



Reinforcement Learning - Lecture 2

Profa. Dra. Esther Luna Colombini
esther@ic.unicamp.br

Prof. Dr. Alexandre Simoes
alexandre.simoes@unesp.br



LaRoCS – Laboratory of Robotics and Cognitive Systems



- Why MDPs?
- Markov Processes
- Markov Reward Processes
 - Discounted Rewards
- Markov Decision Processes
 - How the Markov Decision Process is defined

- These slides were built upon David Silver's Lecture notes on Reinforcement Learning

- **Markov Decision Processes** formally describe an environment for reinforcement learning
- The environment is fully observable
 - The current state completely characterizes the process
- Almost all RL problems can be formalized as MDPs
 - Optimal control primarily deals with continuous MDPs
 - Partially observable problems can be converted into MDPs
 - Bandits are MDPs with one state

- The **Markov property** refers to the **memoryless** property of a stochastic process
 - A stochastic process has the Markov property if
 - the conditional probability distribution $\mathbb{P}[S_{t+1}|S_t]$ of future states S_{t+1} of the process (conditional on both past and present states) depends only upon the **present state**, not on the sequence of events that preceded it
 - 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
 - The state is a sufficient statistic of the future
 - **Stationary Assumption**
 - Transition probabilities are independent of time (t)
 - What if they are not?
 - **A state S_t is Markov** if and only if

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

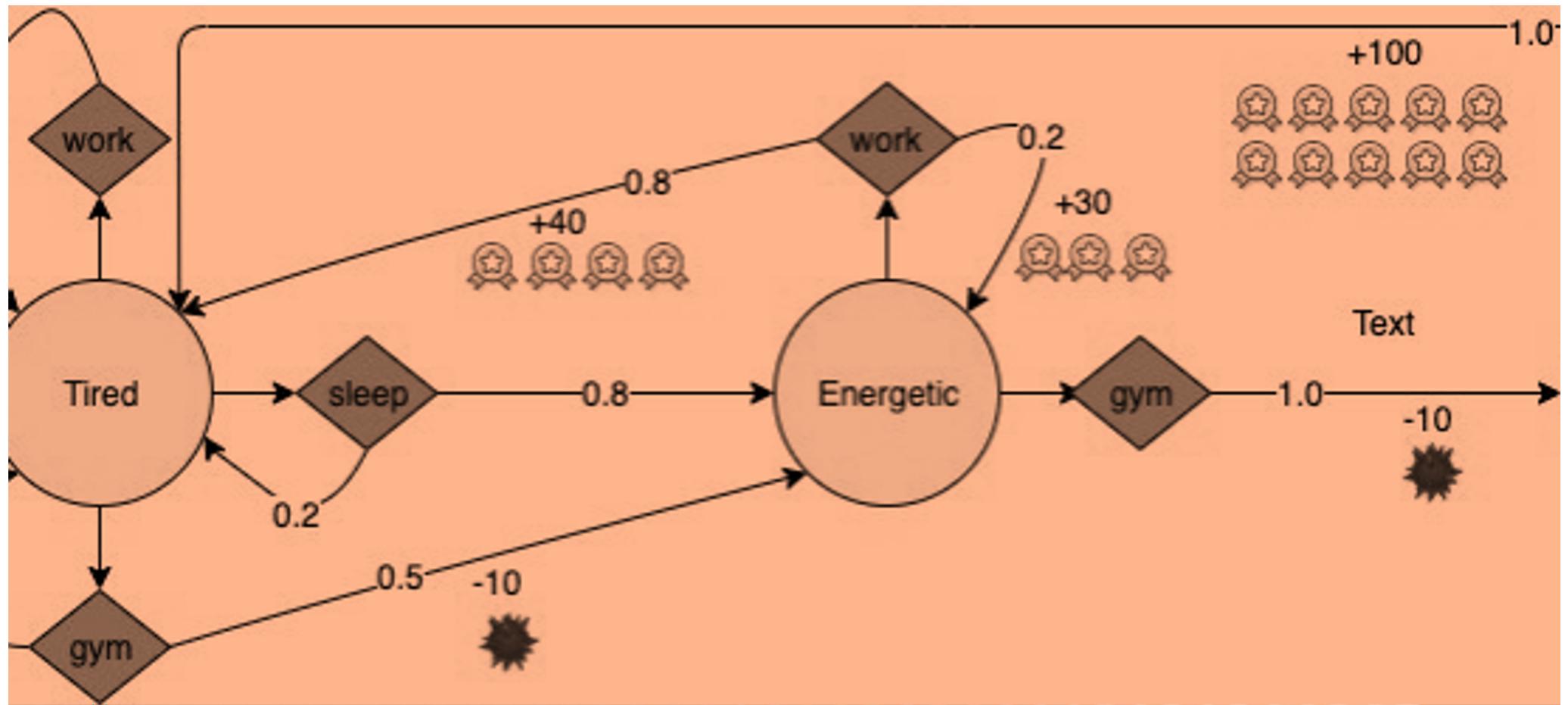
- For a Markov state s and successor state s' , the **state transition probability** is defined by

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

- State transition matrix \mathcal{P} defines transition probabilities from all states s to all successor states s'

$$\mathcal{P} = \begin{matrix} & \text{State to} \\ \text{State from} & \left[\begin{matrix} \mathcal{P}_{11} & \dots & \mathcal{P}_{1n} \\ \vdots & & \vdots \\ \mathcal{P}_{n1} & \dots & \mathcal{P}_{nn} \end{matrix} \right] \end{matrix}$$

With $\sum_{i=1}^n \mathcal{P}_{mi} = 1$, for all m in $[1, n]$



Markov Processes



- A Markov process is a memoryless **random process**, i.e. a sequence of random states S_1, S_2, \dots that are Markov
- A **Markov Process** (or Markov Chain) is a tuple $\langle S, \mathcal{P} \rangle$
 - S is a finite set of states
 - \mathcal{P} is a state transition probability matrix

