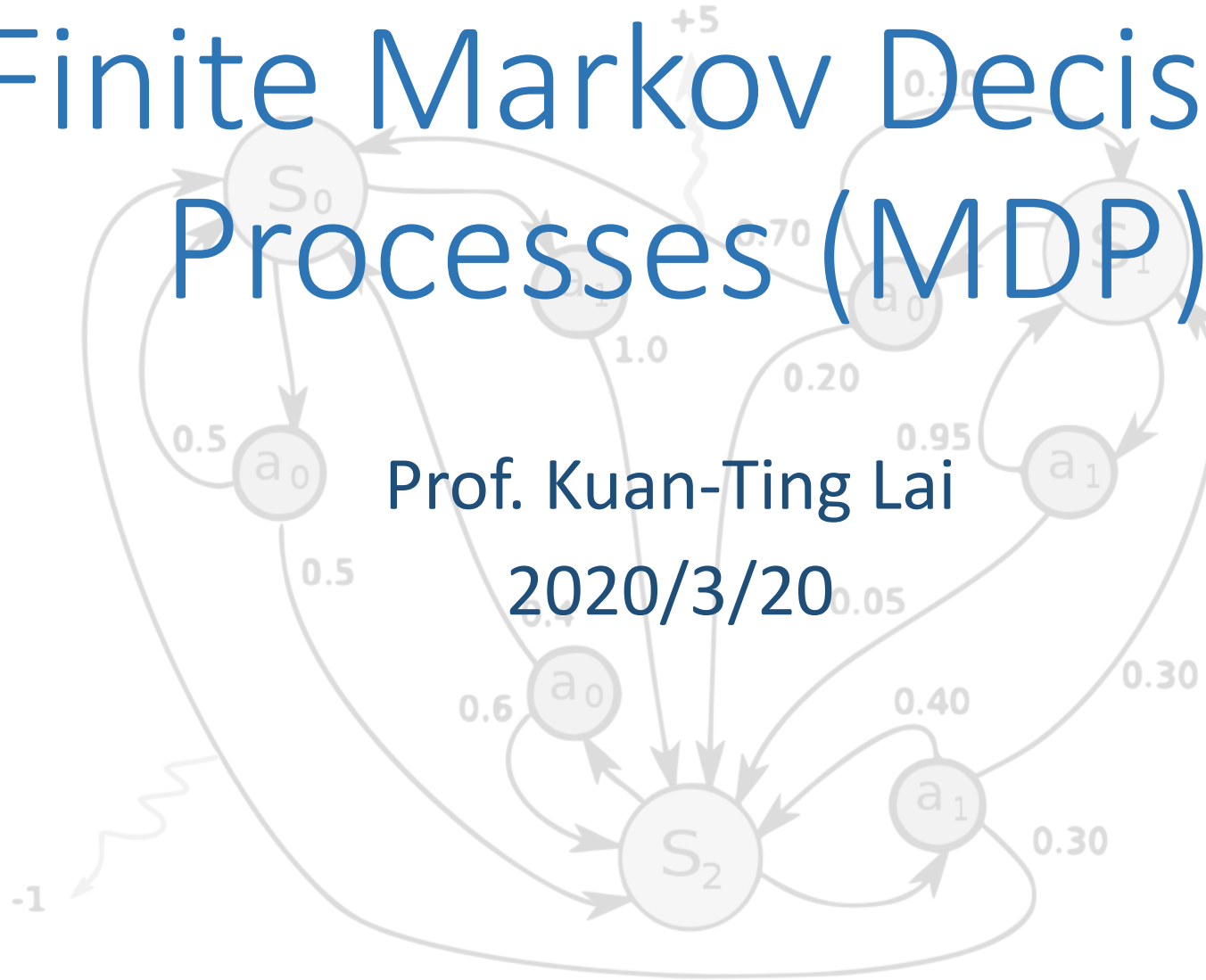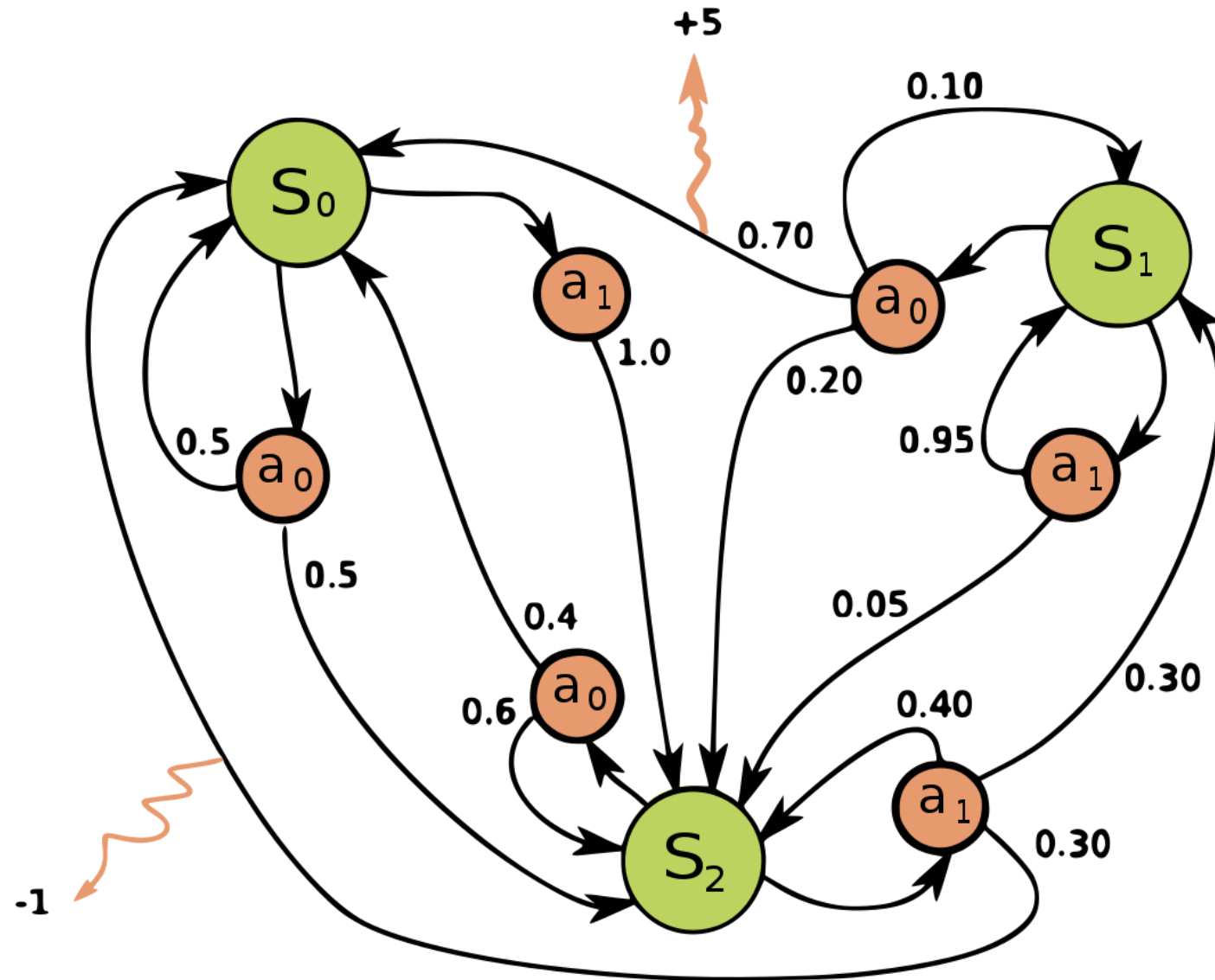# Finite Markov Decision Processes (MDP)

Prof. Kuan-Ting Lai

2020/3/20

# Markov Decision Process (MDP)

# Markov Property

- Current state can represent all information from the past states
- i.e. memoryless
- Let bygones be bygones

