

IS 7033: Artificial Intelligence and Machine Learning

Dr. Peyman Najafirad (Paul Rad)

Associate Professor

Cyber Analytics and AI

210.872.7259

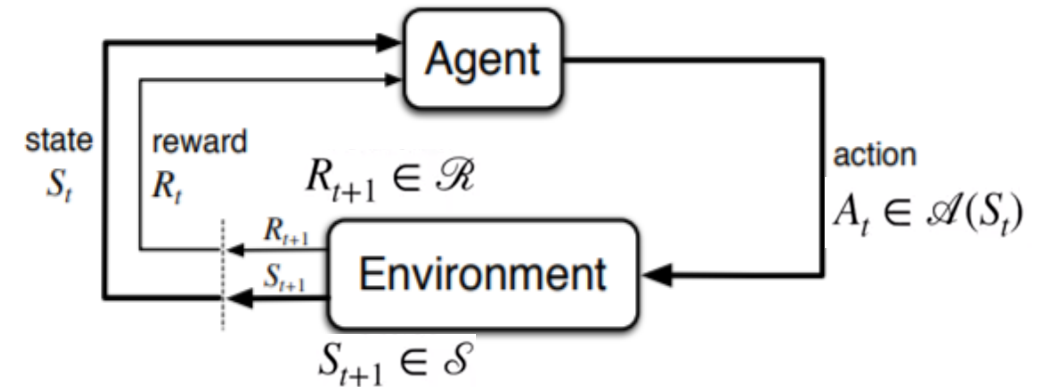
<https://github.com/paulNrad/ProbabilisticGraphModels>

Reinforcement Learning

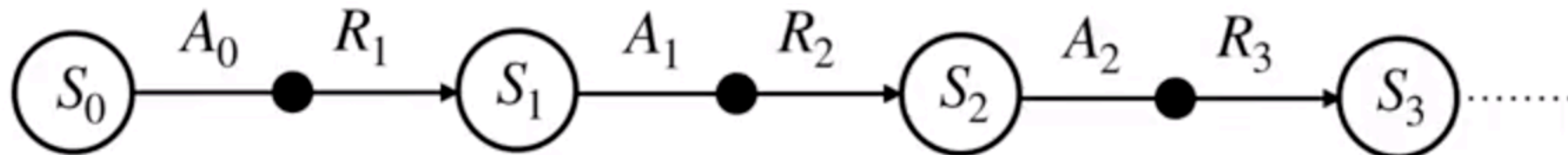
Read Reinforcement Learning by Richard Sutton Chapter 3 and 4

Reinforcement Learning Summary

- The reinforcement learning (RL) framework is characterized by an **agent** learning to interact with its **environment**.
- At each time step, the agent receives the environment's **state** (*the environment present a situation to the agent*), and the agent must choose an appropriate **action** in response. One time step later, the agent receives a **reward** (*the environment indicates whether the agent has responded appropriately to the state*) and a new **state**.
- All agents have the goal to maximize expected **cumulative reward**, or the expected sum of rewards attained over all time steps.



[The agent-environment interaction in reinforcement learning. \(Source: Sutton and Barto, 2017\)](#)



Episodic vs. Continuing Tasks

- **Continuing tasks** are tasks that continue forever, without end.
- **Episodic tasks** are tasks with a well-defined starting and ending point. In this case, we refer to a complete sequence of interaction, from start to finish, as an **episode**. Episodic tasks come to an end whenever the agent reaches a **terminal state**.

Goals, Cumulative Reward, and Discounted Return

- The **return at time step t** is $G_t := R_{t+1} + R_{t+2} + R_{t+3} + \dots$
- The agent selects actions with the goal of maximizing expected (discounted) return.
- The **discounted return at time step t** is $G_t := R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{n=0}^{\infty} \gamma^n R_{t+1+n}$
 - The **discount rate γ** is something that you set, to refine the goal that you have the agent. It must satisfy $0 \leq \gamma \leq 1$.
 - If $\gamma=0$, the agent only cares about the most immediate reward.
 - If $\gamma=1$, the return is not discounted.
 - For larger values of γ , the agent cares more about the distant future. Smaller values of γ result in more extreme discounting, where - in the most extreme case - agent only cares about the most immediate reward.

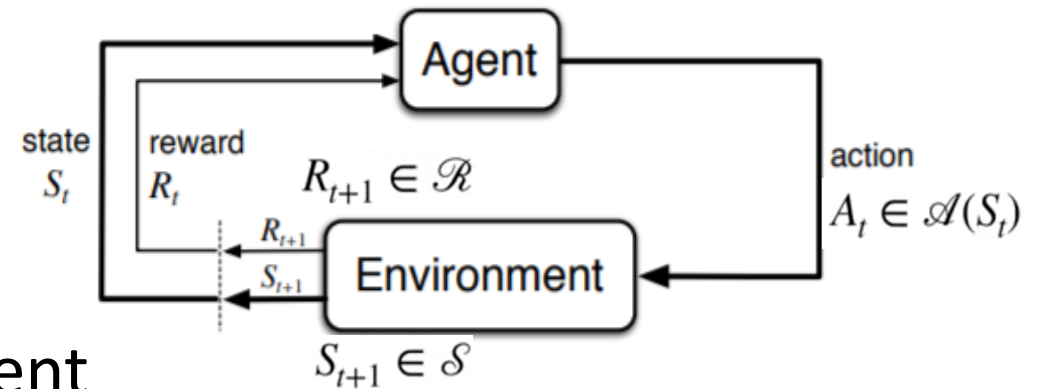
MDPs and One-Step Dynamics

A (finite) Markov Decision Process (MDP) is defined by:

- A (finite) set of states S
- A (finite) set of actions A
- A set of rewards R
- The one-step dynamics of the environment

$$p(s', r | s, \alpha) = P(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = \alpha) \text{ for all } s, s', \alpha, \text{ and } r$$

- A discount rate $\gamma \in [0, 1]$



We can start to think of the solution as a series of actions that need to be learned by the agent towards the pursuit of a goal.