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

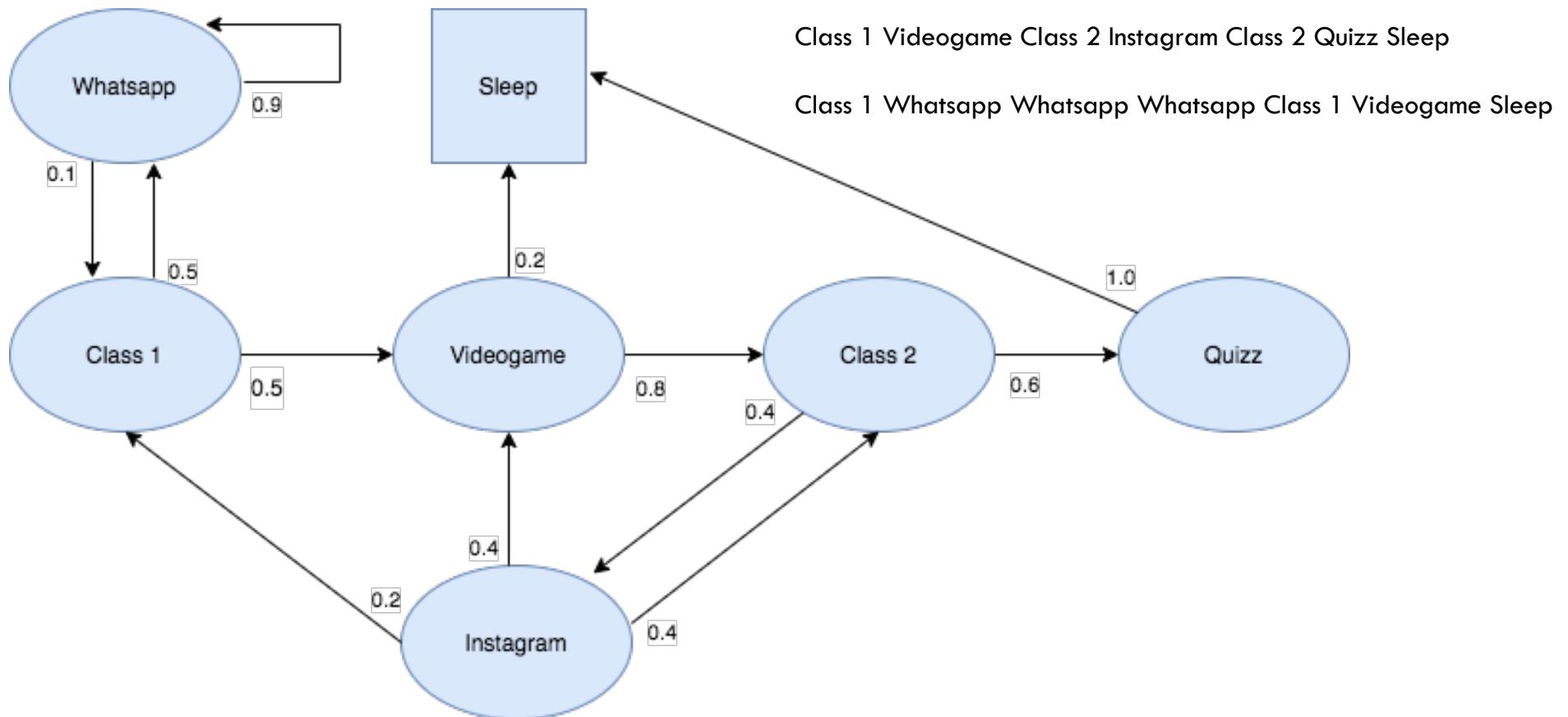


Markov Process

9

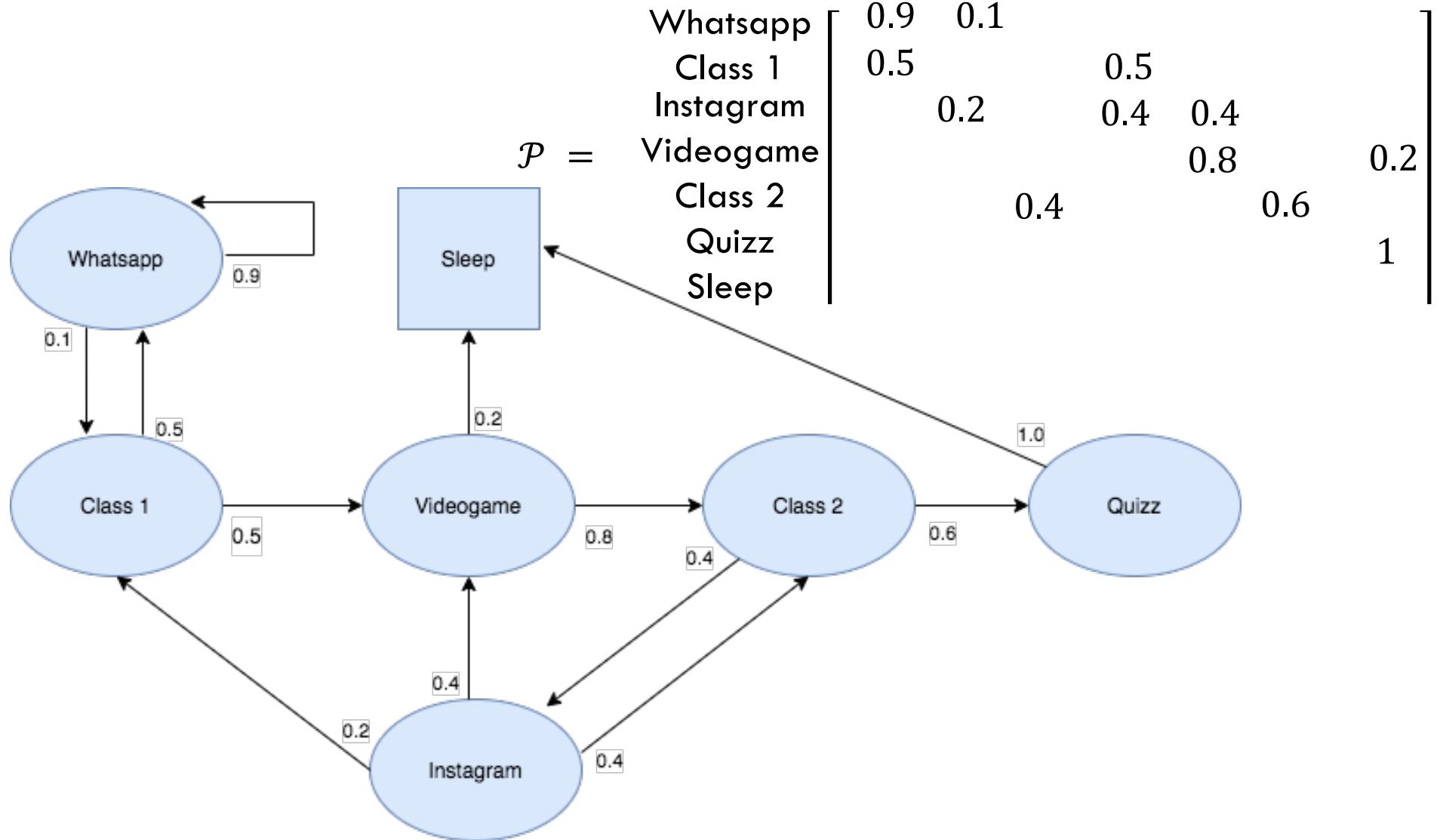
Studying at home Markov Chain

- Let's sample episodes from the Markov Chain starting from $S_1 = \text{Class 1}$



Markov Process

Studying at home Markov Chain

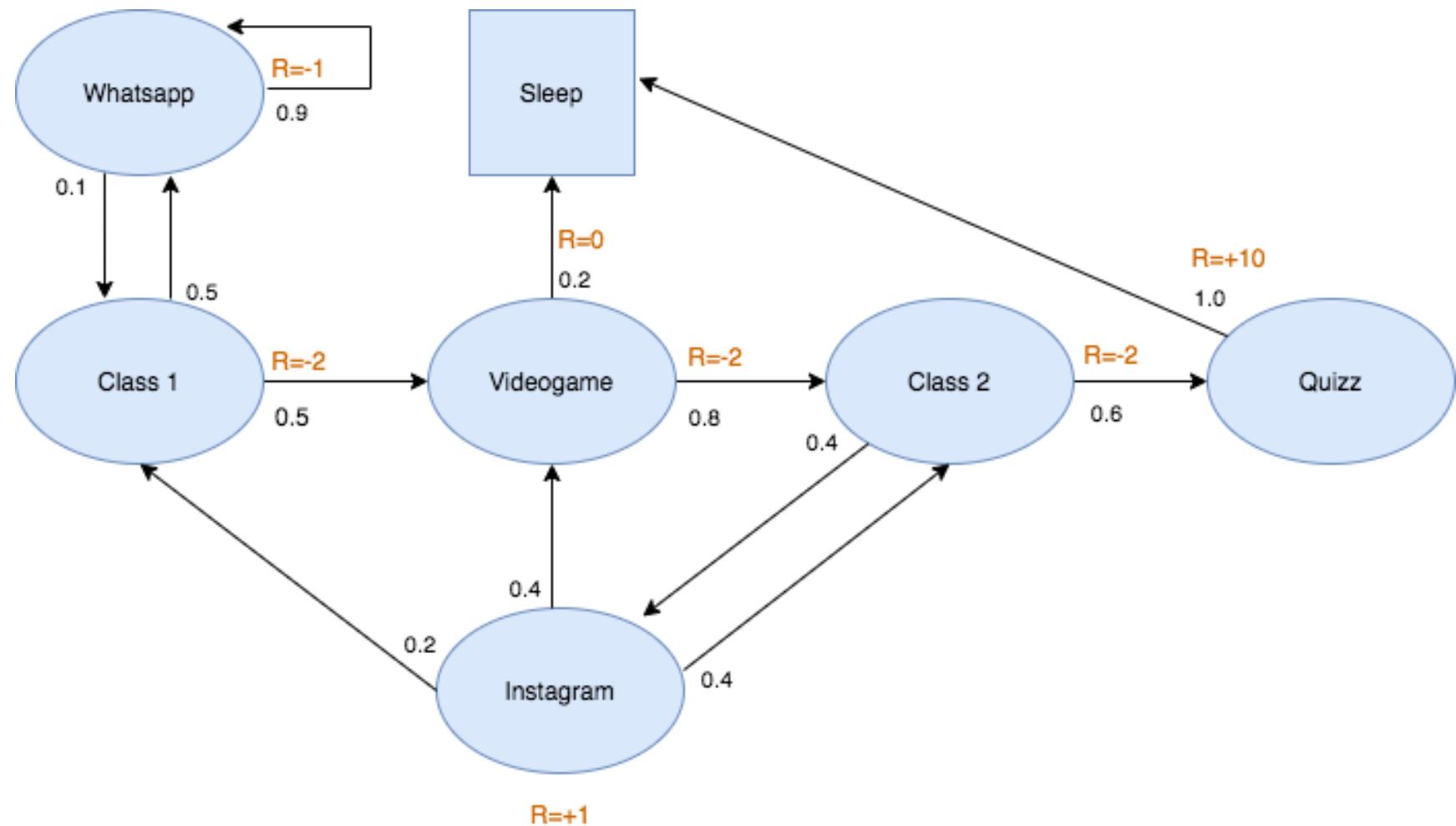


- A Markov reward process is a Markov chain with values
- A **Markov Reward Process** (or MRP) is a tuple $\langle S, \mathcal{P}, \mathcal{R}, \gamma \rangle$
 - S is a finite set of states
 - \mathcal{P} is a state transition probability matrix
$$\mathcal{P}_{ss'} = \mathbb{P}[S_{t+1} = s' | S_t = s]$$
 - \mathcal{R} is a Reward function $\mathcal{R}_S = \mathbb{E}[R_{t+1} | S_t = s]$
 - γ is the discount factor, $\gamma \in [0,1]$
- No actions

Markov Reward Process

12

Studying at home Markov Reward Process



- The **return** G_t is discounted sum of rewards from time step t to horizon

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

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

- Most Markov reward and decision processes are discounted because:
 - 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 behaviour shows preference for immediate reward
 - If episode lengths are always finite, can use $\gamma = 1$

- The state **Value Function** $v(s)$ is the long-term quality of state s , i.e., is the expected return starting from state s

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

- Let's compute the static value for each sample drawn from our MRP
- Start from $S_1 = \text{Class 1}$ and consider $\gamma = 1/2$
 - $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$
 - Class 1 Videogame Class 2 Quizz Sleep
 - $v_1 = -2 + \left(\frac{1}{2} * -2\right) + \left(\frac{1}{4} * -2\right) + \left(\frac{1}{8} * 10\right) = -2.25$
 - Class 1 Videogame Class 2 Instagram Class 2 Quizz Sleep
 - $v_1 = -2 + \left(\frac{1}{2} * -2\right) + \left(\frac{1}{4} * 1\right) + \left(\frac{1}{8} * -2\right) + \left(\frac{1}{16} * -2\right) + \left(\frac{1}{32} * 10\right) = -2.8125$
 - Class 1 Whatsapp Whatsapp Whatsapp Class 1 Videogame Sleep
 - $v_1 = -2 + \left(\frac{1}{2} * -1\right) + \left(\frac{1}{4} * -1\right) + \left(\frac{1}{8} * -1\right) + \left(\frac{1}{16} * -2\right) + \left(\frac{1}{32} * 0\right) = -3$

- Do you think we have defined a good reward function?

- Class 1 Videogame Class 2 Quizz Sleep

- $v_1 = -2 + \left(\frac{1}{2} * -2\right) + \left(\frac{1}{4} * -2\right) + \left(\frac{1}{8} * 10\right) = -2.25$

- Class 1 Videogame Class 2 Instagram Class 2 Quizz Sleep

- $v_1 = -2 + \left(\frac{1}{2} * -2\right) + \left(\frac{1}{4} * 1\right) + \left(\frac{1}{8} * -2\right) + \left(\frac{1}{16} * -2\right) + \left(\frac{1}{32} * 10\right) = -2.8125$

