University at Buffalo
Department of Computer Science and Engineering
School of Engineering and Applied Sciences
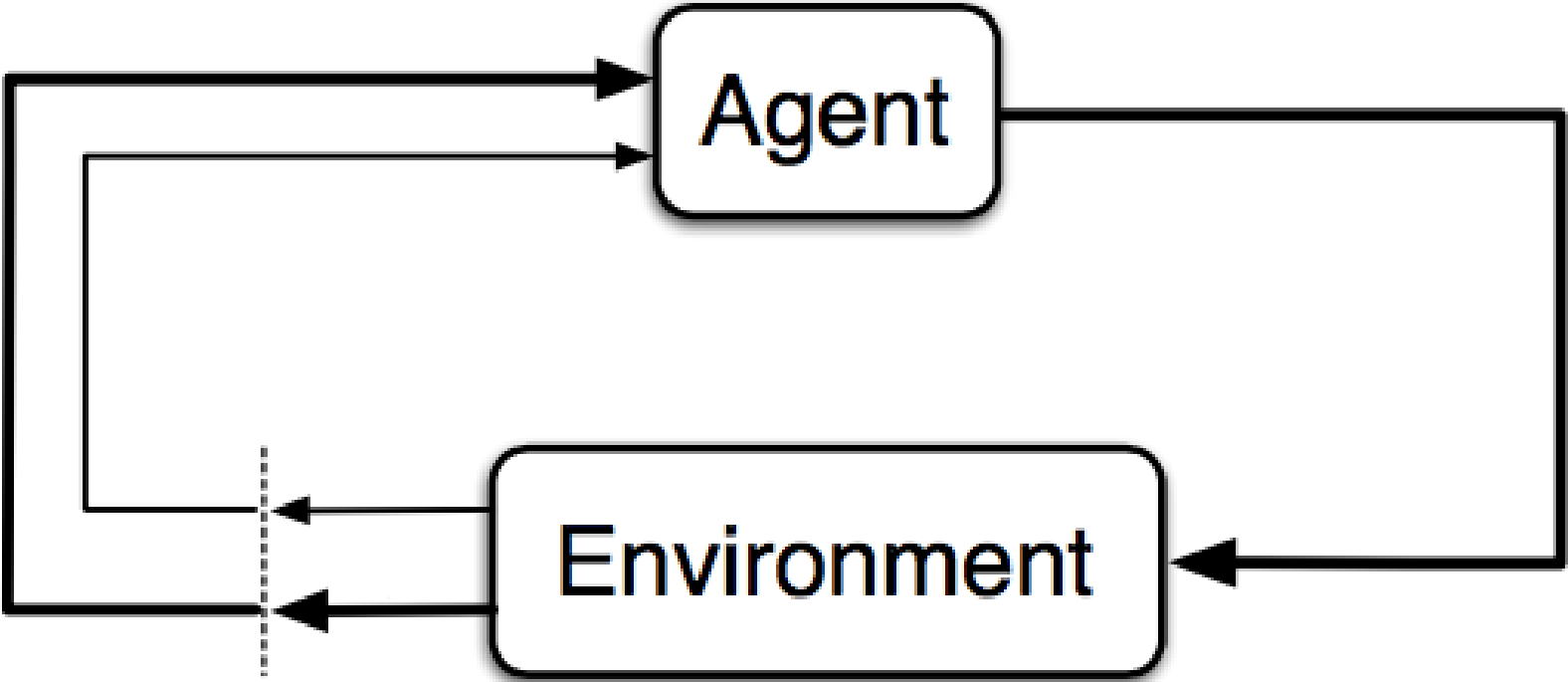
# TABULAR METHODS OVERVIEW

