

3.1: The Agent-Environment Interface

In a finite MDP: sets of states, actions, rewards are all finite!

S_t and R_t have well defined probability distributions that is dependent only on the preceding state and action: s_{t-1}, a_{t-1} . (Markov Property).

The "Markov" part in MDP comes from the Markov property.

Dynamics of the MDP: $p: S \times R \times S \times A \rightarrow [0, 1]$.

$p(s', r | s, a) \doteq \Pr \{ S_t = s', R_t = r | S_{t-1} = s, A_{t-1} = a \}$

$\uparrow \quad \uparrow$
probability of the next state being
The function p defines; s' and its reward being r ,
the dynamics of the ; if we are in state s and take
MDP, aka the ' action a ,
transition function.

Reminder that p is a joint probability distribution for each choice of s and a (not r):

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1 \quad \forall s \in S, a \in A(s).$$

Additional Definitions:

State transition probabilities : $p : S \times S \times R \rightarrow [0, 1]$

$$p(s' | s, a) \doteq \Pr \left\{ S_t = s' \mid S_{t-1} = s, A_{t-1} = a \right\}$$

$$= \sum_{r \in R} p(s', r | s, a)$$

Expected rewards : $r : S \times A \rightarrow \mathbb{R}$

$$r(s, a) \doteq E[R_t | S_{t-1} = s, A_{t-1} = a]$$

$$= \sum_{r \in R} r \sum_{s' \in S} p(s', r | s, a)$$

Expected Reward for state-action-next state : $r : S \times A \times S \rightarrow \mathbb{R}$

$$r(s, a, s') \doteq E[R_t | S_{t-1} = s, A_{t-1} = a, S_t = s']$$

$$= \sum_{r \in R} r \frac{p(s', r | s, a)}{p(s' | s, a)}$$

3.2: Goals and Rewards

- At each time step t , the reward is a simple number $R_t \in \mathbb{R}$.
- The goal of the agent is to maximize cumulative reward in the long run.
- We can use the reward signal to give prior knowledge about how to achieve the task.

3.3: Returns and Episodes

$G_t = \text{Expected Return} = \text{expected sum of all future rewards}$

Episodic case:

- Simple definition of the expected return we want to maximize:
$$G_t = R_t + R_{t+1} + \dots + R_T \quad // \text{sum of ALL future rewards.}$$
- we get a reward at each time step, and T is the final time step.
- This approach makes sense when there is a notion of a final time step, that is, when the agent actions naturally break into sequences, that we call **Episodes**.
- Each episode ends in a special state called the **terminal state**.
- It's important that the start of the episode is independent from the end of the previous episode.
- Tasks with episodes of this kind are called **episodic tasks**. In these, S is the set of all non-terminal states, and S^+ is the set of all states, including terminal ones.

Continuous States

- when the agent-environment interaction does not naturally break into episodes.
- Final timestep is ∞ , so we can't really compute a useful G_t .
- we introduce discounting (to avoid a sum of returns going to ∞):

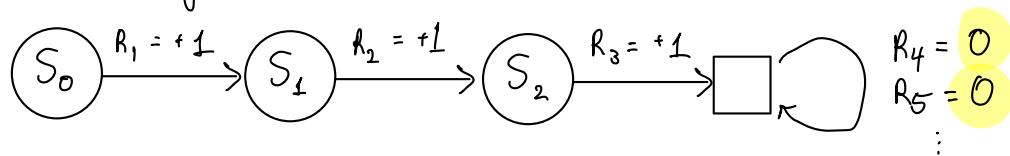
$$\begin{aligned}
 G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\
 &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}, \quad \text{where } \gamma = \text{discount rate, and} \\
 0 \leq \gamma &\leq 1, \quad (\text{for } \gamma = 0, \text{ the agent is myopic}).
 \end{aligned}$$

Recursive form of G_t

$$\begin{aligned}
 G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\
 &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) \\
 &= R_{t+1} + \gamma G_{t+1} \\
 \Rightarrow G_t &= R_{t+1} + \gamma G_{t+1} \quad // \text{notice the recursion.}
 \end{aligned}$$

3.4 : Unified notation for Episodic and Continuing tasks

- s_t = state at timestep t (episodic case).
- The return over a finite # of terms of the episodic case can be treated as the infinite sum by adding an absorbing state:



Unified notation for the return, that includes the possibility of $T = \infty$ or $\gamma = 1$:

$$G_t \stackrel{?}{=} \sum_{k=t+1}^T \gamma^{k-t-1} r_k$$

// can be used for both episodic and continuing tasks.