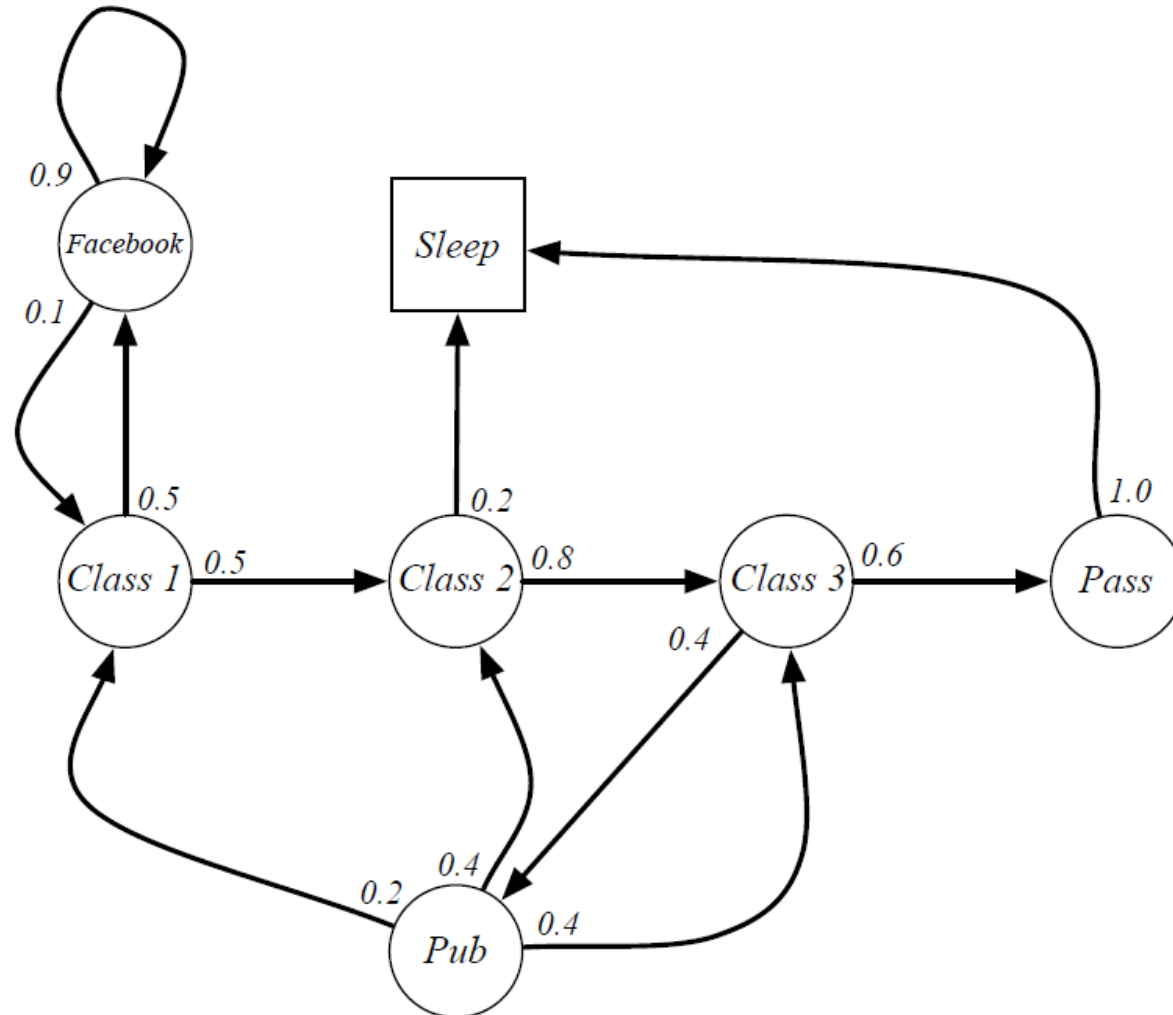## Definition

A state $S_t$ is *Markov* if and only if

$$\mathbb{P}\left[S_{t+1} \mid S_t\right] = \mathbb{P}\left[S_{t+1} \mid S_1, ..., S_t\right]$$

# Markov Process

- A Markov process is a memoryless random process, i.e. a sequence of random states $S_1$, $S_2$, … with Markov property

- Transition probability P(s, s') is the probability of moving from state $s$ to state $s'$

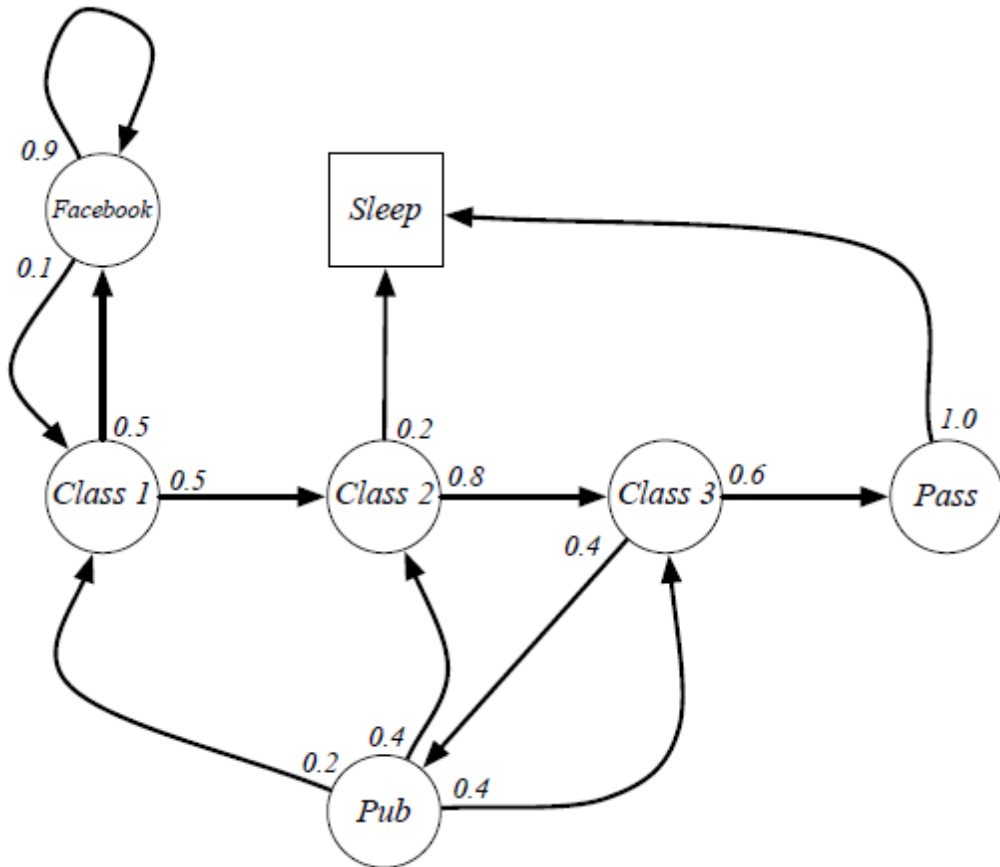$$\mathcal{P}_{ss'} = \mathbb{P}\left[S_{t+1} = s' \mid S_t = s\right]$$