Lecture 7.1
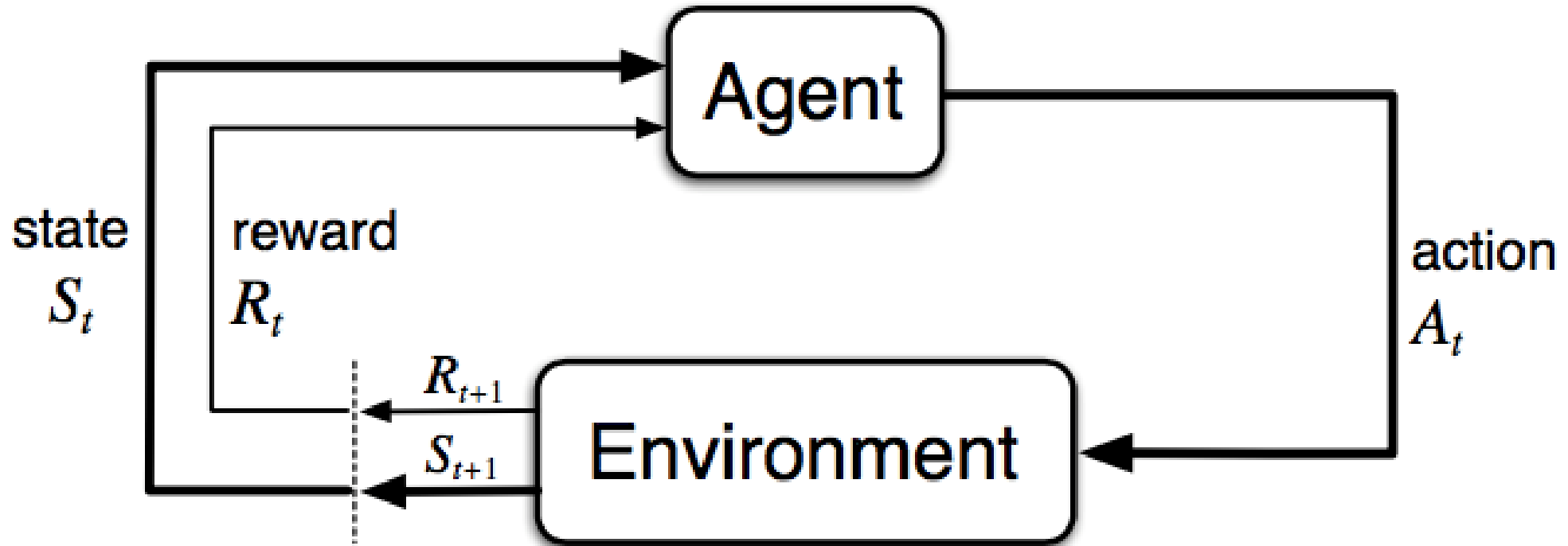
CSE4/510: Reinforcement Learning
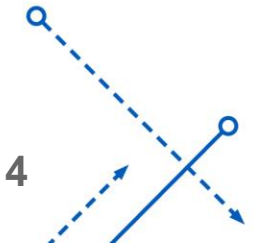
September 17, 2019

# MDP

# MDP



state $S_t$

reward $R_t$

$R_{t+1}$

$S_{t+1}$

Agent

Environment

action $A_t$

**TRUE / FALSE?**

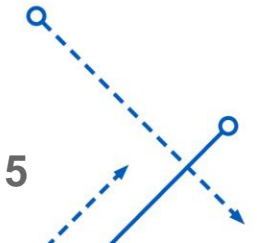Markov Decision Process is defined as:

(s, a, O, P, \gamma)

**FALSE**

Markov Decision Process is defined as:
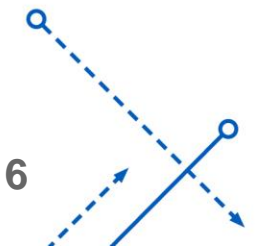
(s, a, O, P, r, \gamma)

# Policy

1. Deterministic

2. Stochastic

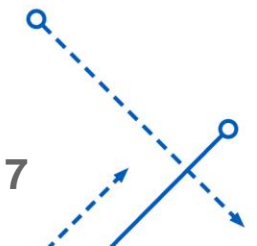A. $\pi(a|s) = \mathbb{P}_\pi[A = a | S = s]$
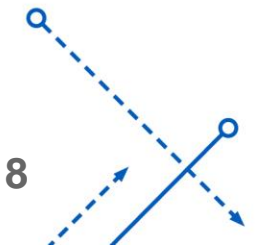
B. $\pi(s) = a$

**ENVIRONMENT**
Deterministic / Stochastic?
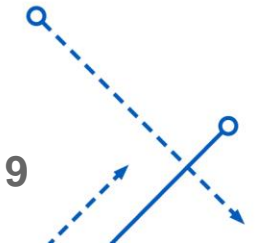

**POLICY**
Deterministic / Stochastic?

# RL agents goal?

# Types of value functions?

# Value Functions

**Types of value functions:**

*State value function* describes the value of a state when following a policy. It is the expected return when starting from state $s$ acting according to our policy $\pi$:

$$V^{\pi}(s) = \mathbb{E}_{\pi}[R_t | S_t = s]$$

*Action value function* tells us the value of taking an action $a$ in state $s$ when following a certain policy $\pi$. It is the expected return given the state and action under $\pi$:
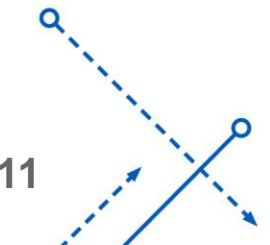
$$Q^{\pi}(s, a) = \mathbb{E}_{\pi}[R_t | S_t = s, A_t = a]$$