3.5: Policies and value functions

Definitions:

- value function: expected value in a state (or a q-state for q-values)
- policy: mapping from states to probabilities of taking each possible action. Determines how the agent behaves.

The value function of a state s under a policy π , denoted as $V_\pi(s)$, is *defined* as the expected return (future rewards) by starting in s and following π thereafter:

$$V_\pi(s) \doteq \underbrace{E_\pi[G_t \mid S_t = s]}_{\substack{\text{value function} \\ \text{of state } s}} = E_\pi \left[\underbrace{\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}}_{\substack{\text{Expected future} \\ \text{rewards if the agent} \\ \text{is in state } s \text{ and} \\ \text{follows policy } \pi.}} \mid S_t = s \right] \quad \forall s \in S$$

$G_t = \text{sum}$
 of future rewards
 $(\text{discounted}).$

V_π = state-value function for policy π .

action-value function q_{π} :

$$q_{\pi}(s, a) \stackrel{?}{=} E_{\pi} \left[G_t \mid S_t = s, A_t = a \right] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

expected future rewards if in
state s , take action a , and
then follow policy π .

Difference between v_{π} and q_{π} is q_{π} conditions on
 $A_t = a$.

$v_{\pi}(s)$ and $q_{\pi}(s)$ can be determined by experience.

If we maintain an average reward received for each state
(or q-state), this average converges to $v_{\pi}(s)$ (resp $q_{\pi}(s)$),
as time passes.

Value functions also have a nice recursive property:
 $V_{\pi}(s) \doteq E_{\pi}[G_t \mid S_t = s]$ // the value of state s under policy π is the expected return (cumulative discounted rewards) starting from state s

$$= E_{\pi} \left[\underbrace{R_{t+1}}_{\substack{\text{immediate} \\ \text{reward at} \\ \text{next timestep}}} + \underbrace{\gamma G_{t+1}}_{\substack{\text{the discounted} \\ \text{future reward from time } t+1 \text{ onward}}} \mid S_t = s \right]$$

$$= \sum_a \underbrace{\pi(a \mid s)}_{\substack{\text{sum over} \\ \text{all available} \\ \text{actions!}}} \sum_{s', r} \underbrace{p(s', r \mid s, a)}_{\substack{\text{probability of} \\ \text{transitioning to state } s' \\ \text{and receiving reward } r \\ \text{given we're in state } s \\ \text{and took action } a}} \underbrace{[r + \gamma E_{\pi}[G_{t+1} \mid S_{t+1} = s']]}_{\substack{\text{immediate} \\ \text{reward} \\ \text{of being} \\ \text{in state} \\ s' \\ \text{value of transition}}}$$

goodness of action.

$$= \boxed{\sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma V_{\pi}(s')]} \quad \forall s \in S$$

// Key insight : $E_{\pi}[G_{t+1} \mid S_{t+1} = s'] \Leftrightarrow V_{\pi}(s')$

The future expected return from state s' is exactly the value function of s' !

If can be viewed as a sum on 3 values:
 a, s', r .

For each triple, we compute the probability $\pi(a \mid s) p(s', r \mid s, a)$, weight the quantity in brackets by this probability, then sum over all the probabilities to get an expected value.

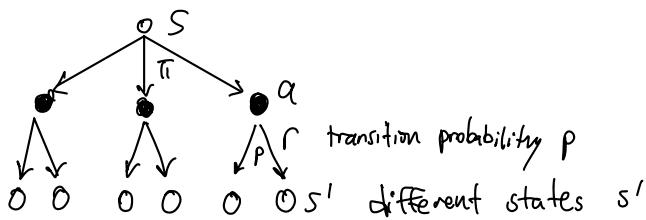
This is the Bellman Eq for V_π : it states a fundamental recursive relationship:

The value of being in a state =

immediate reward you can expect
+ (plus)

discounted value of where you'll end up next.

We can represent this with a back-up diagram:



The value function V_π is the unique solution to its Bellman Eq and this forms the basis of a number of ways to compute, approximate, and learn V_π .