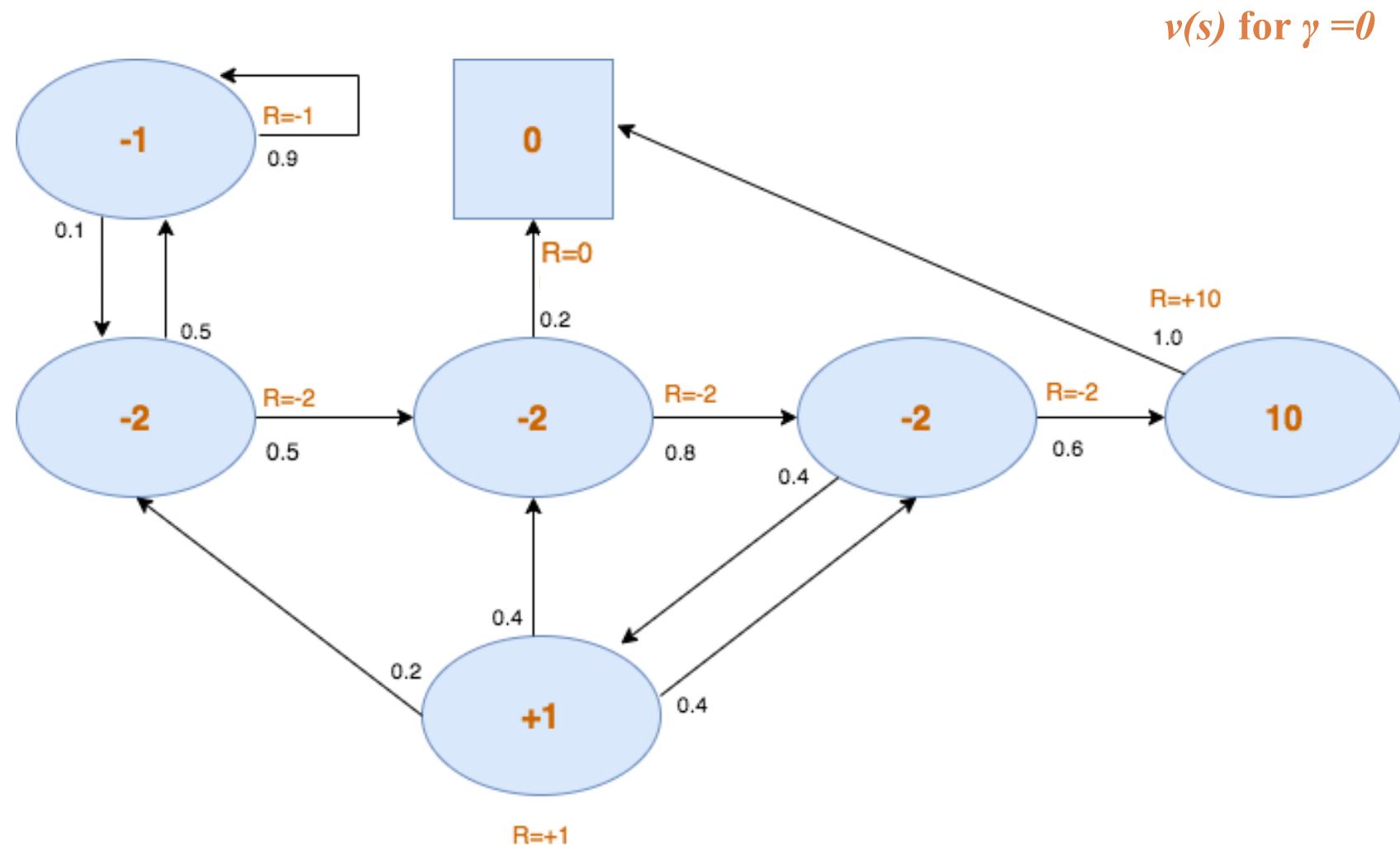
- Class 1 Whatsapp Whatsapp Whatsapp Class 1 Videogame Sleep

- $v_1 = -2 + \left(\frac{1}{2} * -1\right) + \left(\frac{1}{4} * -1\right) + \left(\frac{1}{8} * -1\right) + \left(\frac{1}{16} * -2\right) + \left(\frac{1}{32} * 0\right) = -3$

State-Value Function for Student MRP

18

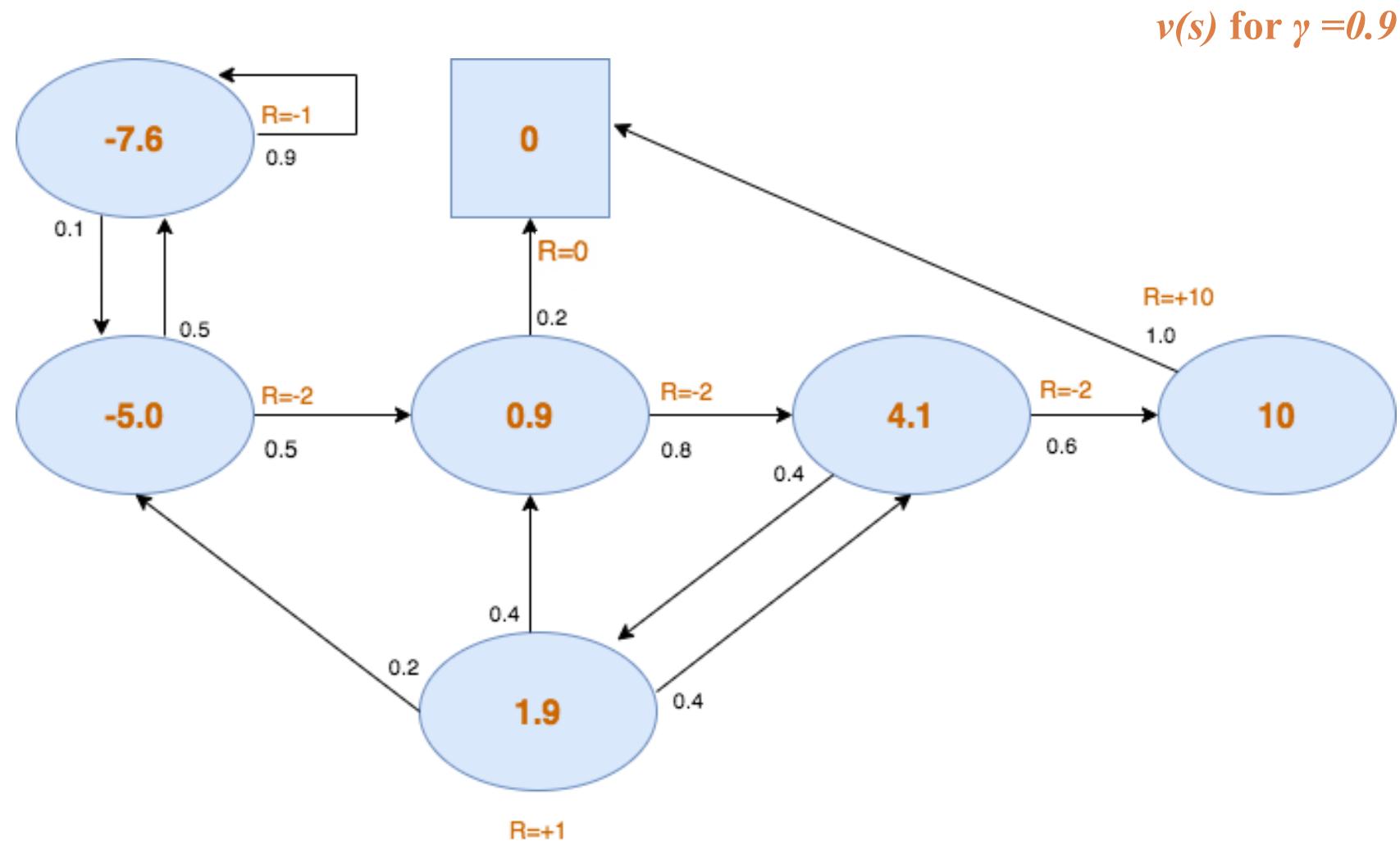
Studying at home MRP



State-Value Function for Student MRP

19

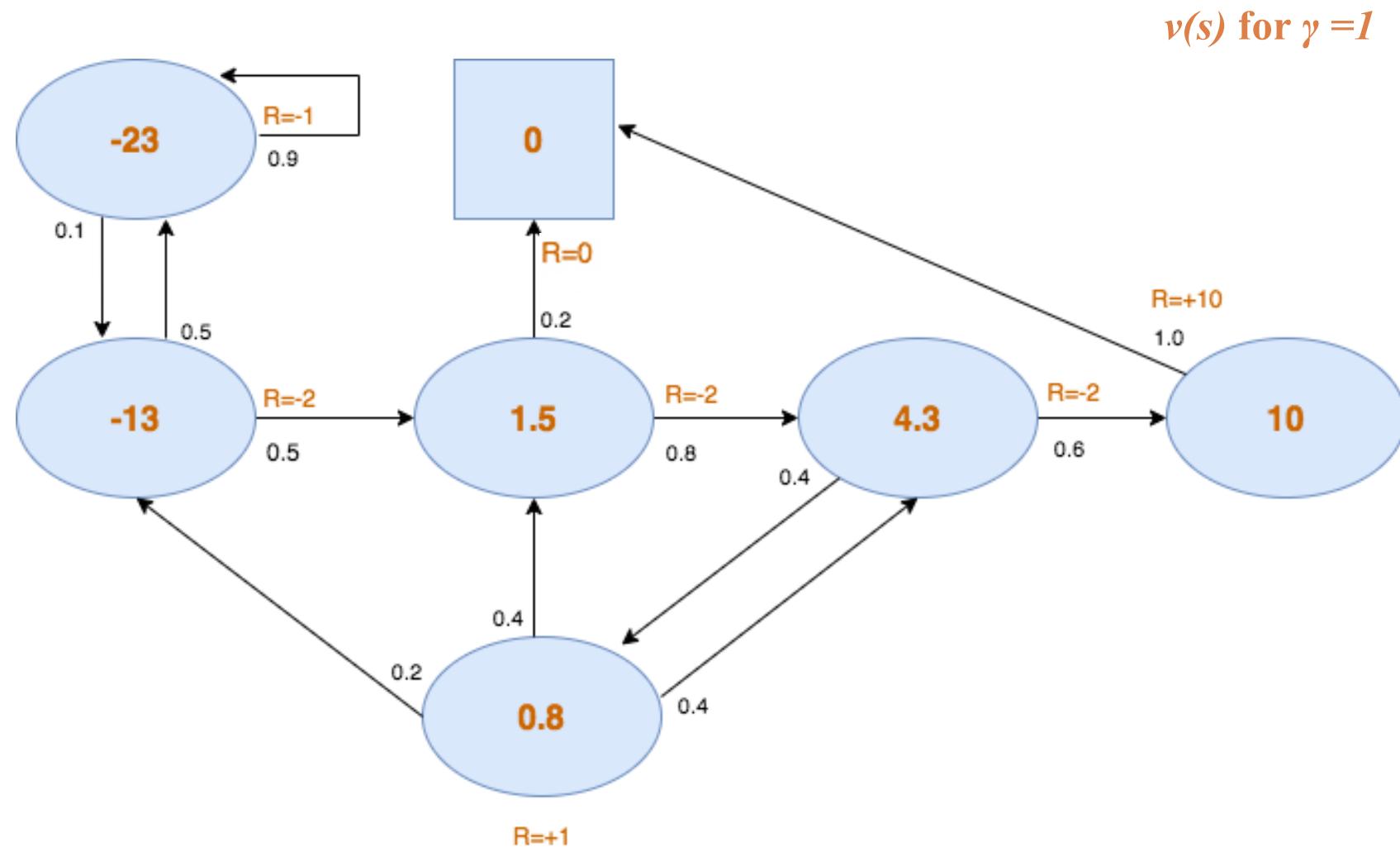
Studying at home MRP



State-Value Function for Student MRP

20

Studying at home MRP

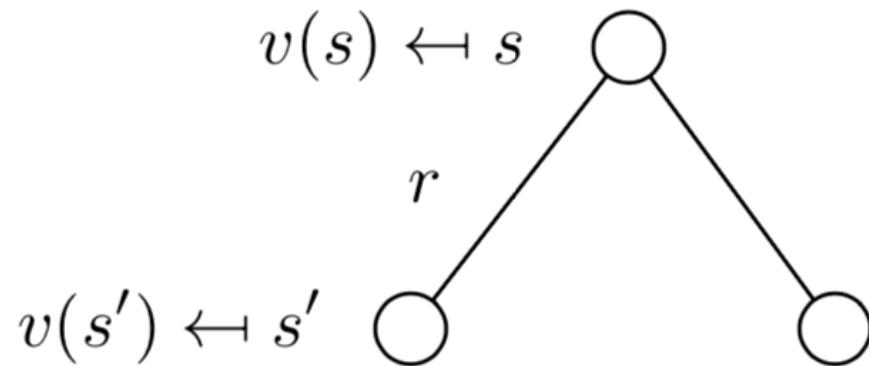


- The value function can be decomposed into two parts:
 - The immediate reward R_{t+1}
 - The discounted value of successor state $\gamma v(S_{t+1})$

