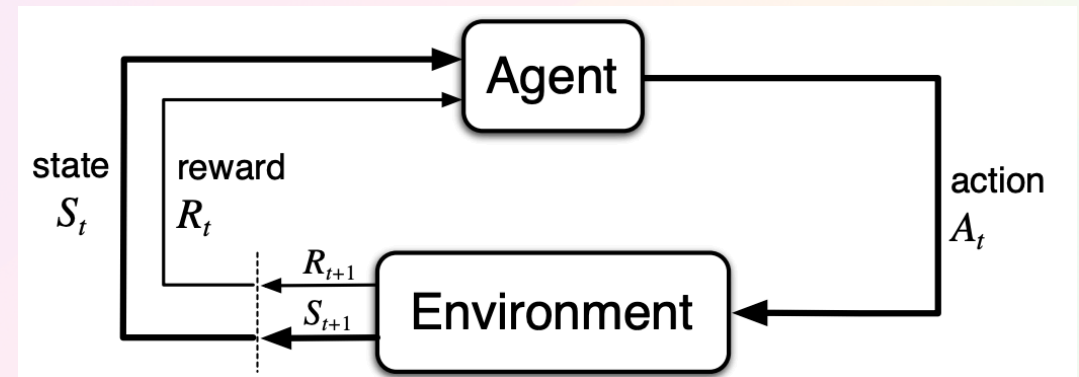


Markov Decision Processes

David I. Inouye

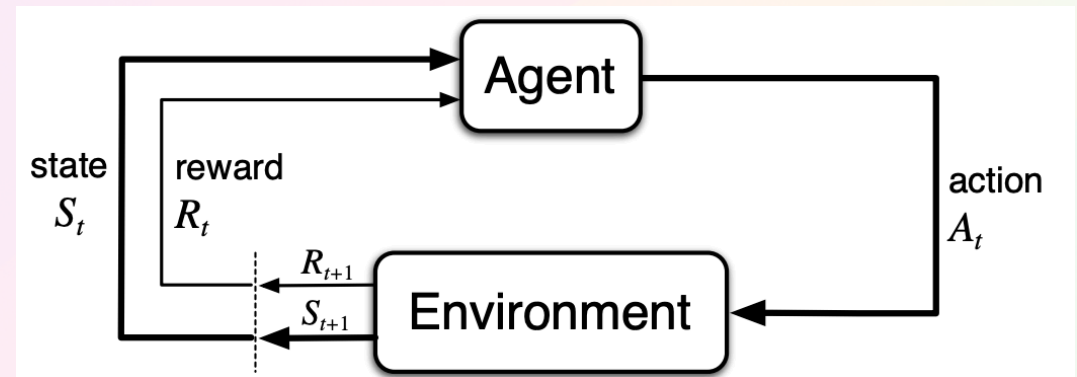
The Agent and Environment Are the Two Key Actors in RL – How Should They Be Defined?

- **Example: Humanoid robot**
 - Is the “agent” the whole robot?
 - What “actions” can the “agent” take? Increase power to left leg motor?
- The agent is probably the “controller” of the robot.
 - In fact, the position of arms and legs might be part of the “environment”.
- Suppose a leg breaks; the controller cannot fix the leg by itself, so the leg is part of the environment.



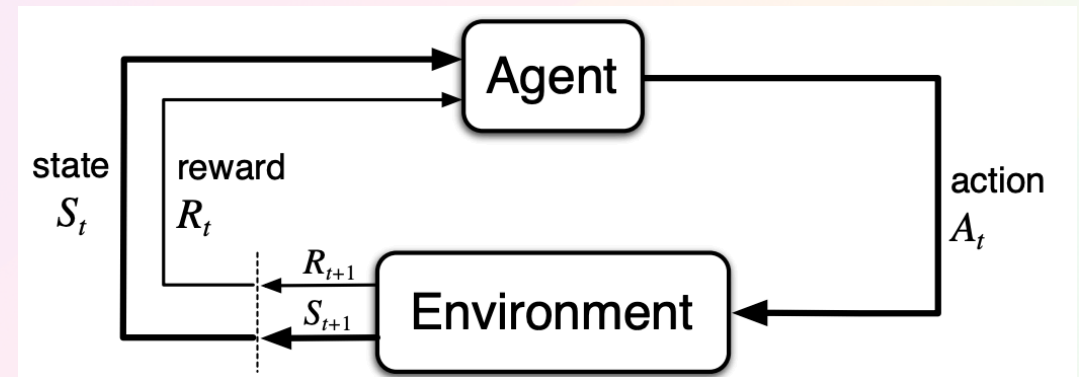
The Agent and Environment Are the Two Key Actors in RL – How Should They Be Defined?

