

# Lecture 5: Mathematical Models of RL

Ziyu Shao

School of Information Science and Technology  
ShanghaiTech University

April 09, 2025

# Outline

- 1 Basic Setting
- 2 Revisit Markov Chain
- 3 Markov Reward Process
- 4 Markov Decision Process
- 5 Summary
- 6 References

# Outline

1 Basic Setting

2 Revisit Markov Chain

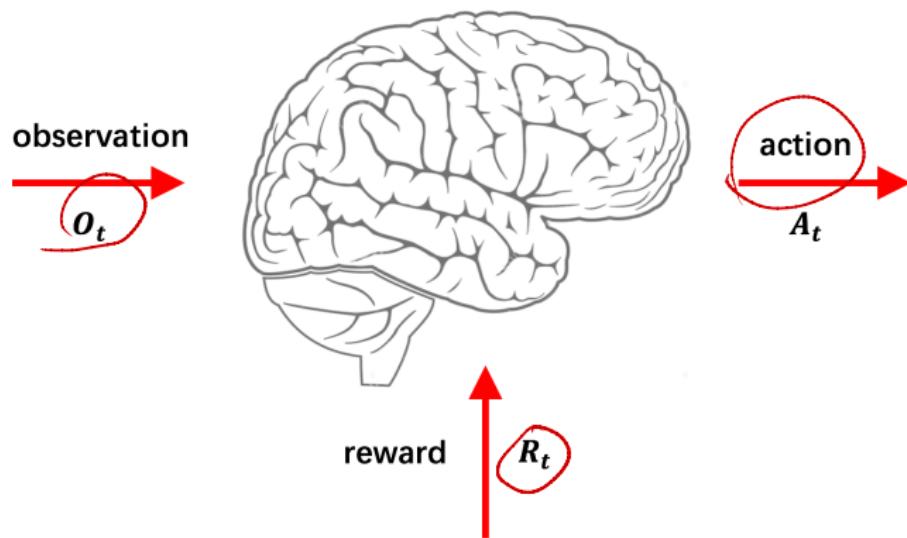
3 Markov Reward Process

4 Markov Decision Process

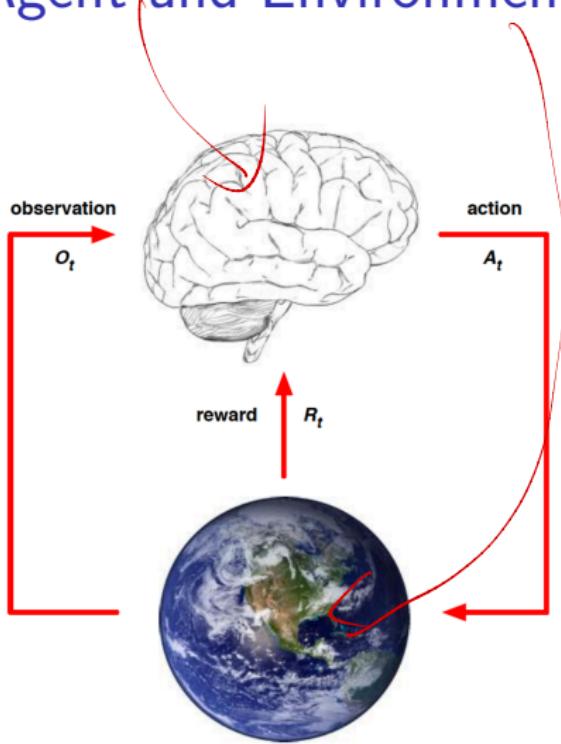
5 Summary

6 References

# Agent and Environment



# Agent and Environment

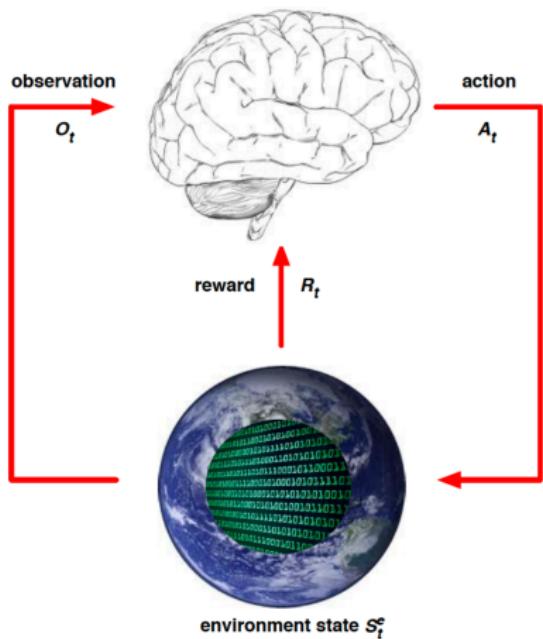


- At each step  $t$  the agent:
  - ▶ Executes action  $A_t$
  - ▶ Receives observation  $O_t$
  - ▶ Receives scalar reward  $R_t$
- The environment:
  - ▶ Receives action  $A_t$
  - ▶ Emits observation  $O_{t+1}$
  - ▶ Emits scalar reward  $R_{t+1}$
- $t$  increments at env. step

# History and State

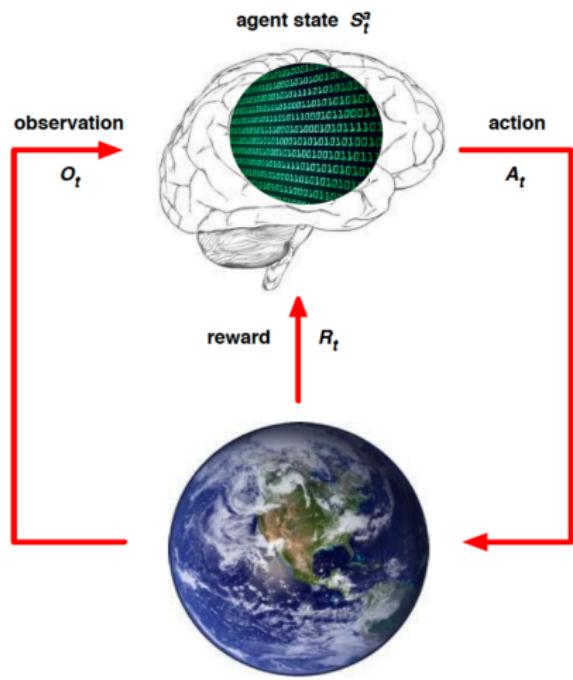
- The **history** is the sequence of observations, actions, rewards  
 $H_t = \underline{O_1}, \underline{R_1}, \underline{A_1}, \dots, \underline{A_{t-1}}, \underline{O_t}, \underline{R_t}$
- i.e. all observable variables up to time  $t$
- i.e. the sensorimotor stream of a robot or embodied agent
- What happens next depends on the history:
  - ▶ The agent selects actions
  - ▶ The environment selects observations/rewards
- **State** is the information used to determine what happens next
- Formally, state is a function of the history:  $S_t = f(H_t)$

# Environment State



- The **environment state  $S_t^e$**  is the environment's private representation
- i.e. whatever data the environment uses to pick the next observation/reward
- The environment state is not usually visible to the agent
- Even if  $S_t^e$  is visible, it may contain irrelevant information

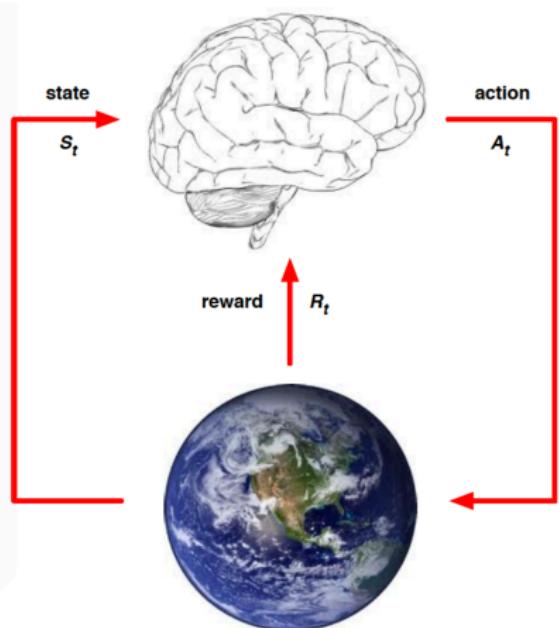
# Agent State



- The **agent state**  $S_t^a$  is the agent's internal representation
- i.e. whatever information the agent uses to pick the next action
- i.e. it is the information used by reinforcement learning algorithms
- It can be any function of history:

$$S_t^a = f(H_t)$$

# Fully Observable Environments



Full observability: agent directly observes environment state

$$O_t = S_t^a = S_t^e$$

- Agent state = environment state
- Formally, this is a **Markov decision process (MDP)**

# Partially Observable Environments

- Partial Observability
  - ▶ A robot with camera vision isn't told its absolute location
  - ▶ A trading agent only observes current prices
  - ▶ A poker playing agent only observes public cards
- Now agent state  $\neq$  environment state
- Formally, this is a partially Observable Markov decision process (POMDP)

# Major Components of An RL Agent

- ① Policy-based RL
- ② Value-function based RL
- ③ Model

An RL agent may include one or more of these components:

- Policy: agent's behavior function
- Value function: how good is each state and/or action
- Model: agent's representation of the environment

# Policy

- A **policy** is the agent's behavior
- It is a map from state to action, e.g.
- Deterministic policy:  $a = \pi(s)$
- Stochastic policy:  $\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$

# Value Function

- Value function is a prediction of future reward
- Used to evaluate the goodness/badness of states
- And therefore to select between actions, e.g.

$$v_{\pi}(s) = \mathbb{E}_{\pi}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$$

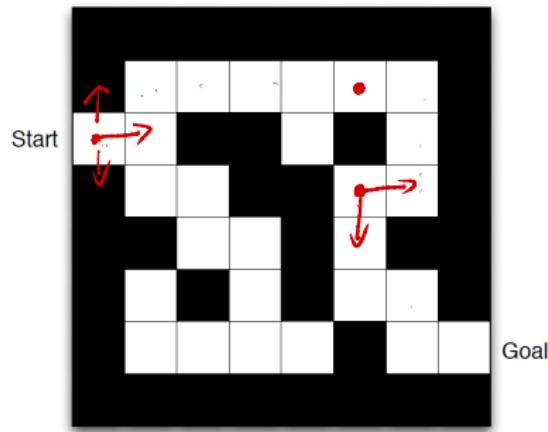
# Model

- A **model** predicts what the environment will do next
- $\mathcal{P}$  predicts the next state
- $\mathcal{R}$  predicts the next (immediate) reward, e.g.

$$\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$$

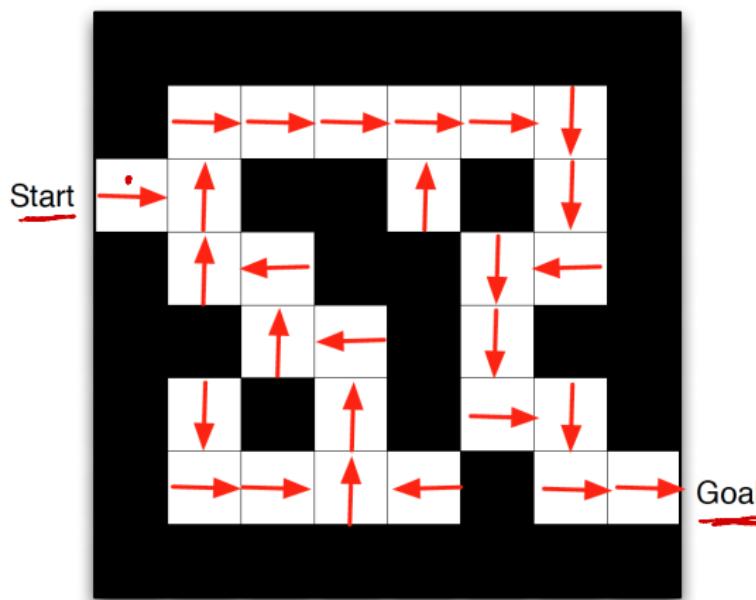
$$\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$$

