Chapter 3: The Reinforcement Learning Problem

Objectives of this chapter:

- describe the RL problem we will be studying for the remainder of the course
- present idealized form of the RL problem for which we have precise theoretical results;
- introduce key components of the mathematics: value functions and Bellman equations;
- describe trade-offs between applicability and mathematical tractability.

The Agent Learns a Policy

Policy at step t, π_t :

a mapping from states to action probabilities $\pi_t(s, a) =$ probability that $a_t = a$ when $s_t = s$

- Reinforcement learning methods specify how the agent changes its policy as a result of experience.
- Roughly, the agent's goal is to get as much reward as it can over the long run.

The Agent-Environment Interface



Agent and environment interact at discrete time steps : t = 0, 1, 2, ...Agent observes state at step t: $s_t \in S$ produces action at step t: $a_t \in A(s_t)$ gets resulting reward : $r_{t+1} \in \Re$ and resulting next state : s_{t+1} $r_{t+1} = \frac{r_{t+2}}{s_t} \frac{r_{t+3}}{s_{t+3}} \frac{r_{t+3}}{s_{t+3}} \frac{r_{t+3}}{s_{t+3}} \cdots$

Getting the Degree of Abstraction Right

2

- Time steps need not refer to fixed intervals of real time.
- Actions can be low level (e.g., voltages to motors), or high level (e.g., accept a job offer), "mental" (e.g., shift in focus of attention), etc.
- States can be low-level "sensations", or they can be abstract, symbolic, based on memory, or subjective (e.g., the state of being "surprised" or "lost").
- An RL agent is not like a whole animal or robot, which consist of many RL agents as well as other components.
- The environment is not necessarily unknown to the agent, only incompletely controllable.
- Reward computation is in the agent's environment because the agent cannot change it arbitrarily.

Goals and Rewards

- Is a scalar reward signal an adequate notion of a goal? maybe not, but it is surprisingly flexible.
- A goal should specify what we want to achieve, not how we want to achieve it.
- A goal must be outside the agent's direct control—thus outside the agent.
- The agent must be able to measure success:
 - explicitly;
 - frequently during its life-span.

Returns

Suppose the sequence of rewards after step t is:

 $r_{t+1}, r_{t+2}, r_{t+3}, \dots$ What do we want to maximize?

In general,

we want to maximize the **expected return**, $E\{R_t\}$, for each step *t*.

Episodic tasks: interaction breaks naturally into episodes, e.g., plays of a game, trips through a maze.

$$R_t = r_{t+1} + r_{t+2} + \dots + r_T,$$

where *T* is a final time step at which a **terminal state** is reached, ending an episode.

Returns for Continuing Tasks

Continuing tasks: interaction does not have natural episodes.

Discounted return:

$$R_{t} = r_{t+1} + \gamma r_{t+2} + \gamma^{2} r_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1}$$

where $\gamma, 0 \le \gamma \le 1$, is the **discount rate**

shortsighted
$$0 \leftarrow \gamma \rightarrow 1$$
 farsighted

Avoid failure: the pole falling beyond a critical angle or the cart hitting end of track. As an episodic task where episode ends upon failure: reward = +1 for each step before failure \Rightarrow return = number of steps before failure As a continuing task with discounted return: reward = -1 upon failure; 0 otherwise \Rightarrow return is related to $-\gamma^k$, for k steps before failure In either case, return is maximized by avoiding failure for as long as possible.