Eq connecting $V_\pi(s)$ with $q_\pi(s, a)$:

A state-value is an expectation over the available action-values under a policy:

$$V_\pi(s) \doteq E_{\pi, p} [G_t \mid S_t = s] \quad // as before$$

$$= E_{\pi, p} [q_\pi(s, A_t) \mid S_t = s]$$

$$= \sum_a \pi(a \mid s) q_\pi(s, a)$$

// this derivation relates $V_\pi(s)$ to $q_\pi(s, a)$.

The following recursive relationship satisfied by $V_\pi(s)$ is called the Bellman Eq for state-values.

The value function $V_\pi(s)$ is the unique solution to its Bellman Eq.

$$\begin{aligned}
 V_\pi(s) &\doteq E_{\pi, p} [G_t \mid S_t = s] \\
 &= E_{\pi, p} [R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\
 &= \sum_a \pi(a \mid s) E_{\pi, p} [R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\
 &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) \cdot \\
 &\quad E_{\pi, p} [r + \gamma G_{t+1} \mid S_t = s, A_t = a, S_{t+1} = s'] \\
 &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) \cdot E_{\pi, p} [r + \gamma G_{t+1} \mid S_{t+1} = s'] \\
 &\quad // \text{ by the Markov Property} \\
 &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) \cdot \left\{ r + \gamma E[G_{t+1} \mid S_{t+1} = s'] \right\} \\
 &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma V_\pi(s')] \quad // V \xrightarrow{\text{to}} V \\
 &\quad \text{relationship}
 \end{aligned}$$

In summary :

$$\begin{aligned}
 V_\pi(s) &= \sum_a \left\{ \pi(a \mid s) q_\pi(s, a) \right\} \quad (\text{action-value perspective}) \\
 &= \sum_a \left\{ \pi(a \mid s) \sum_{s', r} \left[p(s', r \mid s, a) (r + \gamma V_\pi(s')) \right] \right\} \quad (\text{transition perspective})
 \end{aligned}$$

Bellman Eq for $q_{\pi}(s, a)$:

An action-value is an expectation over the possible next state-values under the environment's dynamics.

$$q_{\pi}(s, a) \stackrel{?}{=} E_{p, \pi} [G_t \mid S_t = s, A_t = a]$$

action-value function ; expected ; given that we are in
 for state s , taking ; total return ; state s and take
 action a , and then ; from time t ; action a ,
 following policy π . onward

$$= E_{p, \pi} [R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a]$$

$$= E_{p, \pi} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]$$

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

// this derivation relates $q_{\pi}(s, a)$ to $V_{\pi}(s')$:

$q \rightarrow v$ relationship.

The following recursive relationship satisfied by the Bellman equation for action-values, the action-value function $q_{\pi}(s, a)$ is the unique solution to its Bellman Eq.

$$q_{\pi}(s, a) \stackrel{?}{=} E_{p, \pi} [G_t \mid S_t = s, A_t = a] \quad // \text{as before}$$

$$= E_{p, \pi} [\underbrace{R_{t+1}}_{\text{immediate reward}} + \underbrace{\gamma G_{t+1}}_{\text{future reward, discounted}} \mid S_t = s, A_t = a]$$

total reward

$$\begin{aligned}
&= \sum_{s', r} p(s', r | s, a) E_{\pi, p} [r + \gamma G_{t+1} | S_t = s, A_t = a, S_{t+1} = s'] \\
&= \sum_{s', r} p(s', r | s, a) E_{\pi, p} [r + \gamma G_{t+1} | S_{t+1} = s'] // \text{Markov property} \\
&= \sum_{s', r} p(s', r | s, a) [r + \gamma E_{\pi, p} [G_{t+1} | S_{t+1} = s']] \\
&= \sum_{s', r} p(s', r | s, a) \\
&\quad \left[r + \gamma \sum_{a'} \pi(a' | s') E_{p, \pi} [G_{t+1} | S_{t+1} = s', A_{t+1} = a'] \right] \\
&= \sum_{s', r} p(s', r | s, a) \left[r + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a') \right]
\end{aligned}$$

// this derives the $q \Rightarrow q$ relationship.

In summary:

$$\begin{aligned}
q_{\pi}(s, a) &= \sum_{s', r} \left\{ p(s', r | s, a) [r + \gamma v_{\pi}(s')] \right\} // q-v \\
&= \sum_{s', r} \left\{ p(s', r | s, a) \left[r + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a') \right] \right\} // q-q
\end{aligned}$$

3.6 : Optimal Values and Optimal Policies

The Bellman equation shows us how to compute the value of a state for any policy π .