# Maze Example



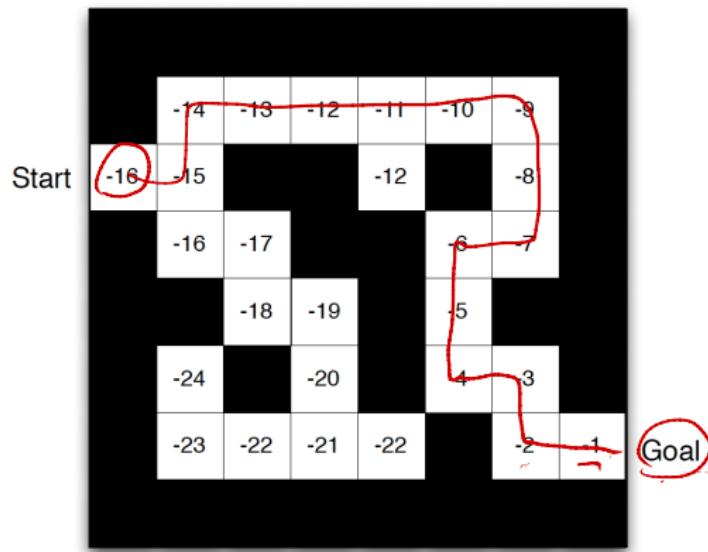
- Rewards: -1 per time-step
- Actions: N, E, S, W
- States: Agent's location

# Maze Example: Policy



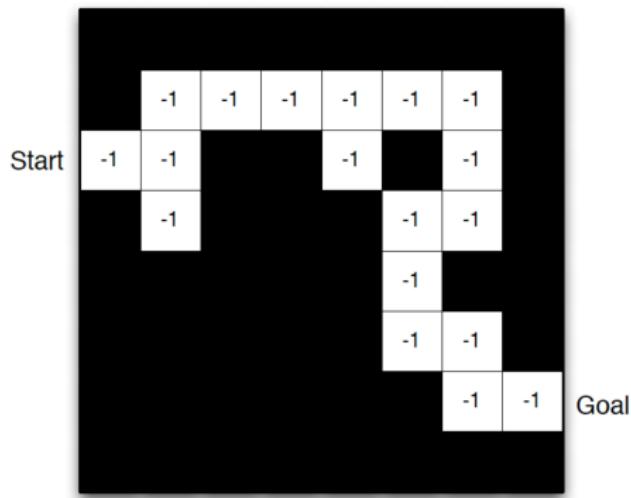
- Arrows represent policy  $\pi(s)$  for each state  $s$

# Maze Example: Value Function



- Numbers represent value  $v_\pi(s)$  of each state  $s$

# Maze Example: Model



- Agent may have an internal model of the environment
- Dynamics: how actions change the state
- Rewards: how much reward from each state
- The model may be imperfect

- Grid layout represents transition model  $\mathcal{P}_{ss'}^a$
- Numbers represents immediate reward  $R_s^a$  from each state  $s$  (same for all  $a$ )

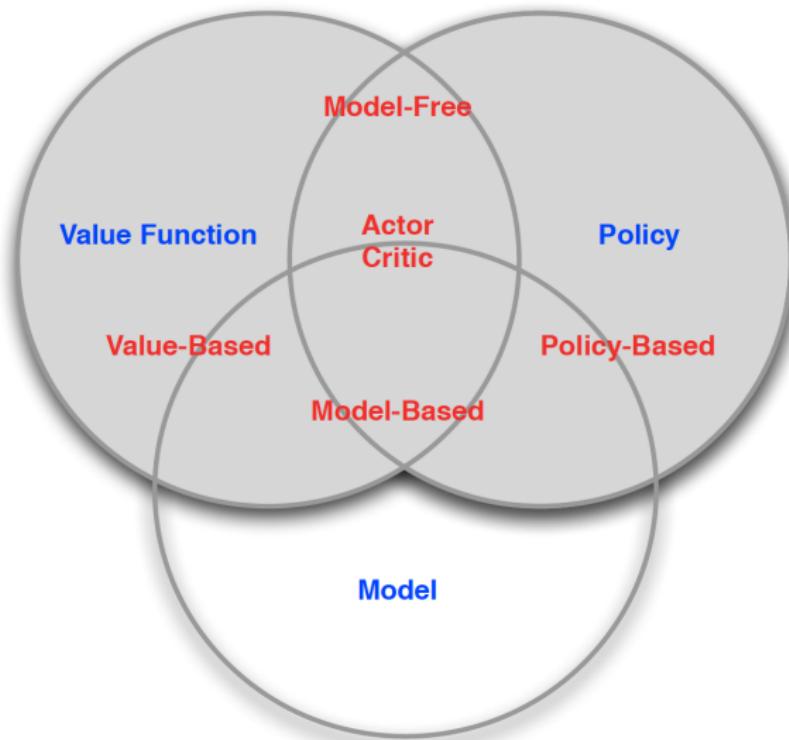
# Categorizing RL Agents: I

- Value based
  - ▶ No Policy (Implicit) DQN
  - ▶ Value Function
- Policy Based
  - ▶ Policy PPO
  - ▶ No Value Function
- Actor Critic
  - ▶ Policy
  - ▶ Value Function

# Categorizing RL Agents: II

- Model Free
  - ▶ Policy and/or Value Function
  - ▶ No Model
- Model Based
  - ▶ Policy and/or Value Function
  - ▶ Model

# RL Agent Taxonomy



# Learning and Planning

Two fundamental problems in sequential decision making

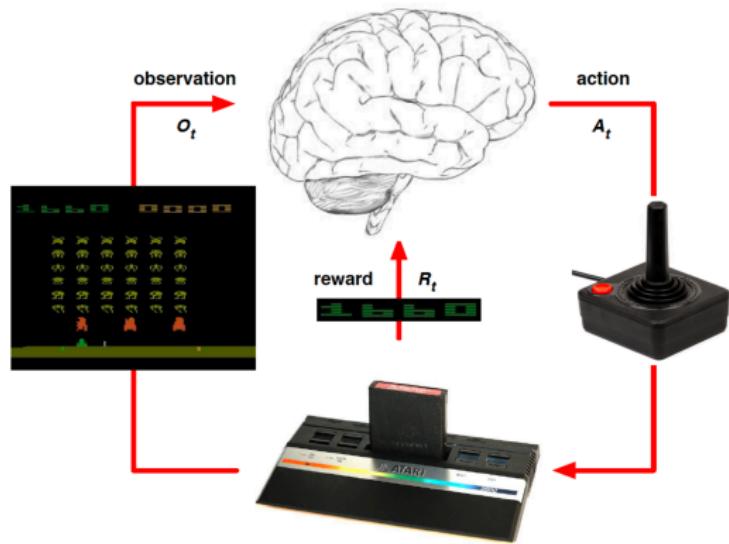
- Reinforcement Learning

- ▶ The environment is initially unknown
- ▶ The agent interacts with the environment
- ▶ The agent improves its policy

- Planning

- ▶ A model of the environment is known
- ▶ The agent performs computations with its model (without any external interaction)
- ▶ The agent improves its policy
- ▶ a.k.a. deliberation, reasoning, introspection, pondering, thought, search

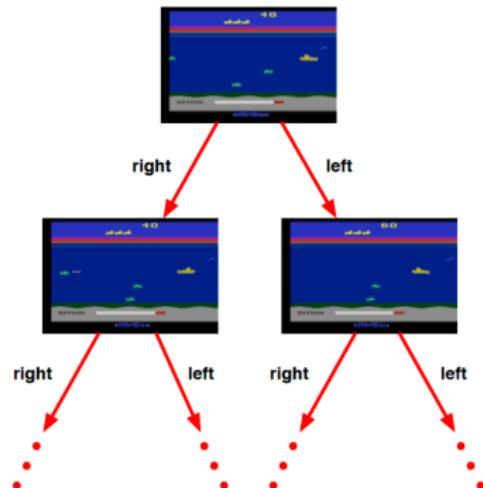
# Atari Example: Reinforcement Learning



- Rules of the game are unknown
- Learn directly from interactive game-play
- Pick actions on joystick, see pixels and scores

# Atari Example: Planning

- Rules of the game are known
- Can query emulator
  - ▶ perfect model inside agent's brain
- If I take action  $a$  from state  $s$ :
  - ▶ what would the next state be?
  - ▶ what would the score be?
- Plan ahead to find optimal policy
  - ▶ e.g. tree search



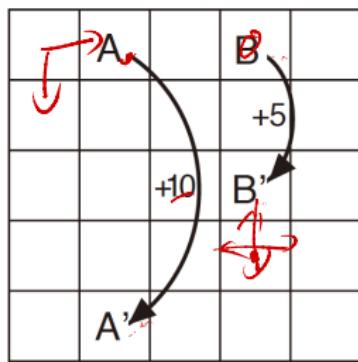
# Prediction and Control

*Estimation*

- **Prediction**: evaluate the future
  - ▶ Given a policy → Value function *estimate*
- **Control**: optimize the future
  - ▶ Find the best policy

# Gridworld Example: Prediction

*estimation*



(a)



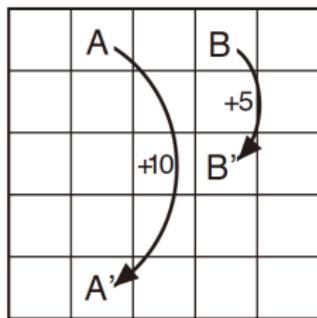
Actions

3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

(b)

What is the value function for the uniform random policy?

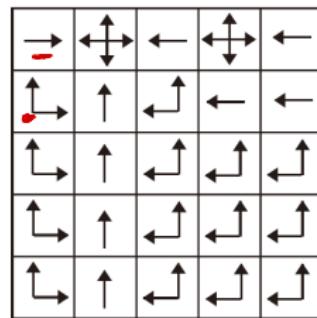
# Gridworld Example: Control



a) gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

b)  $v_*$



c)  $\pi_*$

What is the optimal value function over all possible policies?  
What is the optimal policy?

# Connections with Psychology

- Prediction & Control algorithms in reinforcement learning parallels Classical & Instrumental conditioning in animal learning.
- Environment Models in reinforcement learning parallels Cognitive Maps in animal learning.
  - ▶ they can be learned by supervised learning methods without relying on reward signals
  - ▶ they can be used later to plan behavior
- Model-free & Model-based algorithms in reinforcement learning parallels Habitual & Goal-directed behavior in pyschology.