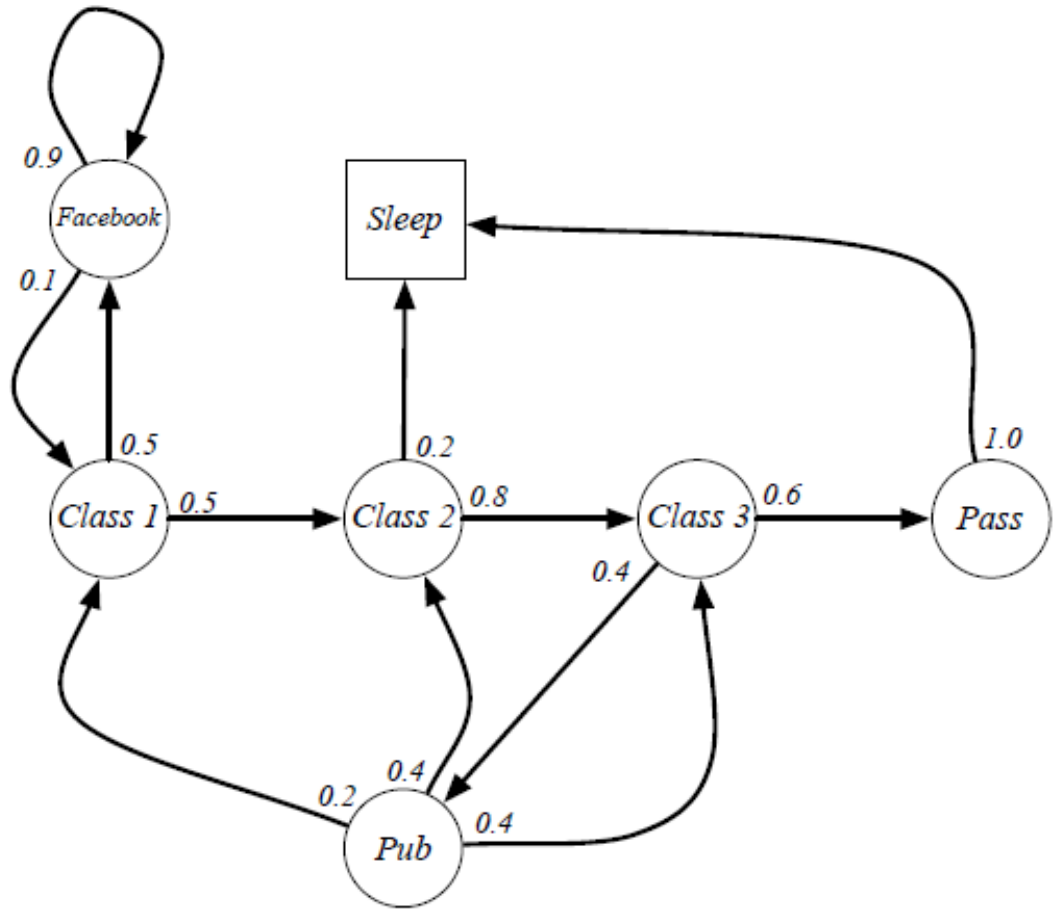# Student Markov Chain

# Student Markov Chain Episodes



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

$$S_1, S_2, ..., S_T$$

- 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



$$\mathcal{P} = \begin{array}{c} \\ C1 \\ C2 \\ C3 \\ Pass \\ Pub \\ FB \\ Sleep \end{array} \begin{array}{ccccccc} C1 & C2 & C3 & Pass & Pub & FB & Sleep \\ & 0.5 & & & & 0.5 & \\ & & 0.8 & & & & 0.2 \\ & & & 0.6 & 0.4 & & \\ & & & & & & 1.0 \\ 0.2 & 0.4 & 0.4 & & & & \\ 0.1 & & & & & 0.9 & \\ & & & & & & 1 \end{array}$$
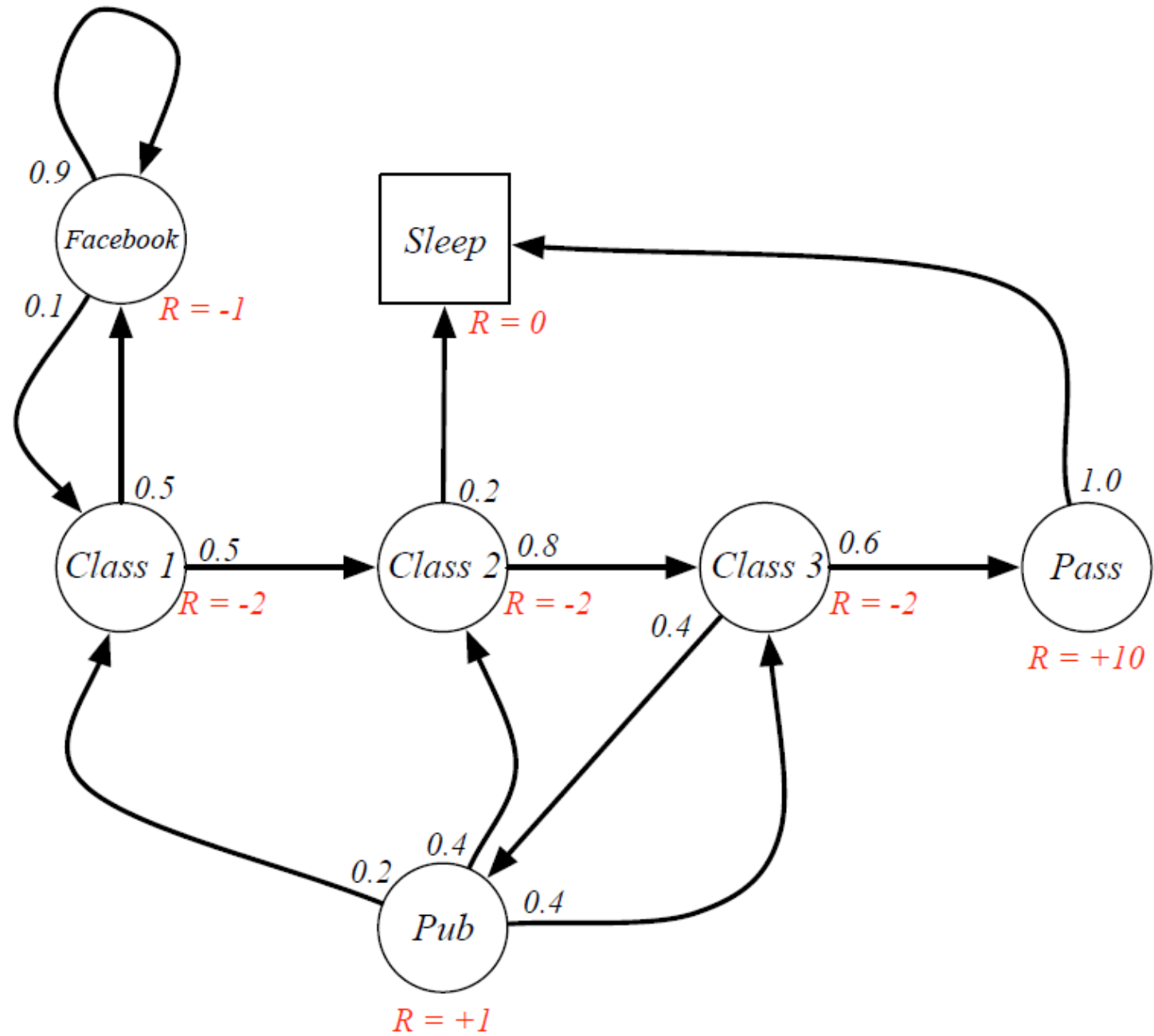
# Adding Reward to Markov Process

- 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,
  $\mathcal{P}_{ss'} = \mathbb{P}\left[S_{t+1} = s' \mid S_t = s\right]$

- $\mathcal{R}$ is a reward function, $\mathcal{R}_s = \mathbb{E}\left[R_{t+1} \mid S_t = s\right]$

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

# Student MRP

# Discounted Future Return $G_t$

- The discount $\gamma \in [0,1]$ is the present value of future rewards
  - $\gamma$ close to 0 leads to "short-sighed" evaluation
  - $\gamma$ close to 1 leads to "far-sighed" evaluation

**Definition**

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

$$G_t = R_{t+1} + \gamma R_{t+2} + \ldots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

# Why add discount factor $\gamma$?

- Uncertainty about the future
- Avoids infinite returns in cyclic Markov processes
- Animal/human behaviour shows preference for immediate reward

# Value Function

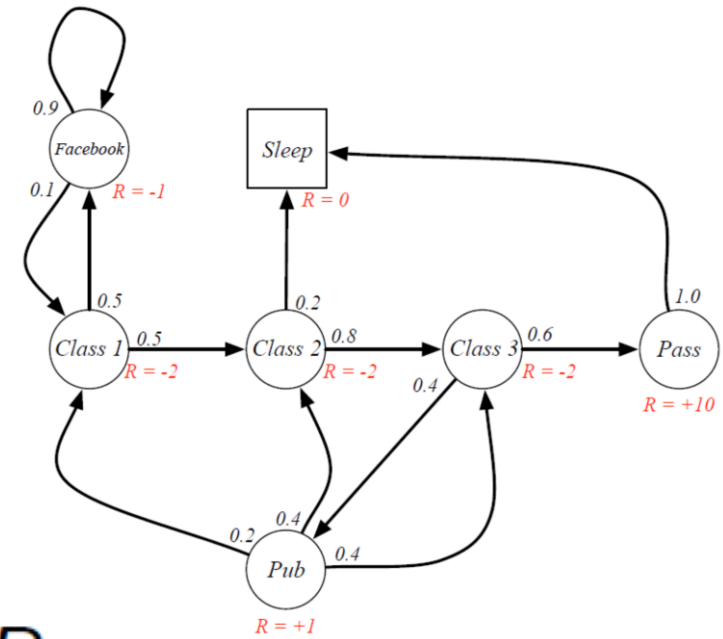- The value function v(s) estimates the long-term value of state s