Policies

A policy determines how an **agent chooses an action in response to the current state**. In other words, it specifies how the agent responds to situations that the environment has presented.

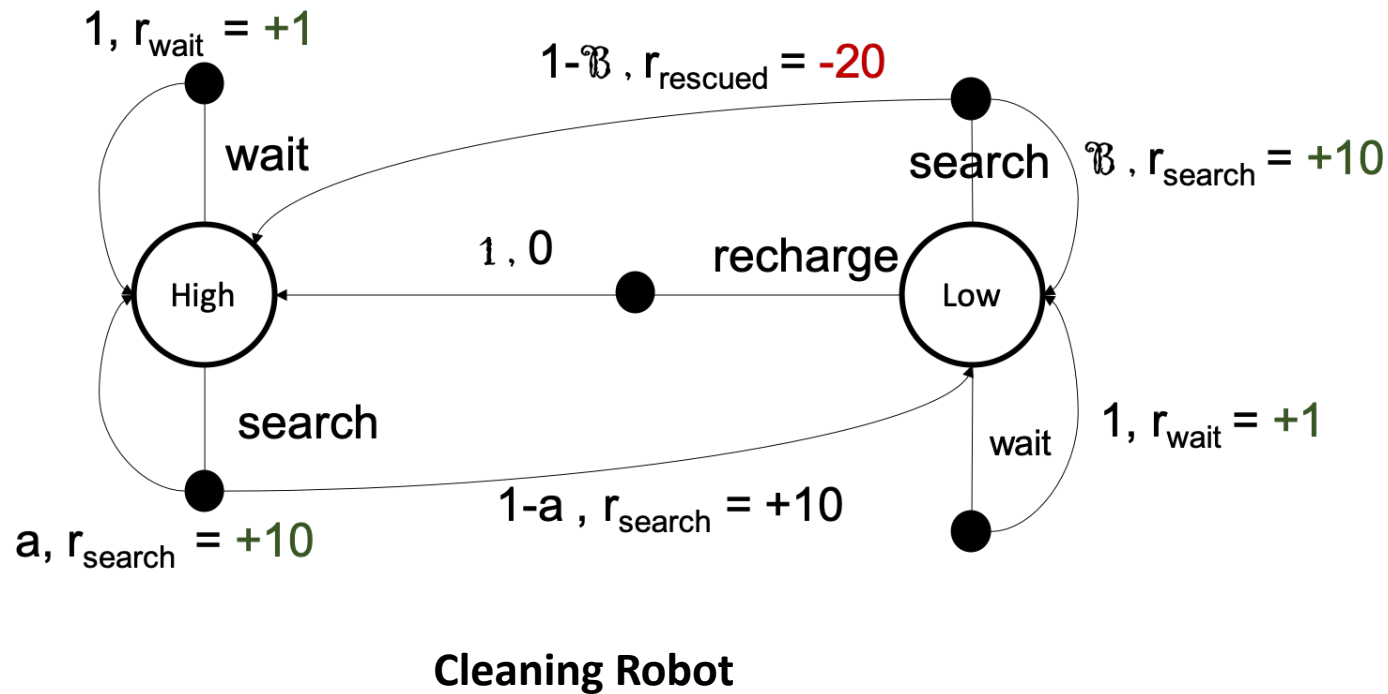
- A **deterministic policy** is mapping from the set of environment states to the set of possible actions:

$$\pi : S \rightarrow A$$

- A **stochastic policy** is a mapping $\pi : S \times A \rightarrow [0,1]$

$$\pi (\alpha \mid s) = P (A_t = \alpha \mid S_t = s)$$

Policies



Stochastic Policy

$$\pi(\text{recharge} \mid \text{low}) = 0.5$$

$$\pi(\text{search} \mid \text{low}) = 0.1$$

$$\pi(\text{wait} \mid \text{low}) = 0.4$$

$$\pi(\text{search} \mid \text{high}) = 0.9$$

$$\pi(\text{wait} \mid \text{high}) = 0.1$$

Deterministic Policy

$$\pi(\text{recharge} \mid \text{low}) = 1$$

$$\pi(\text{search} \mid \text{high}) = 1$$

Now that we know how to establish a policy, what step can we take to make sure the agent's policy is the best one (optimal policy).

Question: Consider a different stochastic policy where:

$$\pi(\text{recharge} \mid \text{low}) = 0.3$$

$$\pi(\text{search} \mid \text{low}) = 0.2$$

$$\pi(\text{wait} \mid \text{low}) = 0.5$$

$$\pi(\text{search} \mid \text{high}) = 0.6$$

$$\pi(\text{wait} \mid \text{high}) = 0.4$$

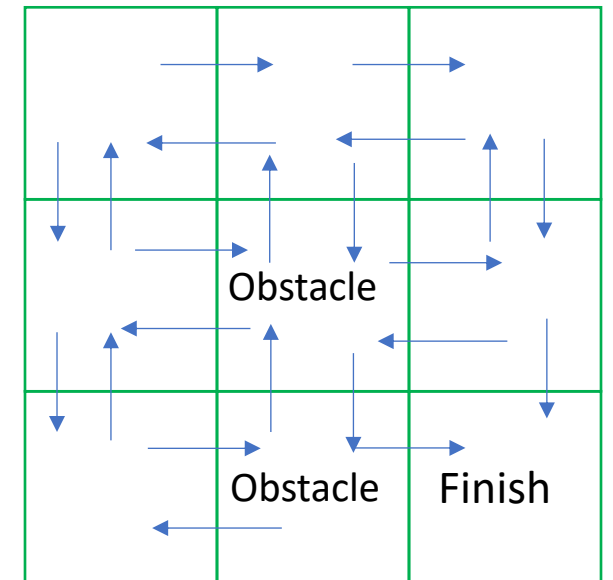
Which of the following statements are true, if the agent follows the policy?

