

Reinforcement Learning II

Reinforcement Learning Roadmap

Knowledge of **Environment**

Perfect knowledge

Known Markov
Decision Process



No knowledge

Must learn from
experience

1

Core concepts in reinforcement learning

Actions, Rewards, Value, Environments, and Policies

2

Markov decision processes

...and Markov chains and Markov reward processes

3

Dynamic Programming

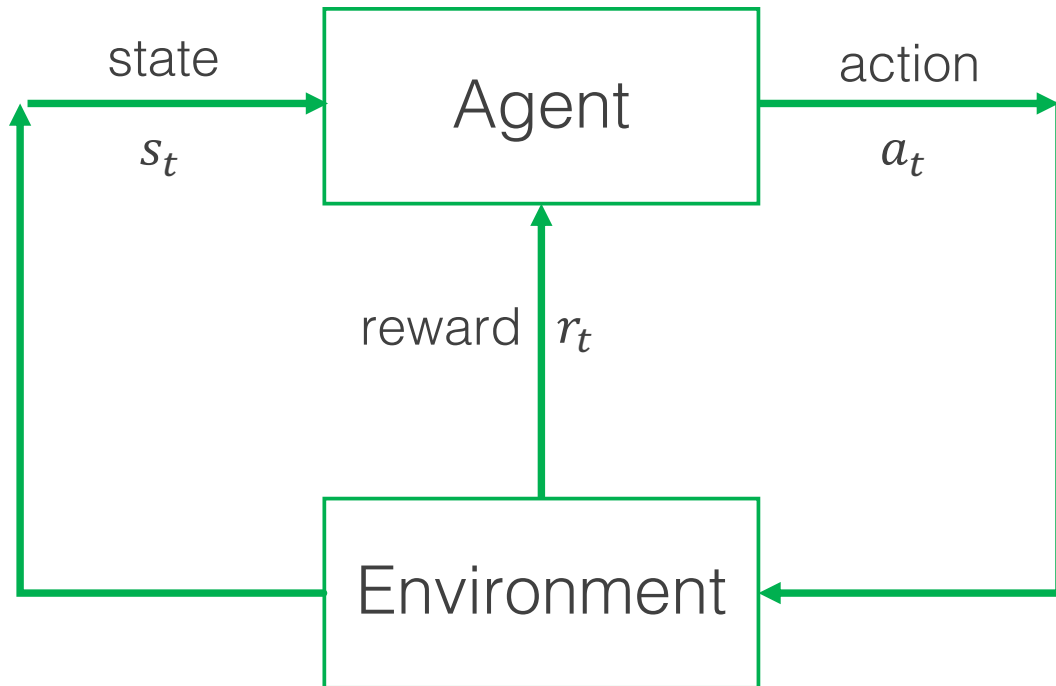
How do we find optimal policies?
(Bellman equations)

4

Monte Carlo Control

How do we estimate our value functions?
How do we use the value functions to choose actions?
How do we learn optimal policies from experience?

Reinforcement Learning Components



Policy (agent behavior), $\pi(s)$

Reward function (the goal), r_t

Value functions (expected returns),
 $v_\pi(s)$ State value
 $q_\pi(s, a)$ Action value

Policy $\pi(s)$

(which actions to take in each state)

Start			
	↓		
→	→	↓	←
↑		↓	
	→	→	↓
	↑		↓
→	↑		↓
Exit			

Reward r_t

(rewards are received after actions are taken)

Start			
	-1		
-1	-1	-1	-1
-1		-1	
	-1	-1	-1
	-1		-1
-1	-1		-1
Exit			

State Value $v_\pi(s)$

(expected cumulative rewards starting from current state **if** we follow the policy)