**Definition**

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

$$v(s) = \mathbb{E}\left[G_t \mid S_t = s\right]$$

# Student MRP Returns



$\bullet\ \gamma = \dfrac{1}{2}$

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

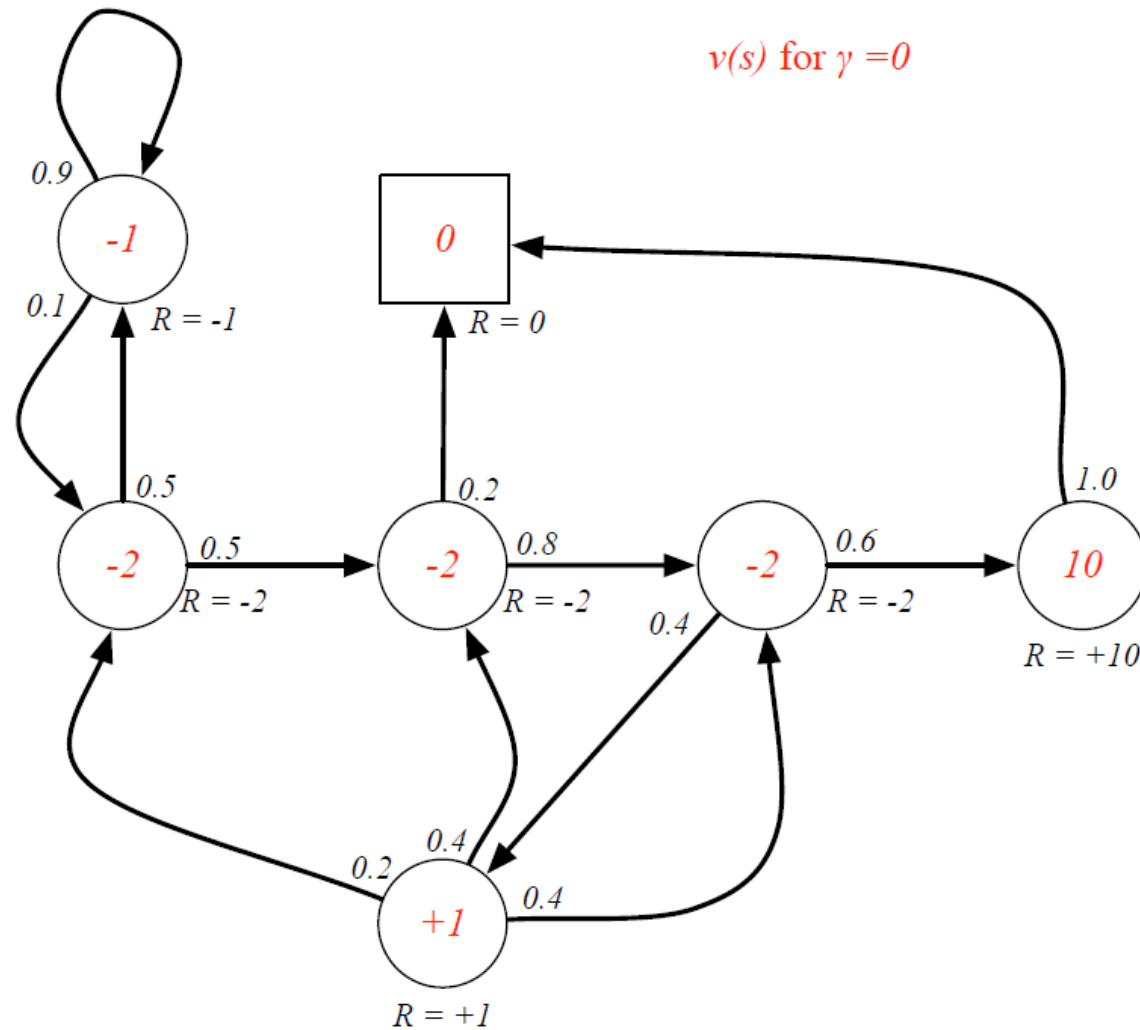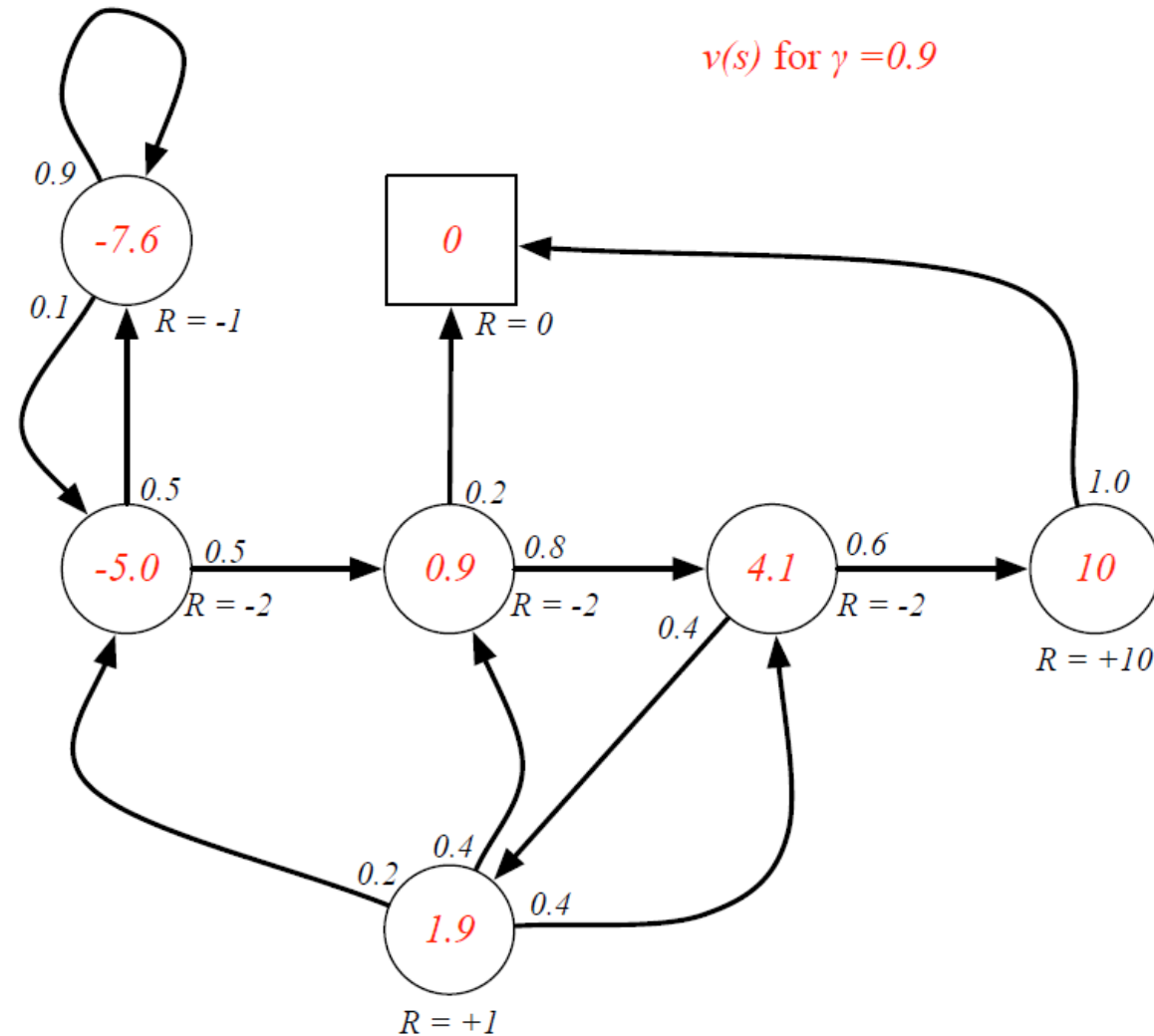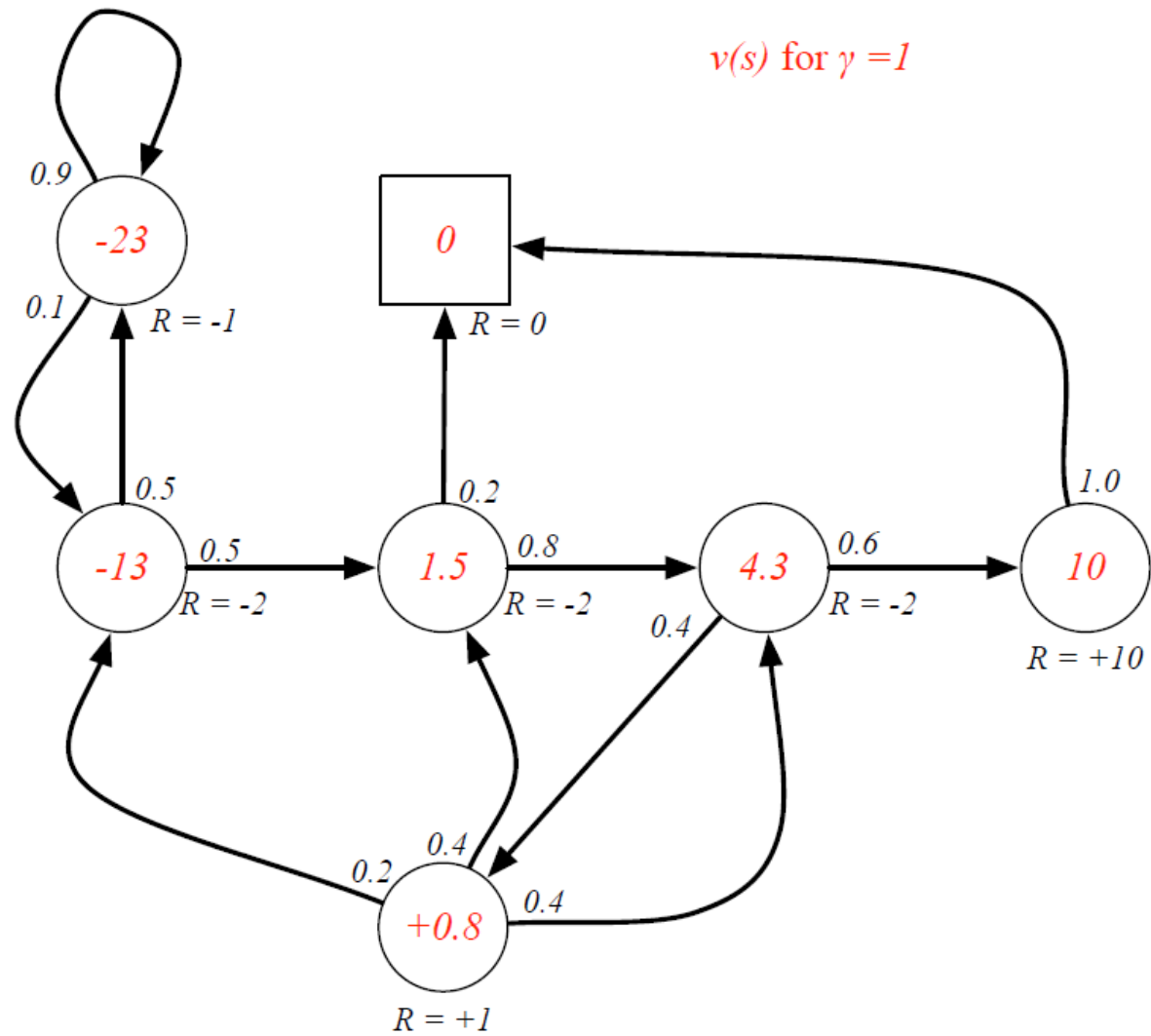| | | |
|---|---|---|
| C1 C2 C3 Pass Sleep | $v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 10 * \frac{1}{8}$ | = −2.25 |
| C1 FB FB C1 C2 Sleep | $v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16}$ | = −3.125 |
| C1 C2 C3 Pub C2 C3 Pass Sleep | $v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 1 * \frac{1}{8} - 2 * \frac{1}{16} ...$ | = −3.41 |
| C1 FB FB C1 C2 C3 Pub C1 ... | $v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16} ...$ | = −3.20 |
| FB FB FB C1 C2 C3 Pub C2 Sleep | | |