- **Example: Oil refinery**
 - Is the “agent” the whole refinery?
 - What “actions” can the “agent” take? Change temperature? Increase electricity?
- The agent is again just the “controller” part of the refinery.
 - Each module may have its own controller/agent.
 - The top-level controller can only “ask” for something (i.e., a command), but it may not happen.



The Agent and Environment Are the Two Key Actors in RL – How Should They Be Defined?

- The **agent** should only be defined as the part that can be **directly and explicitly** controlled via concrete actions.
 - Actions - Select how much power send to each robotic motor.
 - Actions - Change fan speed
 - Actions - Request subcontrollers to produce a certain amount.
 - Note: The boundary should be defined by what can be *controlled* rather than what is *known*.
- The **environment** is anything that is NOT the agent.
 - Actual position of the robot body.
 - Actual temperature in boiler.



Markov Decision Processes (MDP) Mathematically Formalize RL Problems Via Random Variables

- **Sensation / Observation**

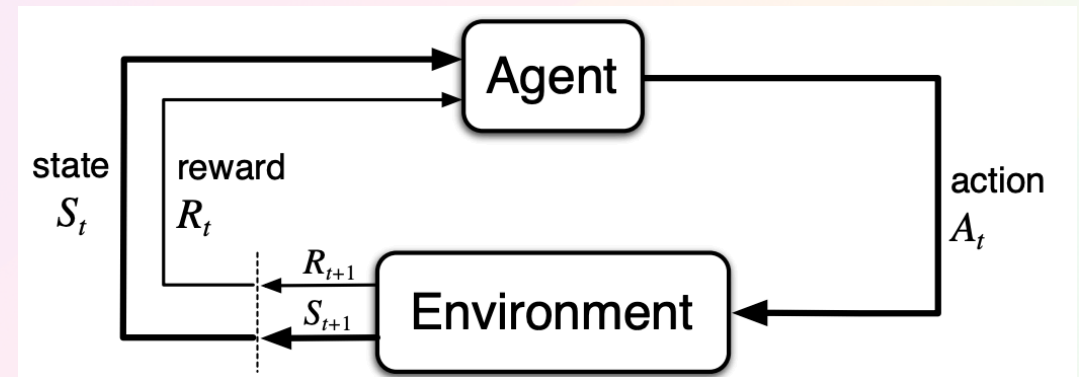
- Observing the environment.
- Encoded as a random variable called state S_t .

- **Action**

- Agent's ability to interact with the environment.
- Encoded as a random variable called action A_t .

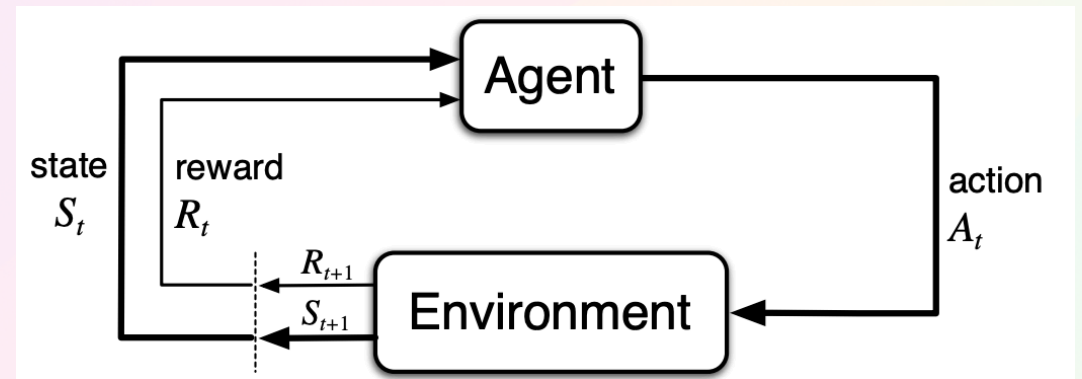
- **Reward**

- Immediate feedback from the environment.
- Encoded as a random variable called reward R_t .



