}, a \in \mathcal{A} : Q^\pi(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} P(s' \mid s, a) \pi(a' \mid s') Q^\pi(s', a').$$

Note that this is a system of linear equations, where the unknown variables are $Q^\pi(s, a)$'s, and the known variables are immediate rewards $r(s, a)$'s and environmental dynamics.

The questions left from the last time are:

1. Does this system have solutions?
2. If so, is the solution unique?

In other words, we would like to find a plausible choice of f , that satisfies:

$$f(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} P(s' \mid s, a) \pi(a' \mid s') f(s', a'). \quad (1)$$

For the first question, we observe that $f(s, a) = Q^\pi(s, a)$ is a valid solution to (1). We will answer the second question in this lecture.

2 Markov Reward Processes (MRPs)

Markov Reward Processes (MRPs) are slight variants of MDPs, where the agent does not have controls on the environment. Formally, an MRP M is described as:

$$M = (\mathcal{S}, \mathbb{P}, r, \gamma, \mu),$$

where, different from those of MDPs,

1. Transition probability $\mathbb{P} : \mathcal{S} \rightarrow \Delta(\mathcal{S})$;
2. Reward function $r : \mathcal{S} \rightarrow [0, 1]$.

Algorithm 1 Agent-environment interaction protocol with an MRP

Initial state $z_0 \sim \mu$
for time steps $t = 0, 1, \dots$, **do**
 Agent receives reward $r_t = r(z_t)$
 Agent observes next state $z_{t+1} \sim \mathbb{P}(\cdot | z_t)$
end for

Similar to MDP setting, we could define the value function in MRP context:

$$V(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s \right].$$

Exercise 1. Give a Bellman consistency equation for V .

Exercise 2. Suppose we execute a stationary policy π in a given MDP M . Can the interaction process be seen as an MRP? Here are two possible answers to this exercise.

Answer 1. MRP has the same state space as MDP, and we could let $z_t = s_t$, and

$$\mathbb{P}(s_{t+1} = s' \mid s_t = s) = \sum_{a \in \mathcal{A}} \pi(a \mid s) \mathbb{P}(s' \mid s, a).$$

Note that there is some subtlety, in that r_t is not a deterministic function of s_t , but still, we could define the reward function in expectation:

$$\mathbb{E}[r_t \mid s_t = s] = \sum_{a \in \mathcal{A}} \pi(a \mid s) r(s, a).$$

Answer 2. MRP has the state space $\mathcal{S} \times \mathcal{A}$, namely $z_t = (s_t, a_t)$, and

$$\mathbb{P}[(s_{t+1}, a_{t+1}) = (s', a') \mid (s_t, a_t) = (s, a)] = \mathbb{P}(s' \mid s, a) \pi(a' \mid s').$$

MRP has the same reward function as MDP.

3 Back to Bellman consistency equation

Recall that,

$$\forall s \in \mathcal{S}, a \in \mathcal{A} : Q^\pi(s, a) = r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} P(s' \mid s, a) \pi(a' \mid s') Q^\pi(s', a').$$

and it helps to look at it from linear algebra perspective.

Denote by $S = |\mathcal{S}|$, and $A = |\mathcal{A}|$,

$$Q^\pi = \begin{bmatrix} \vdots \\ Q^\pi(s, a) \\ \vdots \end{bmatrix}_{SA \times 1},$$
$$r = \begin{bmatrix} \vdots \\ r(s, a) \\ \vdots \end{bmatrix}_{SA \times 1},$$

$$P^\pi = \begin{bmatrix} \vdots & & \\ \pi(a' | s') \cdot P(s' | s, a) & & \\ \vdots & & \end{bmatrix}_{SA \times SA}.$$

Now we are ready to rewrite Bellman consistency equation in the matrix notation,

$$Q^\pi = r + \gamma \cdot P^\pi \cdot Q^\pi.$$

To solve it, rewrite it as:

$$(I - \gamma \cdot P^\pi) \cdot Q^\pi = r.$$

Claim 1. $I - \gamma \cdot P^\pi$ is invertible.

Proof. It suffices to show, $\forall x \neq 0$, $(I - \gamma \cdot P^\pi) \cdot x \neq 0$. To see this, we analyse the ℓ_∞ norm.

$$\begin{aligned} \|(I - \gamma \cdot P^\pi) \cdot x\|_\infty &= \|x - \gamma \cdot P^\pi \cdot x\|_\infty \\ &\geq \|x\|_\infty - \|\gamma \cdot P^\pi \cdot x\|_\infty \text{ (by triangle ineq.)} \end{aligned}$$

Note that each row in P^π sum up to 1, thus we have

$$|(P^\pi \cdot x)_i| = \left| \sum_j w_{ij} x_j \right| \leq \sum_j w_{ij} |x_j| \leq \|x\|_\infty.$$

Therefore,

$$\|(I - \gamma \cdot P^\pi) \cdot x\|_\infty \geq (1 - \gamma) \|x\|_\infty > 0.$$