- $v(s) = \mathbb{E}[G_t | S_t = s]$
 $= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s]$
 $= \mathbb{E}[R_{t+1} + \gamma(R_{t+2} + \dots) | S_t = s]$
 $= \mathbb{E}[R_{t+1} + \gamma G_{t+1} | S_t = s]$
 $= \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) | S_t = s]$

- Let's have a look at the one-step backup diagram

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

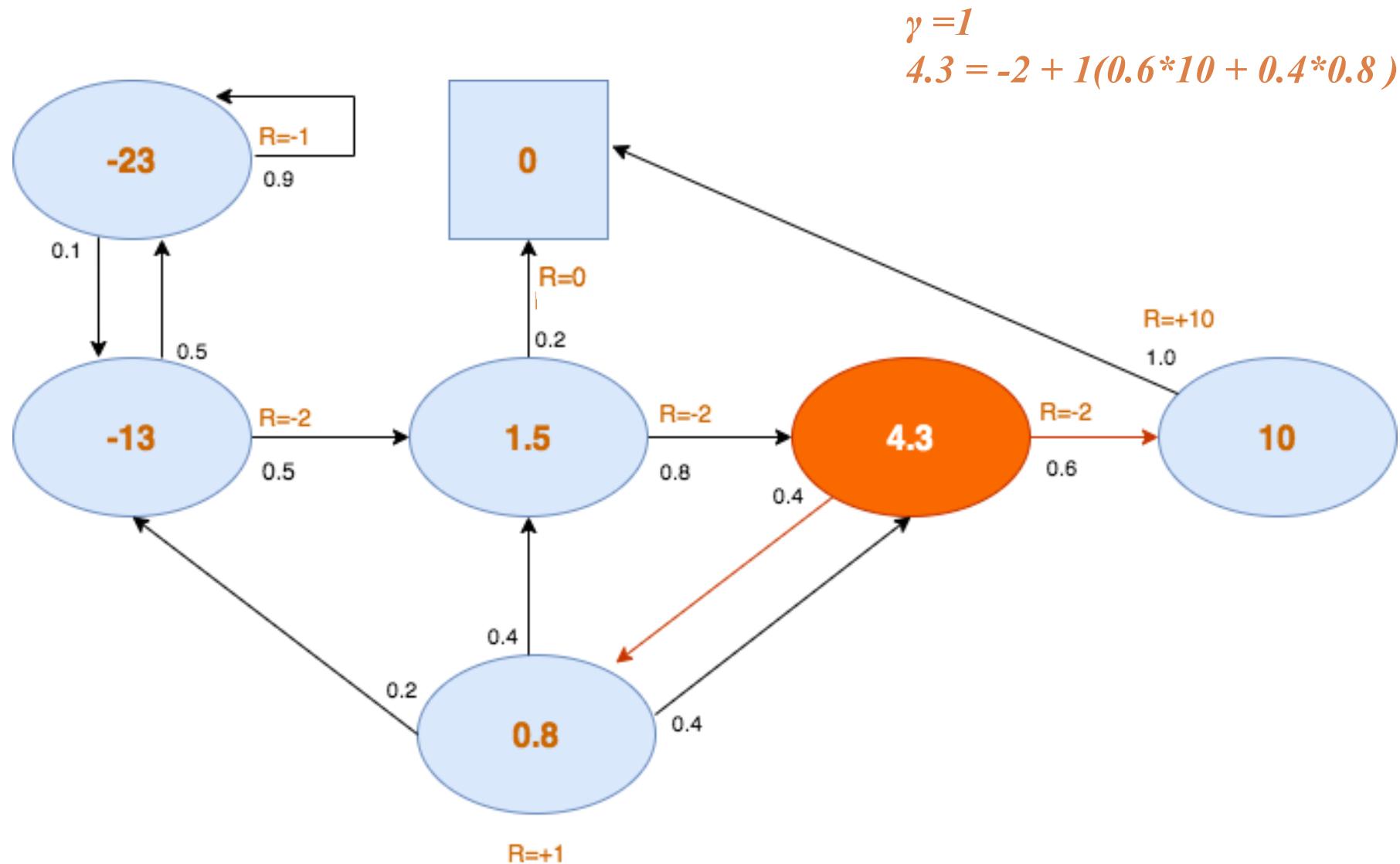


- $v(s) = R_S + \gamma \sum_{s' \in S} \mathcal{P}_{ss'} v(s')$

State-Value Function for Student MRP

23

Studying at home MRP



- The Bellman equation can be expressed concisely using matrices,
 - $v = R + \gamma Pv$
- where v is a column vector with one entry per state

$$\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} \mathcal{P}_{11} & \cdots & \mathcal{P}_{1n} \\ \vdots & & \vdots \\ \mathcal{P}_{n1} & \cdots & \mathcal{P}_{nn} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$$

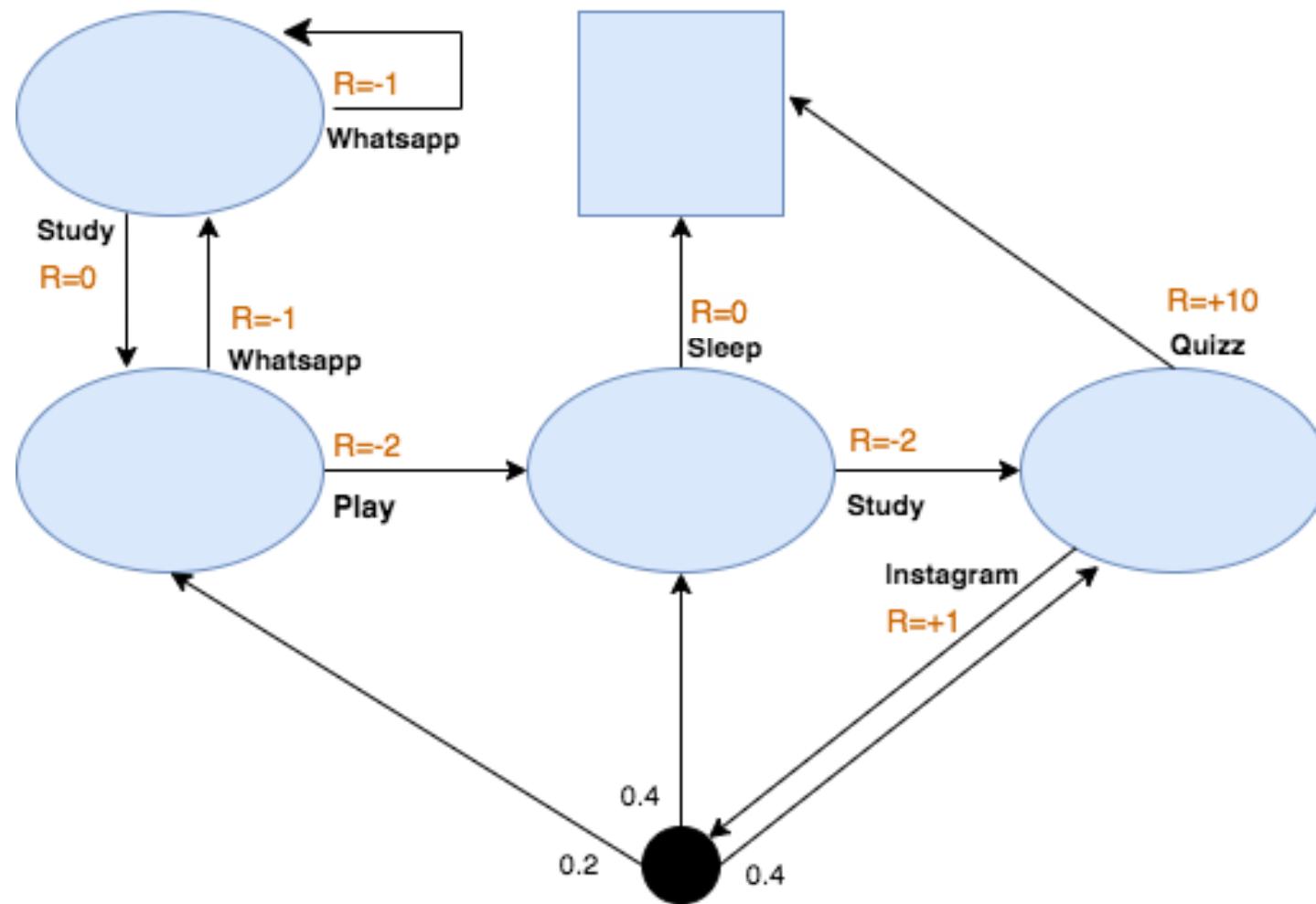
- The Bellman equation is a linear equation It can be solved directly:
 - $v = R + \gamma Pv$
 - $(I - \gamma P)v = R$
 - $v = (I - \gamma P)^{-1} R$
- Computational complexity is $O(n^3)$ for n states
- Direct solution only possible for small MRPs
- There are many iterative methods for large MRPs, e.g.
 - Dynamic programming
 - Monte-Carlo evaluation
 - Temporal-Difference learning