MDPs Model a **Discrete** Sequence of States, Actions, and Rewards

- At each discrete timestep t :
 - Agent receives environment state $S_t \in \mathcal{S}$.
 - Agent decides action based on state, $A_t(S_t) \in \mathcal{A}(S_t)$.
 - Agent receives reward and new state $R_{t+1} \in \mathbb{R}, S_{t+1} \in \mathcal{S}$.
- This creates a sequence: $S_0, A_0, R_1, S_1, A_1, R_2, \dots$



MDPs Model the Joint Distribution of This Sequence by Assuming Conditional Independence

- Using the chain rule, we could define the probability of the following sequence

$$S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \dots$$

$$p(S_0)p(A_0 | S_0)p(S_1, R_1 | S_0, A_0)p(A_1 | S_0, A_0, S_1, R_1)p(S_2, R_2 | S_0, A_0, S_1, R_1, A_1) \dots$$

- **MDPs assume conditional independence:**

- **Action** only depends on **current state**.
- **New state and reward** only depend on **current state and action** (Markov part of MDP).

$$p(S_0) p(A_0 | S_0) p(S_1, R_1 | S_0, A_0) p(A_1 | \cancel{S_0}, \cancel{A_0}, S_1, \cancel{R_1}) p(S_2, R_2 | \cancel{S_0}, \cancel{A_0}, \cancel{R_1}, S_1, A_1) \dots$$

- The **environment's dynamics** are completely represented by:

$$p(S_t, R_t | S_{t-1}, A_{t-1})$$

- The **agent's dynamics** are completely represented by a policy:

$$\pi(A_t | S_t) \triangleq p(A_t | S_t)$$

The Goal or Ultimate Objective Is to Find a Policy That Maximizes the (weighted) Sum of **Future** Rewards

- While rewards are immediate, the goal is to maximize the long-term expected return, i.e., sum of **future rewards**
- **Episodic tasks** have a terminal state.
 - Expected return could be a simple sum: $G_t = \sum_{k=0}^T R_{t+k+1}$
- **Continual tasks** do not have a clear terminal state.
 - A simple sum of rewards would not converge (i.e., infinite)
 - Examples: Controlling an oil refinery or a robot learning to walk
 - Expected return is usually a **weighted sum** of future rewards where γ is the **discount rate**:

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}, 0 < \gamma < 1$$

- If we assume for episodic tasks, $\forall t > T, R_t = 0$, we can simply use the notation for continual tasks but allow $\gamma_t = 1$ if it is episodic.

Returns at the Current Timestep Are Equal to Rewards Plus Discounted Returns at Next Timestep

$$\begin{aligned} G_t &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \\ &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

- Thus, returns are recursively related to each other
 - This simple relationship is fundamental to theoretic development later

The **State-Value Function** Defines the Expected Return of a State **Given** a Particular Policy π

- Remember, rewards are easy. Estimating value is hard.
- Estimating value is the key to most RL algorithms.
- First, we define the value of **states**.
- The **state-value function** for policy π is:

$$v_{\pi}(s) \triangleq \mathbb{E}_{\pi}[G_t | S_t = s] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

- The expectation is based on the given policy π (e.g., random/greedy policy)
- How “good” is this state, given that the agent follows π for all actions afterwards?
- Important: This state-value function will be different for different policies.

The **Bellman Equation** for v_π Is a Special Recursive Equation Whose Solution Is the Value Function

$$\begin{aligned}v_\pi(s) &\triangleq \mathbb{E}_\pi[G_t \mid S_t = s] \\&= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \quad (\text{see previous slide}) \\&= \sum_{a,r,s',a',r',s'' \dots} \pi(a \mid s)p(s', r \mid s, a)\pi(a' \mid s')p(s'', r' \mid s', a') \dots [\dots] \\&= \sum_{a,r,s'} \pi(a \mid s)p(s', r \mid s, a) \left[r + \gamma \sum_{a',r',s'', \dots} \pi(a' \mid s')p(s'', r' \mid s', a')G_{t+1} \right] \\&= \sum_{a,r,s'} \pi(a \mid s)p(s', r \mid s, a) [r + \gamma \mathbb{E}_\pi[G_{t+1} \mid S_{t+1} = s']] \\&= \sum_{a,r,s'} \pi(a \mid s)p(s', r \mid s, a) [r + \gamma v_\pi(s')]\end{aligned}$$

- $v_\pi(s) = \sum_{a,r,s'} \pi(a|s)p(s', r|s, a)[r + \gamma v_\pi(s')]$

Example: Solved Bellman Equation (which is a system of linear equations) to Find the Value Function for a Random Policy

Example 3.5: Gridworld Figure 3.2 (left) shows a rectangular gridworld representation of a simple finite MDP. The cells of the grid correspond to the states of the environment. At each cell, four actions are possible: **north**, **south**, **east**, and **west**, which deterministically cause the agent to move one cell in the respective direction on the grid. Actions that would take the agent off the grid leave its location unchanged, but also result in a reward of -1 . Other actions result in a reward of 0 , except those that move the agent out of the special states A and B. From state A, all four actions yield a reward of $+10$ and take the agent to A'. From state B, all actions yield a reward of $+5$ and take the agent to B'.