# Outline

- 1 Basic Setting
- 2 Revisit Markov Chain
- 3 Markov Reward Process
- 4 Markov Decision Process
- 5 Summary
- 6 References

# Markov Property

## Definition

A state  $S_t$  is Markovian if and only if

$$\mathbb{P}[S_{t+1}|S_t] = \mathbb{P}[S_{t+1}|S_1, \dots, S_t]$$

- “The future is independent of the past given the present”
- The state captures all relevant information from the history
- Once the state is known, the history may be thrown away
- i.e. The state is a sufficient statistic of the future

# Markov Chain

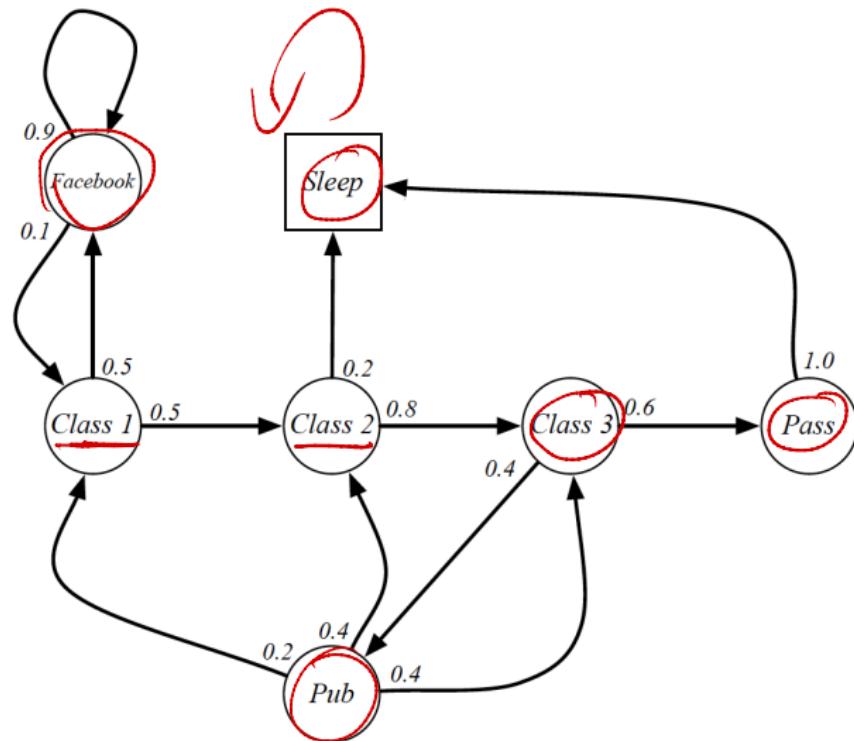
## Definition

A discrete-time Markov chain is a tuple  $\langle \mathcal{S}, \mathcal{P} \rangle$

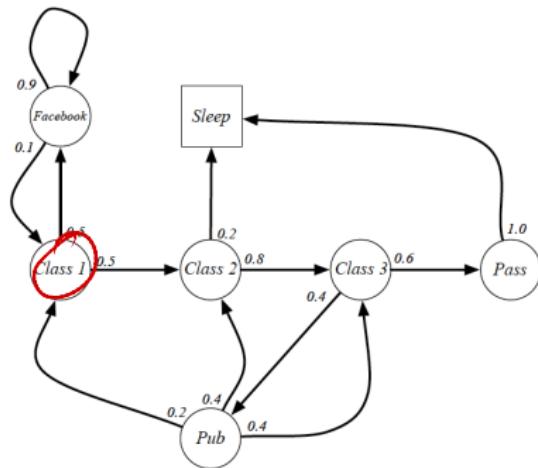
- $\mathcal{S}$  is a (finite) set of states
- $\mathcal{P}$  is a state transition probability matrix

$$\mathcal{P}_{s,s'} = \mathbb{P}[S_{t+1} = s' | S_t = s]$$

# Example: Student Markov Chain



# Example: Student Markov Chain Episodes

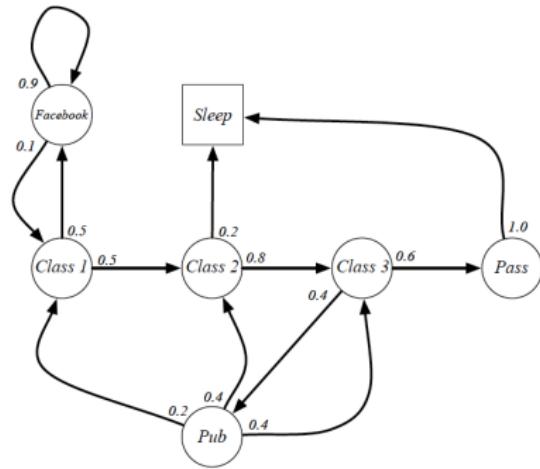


Sample episodes for Student Markov Chain starting from  $S_1 = C1$

$S_1, S_2, \dots, S_T$

- C1 C2 C3 Pass Sleep
- C1 FB FB C1 C2 Sleep
- C1 C2 C3 Pub C2 C3 Pass Sleep
- C1 FB FB C1 C2 C3 Pub C1 FB FB FB C1 C2 C3 Pub C2 Sleep

# Example: Student Markov Chain Transition Matrix



$$\mathcal{P} = \begin{bmatrix} C1 & C2 & C3 & Pass & Pub & FB & Sleep \\ C1 & 0.5 & & & & & 0.2 \\ C2 & & 0.8 & & & & 1.0 \\ C3 & & & 0.6 & 0.4 & & \\ Pass & 0.2 & 0.4 & 0.4 & & & \\ Pub & 0.1 & & & & & \\ FB & & & & & 0.9 & \\ Sleep & & & & & & 1 \end{bmatrix}$$

# Outline

- 1 Basic Setting
- 2 Revisit Markov Chain
- 3 Markov Reward Process
- 4 Markov Decision Process
- 5 Summary
- 6 References

# Markov Reward Process

A Markov reward process is a Markov chain with values.

## Definition

A Markov Reward Process is a tuple  $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

- $\mathcal{S}$  is a (finite) set of states
- $\mathcal{P}$  is a state transition probability matrix

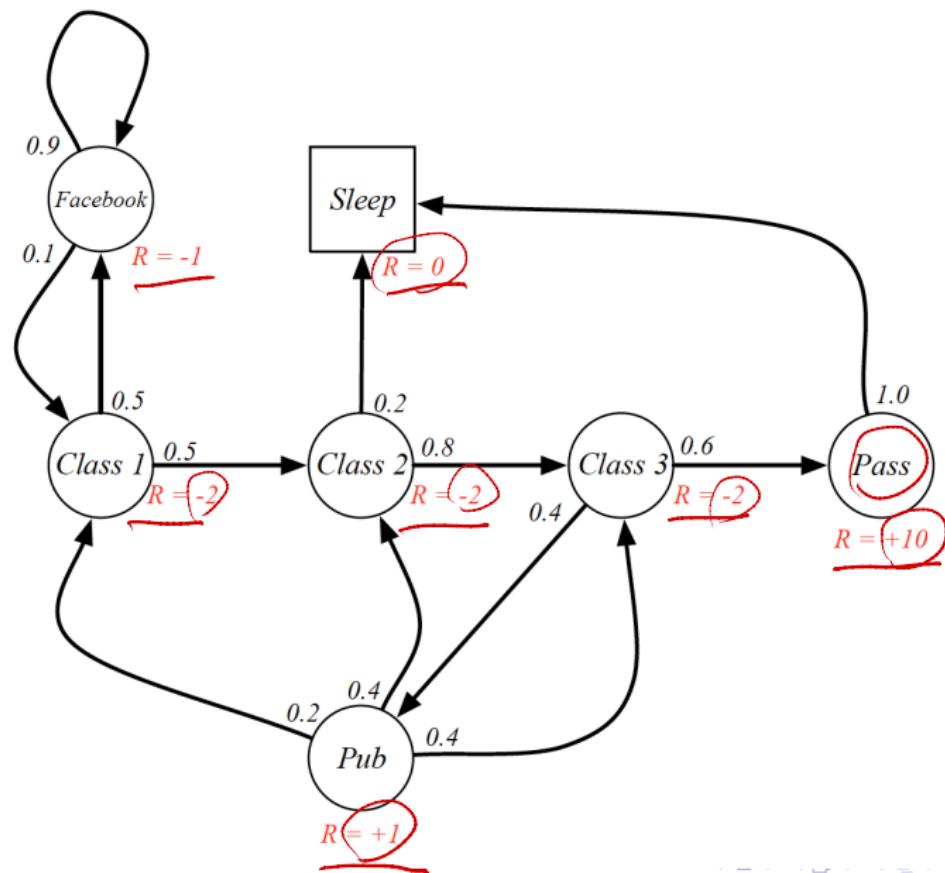
$$\mathcal{P}_{s,s'} = \mathbb{P}[S_{t+1} = s' | S_t = s]$$

- $\mathcal{R}$  is a reward function,

$$\mathcal{R}_s = \mathbb{E}[R_{t+1} | S_t = s]$$

- $\gamma$  is a discount factor,  $\gamma \in [0, 1]$

# Example: Student MRP



# Return

## Definition

The *return*  $G_t$  is the total discounted reward from time-step  $t$ .

$$G_t = \underbrace{R_{t+1} + \gamma R_{t+2} + \dots}_{\text{discounted rewards}} = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

- The *discount*  $\gamma \in [0, 1]$  is the present value of the future rewards
- The value of receiving reward  $R$  after  $k + 1$  time-steps is  $\gamma^k R$ .
- This values immediate reward above delayed reward.
  - ▶  $\gamma$  close to 0 leads to “myopic” evaluation
  - ▶  $\gamma$  close to 1 leads to “far-sighted” evaluation

# Why Discount?

Most Markov reward and decision processes are discounted. Why?

- Mathematically convenient to discount rewards
- Avoids infinite returns in cyclic Markov processes
- Uncertainty about the future may not be fully represented
- If the reward is financial, immediate rewards may earn more interest than delayed rewards
- Animal & human behavior shows preference for immediate reward
- It is sometimes possible to use undiscounted Markov reward processes (i.e.  $\gamma = 1$ ), e.g. if all sequences terminate.

# Value Function

The value function  $v(s)$  gives the long-term value of state  $s$

## Definition

The state value function  $v(s)$  of an MRP is the expected return starting from state  $s$

$$v(s) = \mathbb{E}[G_t | S_t = s]$$

# Example: Student MRP Returns

Sample **returns** for Student MRP:

Starting from  $S_1 = C_1$  with  $\gamma = \frac{1}{2}$

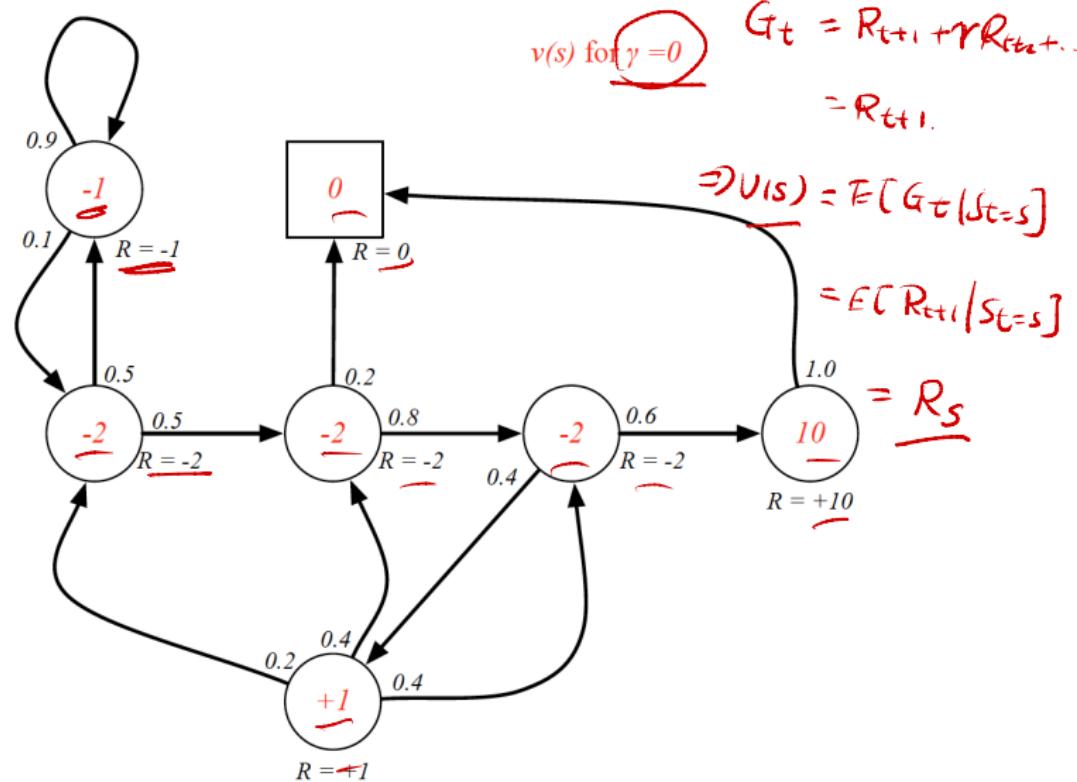
T

$$G_1 = R_2 + \gamma R_3 + \dots + \gamma^{T-2} R_T$$

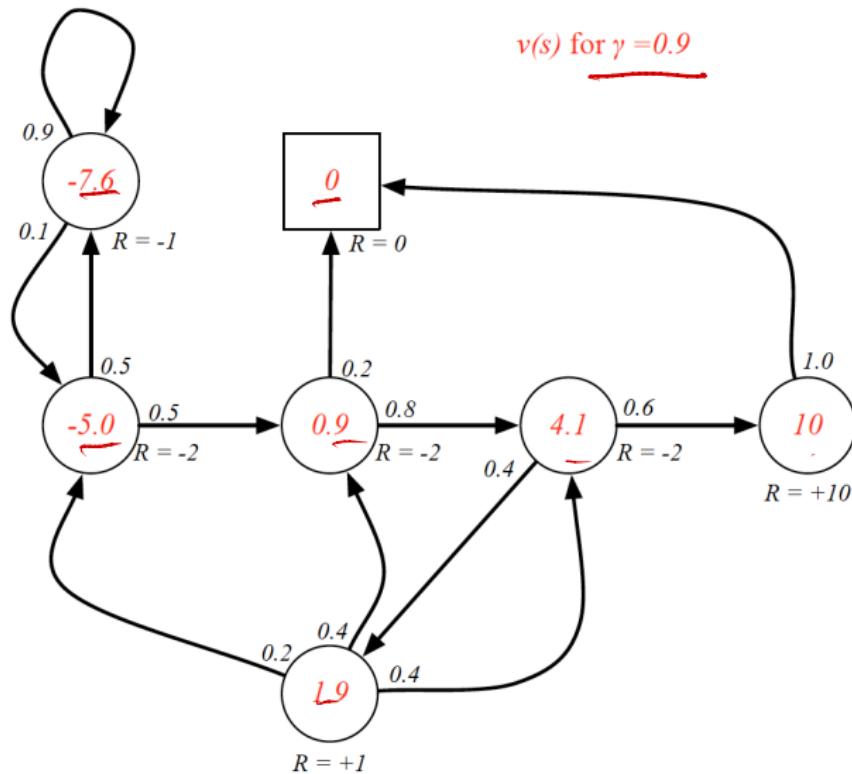
$S_1$	$S_2$	$S_3$	$S_4$	$S_5$	$S_6$	$\dots$
$C_1$	C1	C2	C3	Pass	Sleep	
$C_1$	C1	FB	FB	C1	C2	Sleep
$C_1$	C1	C2	C3	Pub	C2	C3
$C_1$	C1	FB	FB	C1	C2	C3
$C_1$	FB	FB	FB	C1	C2	C3

$$\begin{aligned} v_1 &= -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 10 * \frac{1}{8} &= -2.25 \\ v_1 &= -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16} &= -3.125 \\ v_1 &= -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 1 * \frac{1}{8} - 2 * \frac{1}{16} \dots &= -3.41 \\ v_1 &= -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16} \dots &= -3.20 \end{aligned}$$

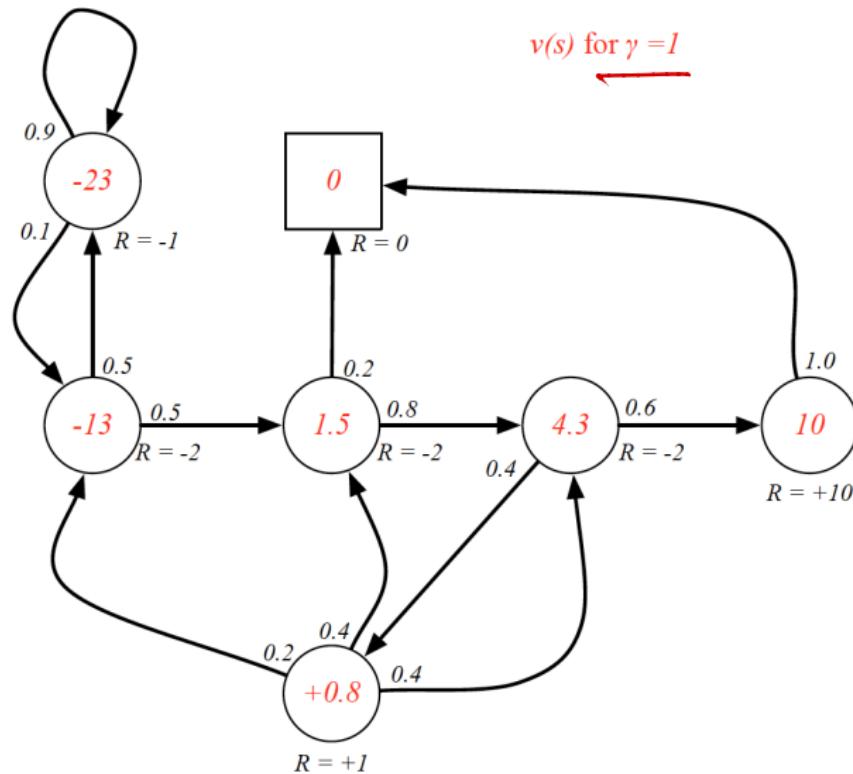
# Example: State-Value Function for Student MRP



# Example: State-Value Function for Student MRP



# Example: State-Value Function for Student MRP



# Bellman Equation for MRPs

$$G_t = R_{t+1} + \gamma G_{t+1}$$

The value function can be decomposed into two parts:

- immediate reward  $R_{t+1}$
- discounted value of successor state  $\gamma v(S_{t+1})$

$$\begin{aligned} v(s) &= \mathbb{E}[G_t \mid S_t = s] \\ &= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \\ &= \mathbb{E}[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) \mid S_t = s] \\ &= \mathbb{E}[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\ &\stackrel{?}{=} \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s] \end{aligned}$$

# Bellman Equation for MRPs

$$1^{\circ} \quad V(s) = E[G_t | S_t = s] ; \quad V(S_t) = E[G_t | S_t] ; \quad \underline{V(S_{t+1}) = E[G_{t+1} | S_{t+1}]}$$

$$2^{\circ} \quad \text{Adam's Law: } E[E[Y|X]] = E[Y] ; \quad \hat{E}(\cdot) = E(\cdot|Z)$$

Adam's Law with Extra Conditionality.

$$\hat{E}(\hat{E}(Y|X)) = \hat{E}(Y)$$

$$\Leftrightarrow E[\underbrace{E[Y|X, Z]}_{|Z} = E[Y|Z]$$

$$3^{\circ} \quad \text{Let } Y = G_{t+1}, X = S_{t+1}, Z = S_t,$$

$$\Rightarrow \underbrace{E[\underbrace{E[G_{t+1}|S_{t+1}, S_t]}_{|S_t}]}_{|S_t} = E[G_{t+1}|S_t]$$

↓ Markov property.

$$\frac{E[\underbrace{E[G_{t+1}|S_{t+1}]}_{|S_t}]}{V(S_{t+1})} \rightarrow \hat{E}[V(S_{t+1})|S_t]$$

$$\Rightarrow E[G_{t+1}|S_t] = E[V(S_{t+1})|S_t] \Rightarrow \text{lets, } \frac{E[G_{t+1}|S_t=s]}{= E[V(S_{t+1})|S_t=s]}$$

# Bellman Equation for MRPs

$$\begin{aligned} 4^{\circ}. \quad & \underline{V(s)} = E[G_t | S_t=s] = E[R_{t+1} + r G_{t+1} | S_t=s] \\ &= E[R_{t+1} | S_t=s] + r \underline{E[G_{t+1} | S_t=s]} \\ &= E[R_{t+1} | S_t=s] + r E[V(S_{t+1}) | S_t=s] \\ &= \underline{E[R_{t+1} + r V(S_{t+1}) | S_t=s]} \end{aligned}$$

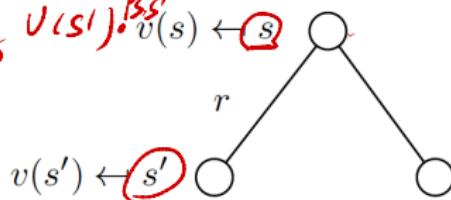
# Bellman Equation for MRPs

$$V(s) = E[R_{t+1} | S_{t+1}=s] + \gamma E[V(S_{t+1}) | S_t=s]$$

NOTE

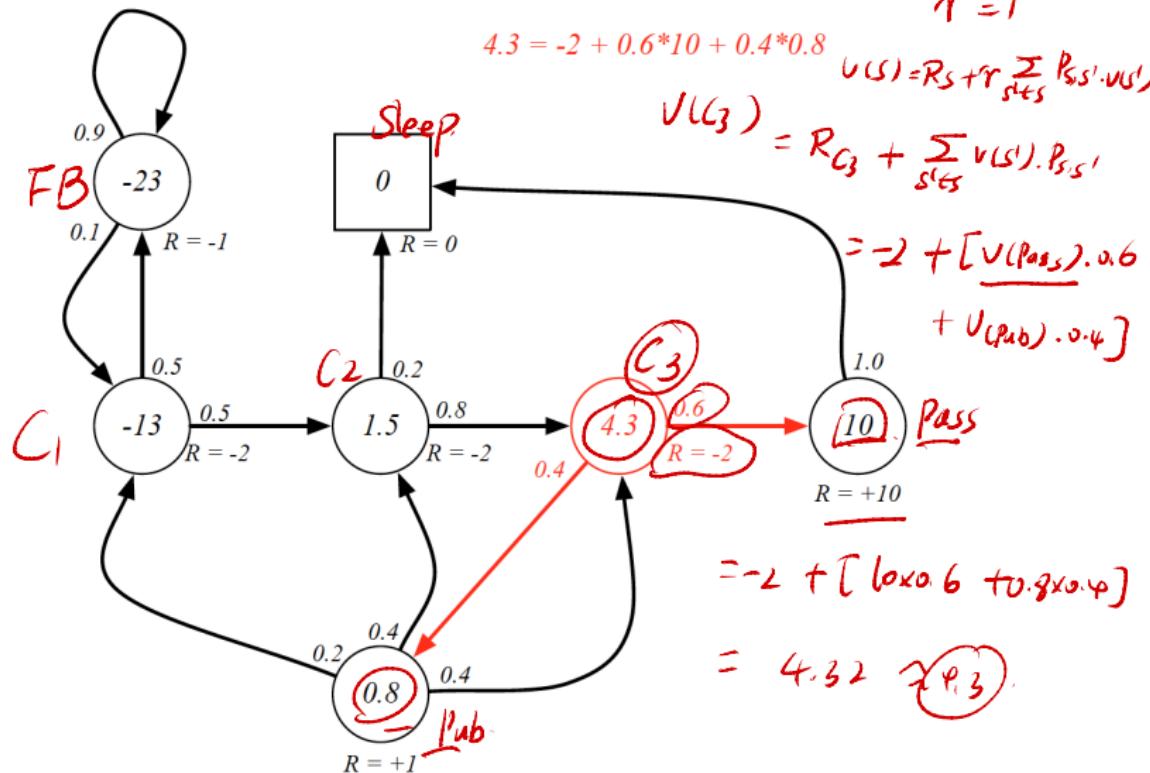
$$= R_s + \gamma \sum_{s' \in S} [E[V(S_{t+1}) | S_t=s, S_{t+1}=s']] \cdot P(S_{t+1}=s' | S_t=s)$$
$$v(s) = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t=s]$$

$$= R_s + \gamma \sum_{s' \in S} v(s') \cdot \mathcal{P}_{ss'}$$



$$v(s) = \underbrace{R_s + \gamma \sum_{s' \in S} \mathcal{P}_{ss'} v(s')}_{\text{Bellman Equation}}$$

# Example: Bellman Equation for Student MRP



# Bellman Equation in Matrix Form

The Bellman equation can be expressed concisely using matrices,

$$\mathbf{v} = \mathcal{R} + \gamma \mathcal{P} \mathbf{v}$$

where  $\mathbf{v}$  is a column vector with one entry per state

$$\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} \mathcal{R}_1 \\ \vdots \\ \mathcal{R}_n \end{bmatrix} + \gamma \underbrace{\begin{bmatrix} \mathcal{P}_{11} & \cdots & \mathcal{P}_{1n} \\ \vdots & \ddots & \vdots \\ \mathcal{P}_{n1} & \cdots & \mathcal{P}_{nn} \end{bmatrix}}_{\mathcal{P}} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$$

# Solving the Bellman Equation

- The Bellman equation is a linear equation
- It can be solved directly:

$$\begin{aligned}\mathbf{v} &= \mathcal{R} + \gamma \mathcal{P} \mathbf{v} \\ (\mathbf{I} - \gamma \mathcal{P}) \mathbf{v} &= \mathcal{R} \\ \mathbf{v} &= (\mathbf{I} - \gamma \mathcal{P})^{-1} \mathcal{R}\end{aligned}$$

- Computational complexity is  $\underline{\mathcal{O}(n^3)}$  for  $n$  states
- Direct solution only possible for small MRPs
- There are many iterative methods for large MRPs, e.g.
  - ▶ Dynamic programming
  - ▶ Monte-Carlo evaluation
  - ▶ Temporal-Different learning

# Outline

- 1 Basic Setting
- 2 Revisit Markov Chain
- 3 Markov Reward Process
- 4 Markov Decision Process
- 5 Summary
- 6 References

M.C.  
M.C + reward  
M.C + reward, action | action-value  
| state-value

# Markov Decision Process

A Markov decision process (MDP) is a Markov reward process with decisions. It is an environment in which all states are Markovian.