# State-Value Function for Student MRP (1)



$v(s)$ for $\gamma = 0$

# State-Value Function for Student MRP (2)



$v(s)$ for $\gamma = 0.9$

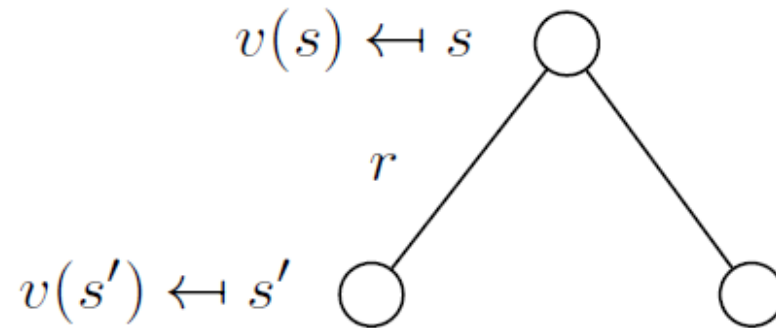# State-Value Function for Student MRP (3)

# Bellman Equation for MRPs

- The value function can be decomposed into two parts:
    - immediate reward $R_{t+1}$
    - discounted value of next state $\gamma$ $v(S_{t+1})$

$$v(s) = \mathbb{E}\left[G_t \mid S_t = s\right]$$

$$= \mathbb{E}\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + ... \mid S_t = s\right]$$

$$= \mathbb{E}\left[R_{t+1} + \gamma\left(R_{t+2} + \gamma R_{t+3} + ...\right) \mid S_t = s\right]$$

$$= \mathbb{E}\left[R_{t+1} + \gamma G_{t+1} \mid S_t = s\right]$$

$$= \mathbb{E}\left[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s\right]$$
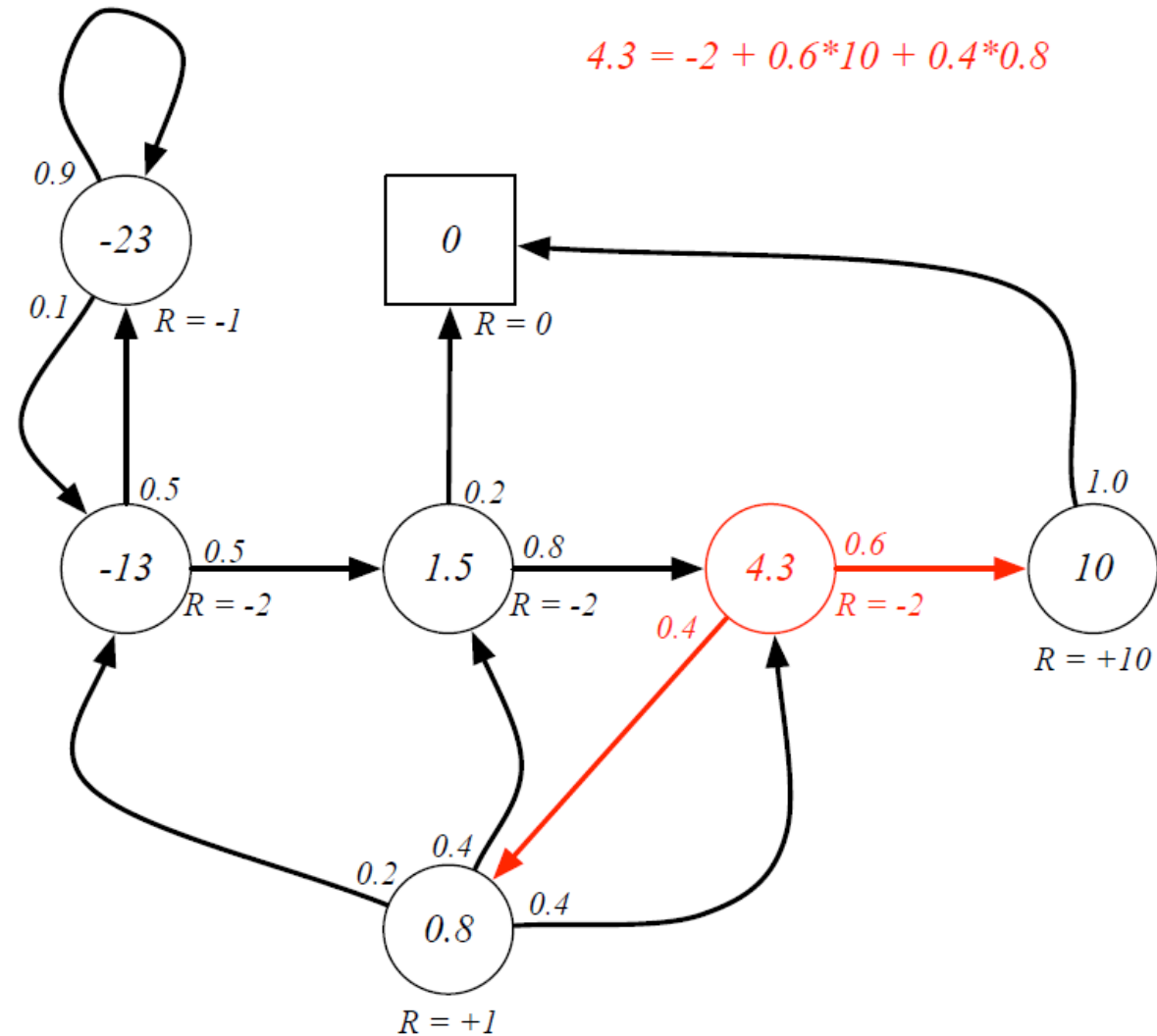
# Backup Diagram for Bellman Equation

$$v(s) = \mathbb{E}\left[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s\right]$$



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

# Calculating Student MDP using Bellman Equation



$4.3 = -2 + 0.6*10 + 0.4*0.8$
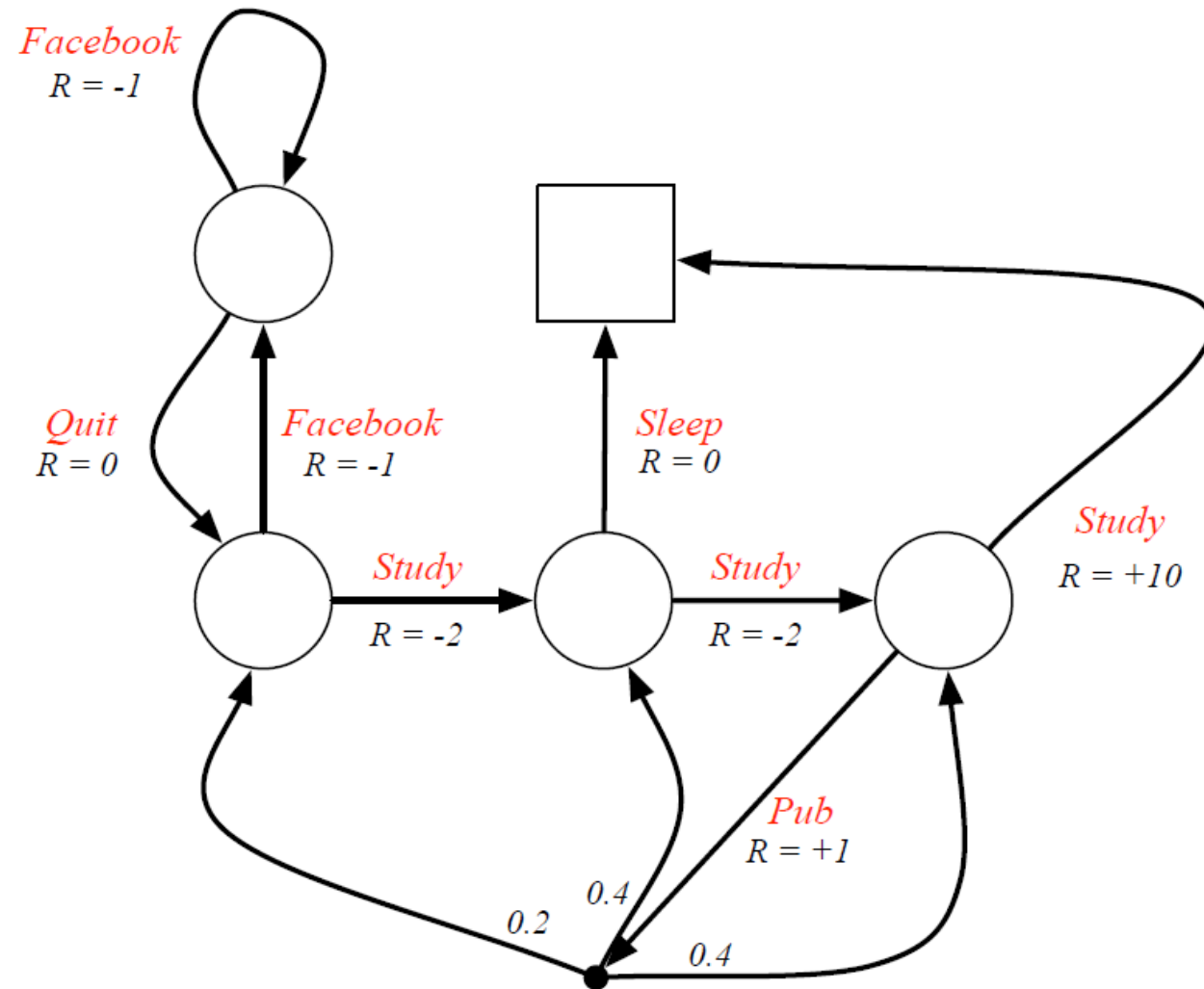
# Markov Decision Process