- A Markov decision process is a Markov chain with values
- A **Markov Decision Process** (or MDP) is a tuple $\langle S, \mathcal{P}, \mathcal{A}, \mathcal{R}, \gamma \rangle$
 - S is a finite set of states
 - A is a finite set of actions
 - \mathcal{P} is a state transition probability matrix
$$\mathcal{P}_{ss'}^{\textcolor{brown}{a}} = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = \textcolor{brown}{a}]$$
 - \mathcal{R} is a Reward function $\mathcal{R}_S = \mathbb{E}[R_{t+1} | S_t = s]$
 - γ is the discount factor, $\gamma \in [0,1]$

Markov Decision Process

27

Studying at home MDP



- A **policy** π is a distribution over actions given states,

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

- A policy fully defines the behaviour of an agent
- MDP policies depend on the current state (not the history)
- i.e. Policies are stationary (time-independent),

$$A_t \sim \pi(\cdot | S_t), \forall t > 0$$

- Given an MDP $\langle S, \mathcal{P}, A, \mathcal{R}, \gamma \rangle$ and a policy π
- The state sequence $S_1, S_2 \dots$ is a Markov Process $\langle S, \mathcal{P}^\pi \rangle$
- The state and reward sequence $S_1, R_1, S_2, R_2 \dots$ is a Markov Reward Process $\langle S, \mathcal{P}^\pi, \mathcal{R}^\pi, \gamma \rangle$
- where

$$\mathcal{P}_{ss'}^\pi, \sum_{\substack{a \in A \\ \infty}}^{\infty} \pi(a|s) \mathcal{P}_{ss'}^a,$$
$$R_s^\pi, \sum_{\substack{a \in A}} \pi(a|s) R_s^a$$

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

$$v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s]$$

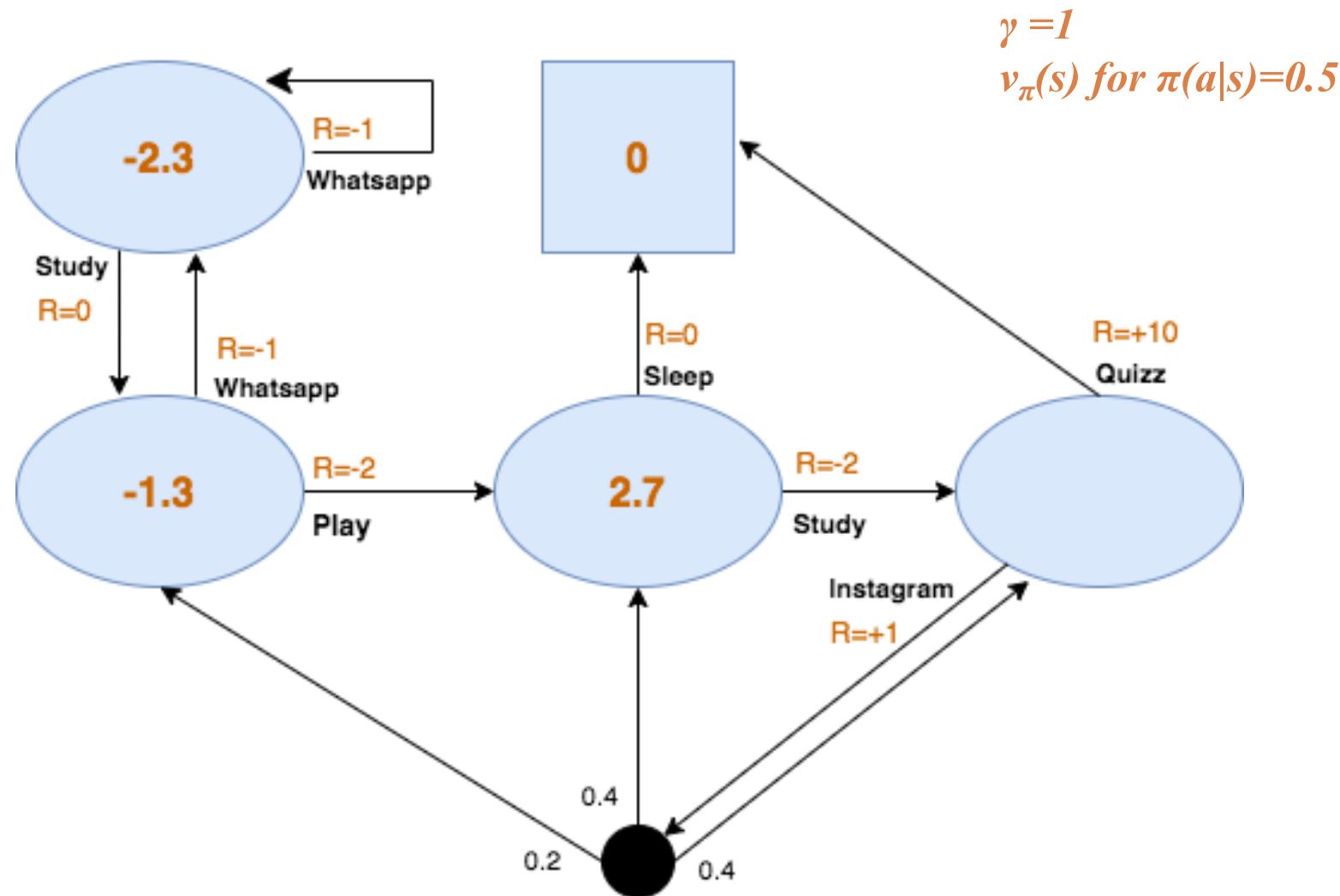
- The **action-value function** $q_\pi(s, a)$ of an MDP is the expected return starting from **state** s , taking **action** a , and then following **policy** π

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

State Value Function with Policy

31

Studying at home MDO



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

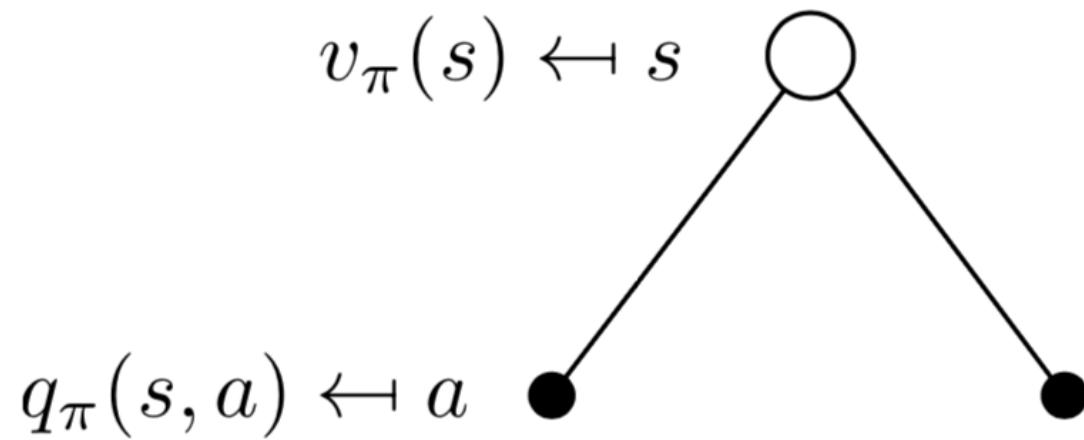
$$v_{\pi}(s) = \mathbb{E}_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s]$$