## Definition

A Markov Decision Process is a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

- $\mathcal{S}$  is a (finite) set of states
- $\mathcal{A}$  is a finite set of actions
- $\mathcal{P}$  is a state transition probability matrix

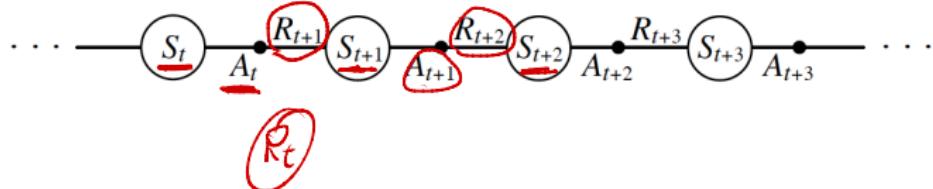
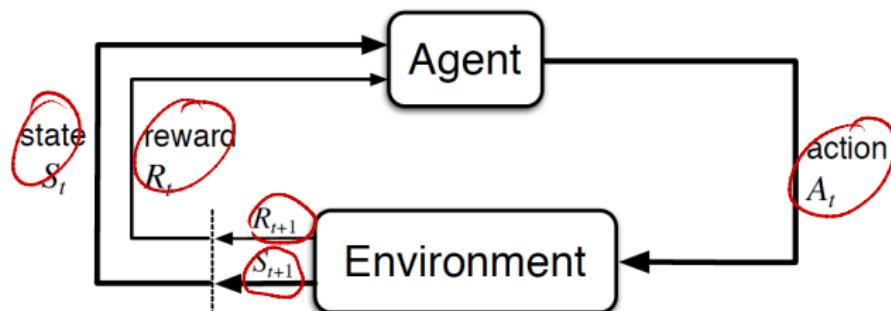
$$\mathcal{P}_{s,s'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$$

- $\mathcal{R}$  is a reward function,

$$\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$$

- $\gamma$  is a discount factor,  $\gamma \in [0, 1]$

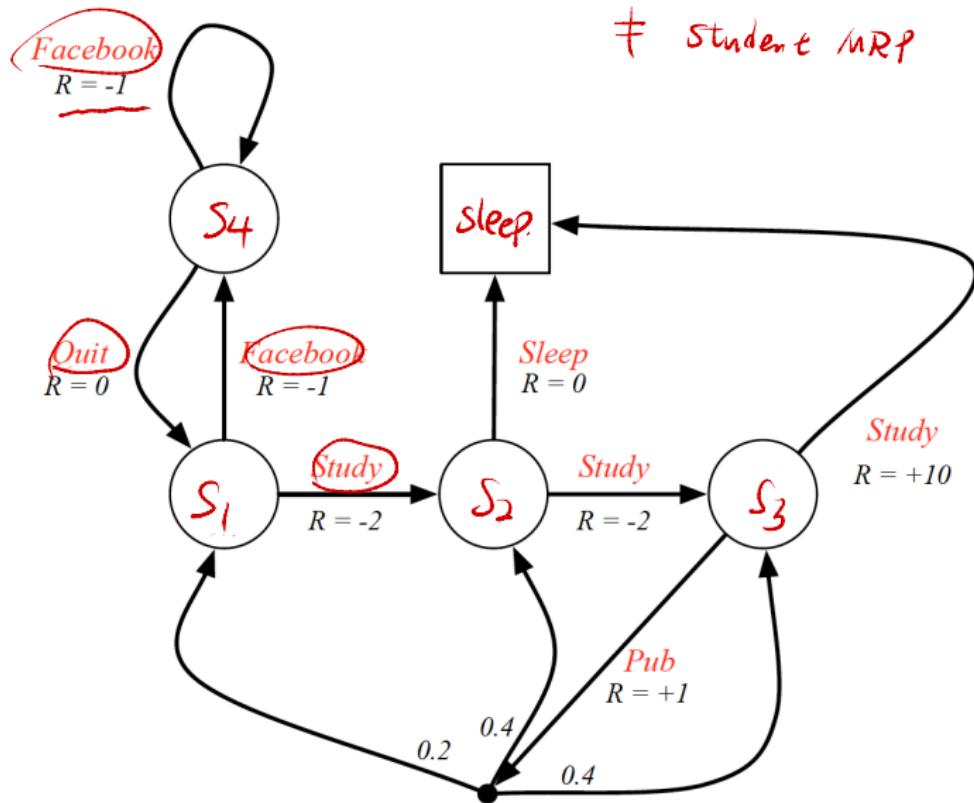
# The Agent-Environment Interface



# Example: Student MDP

$\neq$  Student M.C

$\neq$  Student M.R.P



# Policy

## Definition

A *policy*  $\pi$  is a distribution over actions given states,

$$\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$$

$S_t = s, H_t$

$= (S_0, A_0, R_0, S_1, A_1, R_1, \dots, S_t, R_t)$

- A policy fully defines the behavior of an agent
- MDP policies depend on the current state (not the history)
- i.e. Policies are *stationary* (time-independent),  
 $A_t \sim \pi(\cdot | S_t), \forall t > 0.$

# Policy

$$P_{s,s'}^{\pi} = P^{\pi}[S_{t+1}=s' | S_t=s] \stackrel{\text{Def}}{=} -$$

$$= \sum_{a \in A} P[S_{t+1}=s' | A_t=a, S_t=s] \cdot P(A_t=a | S_t=s)$$

$$= \sum_{a \in A} P_{s,a}^{\pi} \cdot \pi(a|s)$$

- Given an MDP  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$  and a policy  $\pi$
- The state sequence  $S_1, S_2, \dots$  is a  $\langle \mathcal{S}, \mathcal{P}^{\pi} \rangle$
- The state and reward sequence  $S_1, R_2, S_2, \dots$  is a  $\langle \mathcal{S}, \mathcal{P}^{\pi}, \mathcal{R}^{\pi}, \gamma \rangle$

- where  $2^o. R_s^{\pi} = E_{\pi}[R_{t+1} | S_t=s] \stackrel{\text{Def}}{=} \sum_{a \in A} E[R_{t+1} | S_t=s, A_t=a]$

$$P_{s,s'}^{\pi} = \sum_{a \in A} \pi(a|s) P_{s,a}^{\pi} \quad \bullet P(A_t=a | S_t=s)$$

$$R_s^{\pi} = \sum_{a \in A} \pi(a|s) R_s^a \quad = \sum_{a \in A} R_s^a \cdot \pi(a|s)$$

# Value Function

## Definition

The state-value function  $v_\pi(s)$  of an MDP is the expected return starting from state  $s$ , and then following policy  $\pi$

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] \quad \begin{aligned} &= v_\pi(s_t) \\ &= \mathbb{E}_\pi[G_t | S_t = s] \end{aligned}$$

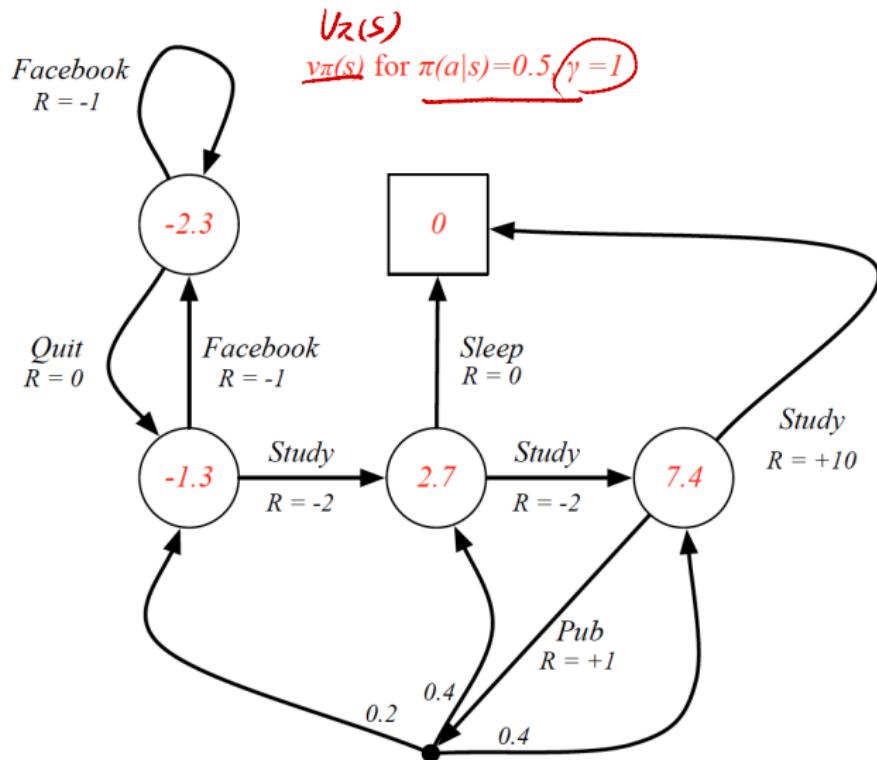
## Definition

The action-value function  $q_\pi(s, a)$  is the expected return starting from state  $s$ , taking action  $a$ , and then following policy  $\pi$

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$

$$q_\pi(s_t, a_t) = \mathbb{E}_\pi[G_t | S_t = s_t, A_t = a_t]$$

# Example: State-Value Function for Student MDP



# Bellman Expectation Equation

The state-value function can again be decomposed into immediate reward plus discounted value of successor state,

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

*Similar to MRP.* Homework

The action-value function can similarly be decomposed,

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]$$

# Bellman Expectation Equation

$$1^{\circ}. \quad q_z(s, a) = E_z [G_t | S_t=s, A_t=a]$$

$$q_z(S_t, A_t) = E_z [G_t | S_t, A_t], \quad q_z(S_{t+1}, A_{t+1}) = E_z [G_{t+1} | S_{t+1}, A_{t+1}]$$

2<sup>o</sup>. Adam's Law with extra conditioning.

$$E[Y|Z] = E[E[Y|X, Z]|Z]$$

Let  $Y = G_{t+1}$ ,  $Z = (S_t, A_t)$ ,  $X = (S_{t+1}, A_{t+1})$ .

$$\begin{aligned} \Rightarrow E[G_{t+1} | S_t, A_t] &= \overset{\text{Adam}}{E}[E[G_{t+1} | \underline{S_{t+1}, A_{t+1}}, S_t, A_t] | S_t, A_t] \\ &\stackrel{\text{Markov}}{=} \underline{E[E[G_{t+1} | S_{t+1}, A_{t+1}] | S_t, A_t]} \\ &= E[q_z(S_{t+1}, A_{t+1}) | S_t, A_t]. \end{aligned}$$

$$\Rightarrow \underline{E[G_{t+1} | S_t=s, A_t=a]} = E_z [q_z(S_{t+1}, A_{t+1}) | S_t=s, A_t=a]$$

by  $S, a$ .