A Unified Notation

- In episodic tasks, we number the time steps of each episode starting from zero.
- We usually do not have to distinguish between episodes, so we write s_t instead of $s_{t,j}$ for the state at step t of episode j.
- Think of each episode as ending in an absorbing state that always produces reward of zero:

$$(s_0 \xrightarrow{r_1 = +1} (s_1 \xrightarrow{r_2 = +1} (s_2 \xrightarrow{r_3 = +1} (s_2 \xrightarrow{r_4 = } r_5 \xrightarrow{r_4 \xrightarrow{r_4 = } r_5 \xrightarrow{r_4 = } r_5 \xrightarrow{r_4 \xrightarrow{r_4 = } r_5 \xrightarrow{r_4 \xrightarrow{r_4 \to } r_5 \xrightarrow{r_4 \xrightarrow{r_4 \to } r_5 \xrightarrow{r_4 \xrightarrow{r_4 \to } r_5 \xrightarrow{r_4 \xrightarrow{r_4 \xrightarrow{r_4 \to } r_5 \xrightarrow{r_4 xr_4 \xrightarrow{r_4 \xrightarrow{r_4 xr_4 \xrightarrow{r_4 \xrightarrow{r_4 xr_4 \xrightarrow{r_4$$

• We can cover all cases by writing



where γ can be 1 only if a zero reward absorbing state is always reached.

The Markov Property

- By "the state" at step *t*, the book means whatever information is available to the agent at step *t* about its environment.
- The state can include immediate "sensations," highly processed sensations, and structures built up over time from sequences of sensations.
- Ideally, a state should summarize past sensations so as to retain all "essential" information, i.e., it should have the **Markov Property**:

$$\Pr\{s_{t+1} = s', r_{t+1} = r \mid s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_1, s_0, a_0\} = \Pr\{s_{t+1} = s', r_{t+1} = r \mid s_t, a_t\}$$
for all s', r, and histories $s_t, a_t, r_t, s_{t-1}, a_{t-1}, \dots, r_t, s_0, a_0$.

II

Markov Decision Processes

- If a reinforcement learning task has the Markov Property, it is basically a **Markov Decision Process (MDP)**.
- If state and action sets are finite, it is a finite MDP.
- To define a finite MDP, you need to give:
 - state and action sets
 - one-step "dynamics" defined by transition probabilities

$$P^a_{ss'} = \Pr\left\{s_{t+1} = s' \mid s_t = s, a_t = a\right\} \text{ for all } s, s' \in S, a \in A(s).$$

• reward expectations:

$$R_{ss'}^{a} = E\left\{r_{t+1} \mid s_{t} = s, a_{t} = a, s_{t+1} = s'\right\} \text{ for all } s, s' \in S, \ a \in A(s).$$

An Example Finite MDP

Recycling Robot

- At each step, robot has to decide whether it should (1) actively search for a can, (2) wait for someone to bring it a can, or (3) go to home base and recharge.
- Searching is better but runs down the battery; if runs out of power while searching, has to be rescued (which is bad).
- Decisions made on basis of current energy level: high, low.

13



1-a, R search

14

16

14

Value Functions

• The value of a state is the expected return starting from that state; depends on the agent's policy:

State - value function for policy π :

$$V^{\pi}(s) = E_{\pi} \left\{ R_{t} \mid s_{t} = s \right\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} \mid s_{t} = s \right\}$$

• The value of taking an action in a state under policy π is the expected return starting from that state, taking that action, and thereafter following π :

Action - value function for policy π :

$$Q^{\pi}(s,a) = E_{\pi} \left\{ R_{t} \mid s_{t} = s, a_{t} = a \right\} = E_{\pi} \left\{ \sum_{k=0}^{\infty} \gamma^{k} r_{t+k+1} \mid s_{t} = s, a_{t} = a \right\}$$
15

Bellman Equation for a Policy π

The basic idea:

 $\alpha, \mathcal{R}^{\text{search}}$

13

$$R_{t} = r_{t+1} + \gamma r_{t+2} + \gamma^{2} r_{t+3} + \gamma^{3} r_{t+4} \cdots$$
$$= r_{t+1} + \gamma \left(r_{t+2} + \gamma r_{t+3} + \gamma^{2} r_{t+4} \cdots \right)$$
$$= r_{t+1} + \gamma R_{t+1}$$

So:

$$V^{\pi}(s) = E_{\pi} \{ R_{t} | s_{t} = s \}$$

$$= E_{\pi} \{ r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_{t} = s \}$$