□

Since $I - \gamma \cdot P^\pi$ is invertible, we are able to claim that

$$Q^\pi = (I - \gamma \cdot P^\pi)^{-1} \cdot r.$$

Now we take a closer look at entries of $(I - \gamma \cdot P^\pi)^{-1}$. Analogous to the algebraic fact that $\frac{1}{1-\gamma} = \sum_{t=0}^{\infty} \gamma^t$, we also have,

$$\begin{aligned} (I - \gamma \cdot P^\pi)^{-1}_{(s,a),(s',a')} &= \sum_{t=0}^{\infty} \gamma^t (P^\pi)^t_{(s,a),(s',a')} \\ &= \sum_{t=0}^{\infty} \gamma^t P[(s_t, a_t) = (s', a') \mid (s_0, a_0) = (s, a)] \geq 0. \end{aligned}$$

This is called discounted state-action visitation (occupancy).

To see the second equality, we analyse the case where $t = 2$, as an example; this can be easily generalized to interpret $(P^\pi)^t_{(s,a),(s',a')}$ as $P[(s_t, a_t) = (s', a') \mid (s_0, a_0) = (s, a)]$.

$$\begin{aligned} (P^\pi)^2_{(s,a),(s',a')} &= \sum_{\tilde{s}, \tilde{a}} (P^\pi)_{(s,a),(\tilde{s},\tilde{a})} \cdot (P^\pi)_{(\tilde{s},\tilde{a}), (s',a')} \\ &= \sum_{\tilde{s}, \tilde{a}} P[(s_1, a_1) = (\tilde{s}, \tilde{a}) \mid (s_0, a_0) = (s, a)] \cdot P[(s_2, a_2) = (s', a') \mid (s_1, a_1) = (\tilde{s}, \tilde{a})] \\ &= P[(s_2, a_2) = (s', a') \mid (s_0, a_0) = (s, a)] \end{aligned}$$

where the first equality is by matrix multiplication formula, the second equality is by the definition of P^π , and the third is by Markov property and total law of probability.

4 Optimal value function

Optimal value function could be thought of as a more objective measure of how advantage a state is. Previously, we defined the value function of a state, but that depends on the policy the agent executes. For example, think about a state that is very close to the goal state, and it is promising to get high reward. But the policy may not be performing so well, and it is actually directing away from the goal state. In this case, the value function of the state may be bad, although the state itself is good. Thus we would like a more objective measure of how "good" a state is.

The way we formulate this idea is to define the optimal value function and optimal action-value function as follows:

$$V^*(s) = \max_{\pi \in \Pi} V^\pi(s),$$

$$Q^*(s, a) = \max_{\pi \in \Pi} Q^\pi(s, a),$$

here, Π is the collection of all policies (including nonstationary ones).

Representing V^* with Q^* :

$$\begin{aligned} V^*(s) &= \max_{\pi} \mathbb{E}[G_0 \mid s_0 = s, \pi] \\ &= \max_{\pi} \sum_{a \in \mathcal{A}} \pi(a \mid s) \mathbb{E}[G_0 \mid s_0 = s, a_0 = a, \pi] \\ &\leq \max_{\pi} \sum_{a \in \mathcal{A}} \pi(a \mid s) \max_{\pi'} \mathbb{E}[G_0 \mid s_0 = s, a_0 = a, \pi'] \\ &= \max_{\pi} \sum_{a \in \mathcal{A}} \pi(a \mid s) Q^*(s, a). \end{aligned} \tag{2}$$

Note that the optimal π will be: $\pi(a \mid s) = \begin{cases} 1 & a = \operatorname{argmax}_{a'} Q^*(s, a') \\ 0 & \text{otherwise} \end{cases}$,

then we further have,

$$V^*(s) \leq \max_a Q^*(s, a).$$

Representing Q^* using V^* :

$$\begin{aligned} Q^*(s, a) &= \max_{\pi} \mathbb{E}[G_0 \mid s_0 = s, a_0 = a, \pi] \\ &= r(s, a) + \gamma \cdot \max_{\pi} \mathbb{E}[G_1 \mid s_0 = s, a_0 = a, \pi] \\ &= r(s, a) + \gamma \cdot \max_{\pi} \sum_{s' \in \mathcal{S}} P(s' \mid s, a) \mathbb{E}[G_1 \mid s_0 = s, a_0 = a, s_1 = s', \pi] \\ &\leq r(s, a) + \gamma \cdot \sum_{s' \in \mathcal{S}} P(s' \mid s, a) \max_{\pi} \mathbb{E}[G_1 \mid s_0 = s, a_0 = a, s_1 = s', \pi] \\ &\leq r(s, a) + \gamma \cdot \sum_{s' \in \mathcal{S}} P(s' \mid s, a) V^*(s') \text{ (Next class)} \end{aligned} \tag{3}$$

We will show in next class that, combining equations (2) and (3), we'll get:

$$Q^*(s, a) \leq r(s, a) + \gamma \cdot \sum_{s' \in \mathcal{S}} P(s' \mid s, a) \max_{a'} Q^*(s', a').$$