# Bellman Expectation Equation

$$G_t = R_{t+1} + r G_{t+1}$$

$$3^{\circ}. \quad q_z(s, a) \triangleq E_z[G_t | S_t=s, A_t=a]$$

$$= E_z[R_{t+1} + r G_{t+1} | S_t=s, A_t=a]$$

$$\begin{aligned} &= \underbrace{E_z[R_{t+1} | S_t=s, A_t=a]}_{\downarrow} + r \underbrace{E_z[G_{t+1} | S_t=s, A_t=a]}_{+ r E_z[q_z(S_{t+1}, A_{t+1}) | S_t=s, A_t=a]} \\ &= \end{aligned}$$

$$= E_z[R_{t+1} + r q_z(S_{t+1}, A_{t+1}) | S_t=s, A_t=a]$$

# Bellman Expectation Equation for $V^\pi$

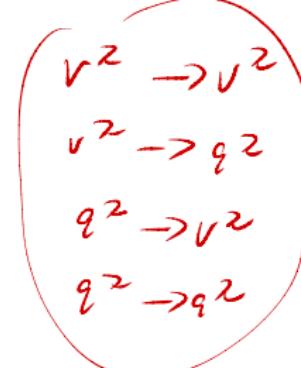
$$V_\pi(s) = E_\pi[G_t | S_t=s]$$

$$\stackrel{\text{NOTE}}{=} \sum_{a \in A} \underbrace{E_\pi[G_t | S_t=s, A_t=a]}_{q_\pi(s,a)} \cdot \underbrace{P(A_t=a | S_t=s)}_{\pi(a|s)}$$

$$= \sum_{a \in A} q_\pi(s, a) \cdot \pi(a|s)$$

$$v_\pi(s) \leftarrow s$$

$$q_\pi(s, a) \leftarrow a$$



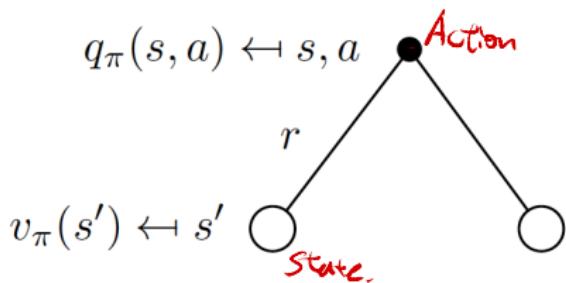
$$v_\pi(s) = \sum_{a \in A} \pi(a|s) q_\pi(s, a)$$

$$E_\pi[q_\pi(s, a)]$$

(A) action

$\sim z(a|s)$

# Bellman Expectation Equation for $Q^\pi$



$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s')$$

# Bellman Expectation Equation for $Q^\pi$

1°.  $E[\pi](q_\pi(s_{t+1}, a_{t+1}) | s_{t+1} = s', s_t = s, a_t = a]$

NOTE

$$= \sum_{a' \in A} E[\pi](q_\pi(s_{t+1}, a_{t+1}) | s_{t+1} = s', a_{t+1} = a', s_t = s, a_t = a) \cdot P(a_{t+1} = a')$$

$$= \sum_{a' \in A} E[\pi](q_\pi(s_{t+1}, a_{t+1}) | s_{t+1} = s', a_{t+1} = a') \cdot P(a_{t+1} = a' | s_{t+1} = s')$$
$$= \left( \sum_{a' \in A} q_\pi(s_t, a') \cdot \pi(a'|s_t) \right) = v_\pi(s_t)$$

2°.  $E[\pi](q_\pi(s_{t+1}, a_{t+1}) | s_t = s, a_t = a)$

NOTE

$$= \sum_{s' \in S} E[\pi](q_\pi(s_{t+1}, a_{t+1}) | s_t = s, a_t = a, s_{t+1} = s') \cdot P(s_{t+1} = s' | s_t = s, a_t = a)$$

$$= \left( \sum_{s' \in S} v_\pi(s') \cdot p_{s, s'}^a \right)$$

# Bellman Expectation Equation for $Q^\pi$

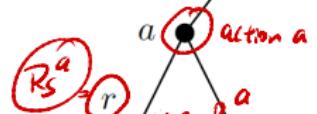
$$\begin{aligned} 3^{\circ}. \quad q_\pi(s, a) &= E_\pi [R_{t+1} + \gamma q_\pi(s_{t+1}, a_{t+1}) \mid s_t=s, a_t=a] \\ &= \frac{E_\pi [R_{t+1} \mid s_t=s, a_t=a]}{R_s^a + \gamma \sum_{s' \neq s} P_{s,s'}^a \cdot q_\pi(s')} \underbrace{E_\pi [q_\pi(s_{t+1}, a_{t+1}) \mid s_t=s, a_t=a]}_{\text{Bellman Expectation Equation}} \end{aligned}$$

# Bellman Expectation Equation for $v_\pi$

fixed point  
 $(fix)=x$

$$v_\pi(s) \leftarrow s$$

W.p.  $\pi(a|s)$

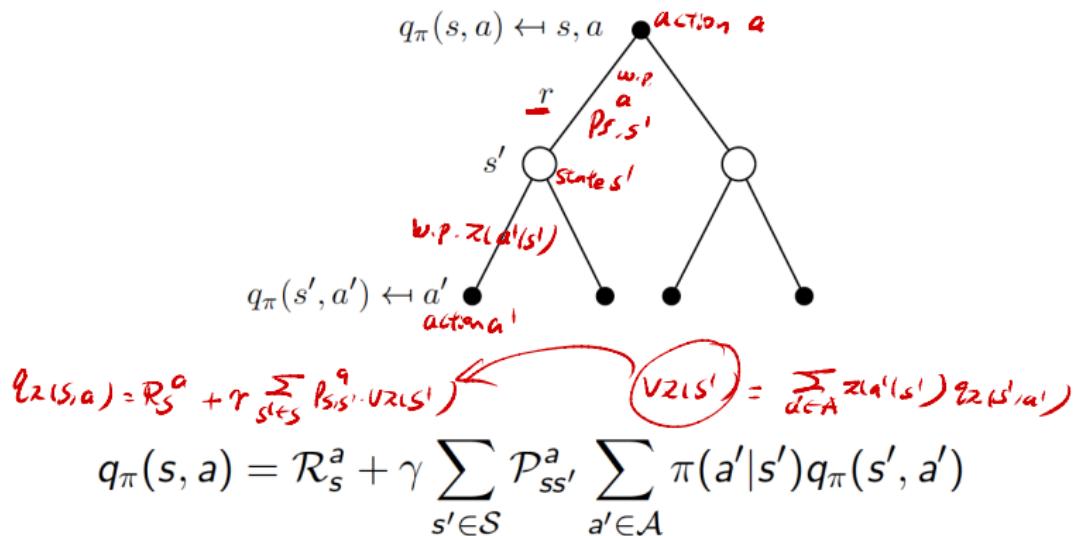


State  $s'$

$$V_\pi(s) = \underbrace{\sum_{a \in A} \pi(a|s) q_\pi(s, a)}_{q_\pi(s, a) = R_s^a + r \sum_{s' \in S} P_{s,s'}^a v_\pi(s')}$$

$$v_\pi(s) = \sum_{a \in A} \pi(a|s) \left( R_s^a + \gamma \sum_{s' \in S} P_{s,s'}^a v_\pi(s') \right)$$

# Bellman Expectation Equation for $q_\pi$



# Example: Bellman Expectation Equation in Student MDP

Policy  $\pi$ : equal prob for all options. , discount  $\gamma = 1$

$$\text{Facebook} \quad R = -1$$

$$7.4 = 0.5 * (1 + 0.2 * -1.3 + 0.4 * 2.7 + 0.4 * 7.4) + 0.5 * 10$$

$$\pi(\text{Facebook} | s_4)$$

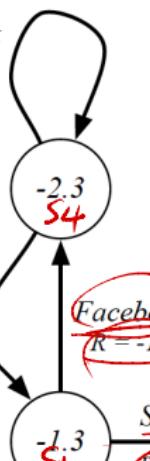
$$= \pi(\text{Quit} | s_4)$$

$$= 0.5$$

$$\text{Quit} \quad R = 0$$

$$\pi(\text{Study} | s_1) = 0.5$$

$$\pi(\text{Facebook} | s_1) = 0.5$$



$$P_{s_1, s_2}^{\text{Study}} = 1$$

$$7.4$$

$$+ 0.5 * 10$$



$$\pi(\text{Sleep} | s_2) = 0.5$$

$$\pi(\text{Study} | s_2) = 0.5$$

$$2.7$$

$$R = -2$$

$$7.4$$

$$R = +10$$

$$\pi(\text{Study} | s_3) = \pi(\text{Pub} | s_3) = 0.5$$

$$= 0.5$$

$$7.4$$

$$R = -2$$

$$7.4$$

$$R = +10$$

$$P_{s_3, s_1}^{\text{Pub}} = 0.2$$

$$P_{s_3, s_2}^{\text{Pub}} = 0.4$$

$$P_{s_3, s_3}^{\text{Pub}} = 0.4$$

# Example: Bellman Expectation Equation in Student MDP

$$V_j \triangleq V_Z(S_j)$$

$$1^{\circ}. \quad V_Z(S) = \sum_{a \in A} \pi(a|S) [R_S^a + r \sum_{S' \in S} p_{S,S'} V_Z(S')] ;$$

$$\begin{aligned} 2^{\circ}. \quad V_1 &\triangleq V_Z(S_1) = 0.5 [R_{S_1}^{\text{Study}} + 1 \cdot 1 \cdot V_Z(S_2)] \\ &\quad + 0.5 [R_{S_1}^{\text{Facebook}} + 1 \cdot 1 \cdot V_Z(S_4)] \\ &= 0.5 (-2 + V_2) + 0.5 (-1 + V_4) \end{aligned}$$

$$V_2 \triangleq V_Z(S_2) = 0.5 (-2 + V_3) + 0.5 (0 + 0)$$

$$\begin{aligned} V_3 &\triangleq V_Z(S_3) = 0.5 (1 + 0.2 V_1 + 0.4 V_2 + 0.4 V_3) \\ &\quad + 0.5 (0 + 0) \end{aligned}$$

$$V_4 \triangleq V_Z(S_4) = 0.5 (0 + V_1) + 0.5 (-1 + V_4)$$

# Example: Bellman Expectation Equation in Student MDP

3°.  $V_1 = -1.3, V_2 = 2.7, V_3 = 7.4, V_4 = -2.3$

4°.  $q_2(s, a) = R_s^a + \gamma \sum_{s'} p_{s,s'}^a V_2(s')$

$\Rightarrow q_2(s_1, \text{study}) = -2 + 1 \cdot V_2 = -2 + 2.7 = 0.7$

$q_2(s_1, \text{Facebook}) = -1 + 1 \cdot V_4 = -1 - 2.3 = -3.3$

$q_2(s_1, \text{Sleep}) = 0 + 0 = 0$

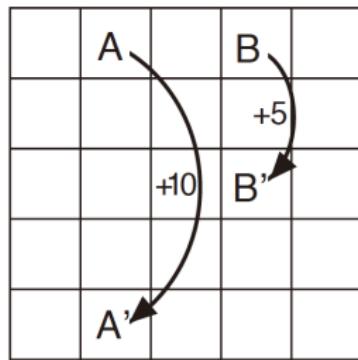
$q_2(s_2, \text{Study}) = -2 + 1 \cdot V_3 = -2 + 7.4 = 5.4$

$q_2(s_3, \text{Study}) = 10, q_2(s_3, \text{Pub}) = 4.78$

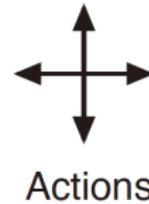
$q_2(s_4, \text{Facebook}) = -3.3; q_2(s_4, \text{Quit}) = -1.3$

# Example: Bellman Expectation Equation in Student MDP

# Homework: Bellman Expectation Equation in Gridworld



(a)



3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

(b)

What is the value function for the uniform random policy?