- a) If the battery is low, the agent will always decide to wait for cans.
- b) If the battery level is high, the agent chooses to search for a can with 60% probability, and otherwise waits for a can.
- c) If the battery level is low, the agent is most likely to decide to wait for cans.

Grid world

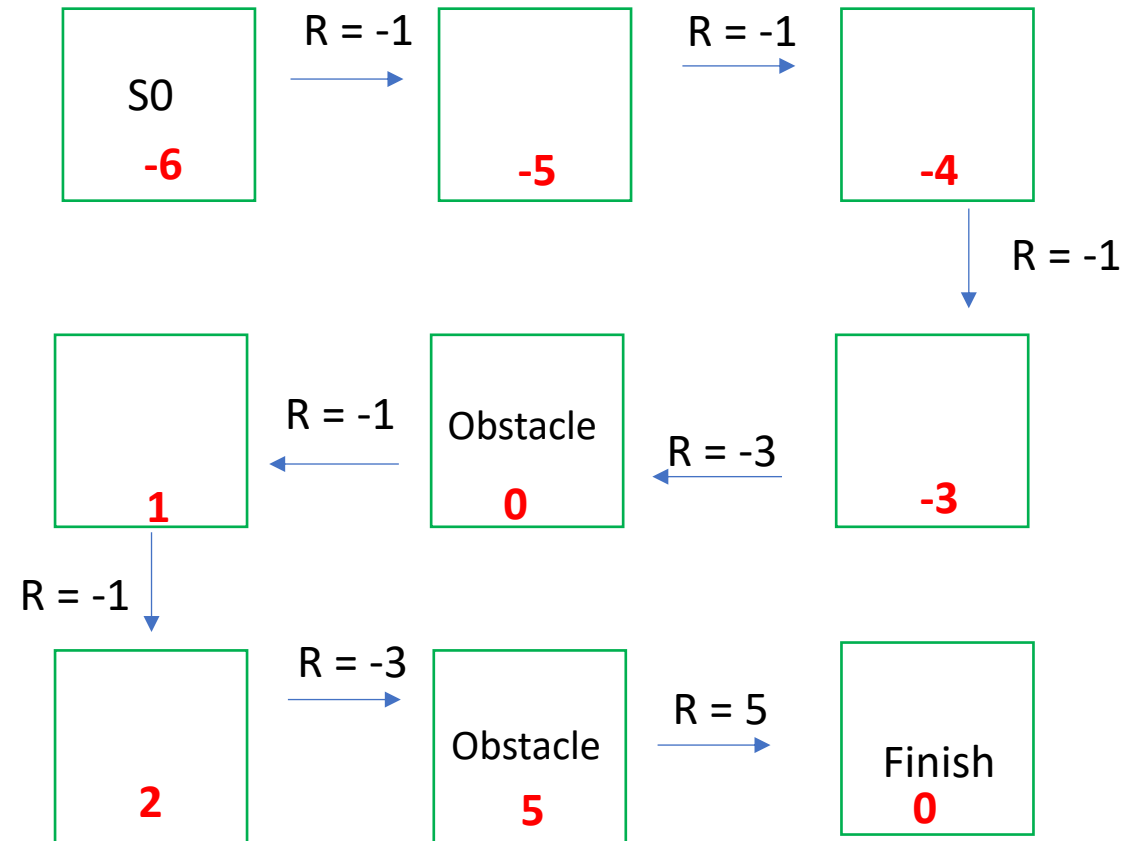
Let's consider:

- The agent can only move up, down, left or right, and can only take actions that lead in to not fall off the grid.
- The goal of the agent to get to the bottom right hand corner of the world as quickly as possible.
- The agent receives a rewards of negative one for most transitions. If it leads to a mountain rewards of 3 and the finish state rewards of 5.



MDP and Agent Policy

Let's choose a policy that the agent visits every state in a roundabout manner.



Cumulative Return for S0 $\rightarrow (-1) + (-1) + (-1) + (-3) + (-1) + (-1) + (-3) + (5) = -6$

State-Value Function

For each state, the **state-value function** yields the **expected return**, if the agent started in that state, and then followed the policy for all times steps.

-6	-5	-4
1	0	-3
2	5	0

We call v_{π} the state-value function for policy π

The value of state s under a policy π is:

$$v_{\pi}(s) = E_{\pi}[G_t | S_t = s]$$

For each state s

It yields the **expected return** If the agent starts in **state s** and then uses **the policy** to choose its actions for **all time steps**.