Or, without the expectation operator:

$$V^{\pi}(s) = \sum_{a} \pi(s, a) \sum_{s'} P^{a}_{ss'} \Big[R^{a}_{ss'} + \gamma V^{\pi}(s') \Big]$$
16

More on the Bellman Equation

$$V^{\pi}(s) = \sum_{a} \pi(s, a) \sum_{s'} P^{a}_{ss'} \Big[R^{a}_{ss'} + \gamma V^{\pi}(s') \Big]$$

This is a set of equations (in fact, linear), one for each state. The value function for π is its unique solution.

Backup diagrams:





Gridworld

- Actions: north, south, east, west; deterministic.
- If would take agent off the grid: no move but reward = -1
- Other actions produce reward = 0, except actions that move agent out of special states A and B as shown.



Optimal Value Functions

- For finite MDPs, policies can be **partially ordered**: $\pi \ge \pi'$ if and only if $V^{\pi}(s) \ge V^{\pi'}(s)$ for all $s \in S$
- There is always at least one (and possibly many) policies that is better than or equal to all the others. This is an **optimal policy**. We denote them all π *.
- Optimal policies share the same **optimal state-value function**:

$$V^*(s) = \max_{\pi} V^{\pi}(s) \text{ for all } s \in S$$

• Optimal policies also share the same **optimal action-value function**:

 $Q^*(s, a) = \max Q^{\pi}(s, a)$ for all $s \in S$ and $a \in A(s)$

This is the expected return for taking action a in state s and thereafter following an optimal policy.



Bellman Optimality Equation for V*

The value of a state under an optimal policy must equal the expected return for the best action from that state:



Bellman Optimality Equation for Q^* $Q^*(s, a) = E\left\{r_{t+1} + \gamma \max_{a'} Q^*(s', a') | s_t = s, a_t = a\right\}$ $= \sum_{s'} P^a_{ss'} \left[R^a_{ss'} + \gamma \max_{a'} Q^*(s', a')\right]$ The relevant backup diagram: Q^* is the unique solution of this system of nonlinear equations. Why Optimal State-Value Functions are Useful Any policy that is greedy with respect to V^* is an optimal policy. Therefore, given V^* , one-step-ahead search produces the long-term optimal actions. E.g., back to the gridworld: $\underbrace{\square A + B + 5 - 198 + 178 + 160 + 144 + 130 - 178 + 160 + 144 + 130 - 178 + 160 + 144 + 130 - 1178 + 160 + 144 + 130 - 1178 + 160 + 144 + 130 + 178 + 18$



 Agent-environment interaction States Actions Rewards Policy: stochastic rule for selecting actions Return: the function of future rewards agent tries to maximize Episodic and continuing tasks Markov Property Markov Decision Process Transition probabilities Expected rewards Value functions State-value function for a policy Action-value function for a policy Optimal state-value function Optimal action-value functions Optimal value functions Optimal policies Bellman Equations The need for approximation 	Summary	
27	 Agent-environment interaction States Actions Rewards Policy: stochastic rule for selecting actions Return: the function of future rewards agent tries to maximize Episodic and continuing tasks Markov Property Markov Decision Process Transition probabilities Expected rewards 	 Value functions State-value function for a policy Action-value function for a policy Optimal state-value function Optimal action-value functions Optimal value functions Optimal policies Bellman Equations The need for approximation

Solving the Bellman Optimality Equation

- Finding an optimal policy by solving the Bellman Optimality Equation requires the following:
 - accurate knowledge of environment dynamics;
 - we have enough space and time to do the computation;
 - the Markov Property.
- How much space and time do we need?
 - polynomial in number of states (via dynamic programming methods; Chapter 4),
 - BUT, number of states is often huge (e.g., backgammon has about 10**20 states).
- We usually have to settle for approximations.
- Many RL methods can be understood as approximately solving the Bellman Optimality Equation.