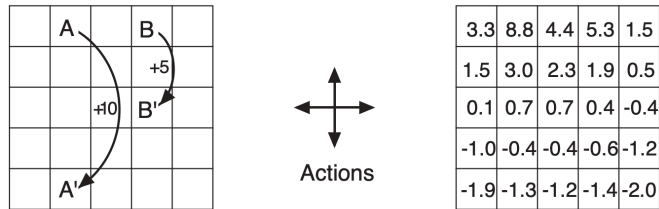


Figure 3.2: Gridworld example: exceptional reward dynamics (left) and state-value function for the equiprobable random policy (right).

- $$v_{\pi}(s) = \sum_{a,r,s'} \pi(a | s) p(s', r | s, a) [r + \gamma v_{\pi}(s')]$$
- $$v_{\pi}([1, 1]) = N + S + E + W$$
- $$= \frac{1}{4} \left((-1 + \gamma \cdot 3.3) + (0 + \gamma \cdot 1.5) + (0 + \gamma \cdot 8.8) + (-1 + \gamma \cdot 3.3) \right)$$
- $$= \frac{1}{4} \left((-1 + 0.9 \cdot 3.3) + (0 + 0.9 \cdot 1.5) + (0 + 0.9 \cdot 8.8) + (-1 + 0.9 \cdot 3.3) \right)$$
- $$= 3.3 = v_{\pi}([1, 1])$$

Example copied verbatim from page 60 of Sutton, R. S., & Barto, A. G. (2018).

Reinforcement learning: An introduction. MIT Press.

The **Action-Value Function** Defines the Expected Return of Taking an Action Given a Current State and a Particular Policy π

- We now define the value of an **action** given a current state and policy π .
- The **action-value function** for policy π is defined as:

$$q_{\pi}(s, a) \triangleq \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

- The key difference is that now we condition on the next action a
- The expectation is based on the given policy π (e.g., random/greedy policy)
- How “good” is the **action, given** that the agent first takes action a and then follows π afterwards?
- **Important:** This action-value function will be different for different policies.

Optimal Policies and Optimal Value Functions Can Be Defined and Exist Theoretically

- How should we define an optimal policy?
 - We want it to mean that this policy will produce the most expected return (i.e., long-term reward).
- One policy π is better than another policy π' if its value function is better for all possible states:

$$\pi \geq \pi' \iff v_\pi(s) \geq v_{\pi'}(s), \quad \forall s \in \mathcal{S}$$

- There exists at least one policy that is better or equal to all others called an **optimal policy** π_* .

$$\pi_* \geq \pi, \quad \forall \pi$$

- All optimal policies share the same **optimal state-value function**:

$$v_*(s) \triangleq \max_{\pi} v_\pi(s), \quad \forall s \in \mathcal{S}$$

- Similarly, they share the same **optimal action-value function**:

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

The **Bellman Optimality Equation** for v_* Enables Solving for the **Optimal** State-Value Function

$$\begin{aligned}v_*(s) &= \max_a q_{\pi_*}(s, a) \\&= \max_a \mathbb{E}[G_t | S_t = s, A_t = a] \\&= \max_a \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\&= \max_a \mathbb{E}[R_{t+1} + \gamma v_*(S_{t+1}) | S_t = s, A_t = a] \\&= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]\end{aligned}$$

- No optimal policy here! Just the model part since we are taking maximum over actions.
- For finite MDPs, this can be solved via a system of non-linear equations.

The **Bellman Optimality Equation** for q_* Enables Solving for the **Optimal** Action-Value Function

$$\begin{aligned}q_*(s, a) &= \mathbb{E}[R_{t+1} + \gamma \max_a q_*(s, a) | S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} q_*(s', a')]\end{aligned}$$

- Again, no policy here, just the q function! This can be solved as well.
- This can be solved as well.

The Optimal Value Functions Can Be Used to Construct an Optimal Policy!

- If we have $v_*(s)$, we can simply choose the optimal action that will maximize the value after taking one action :

$$a_t^* = \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')]$$

- Called one-step search or **greedy search** with respect to state-value function $v_*(s)$ (rather than greedy with respect to the expected immediate reward $\mathbb{E}[R_t | A_t = a]$)
 - **Greedy/local** w.r.t. v_* is **globally** optimal.
- If we have $q_*(s, a)$, then the optimal action is even simpler:

$$a_t^* = \arg \max_a q_*(s, a)$$

- The action-value function “caches” the one-step search values.