Bellman Equations

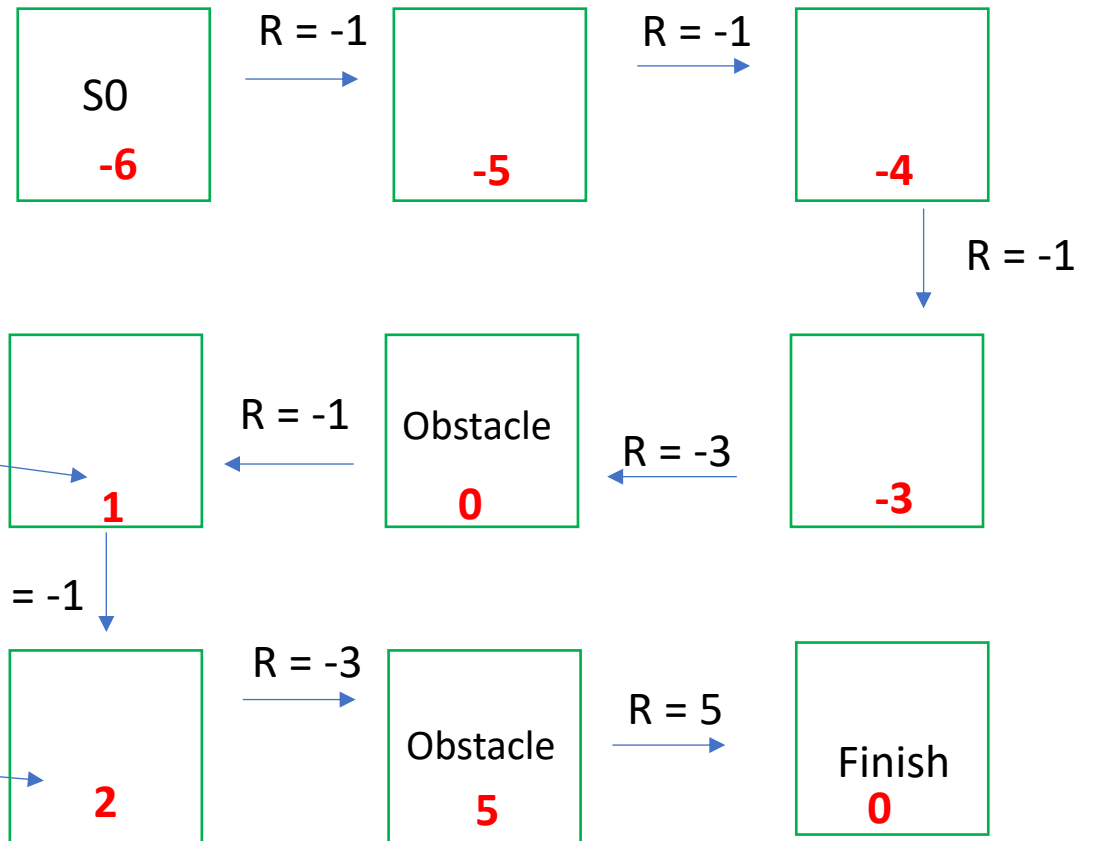
The value of any state (1)
=

the immediate reward (-1)

+

The value of of the state that follows (2)

Discount rate is 1



Bellman Expected Equation

$$v_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s]$$

The value of any state s = the expected value of the immediate rewards
+

The discounted value of the state that follows under that policy.

Question: State-Value Functions for π'