The action-value function can similarly be decomposed,

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

 Bellman Expectation Equation for V^π

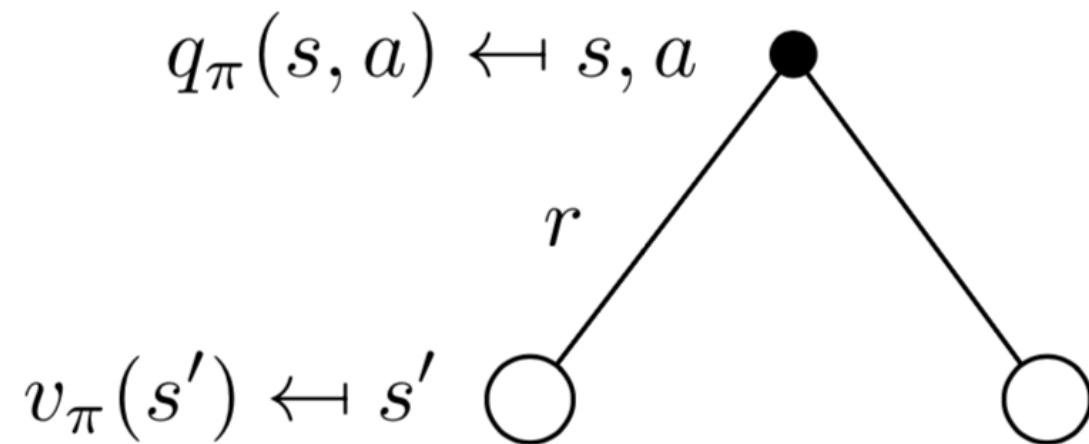
33



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

●●●●● Bellman Expectation Equation for Q^π

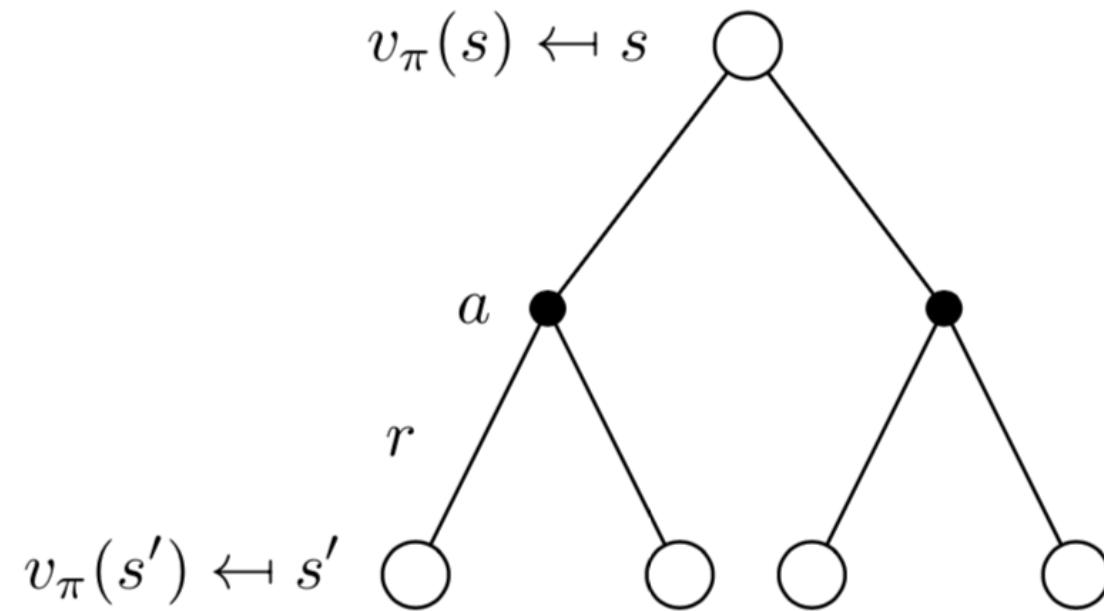
34



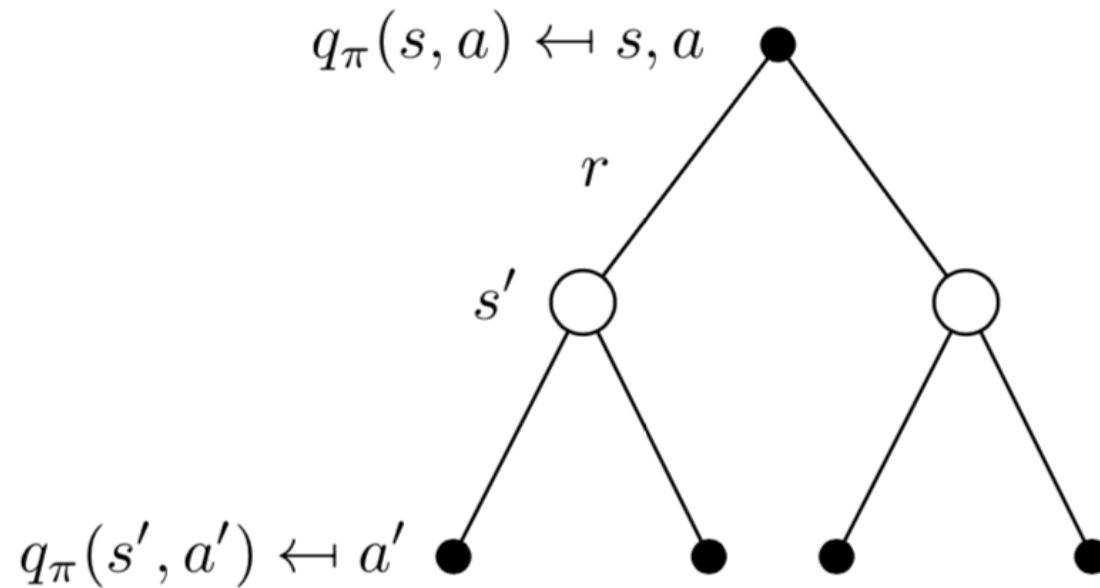
$$q_\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} \mathcal{P}_{ss'}^a v_\pi(s')$$

●●●●● Bellman Expectation Equation for V^π

35



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

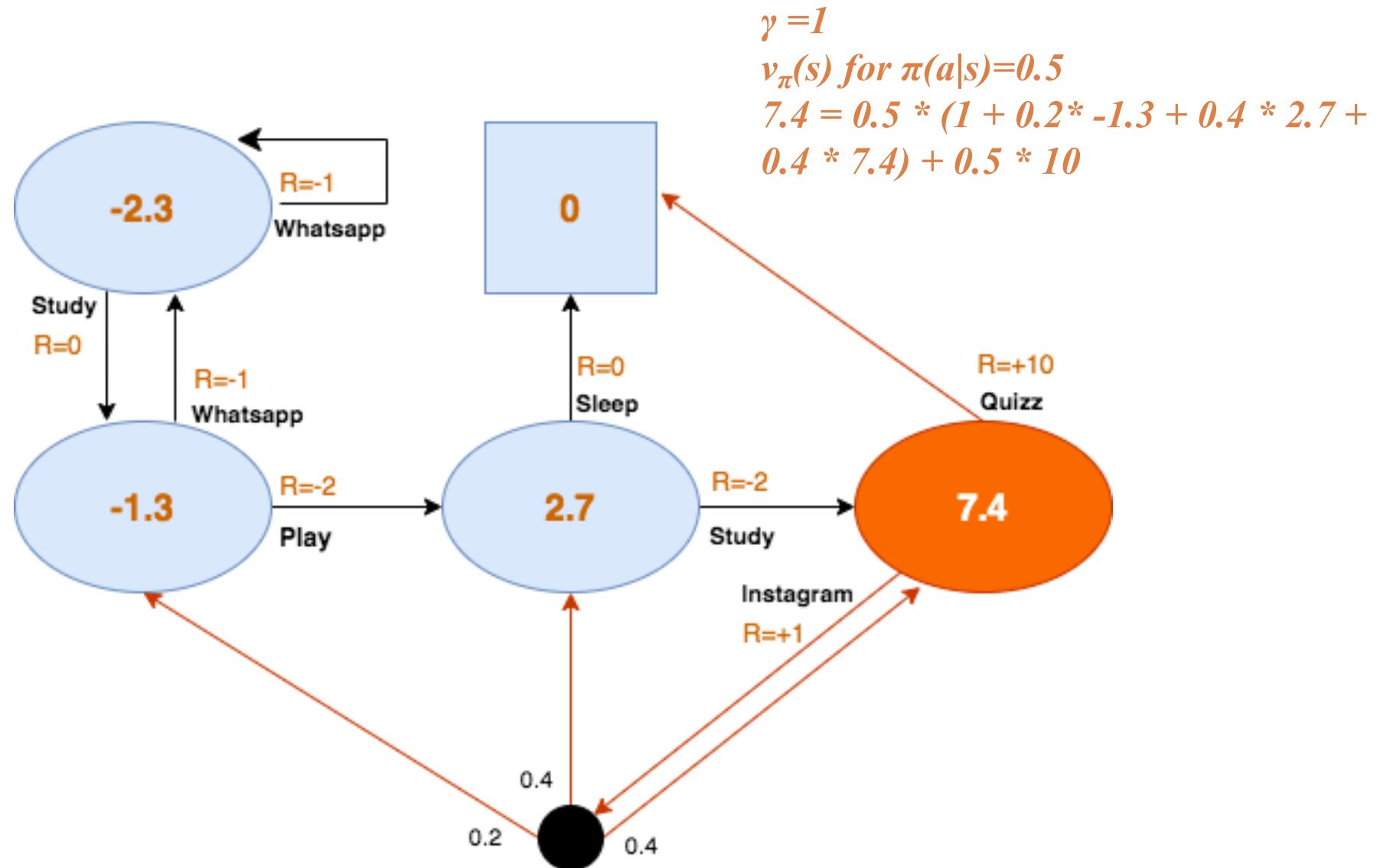


$$q_\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} \mathcal{P}_{ss'}^a \sum_{a' \in A} \pi(a'|s') q_\pi(s', a')$$

Markov Decision Process

37

Studying at home MDP



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

$$v_\pi(s) = R^\pi + \gamma P^\pi$$

- with direct solution

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

- 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

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

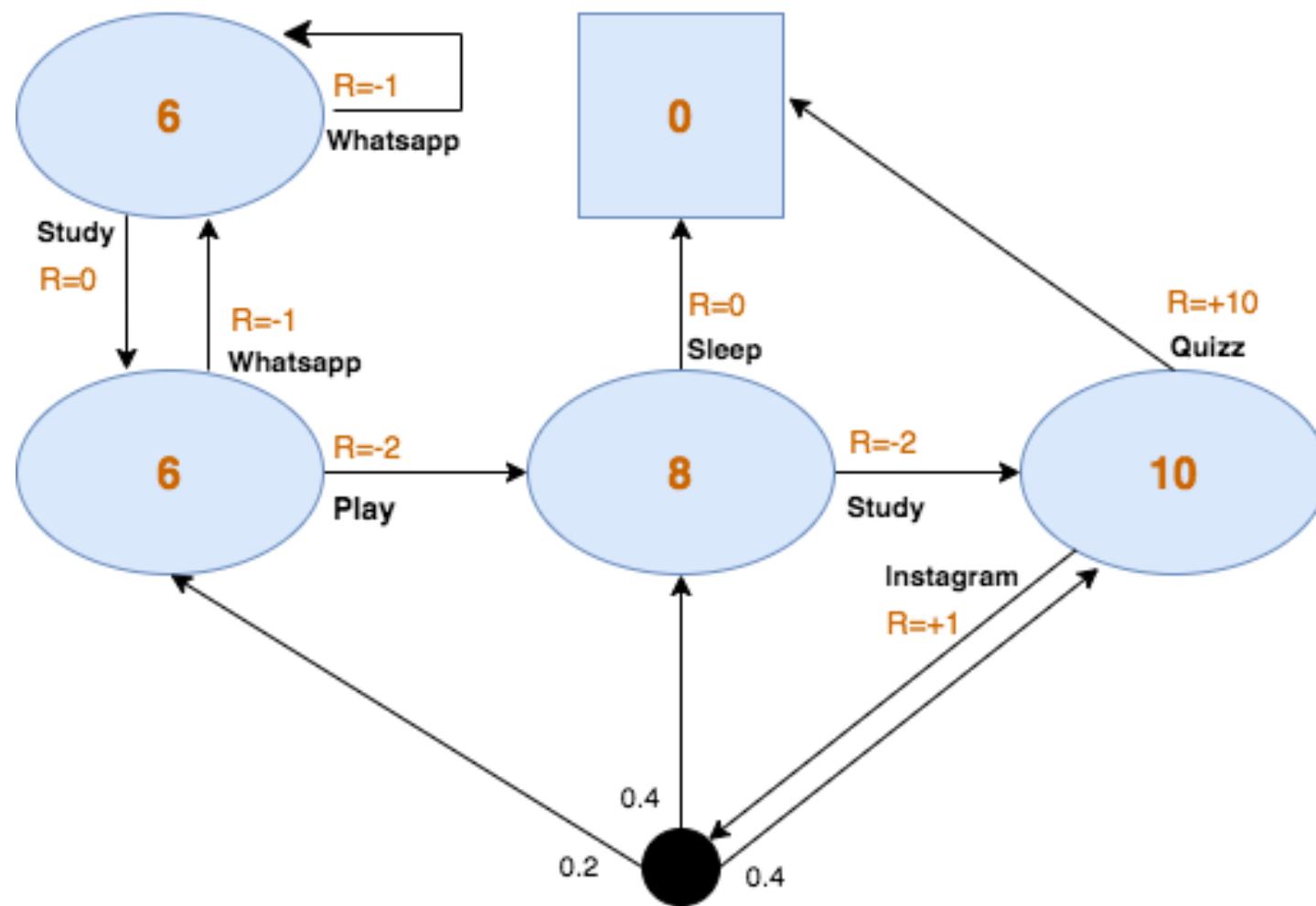
- The optimal value function specifies the best possible performance in the MDP.
- An MDP is “solved” when we know the optimal value function
 -