Example: With the Gridworld Optimal Value Function, We Can Define the Optimal Policy

Example 3.8: Solving the Gridworld Suppose we solve the Bellman equation for v_* for the simple grid task introduced in Example 3.5 and shown again in Figure 3.5 (left). Recall that state A is followed by a reward of +10 and transition to state A', while state B is followed by a reward of +5 and transition to state B'. Figure 3.5 (middle) shows the optimal value function, and Figure 3.5 (right) shows the corresponding optimal policies. Where there are multiple arrows in a cell, all of the corresponding actions are optimal.

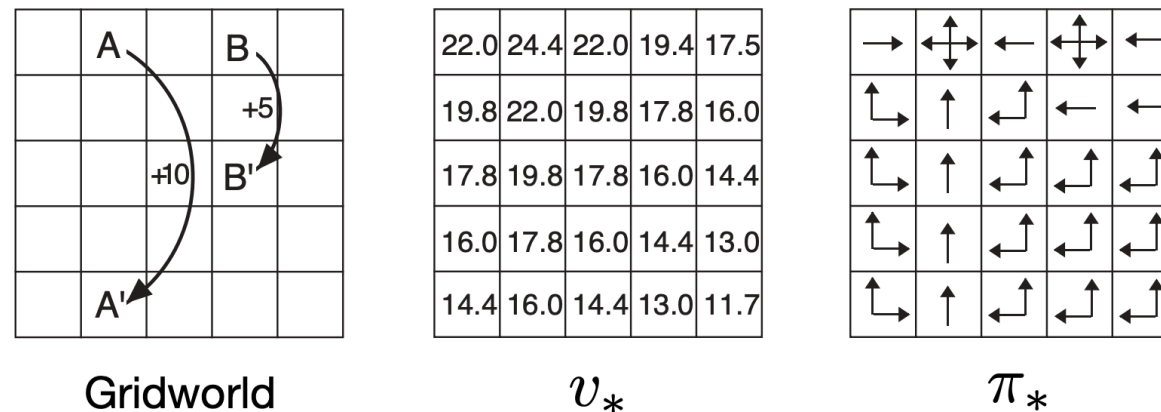


Figure 3.5: Optimal solutions to the gridworld example. ■

Example copied verbatim from page 65 of Sutton, R. S., & Barto, A. G. (2018).

Reinforcement learning: An introduction. MIT Press.

MDPs solved!?! No... Solution relies on several difficult assumptions

1. The dynamics of the environment are known
(i.e., $p(S_t, R_t \mid S_{t-1}, A_{t-1})$ is known perfectly)
 - Except in simple simulations like grid world, these dynamics are rarely known
 - Even for simple games like chess, the opponent's strategy is unknown so the dynamics are unknown
2. Computational resources are sufficient for this calculation
 - For example, even a simple backgammon game has 10^{20} states
3. The states have the Markov property
 - All relevant information of the past must be represented in the environment state
 - It is often difficult to ensure this property in real-world scenarios

RL is about approximating solutions to these MDPs

- While optimal solutions are almost never possible, approximations can still be quite useful
- Most RL algorithms focus on estimating the value functions in some way
- RL algorithms often substitute knowledge of the environment with actual experience in place of knowledge (like in bandits)
- One saving grace in practice is that not all states are equally likely
 - This is akin to the notion that real-world data (e.g., images) is much simpler than the high-dimensional space it lives in
 - Thus, generalization beyond the theoretic “worst case” is possible in practice
- While there are other algorithms for solving MDPs, **the online nature of RL distinguishes it from other approaches to MDPs**
 - It can collect more information on the “common” states and thus do well

Summary of MDPs

- Agent and environment mappings for problems
 - The agent is defined as what is **directly controllable** (NOT mere knowledge)
 - The environment is everything else
- MDPs model this interaction with **states**, **actions** and **rewards**
 - The environment is defined by the transition distribution
$$p(S_t, R_t \mid S_{t-1}, A_{t-1})$$
 - The agent is defined by a policy $\pi(A_t \mid S_t)$
 - Together, they define a joint distribution over sequences of states, actions and rewards
- The agent attempts to optimize the **expected return**, which is the discounted sum of future rewards
- The **state-value and action-value functions** represent the long-term value of states or actions
- An **optimal policy** can be constructed from optimal value functions
- In practice, RL algorithms usually approximate these value functions in an **online manner** even if the environment is unknown

Reference

- Based on the excellent RL book by Sutton and Barto
 - <http://incompleteideas.net/book/the-book-2nd.html>