You will calculate the value function corresponding to a particular policy

Deterministic policy

$\pi(s_1)=\text{right}$

$\pi(s_2)=\text{right}$

$\pi(s_3)=\text{down}$

$\pi(s_4)=\text{up}$

$\pi(s_5)=\text{right}$

$\pi(s_6)=\text{down}$

$\pi(s_7)=\text{right}$

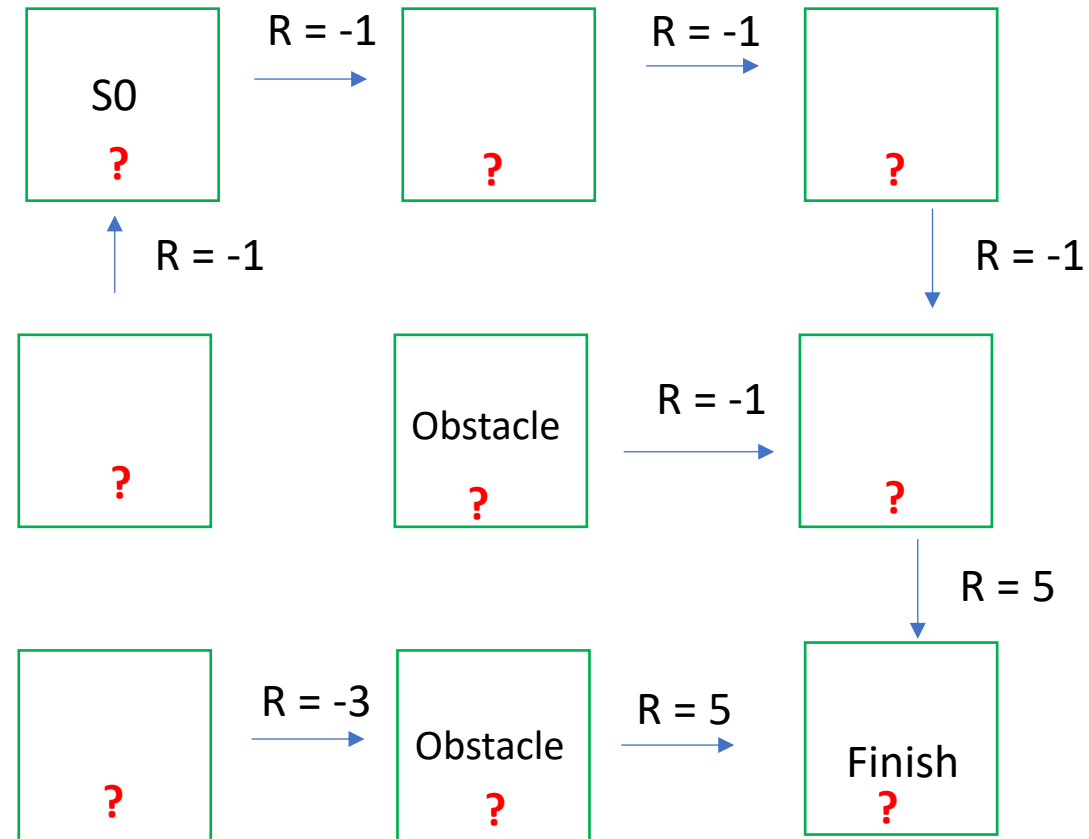
$\pi(s_8)=\text{right}$

Assume $\gamma=1$

Questions:

What is $v_{\pi}(s_4)$?

What is $v_{\pi}(s_1)$?



State value functions for different policies

State value function for policy π

-6	-5	-4
1	0	-3
2	5	0

State value function for policy π'

2	3	4
1	4	5
2	5	0

$\pi' \geq \pi$ if and only if $v_{\pi'}(s) \geq v_{\pi}(s)$ for all $s \in S$

Optimal policy

Definition

$\pi' \geq \pi$ if and only if $v_{\pi'}(s) \geq v_{\pi}(s)$ for all $s \in S$

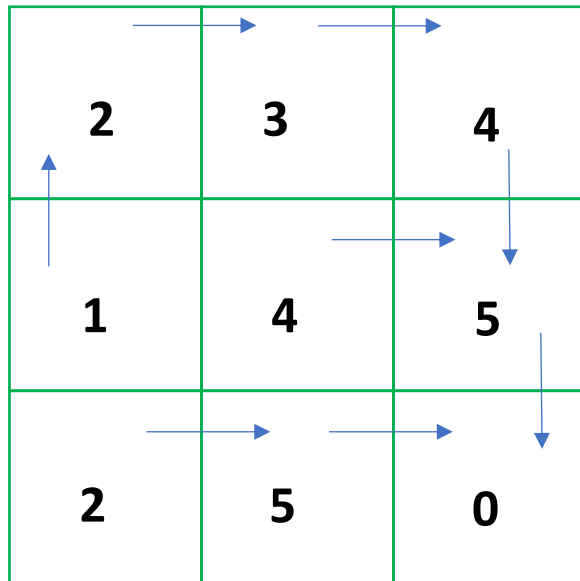
Definition

An optimal policy **π^*** satisfies **$\pi^* \geq \pi$** for all **π**

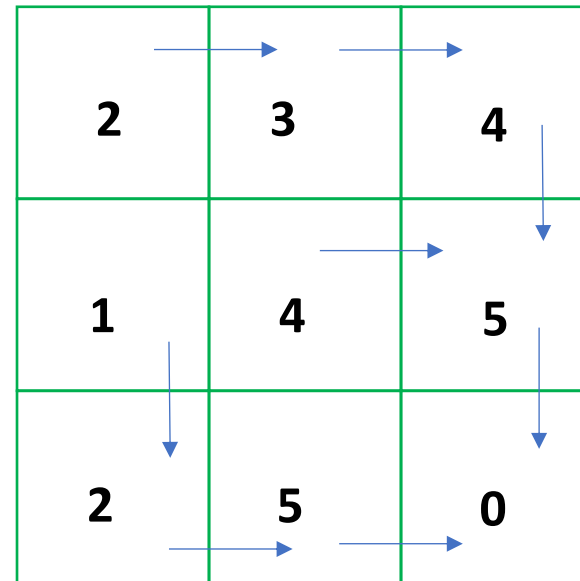
The optimal state value function is denoted v^*

Both are Optimal Policies

State value function for policy π'



State value function for policy π''



Action Value Function $q_{\pi}(s, a)$

We call v_{π} the state-value function for policy π

We call q_{π} the action-value function for policy π

The value of state s under a policy π is:

$$v_{\pi}(s) = E_{\pi}[G_t | S_t = s]$$

The value of taking an action under a policy π is:

$$q_{\pi}(s, \alpha) = E_{\pi}[G_t | S_t = s, A_t = \alpha]$$

For each state s

It yields the **expected return** If the agent starts in **state s** and then uses **the policy** to choose its actions for all future time steps.

For each **state s and action a**

It yields the **expected return** If the agent starts in **state s** and then chooses **action a** and then uses **the policy** to choose its actions for all future time steps.

Action Value for Policy π

Action value function for policy π'

<div>02</div>	<div>13</div>	<div>24</div>
<div>11</div>	<div>024</div>	<div>135</div>
<div>02</div>	<div>115</div>	<div>1</div>

Optimal Policy

- The agent interacts with the environment. From that interactions, it estimates the optimal action value function. Then the agent uses that value function to get the optimal policy.

Interaction $\rightarrow q_* \rightarrow \pi_*$

Action value function for policy π'

<div>02</div>	<div>13</div>	<div>24</div>
<div>11</div>	<div>02</div>	<div>13</div>
<div>02</div>	<div>15</div>	<div>1</div>

MDP Summary

- So far, we have learned how to specify a solution to the reinforcement learning problem. In particular, the **optimal policy** π^* specifies - for each environment state - how the agent should select an action towards its goal of maximizing expected rewards.
- We learned that the agent could structure its search for an optimal policy by first estimating the **optimal action-value function** q^* ; then, once q^* is known, π^* is quickly obtained.