Markov Decision Process

40

Studying at home MDP

$v_*(s)$ for $\gamma = 1$

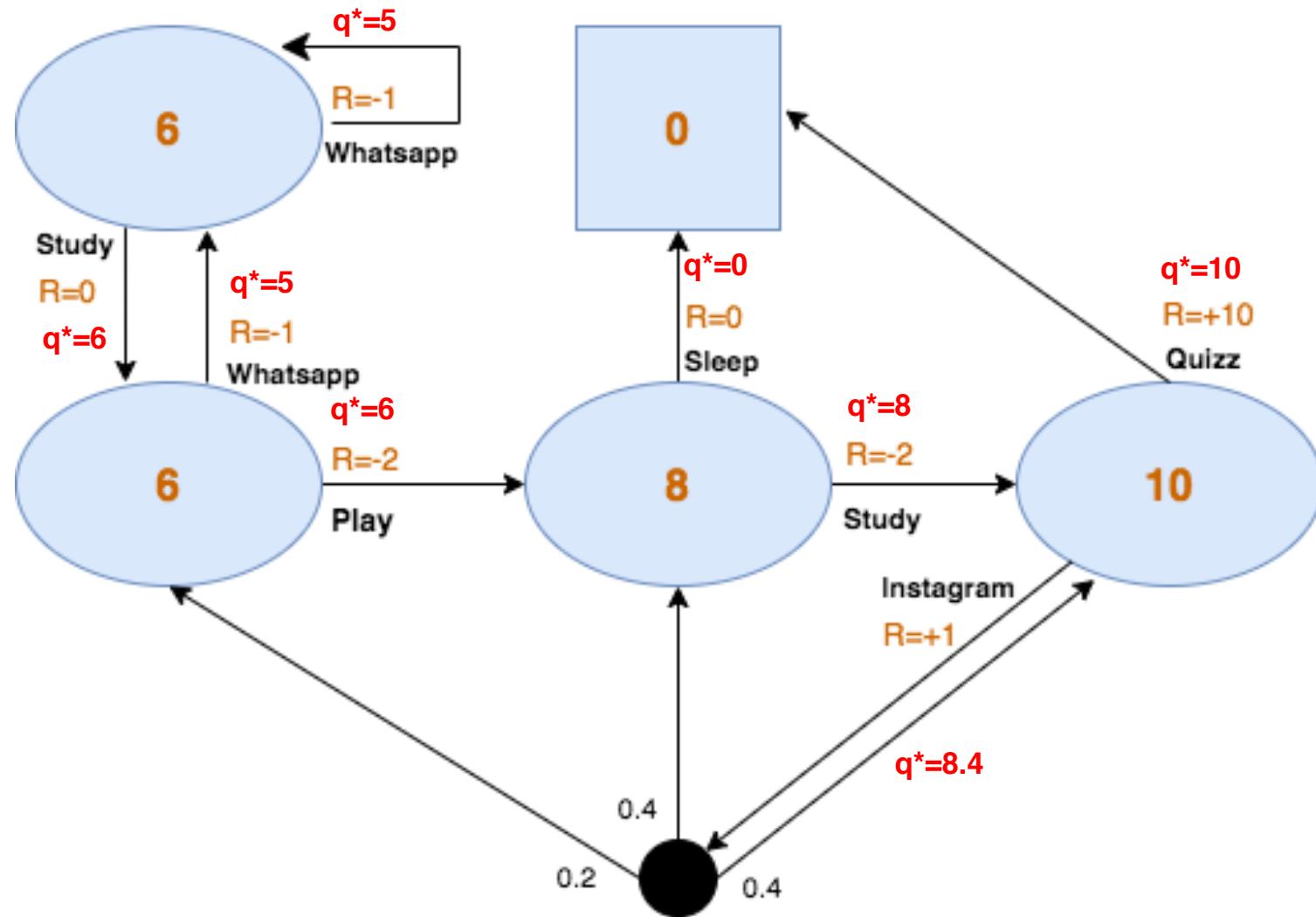


Markov Decision Process

41

Studying at home MDP

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



- Define a partial ordering over policies

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

- Theorem
- For any Markov Decision Process
 - There exists an optimal policy π_* 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,

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

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

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

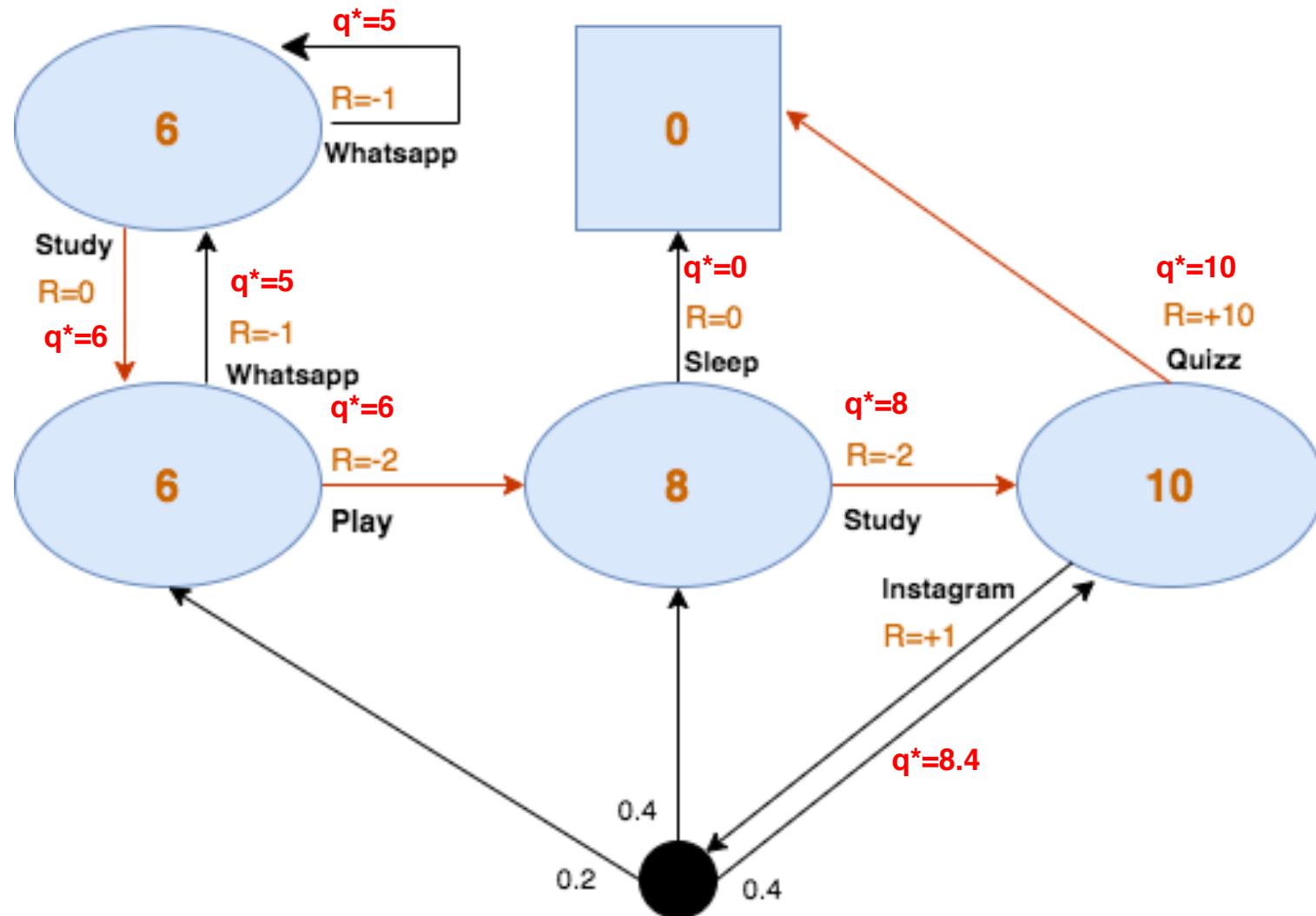
- There is always a deterministic optimal policy for any MDP
- If we know $q_*(s, a)$, we immediately have the optimal policy

Optimal Policy for our MDP

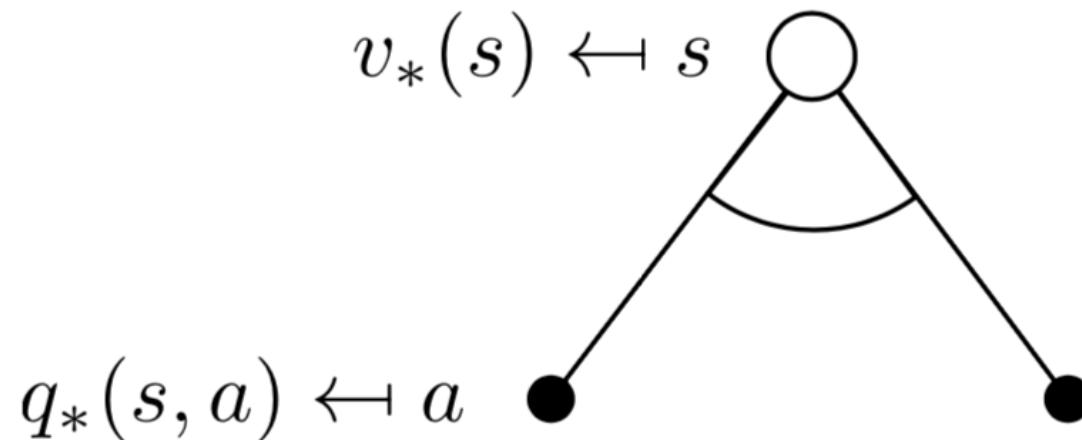
44

Studying at home MDP

$\pi_*(s)$ for $\gamma = 1$

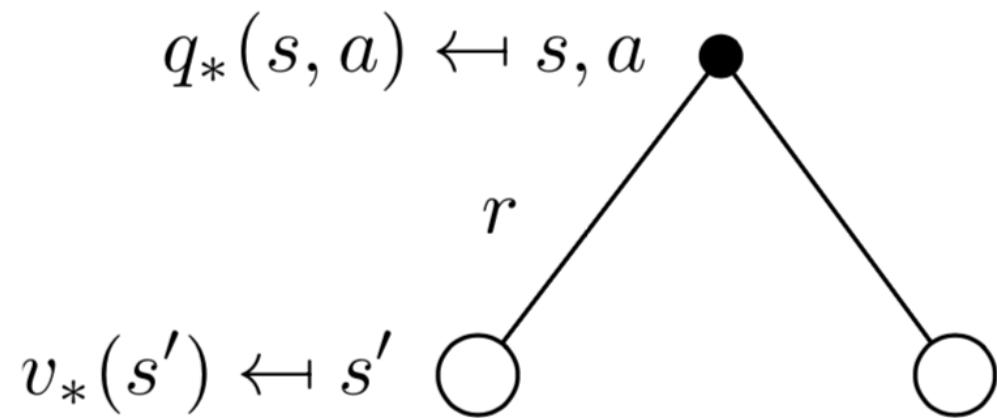


- The optimal value functions are recursively related by the Bellman optimality equations:



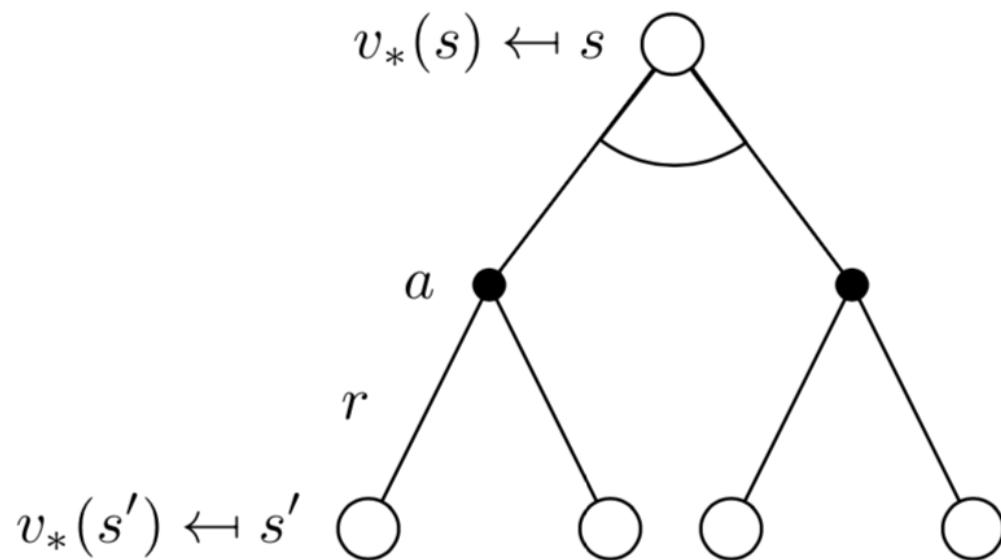
$$v_*(s) = \max_a q_*(s, a)$$

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



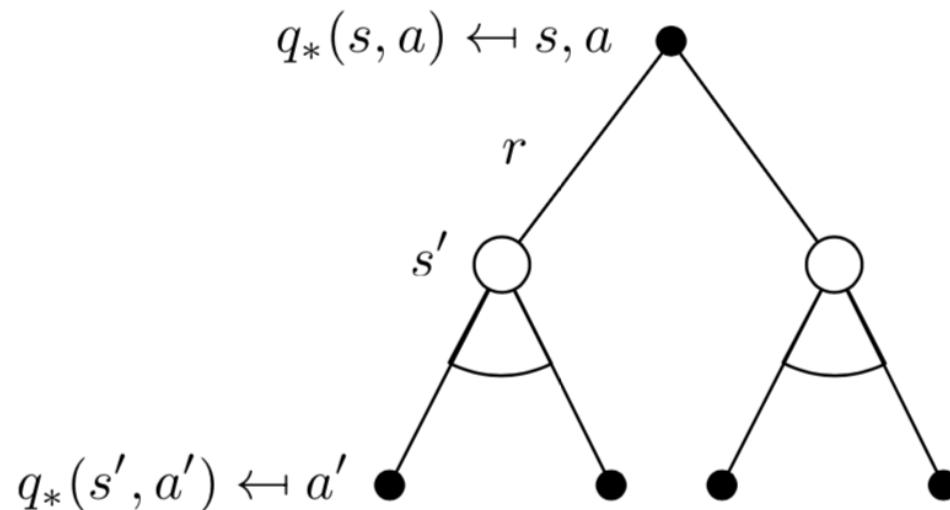
$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in S} \mathcal{P}_{ss'}^a v_*(s')$$

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



$$v_*(s) = \max_a \left(R_s^a + \gamma \sum_{s' \in S} \mathcal{P}_{ss'}^a v_*(s') \right)$$

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



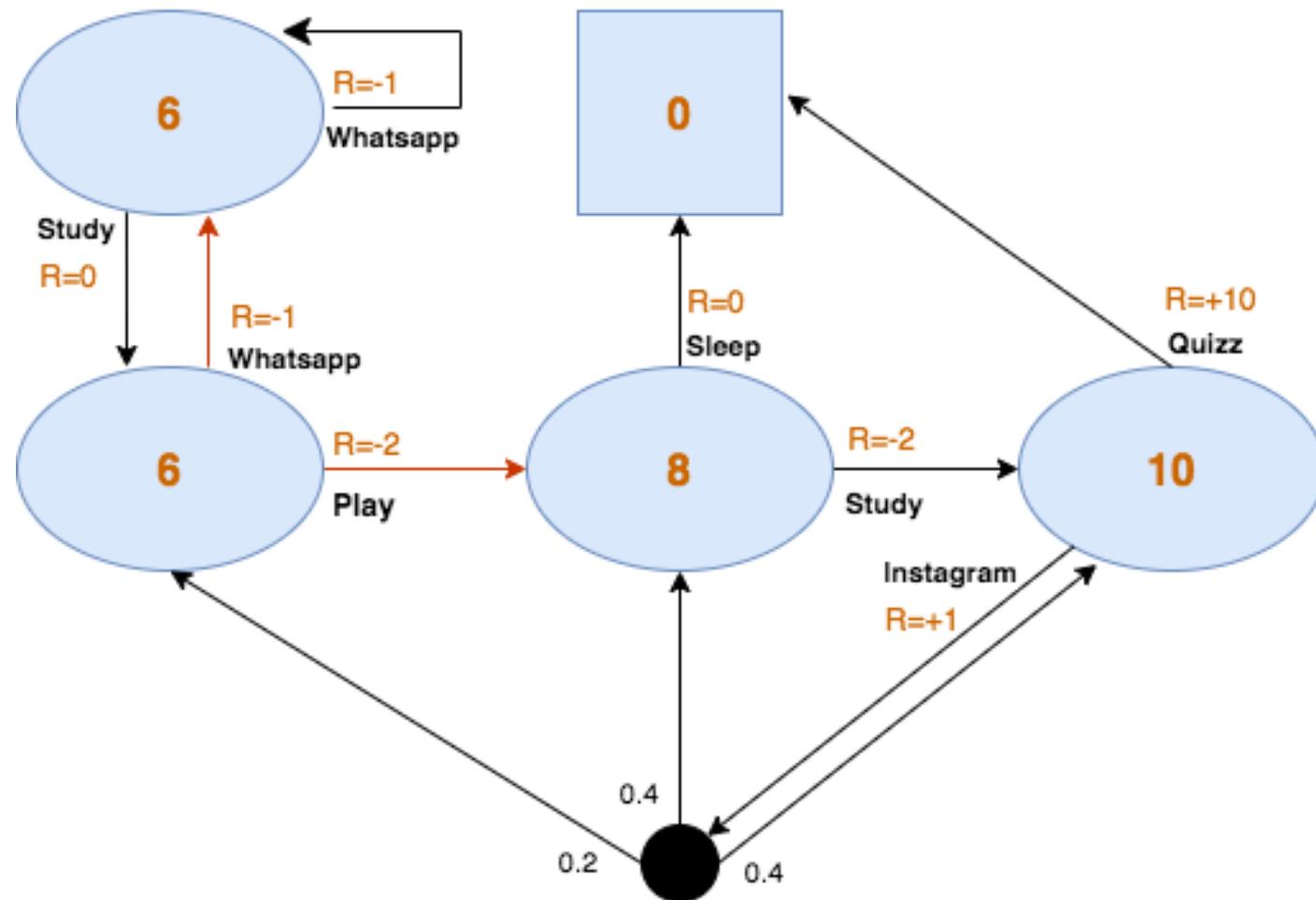
$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \max_{a'} q_*(s', a')$$

Markov Decision Process

49

Studying at home Markov Chain

$$6 = \max \{-2 + 8, -1 + 6\}$$



- Bellman Optimality Equation is non-linear
- No closed form solution (in general)
- Many iterative solution methods
 - Value Iteration
 - Policy Iteration
 - Q-learning
 - Sarsa

- Infinite and continuous MDPs
- Partially observable MDPs
- Undiscounted, average reward MDPs

- What happens if we do not have full MDP?
 - That is, we need to learn the associated rewards
 - Well .. We know about states and actions
 - We do not know about the system model (transition function) or the reward function
 - We can learn from experience and carry out actions to generate such experiences.
- This is the main goal of reinforcement learning...

- What is the Markov property?
- Why do we need the Markov Decision Process?
- When do we prefer immediate rewards?
- What is the use of the discount factor?
- Why do we use the Bellman function?
- How would you derive the Bellman equation for a q function?
- How are the value function and Q function related?
- How could we solve iteratively an MDP?
 - ▣ Next class

Lecture 2

- **Reading:**

- RUSSELL, S. NORVIG, P. Artificial Intelligence.
3a edição. Chapter 21.
- BARTO, A., SUTTON, R. Reinforcement
Learning: An Introduction. Second Edition.
Freely Available at:
<https://drive.google.com/file/d/1xeUDVGWGUUv1-ccUMAZHJLej2C7aAFWY/view?usp=sharing>

Lecture 2

- BARTO, A., SUTTON, R. Reinforcement Learning: An Introduction. Second Edition.
- MURPHY, R. R. Introduction to AI robotics. MIT Press, 2002.
- Lex Fridman, MIT Deep Learning Course, MIT, 2019.
- DUDEK, G.; JENKIN, M. Computational Principles of mobile robotics. Cambridge Press, 2000.
- ROMERO, R. A. F.; PRESTES, E.; OSÓRIO, F.; WOLF, D. (Orgs) Robótica móvel. LTC, 2014.
- BROOKS, R. Intelligence without representation. *Artificial Intelligence*, 47:139-159, 1991.
- RUSSEL, S. NORVIG, P. Artificial Intelligence: a modern approach. Prentice Hall, 2002.
- BRATKO, I. PROLOG: programming for artificial intelligence. Addison Wesley, 2nd edition, 1990.

This material is part of the Machine Learning Course
By Esther Colombini and Alexandre Simões