# Value Functions

V(s) can also be interpreted, as the cumulative future reward

**Are we missing something?**

V(s) can also be interpreted, as the <span style="color:red">expected</span> cumulative future <span style="color:red">discounted</span> reward

**1. Evaluate**

**A.** $V_\pi(s) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \ldots | S_t = s]$

**2. Improve**

**B.** $\pi' = greedy(V_\pi)$

**Evaluate** $\qquad V_\pi(s) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \ldots | S_t = s]$

**Improve** $\qquad \pi' = greedy(V_\pi)$

# Dynamic Programming
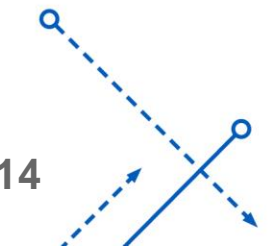
Given a policy $\pi$

- **Evaluate** the policy $\pi$

$$V_\pi(s) = \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \dots | S_t = s]$$

- **Improve** the policy by acting greedily with respect to $v_\pi$

$$\pi' = \text{greedy}(V_\pi)$$

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} V_{\pi_2} \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} V_*$$

starting $V$ $\pi$

$V = V\pi$

$\pi = greedy(V)$

$V^*$
$\pi^*$

# Dynamic Programming

**1. Distribution model**

**A. Produce a single outcome taken according to its probability of occurring**

**2. Sample model**

**B. List all possible outcomes and their probabilities**

**Model-free RL**

Model

Model learning

Simulation

Planning

Interaction with Environment

Experience

Value function

Direct RL methods

Greedification

Policy

**Planning (or model-based RL)**

Model learning

Simulation

Planning

Interaction with Environment

Experience

Direct RL methods

Value function

Greedification

Policy

# MC / DP / TD ?

Bootstrapping

MC
DP
TD

Bootstrapping

MC
✅ DP
✅ TD

Sampling

MC
DP
TD

Sampling

✅ MC

DP

✅ TD

**Dynamic Programming**  A. $V(S_t) \leftarrow V(S_t) + \alpha \left(R_{t+1} + \gamma V(S_{t+1}) - V(S_t)\right)$

**Monte Carlo**  B. $V(S_t) \leftarrow E_\pi \left[R_{t+1} + \gamma V(S_{t+1})\right] = \sum_a \pi(a|S_t) \sum_{s',r} p(s', r|S_t, a)[r + \gamma V(s')]$

**Temporal Difference**  C. $V(S_t) \leftarrow V(S_t) + \alpha \left(G_t - V(S_t)\right)$

# Value Based RL

**Dynamic Programming**
$$V(S_t) \leftarrow E_\pi \left[ R_{t+1} + \gamma V(S_{t+1}) \right] = \sum_a \pi(a|S_t) \sum_{s',r} p(s', r|S_t, a)[r + \gamma V(s')]$$

**Monte Carlo**
$$V(S_t) \leftarrow V(S_t) + \alpha \left( {\color{red}G_t} - V(S_t) \right)$$

**Temporal Difference**
$$V(S_t) \leftarrow V(S_t) + \alpha \left( {\color{red}R_{t+1} + \gamma V(S_{t+1})} - V(S_t) \right)$$

# Dynamic Programming

**Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$**

1. Initialization
   $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation
   Loop:
   $\quad\quad \Delta \leftarrow 0$
   $\quad\quad$ Loop for each $s \in \mathcal{S}$:
   $\quad\quad\quad v \leftarrow V(s)$
   $\quad\quad\quad V(s) \leftarrow \sum_{s',r} p(s',r\,|\,s,\pi(s))\big[r + \gamma V(s')\big]$
   $\quad\quad\quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
   $\quad$ until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement
   $policy\text{-}stable \leftarrow true$
   For each $s \in \mathcal{S}$:
   $\quad\quad old\text{-}action \leftarrow \pi(s)$
   $\quad\quad \pi(s) \leftarrow \arg\max_a \sum_{s',r} p(s',r\,|\,s,a)\big[r + \gamma V(s')\big]$
   $\quad\quad$ If $old\text{-}action \neq \pi(s)$, then $policy\text{-}stable \leftarrow false$
   If $policy\text{-}stable$, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

## TD / Monte Carlo ?