The Bellman optimality equations focus on optimal policy (the policy with the most reward over the long run).
Mathematically:

- $\pi \geq \pi'$ if $V_{\pi}(s) \geq V_{\pi'}(s) \quad \forall s \in S$
- the optimal policy is \geq all other policies.
- the optimal functions are called v_* and q_*

\Rightarrow Optimal Policy:

$$\pi_*(s) = \operatorname{argmax}_a q_*(s, a) \quad \forall s \in S$$

Optimal state-value function:

$$v_*(s) = \max_{\pi} V_{\pi}(s) \quad \forall s \in S$$

Optimal action value function:

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a) \quad \forall s \in S, a \in A(s)$$

For a state-action pair, the optimal function gives the expected return for taking action a in state s and thereafter following an optimal policy.

$$q_*(s, a) = E[R_{t+1} + \gamma V_*(S_{t+1}) | S_t = s, A_t = a]$$

Bellman Optimality Equations

It expresses the fact that the value of a state under the optimal policy must equal the expected return for the best action from that state:

$$\begin{aligned}
 V^*(s) &= \max_{a \in A(s)} q_{\pi^*}(s, a) \quad \forall s \in S \\
 &= \max_a E_{\pi^*} [G_t \mid S_t = s, A_t = a] \\
 &= \max_a E_{\pi^*} [R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\
 &= \max_a E_{\pi^*} [R_{t+1} + \gamma V^*(S_{t+1}) \mid S_t = s, A_t = a] \\
 V^*(s) &= \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma V^*(s')] \quad \text{recursion!}
 \end{aligned}$$

The last 2 equations are the form of Bellman optimality eq for V^* .

For q^* , we have:

$$\begin{aligned}
 q^*(s, a) &= E[R_{t+1} + \gamma \max_{a'} q^*(s_{t+1}, a') \mid S_t = s, A_t = a] \\
 &= \sum_{s', r} p(s', r \mid s, a) [r + \gamma \underbrace{\max_{a'} q^*(s', a')}_{V^*(s')}]
 \end{aligned}$$

? q^* is "better" than V^* , because q^* is a 2-step planning process. We know the values as well as the values of the actions in the next state.

$$\text{since } \pi_*(s) = \arg\max_a q_*(s, a)$$

$$\Rightarrow \pi_*(s) = \arg\max_a \sum_{s', r} p(s', r | s, a) [r + \gamma \underbrace{\max_{a'} q_*(s', a')}_{V_*(s')}]$$

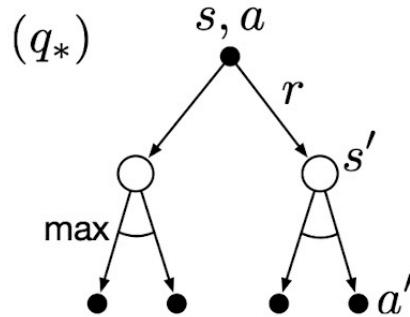
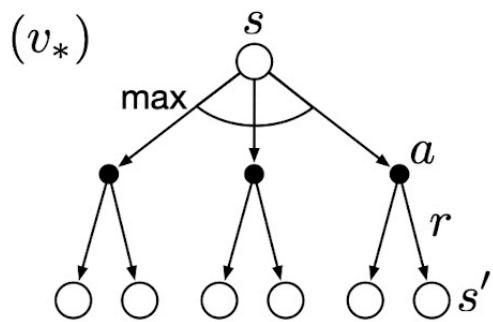


Figure 3.4: Backup diagrams for v_* and q_*

- For finite MDPs, the Bellman optimality Eq for V_* has a unique solution.
It's a system of equations, one for each state (one unknown for each state).
- Once we have V_* , it's easy to find the optimal policy (greedy wrt the optimal evaluation function V_*).

This solution relies on 3 assumptions which may not be true in reality;

- 1) we know the dynamics of the env.
- 2) we have enough computational resources.
- 3) Markov property.

3.7: Optimality and approximation

In a practical setting, we almost always have to settle for an approximation, we can do this cleverly:

make approximately optimal policies in states that have high probability of occurring, at the expense of making poor decisions in the states that have low probability of occurring.

3.8 : Summary

The optimal value function is unique, the optimal policy isn't.

3.10 : Partially Observable MDPs