- A Markov decision process (MDP) is a Markov reward process with decisions.

**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,
  $\mathcal{P}_{ss'}^a = \mathbb{P}\left[S_{t+1} = s' \mid S_t = s, A_t = a\right]$
- $\mathcal{R}$ is a reward function, $\mathcal{R}_s^a = \mathbb{E}\left[R_{t+1} \mid S_t = s, A_t = a\right]$
- $\gamma$ is a discount factor $\gamma \in [0, 1]$.

# Student MDP with Actions

# Policy

- MDP Policies only depend on the current state, i.e. stationary

## Definition

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

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

# Policies

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

$$\mathcal{P}^\pi_{s,s'} = \sum_{a \in \mathcal{A}} \pi(a|s) \mathcal{P}^a_{ss'}$$

$$\mathcal{R}^\pi_s = \sum_{a \in \mathcal{A}} \pi(a|s) \mathcal{R}^a_s$$

# Value Function

The *state-value function* $v_\pi(s)$ of an MDP is the expected return starting from state $s$, and then following policy $\pi$
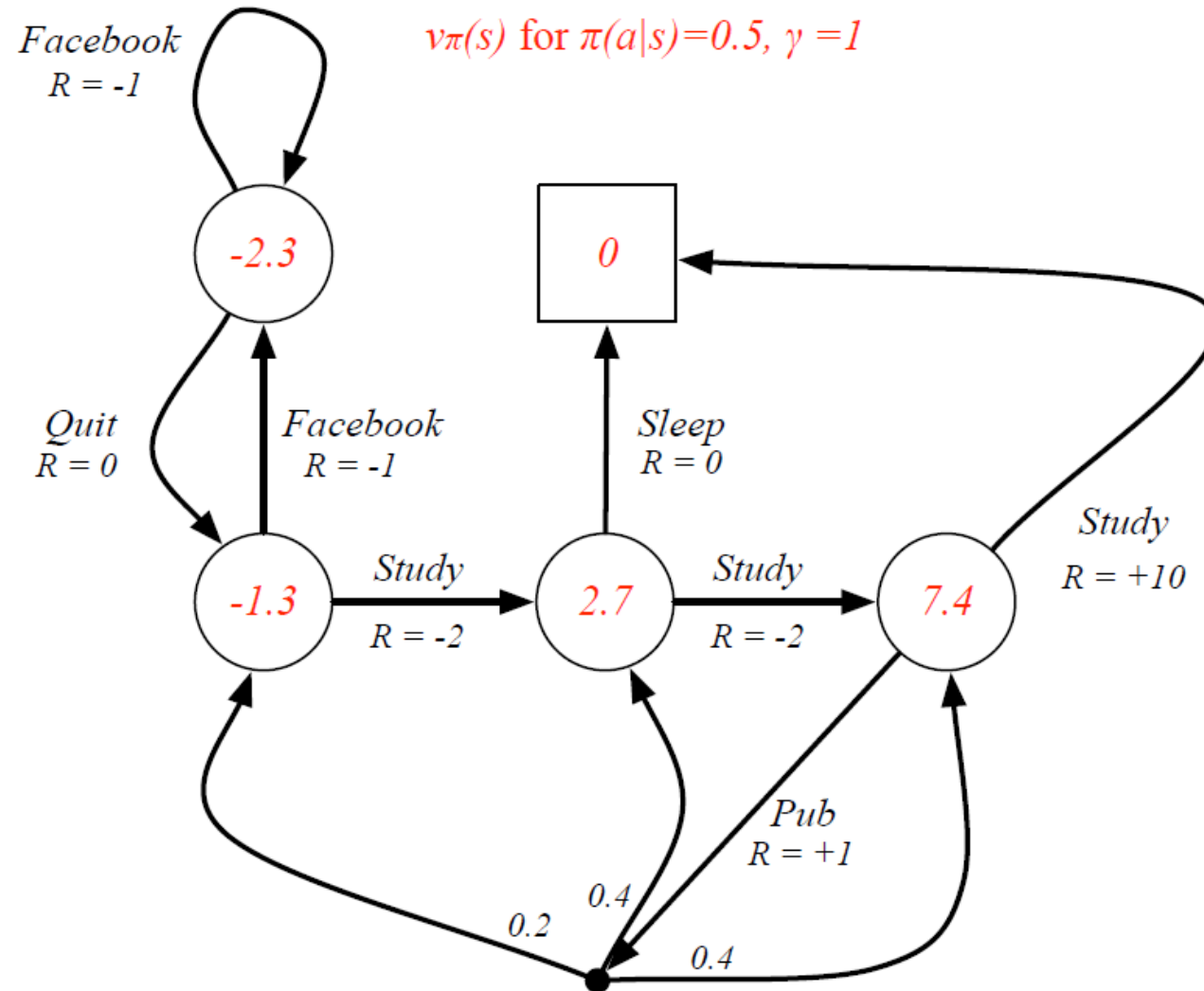
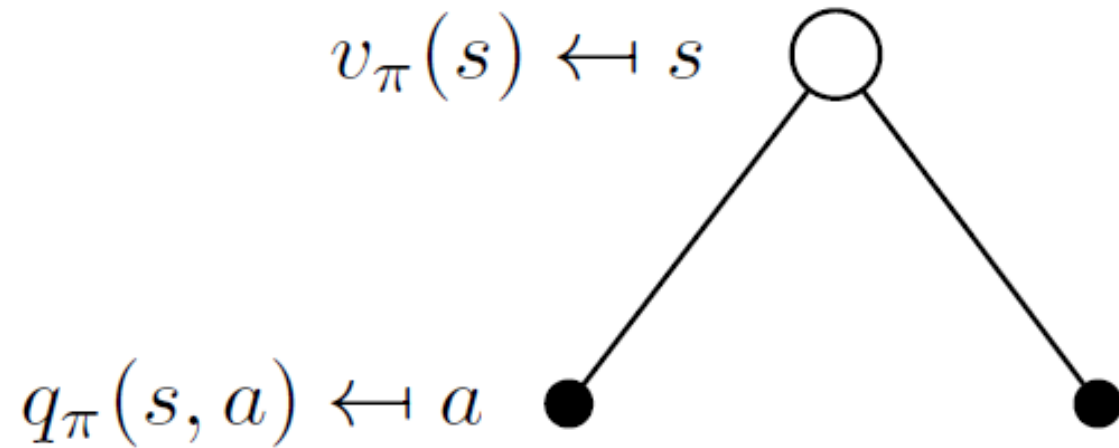$$v_\pi(s) = \mathbb{E}_\pi\left[G_t \mid S_t = s\right]$$