Bellman Equation

$$p(s', r | s, a)$$

$$E_{\pi}[R(t+1) | S_t = s] = \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) r$$

Then for Bellman Equation:

$$v_{\pi}(s) = E_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s]$$

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) \{r + \gamma v_{\pi}(s')\}$$

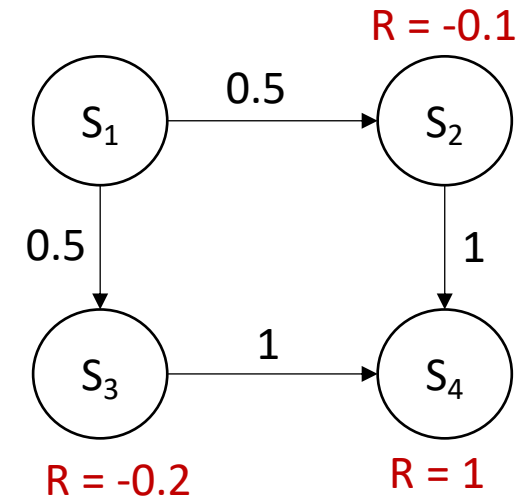
Bellman Equation by Example

- Start at S_1
- Discount factor $\gamma = 0.9$
- Terminal State $S_4 \rightarrow V(s_4) = 0$
- Agent Policy $\pi(a|s)$
- My reward and next state $p(s', r | s, a)$

$$V(s_2) = 1$$

$$V(s_3) = 1$$

$$V(s_1) = 0.5[-0.1 + 0.9] + 0.5[-0.2 + 0.9] = 0.75$$



Bellman Equation by Example

Discount factor $\gamma = 0.9$

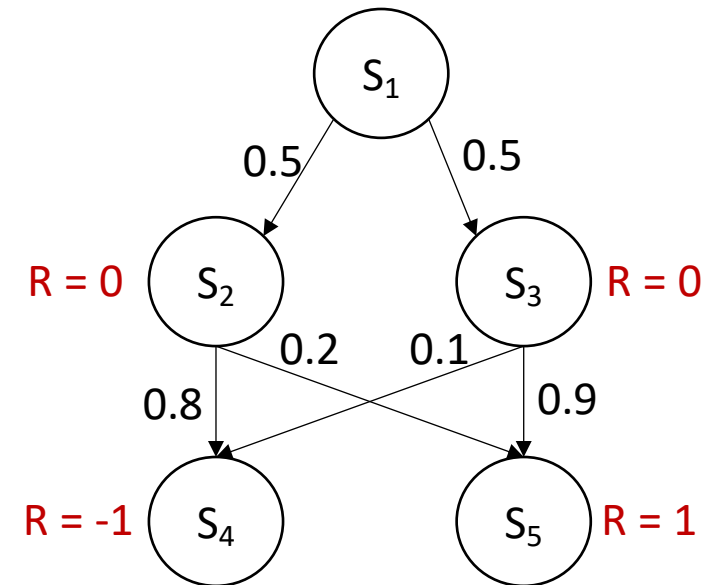
$V(S_4)$ and $V(S_5) = 0$ Terminal State

$$V(S_3) = 0.9 * 1 + 0.1 * (-1) = 0.8$$

$$V(S_2) = 0.2 * 1 + 0.8 * (-1) = -0.6$$

$$\begin{aligned} V(S_1) &= 0.5 * (0 + \gamma V(s_2)) + 0.5 * (0 + \gamma V(s_3)) \\ &= 0.5 * (-0.6 * 0.9) + 0.5 * (0.8 * 0.9) \end{aligned}$$

$$V(S_1) = 0.09$$



Bellman Equation by Example

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

Discount factor $\gamma = 0.9$

$$V(s_1) = 0.3 (-0.1 + \gamma V(s_1)) + 0.7 (-0.1 + \gamma V(s_2))$$

$$V(s_1) = 0.3 (-0.1 + 0.9 V(s_1)) + 0.7 (-0.1 + 0.9 V(s_2))$$

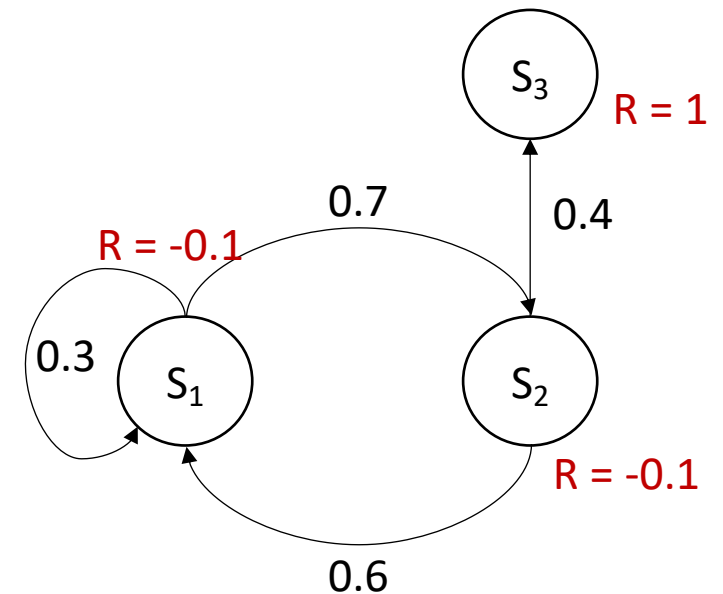
$$V(s_1) = -0.1 + 0.27 V(s_1) + 0.63 V(s_2)$$