# Bellman Expectation Equation (Matrix Form)

The Bellman expectation equation can be expressed concisely using the induced MRP,

$$v_\pi = \underbrace{\mathcal{R}^\pi}_{\text{---}} + \underbrace{\gamma \mathcal{P}^\pi v_\pi}_{\text{---}}$$

with direct solution

$$v_\pi = (I - \gamma \mathcal{P}^\pi)^{-1} \mathcal{R}^\pi$$

# Optimal Value Function

## Definition

The optimal state-value function  $v_*(s)$  is the maximum value function over all policies

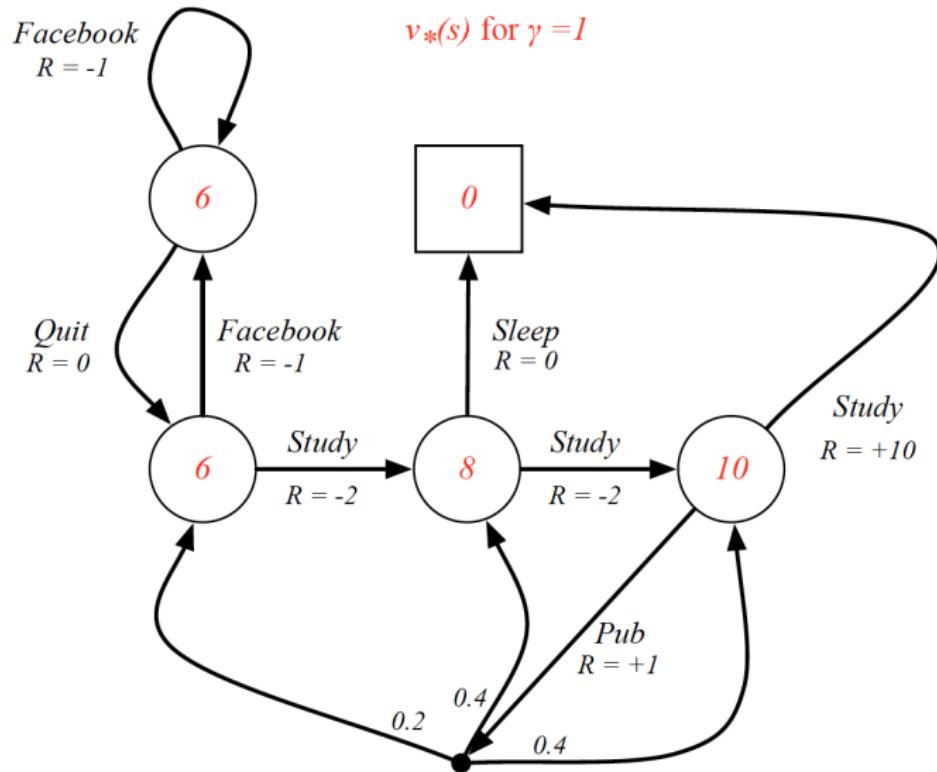
$$v_*(s) = \max_{\pi} v_{\pi}(s)$$

The optimal action-value function  $q_*(s, a)$  is the maximum action-value function over all policies

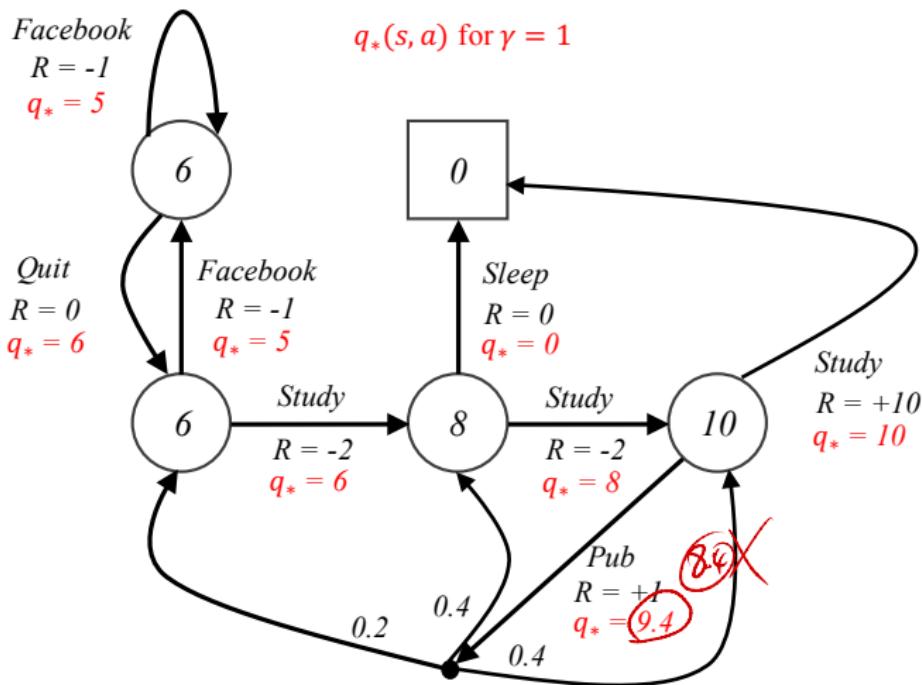
$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

- The optimal value function specifies the best possible performance in the MDP.
- An MDP is “solved” when we know the optimal value function.

# Example: Optimal Value Function for Student MDP



# Example: Optimal Action-Value Function for Student MDP



# Optimal Policy

Define a partial ordering over policies

$$\pi \geq \pi' \text{ if } v_\pi(s) \geq v_{\pi'}(s), \forall s$$

## Theorem

For any Markov Decision Process

- There exists an optimal policy  $\pi_*$  that is better than or equal to all other policies,  $\pi_* \geq \pi, \forall \pi$
- All optimal policies achieve the optimal value function,  
 $v_{\pi_*}(s) = v_*(s)$
- All optimal policies achieve the optimal action-value function,  
 $q_{\pi_*}(s, a) = q_*(s, a)$

# Finding an Optimal Policy

(Planning / RL)

Algorithm

(State)

① optimal value function + implicitly policy  $\Rightarrow$

Value function based

②

optimal policy (explicitly) + implicitly

policy based

optimal Value function

An optimal policy can be found by maximizing over  $q_*(s, a)$ ,

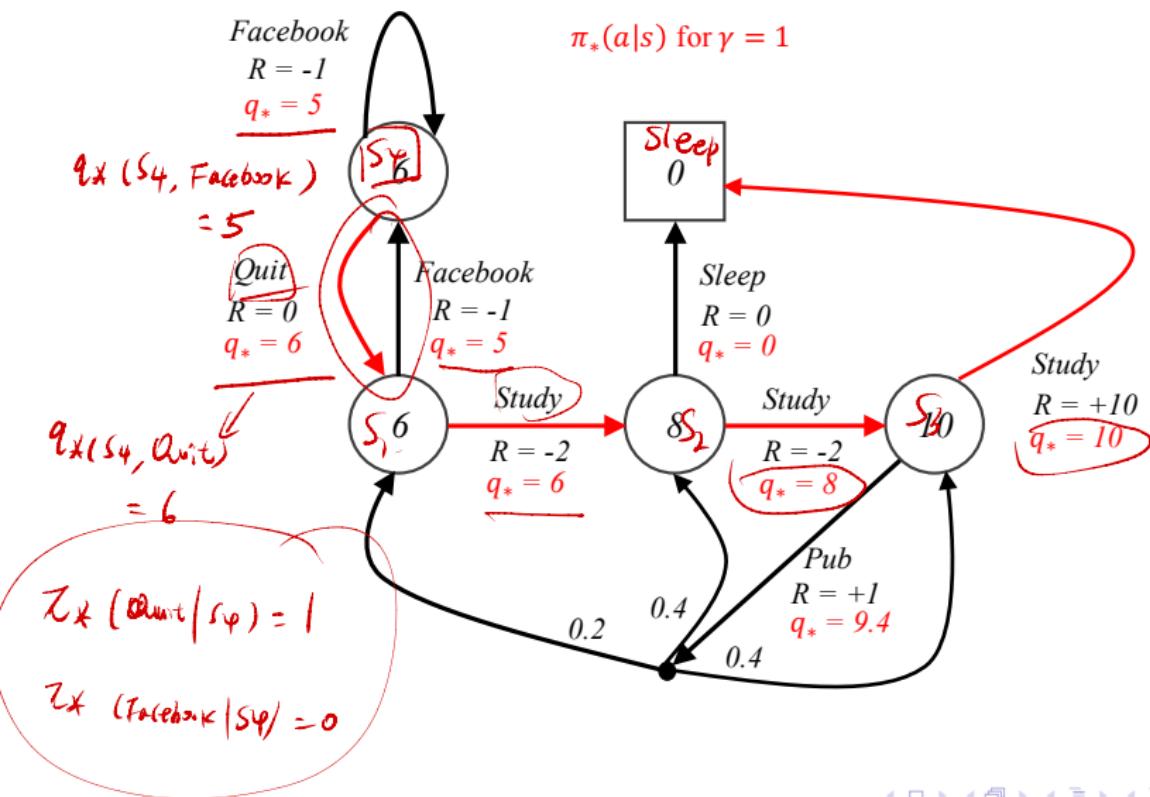
① + ②  $\Rightarrow$  Actor-critic

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } \underset{a \in A}{\arg \max} q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

- There is always a deterministic optimal policy for any MDP
- If we know  $q_*(s, a)$ , we immediately have the optimal policy

$U^*(s)$

## Example: Optimal Policy for Student MDP



# Bellman Optimality Equation for $v_*$

$U^*$	$\varnothing$
$q^*$	$U^*$
$U^*$	$U^*$
$q^*$	$q^*$

The optimal value functions are recursively related by the Bellman optimality equations:

Linear Programming (LP)

$$\max \theta_1 x_1 + \theta_2 x_2$$

$$\theta = (\theta_1, \theta_2)$$

$$0 \leq \theta_1, \theta_2 \leq 1$$

$$\theta_1 + \theta_2 = 1$$

$$\theta_1^* = 0, \theta_2^* = 1, x_2$$

$$\theta_1^* = 1, \theta_2^* = 0, x_1$$

$$\max(x_1, x_2)$$

$$\begin{array}{l} \text{const:} \\ (x_1, x_2) \\ x_1 \leq x_2 \end{array}$$

$$v_*(s) \leftarrow s$$

$$q_*(s, a) \leftarrow a$$

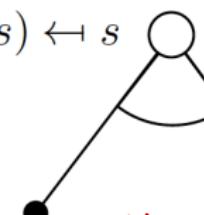
$$v_*(s) = \max_a q_*(s, a)$$

$$U_{LP}(s) = \max_a U_{LP}(s)$$

$$V_{LP}(s) = \sum_{a \in A} \pi(a|s) q_*(s, a)$$

$$\sum_{a \in A} \pi(a|s) = 1$$

$$0 \leq \pi(a|s) \leq 1, \forall a \in A$$



# Bellman Optimality Equation for $Q^*$

3°.  $E[R_{t+1} | S_t = s, A_t = a] = R_s^a ; E[V^*(S_{t+1}) | S_t = s, A_t = a] \stackrel{\text{LOTE}}{=} \sum_{s' \in S} E[V_A(S_{t+1}) | S_{t+1} = s', S_t = s, A_t = a]$

$q_*(s, a) = R_s^a + r \sum_{s' \in S} P_{s,s'}^a \cdot V_A(s')$

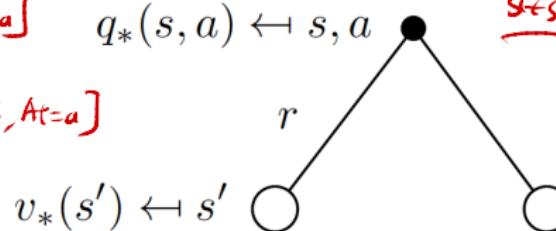
$q_*(s, a) = E[R_{t+1} | S_t = s, A_t = a] + r E[V^*(S_{t+1}) | S_t = s, A_t = a]$

$+ r E[V_A(S_{t+1}) | S_t = s, A_t = a]$

$$q_*(s, a) \leftarrow s, a$$

$$= \sum_{s' \in S} V_A(s') \cdot P_{s,s'}^a$$

$= E[R_{t+1} + r V_A(S_{t+1}) | S_t = s, A_t = a]$



1°. Bellman

Expectation  
Evaluation.