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

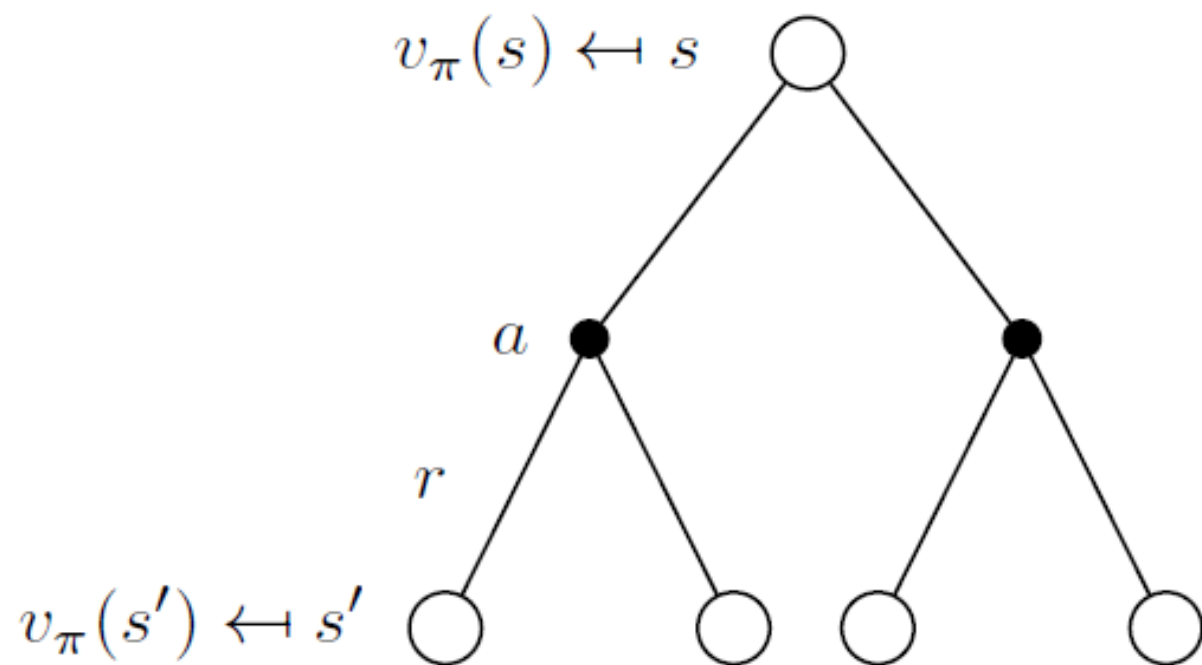$$q_\pi(s, a) = \mathbb{E}_\pi\left[G_t \mid S_t = s, A_t = a\right]$$
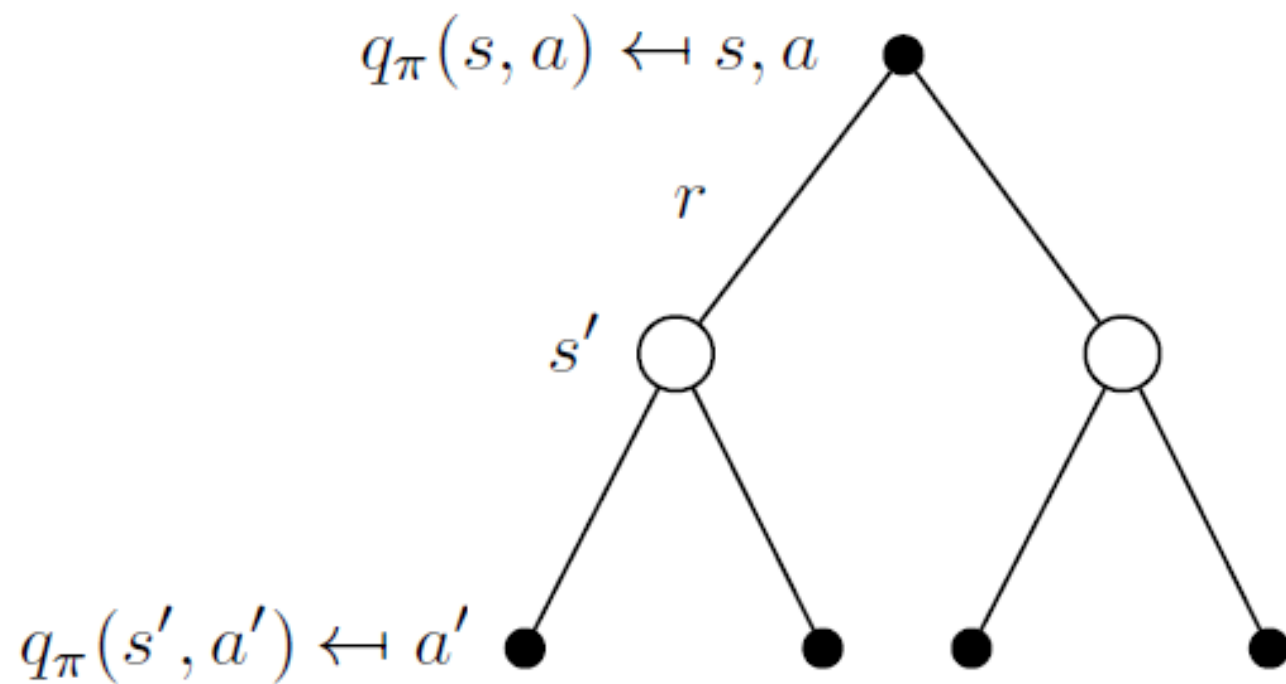
# State-Value Function for Student MDP