$$V(s_2) = 0.4 (1 + \gamma V(s_3)) + 0.6 (-0.1 + \gamma V(s_1))$$

$$V(s_2) = 0.4 (1 + 0.9 V(s_3)) + 0.6 (-0.1 + 0.9 V(s_1))$$

$$V(s_2) = 0.34 + 0.54 V(s_1) + 0.36 V(s_3)$$

$V(s_3) = 0$ Terminal state



$$V(s_1) = 0.293$$

$$V(s_2) = 0.498$$

$$V(s_3) = 0$$

Relationship between V and Q

$$V_*(s) = \text{Max}_{\alpha} \{Q_*(s, a)\}$$

Implementing the Optimal Policy

- Key point: value function already takes future rewards into account
- Just greedily choose the action that yields the best next-state value $V(s')$ requires look-ahead search
- If we have $Q(s, a)$, no need to look ahead, simply choose argmax

The Bellman Equation

$$v_{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) \{r + \gamma v_{\pi}(s')\}$$

Using iterative policy evaluation

Monte Carlo Method

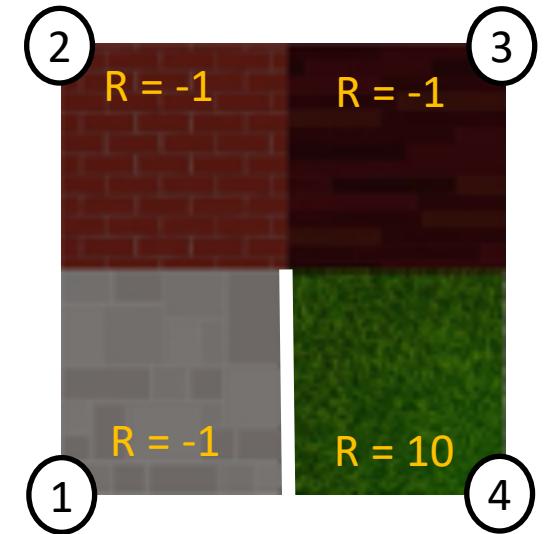
Monte Carlo methods, or **Monte Carlo** experiments, are a broad class of computational algorithms that rely on repeated random sampling to obtain numerical results.

The underlying concept is to use randomness to solve problems that might be deterministic in principle.

Monte Carlo methods are mainly used in three problem classes: optimization, numerical integration, and generating draws from a probability distribution.

Monte Carlo Methods

- The agent is looking for the optimal policy π_*
 - for each state which action is the best?
- Assume environment's dynamics are stochastic (assume 70% moves in the direction takes action)
- Since environment's dynamics are stochastic, to truly understand the environment, the agent needs more episodes.
- Agent collects episodes, randomly conducts actions in each state, and consolidates the information in a table, and then uses the table to come up with a better policy.

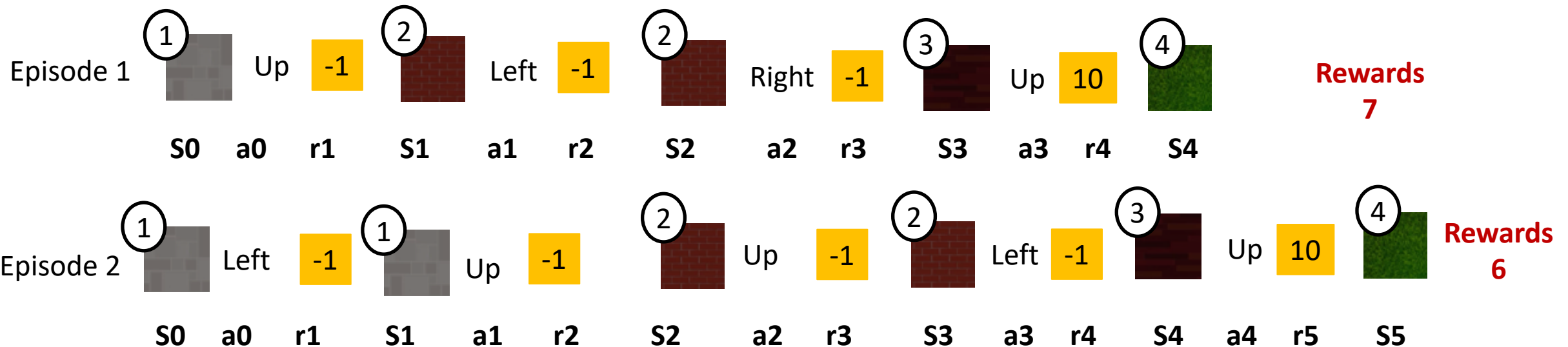
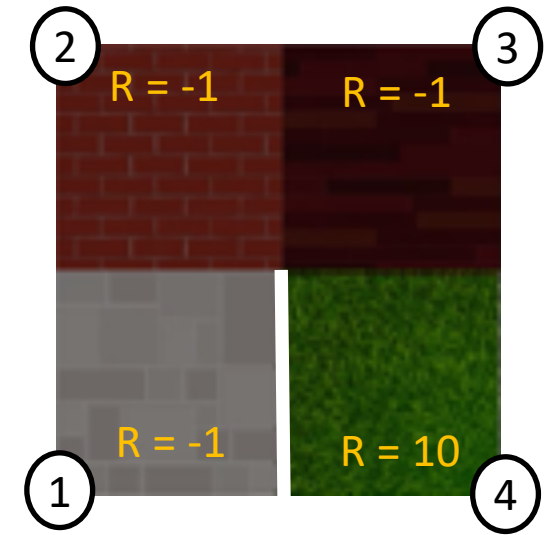