$q_*(s, a) = R_s^a + r \sum_{s' \in S} P_{s,s'}^a \cdot V_A(s') ; \Rightarrow q_*(s, a) = \max_z q_z(s, a)$

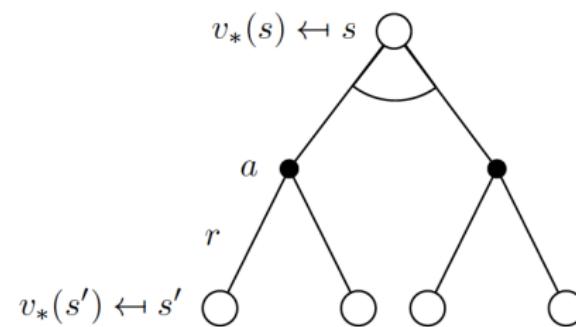
$q_*(s, a) = R_s^a + r \sum_{s' \in S} P_{s,s'}^a \cdot (\max_z V_z(s'))$

$q_*(s, a) = R_s^a + r \sum_{s' \in S} P_{s,s'}^a \cdot V^*(s')$

# Bellman Optimality Equation for $V^*$

$$2^o. V_*(s) = \max_a q_*(s, a); \quad q_*(s, a) = E[R_{t+1} + rV_*(s_{t+1}) | S_t=s, A_t=a]$$

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



$$1^o. V_*(s) = \max_a q_*(s, a); \quad q_*(s, a) = R_s^a + r \sum_{s' \in S} p_{s, s'}^a V_*(s')$$

$$v_*(s) = \max_a (R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s'))$$

# Bellman Optimality Equation for $Q^*$

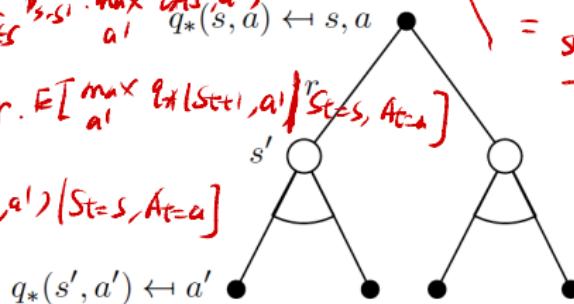
$$2^{\circ}. E[R_{t+1} | S_t=s, A_t=a] = R_s^a ; \quad E[\max_{a'} q_*(s_{t+1}, a') | S_t=s, A_t=a] \stackrel{LOTE}{=} \underline{\sum_{s' \in S} E[\max_{a'} q_*(s_{t+1}, a') | S_{t+1}=s', S_t=s, A_t=a]} .$$

$$q_*(s, a) = \mathbb{E} \left[ R_{t+1} + \gamma \max_{a'} q_*(s_{t+1}, a') | S_t = s, A_t = a \right] P[S_{t+1}=s' | S_t=s, A_t=a]$$

$$3^{\circ}. Q_*(s, a) = R_s^a + \gamma \sum_{S \in S} p_{s, s'}^a \cdot \max_{a'} q_*(s', a') \leftarrow s, a = \sum_{S \in S} \max_{a'} q_*(s', a') \cdot p_{s, s'}^a$$

$$= E[R_{t+1} | S_t=s, A_t=a] + \gamma \cdot E[\max_{a'} q_*(s_{t+1}, a') | S_t=s, A_t=a]$$

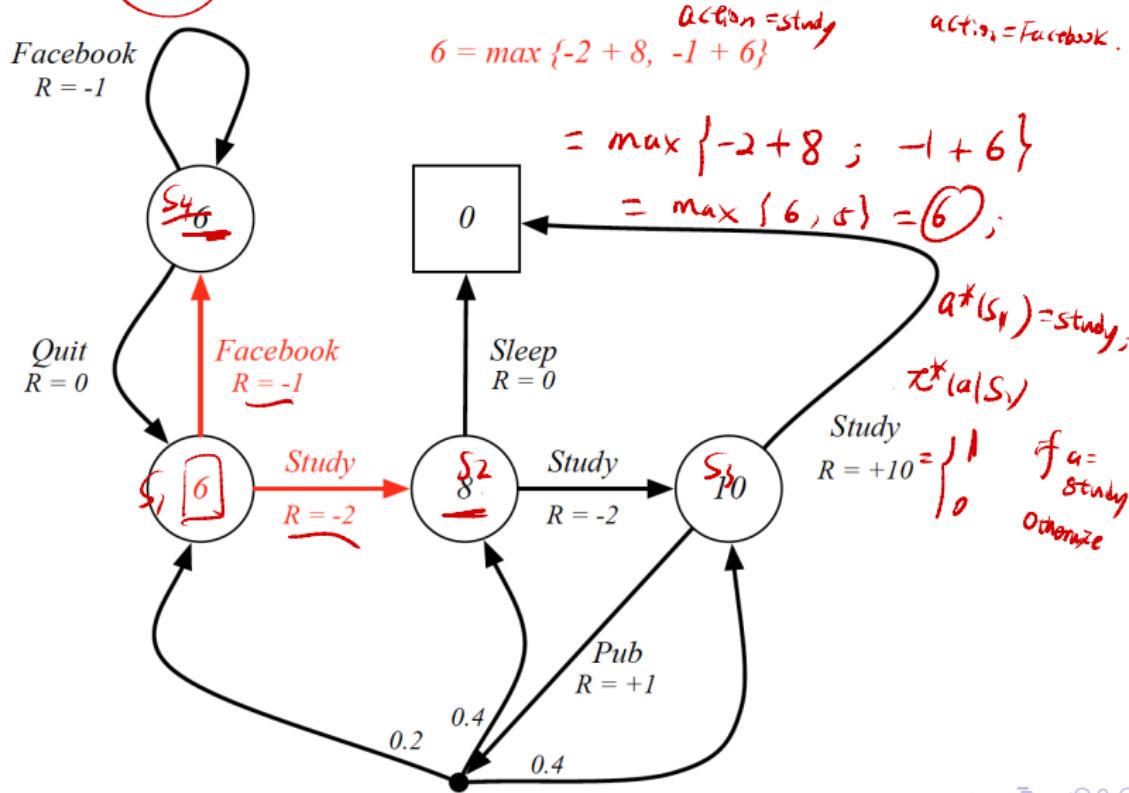
$$= E[R_{t+1} + \gamma \max_{a'} q_*(s_{t+1}, a') | S_t=s, A_t=a]$$



$$1^{\circ}. Q_*(s, a) = R_s^a + \gamma \sum_{S \in S} p_{s, s'}^a V_*(s') ; \quad V_*(s') = \max_{a'} q_*(s', a')$$

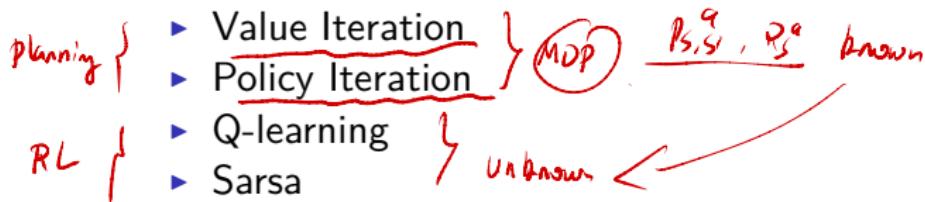
$$q_*(s, a) = R_s^a + \gamma \sum_{S \in S} p_{s, s'}^a \max_{a'} q_*(s', a')$$

# Example: Bellman Optimality Equation in Student MDP



# Solving the Bellman Optimality Equation in General

- Bellman Optimality Equation is non-linear
- No closed form solution (in general)
- Many iterative solution methods



# 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
  - ▶ BUT, number of states is often huge
  - ▶ So exhaustive sweeps of the state space are not possible

# Approximation and Reinforcement Learning

ADP | OR / control / AI

- RL methods: Approximating Bellman optimality equations
- Balancing reward accumulation and system identification (model learning) in case of unknown dynamics
- The on-line nature of reinforcement learning makes it possible to approximate optimal policies in ways that put more effort into learning to make good decisions for frequently encountered states, at the expense of less effort for infrequently encountered states.

# Outline

1 Basic Setting

2 Revisit Markov Chain

3 Markov Reward Process

4 Markov Decision Process

5 Summary

6 References

# Markov Decision Process

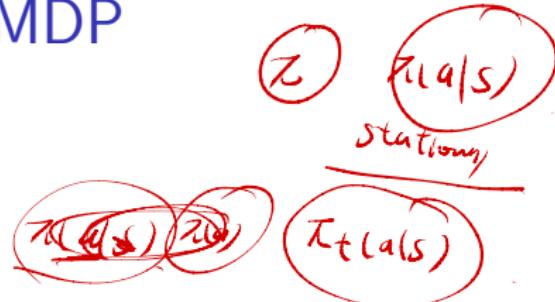
- Markov decision processes formally describe an environment for reinforcement learning
- Where the environment is fully observable
- i.e. The current state completely characterizes the process
- Almost all RL problems can be formalized as MDPs, e.g.
  - ▶ Optimal control primarily deals with continuous MDPs HJB
  - ▶ Partially observable problems can be converted into MDPs
  - ▶ Bandits are MDPs with one state

# Markov Models in General

Markov Models	Do we have control over the <u>state transitions</u> ?	
	No	Yes
Are the states completely observable?	Yes	<u>Markov Chain</u>
	No	<u>HMM</u> Hidden Markov Model
		<u>MDP</u> Markov Decision Process
		<u>POMDP</u> Partially Observable Markov Decision Process

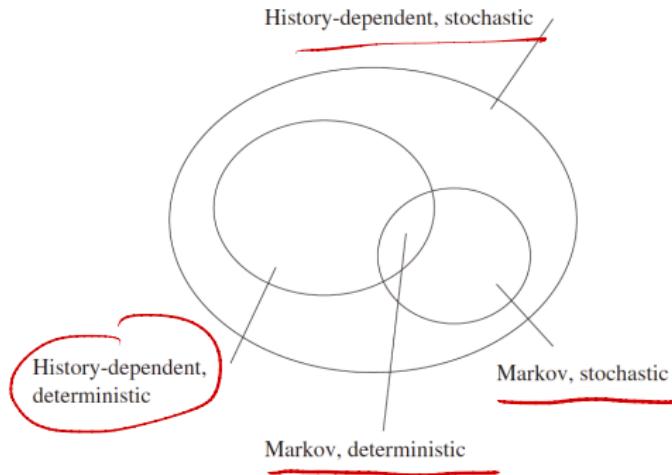
# Different Policy Families for MDP

$$\begin{array}{lll} \text{MAB} & \underline{\pi_t(a)} & : \underline{\pi_t(a|s_0)} \\ \text{MDP} & \underline{\pi_t(a|s)} & \end{array}$$



Policy $\pi_t$	Deterministic	Stochastic
Markov	$s_t \rightarrow a_t$	$a_t, s_t \rightarrow [0, 1]$
History-dependent	$h_t \rightarrow a_t$	$h_t, s_t \rightarrow [0, 1]$

# Different Policy Families for MDP



# Main Performance Criteria of MDP

- The finite criterion

$$\mathbb{E}[r_0 + r_1 + \dots + r_{N-1} | s_0]$$

- The discounted criterion

$$\mathbb{E}[r_0 + \gamma r_1 + \cancel{\gamma^2} r_2 + \dots + \cancel{\gamma^t} r_t + \dots | s_0]$$

- The total reward criterion

$$\mathbb{E}[r_0 + r_1 + \dots + r_t + \dots | s_0]$$

- The average criterion

$$\lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E}[r_0 + r_1 + \dots + r_{n-1} | s_0]$$

# Outline

1 Basic Setting

2 Revisit Markov Chain

3 Markov Reward Process

4 Markov Decision Process

5 Summary

6 References

# Main References

- Reinforcement Learning: An Introduction (second edition), R. Sutton & A. Barto, 2018.
- RL course slides from Richard Sutton, University of Alberta.
- RL course slides from David Silver, University College London.