# Backup Diagram for $v_\pi$ and $q_\pi$

$$v_\pi(s) \leftarrowtail s$$

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

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

$$v_\pi(s) \leftharpoondown s$$
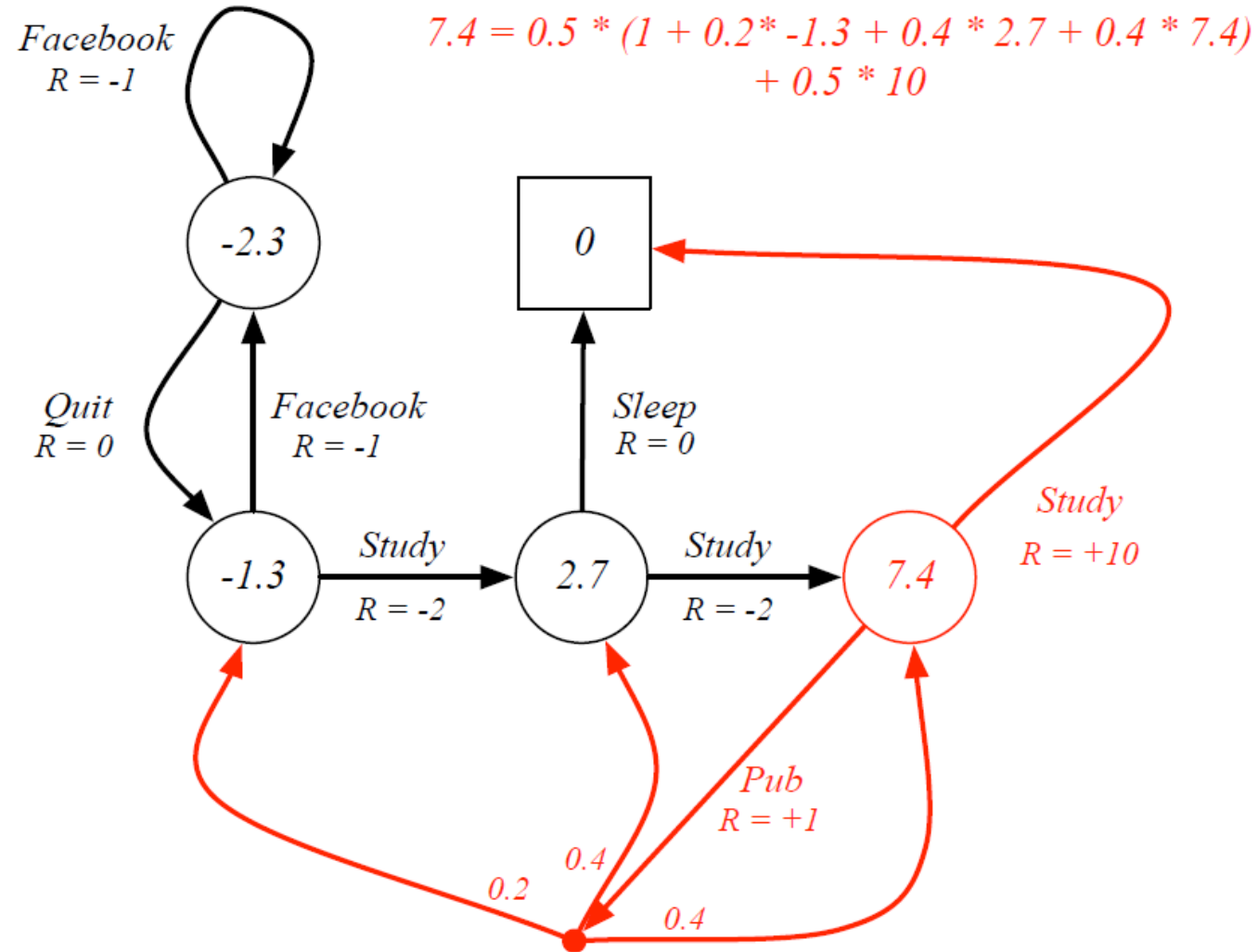
$$a$$

$$r$$

$$v_\pi(s') \leftharpoondown s'$$

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

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

# Bellman Expectation Equation for Student MDP



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

Facebook
R = -1

-2.3

0

Quit
R = 0

Facebook
R = -1

Sleep
R = 0

Study
R = +10

-1.3

Study

2.7

Study

7.4

R = -2

R = -2

Pub
R = +1

0.4

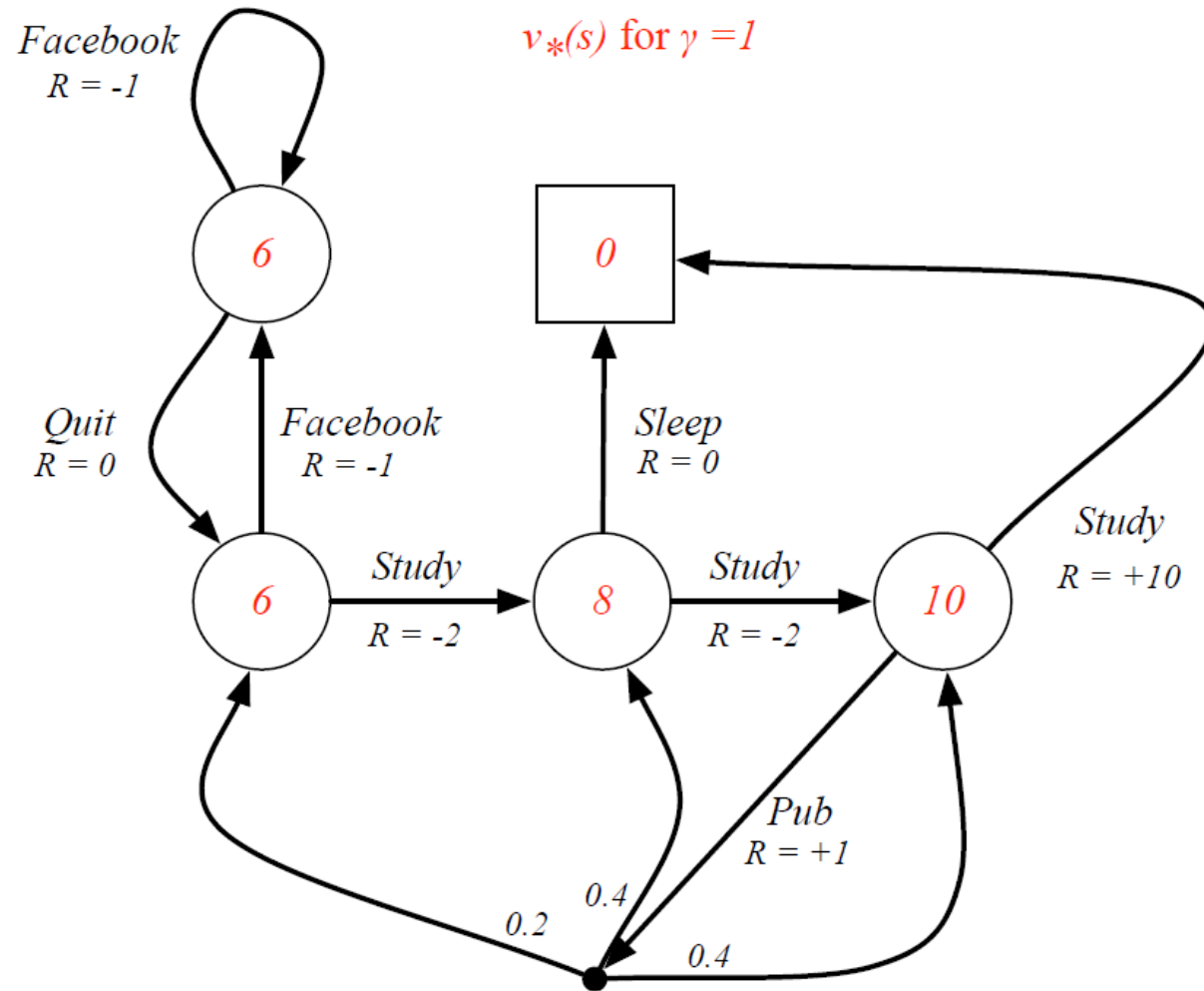0.2

0.4

# Optimal Value Function

## Definition

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

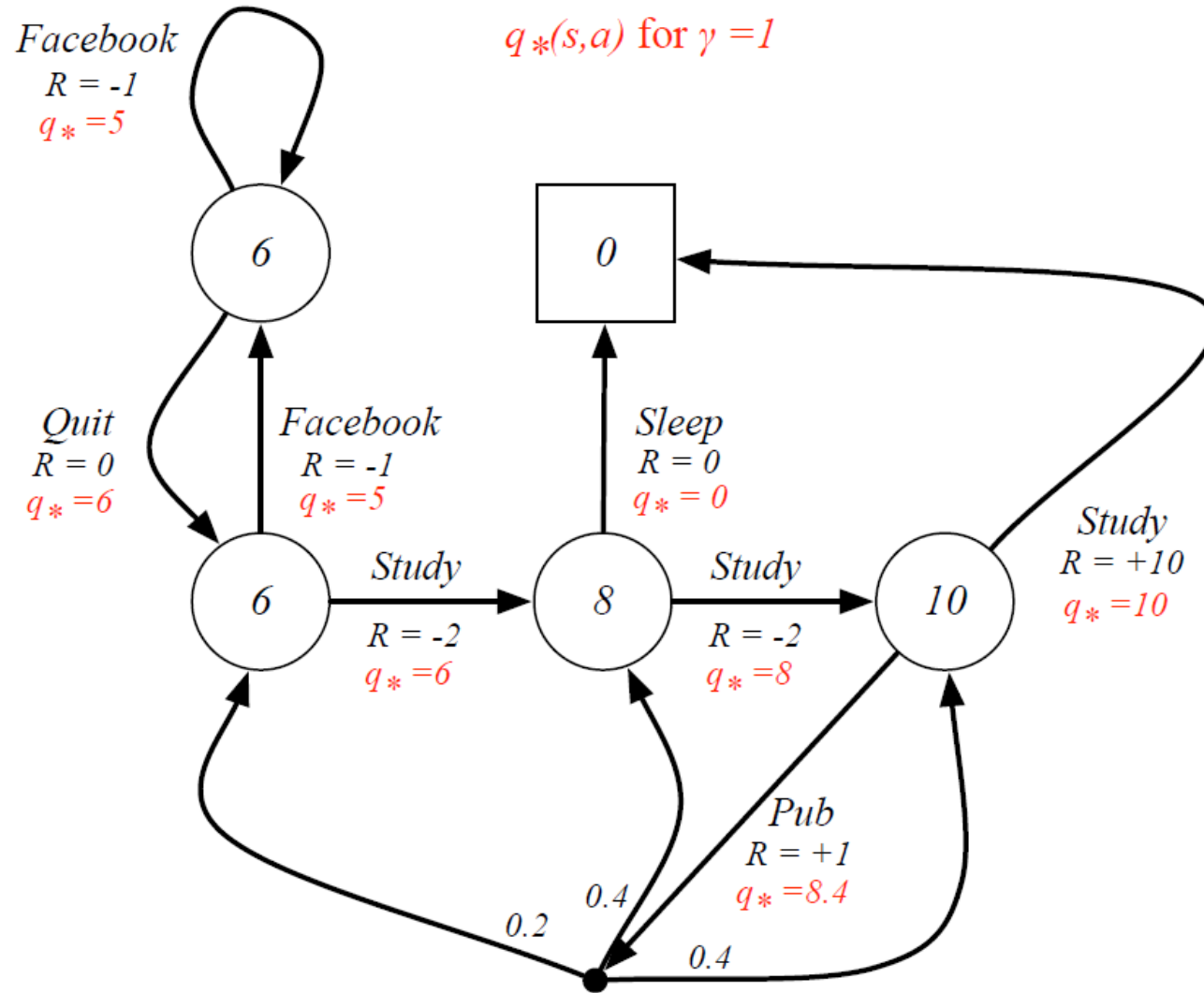$$v_*(s) = \max_{\pi} v_{\pi}(s)$$

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

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

# Optimal Value Function for Student MDP

# Optimal Action-Value Function for Student MDP



Facebook
R = -1
$q_* = 5$

$q_*(s,a)$ for $\gamma = 1$

0

6

Quit
R = 0
$q_* = 6$

Facebook
R = -1
$q_* = 5$

Sleep
R = 0
$q_* = 0$

Study
R = +10
$q_* = 10$

6

Study
R = -2
$q_* = 6$

8

Study
R = -2
$q_* = 8$

10

Pub
R = +1
$q_* = 8.4$

0.4

0.2

0.4

# Reference

- Davlid Silver, Lecture 2: Markov Decision Processes, Reinforcement Learning (https://www.youtube.com/watch?v=lfHX2hHRMVQ&list=PLqYmG7hTraZDM-OYHWgPebj2MfCFzFObQ&index=2)

- Chapter 3, Richard S. Sutton and Andrew G. Barto, "Reinforcement Learning: An Introduction," 2$^{nd}$ edition, Nov. 2018