Actions

Up , Down, Left, Right

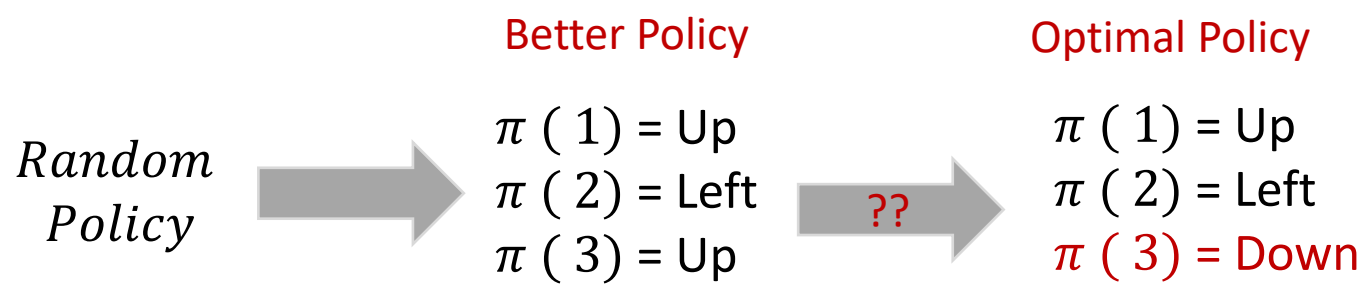
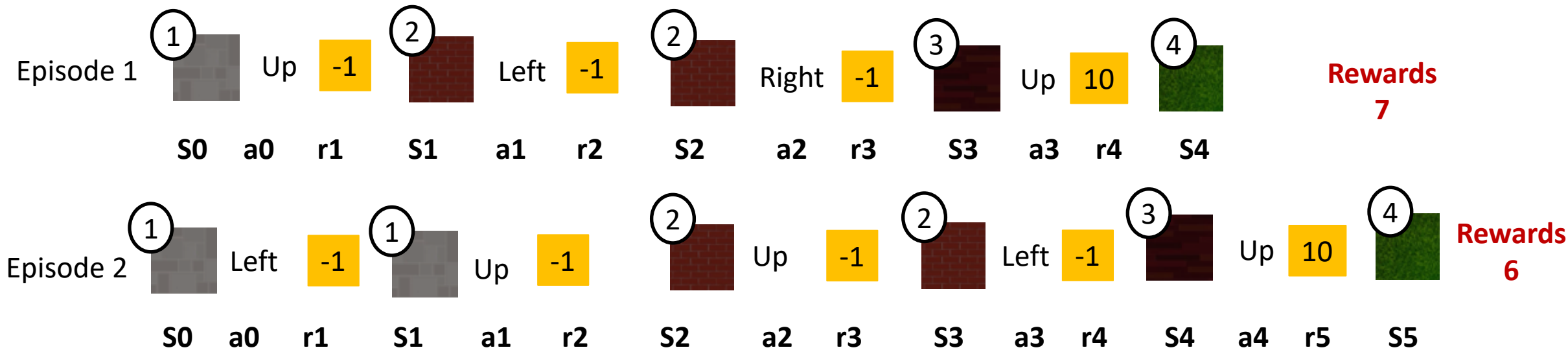
Discount Rate = 1

Monte Carlo Methods



The Agent is looking for the optimal policy? How can the agent use it to improve its strategy?

Monte Carlo Prediction



	Up	Down	Left	Right
1	+7	+6	+5	+6
2	+8	+7	+8	+9
3	+10	+8	+9	+9

Q Table

Action value function

MC Prediction Algorithm

- So we have discussed how the agent can take a **policy**, like the random policy, use it to collect some episodes, and then consolidate the results to arrive at a **better policy** by estimating the action-value function with a Q-table.
- Note that each entry in the Q-table corresponds to a particular state and action. To populate an entry, we use the return that followed when the agent was in that state, and chose the action. In the event that the agent has selected the same action many times from the same state, we need only average the returns.

Option 1: Every-visit MC Prediction

Average the returns following all visits to each state-action pair, in all episodes.

Option 2: First-visit MC Prediction

For each episode, we only consider the first visit to the state-action pair. The pseudocode for this option can be found below.

First Visit Monte Carlo Prediction

- **Q-table**, with a row for each state and a column for each action. The entry corresponding to state s and action a is denoted $Q(s,a)$.
- **N** - table that keeps track of the number of first visits we have made to each state-action pair.
- **returns_sum** - table that keeps track of the sum of the rewards obtained after first visits to each state-action pair.

Algorithm 9: First-Visit MC Prediction (*for action values*)

Input: policy π , positive integer $num_episodes$

Output: value function Q ($\approx q_\pi$ if $num_episodes$ is large enough)

Initialize $N(s,a) = 0$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Initialize $returns_sum(s,a) = 0$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

for $i \leftarrow 1$ **to** $num_episodes$ **do**

 Generate an episode $S_0, A_0, R_1, \dots, S_T$ using π

for $t \leftarrow 0$ **to** $T - 1$ **do**

if (S_t, A_t) is a first visit (with return G_t) **then**

$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$

$returns_sum(S_t, A_t) \leftarrow returns_sum(S_t, A_t) + G_t$

end

end

$Q(s,a) \leftarrow returns_sum(s,a)/N(s,a)$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

return Q

MC Control