Start			
	-8		
-8	-7	-6	-7
-9		-5	
	-5	-4	-3
	-6		-2
-8	-7		-1
Exit			

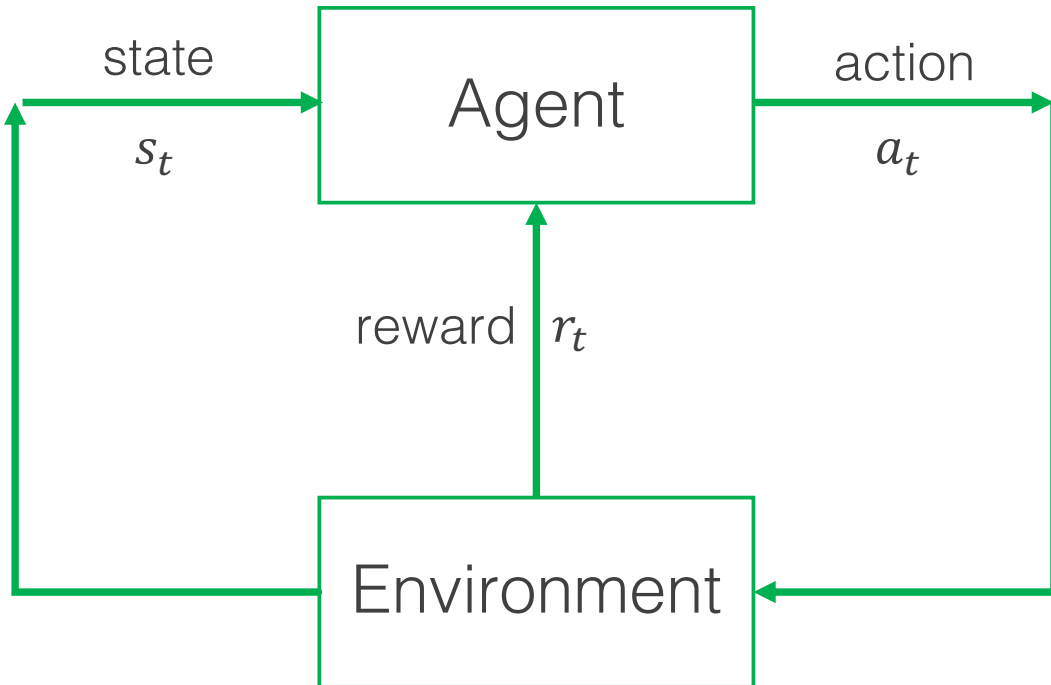
Action Value $q_\pi(s, a)$

(expected cumulative rewards starting from current state **if** we take action a then follow the policy)

↑	-9
→	-7
←	-9

↑	-4
↓	-2

Value functions



State Value function, $v_\pi(s)$

- How “good” is it to be in a state, s_t then follow policy π to choose actions
- Total expected rewards

$$v_\pi(s) = E_\pi[G_t | s_t = s]$$

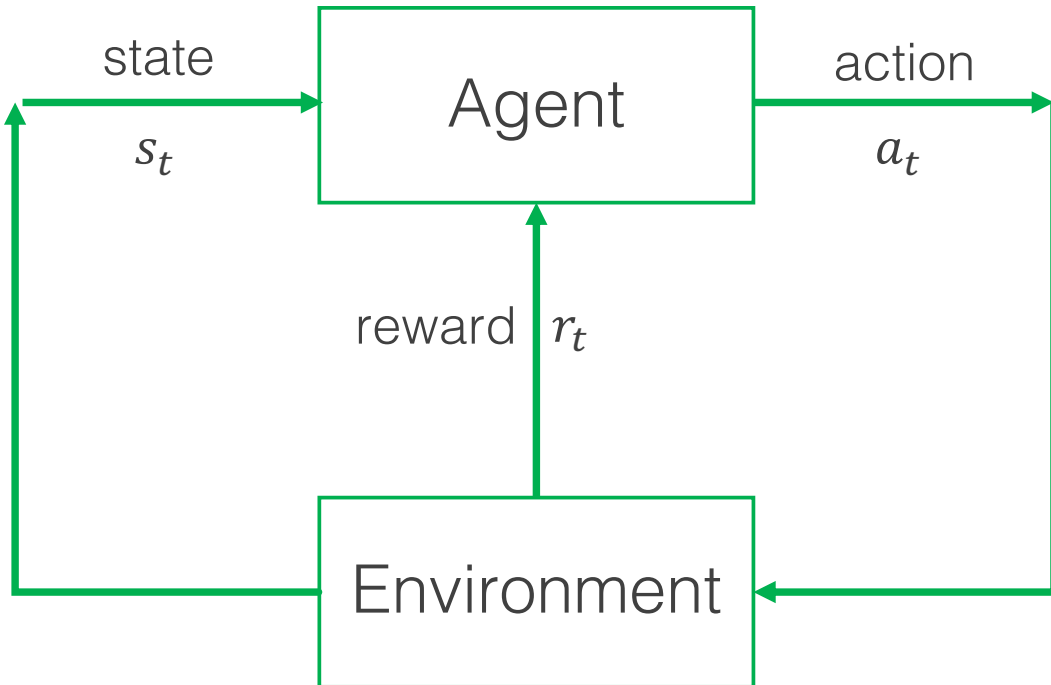
Action Value function, $q_\pi(s, a)$

- How “good” is it to be in a state, s , take action a , then follow policy π to choose actions
- Total expected rewards

$$q_\pi(s, a) = E_\pi[G_t | s_t = s, a_t = a]$$

Where
$$G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

Policy



Policy, $\pi(s)$

- Selects an action to choose based on the state
- Determines an agent's “behavior”

Deterministic policy:

$$a = \pi(s)$$

Stochastic policy:

$$\pi(a|s) = P(A_t = a|S_t = s)$$

Helps us “explore” the state space

RL tries to learn the “best” policy

Returns / cumulative reward

Episodic tasks (finite number, T , of steps, then reset)

$$G_t = r_{t+1} + r_{t+2} + \dots + r_T$$

Continuing tasks with discounting ($T \rightarrow \infty$)

$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

where $0 \leq \gamma \leq 1$ is the discount rate

γ makes the agent care more/less about immediate rewards

Building blocks for the full RL problem

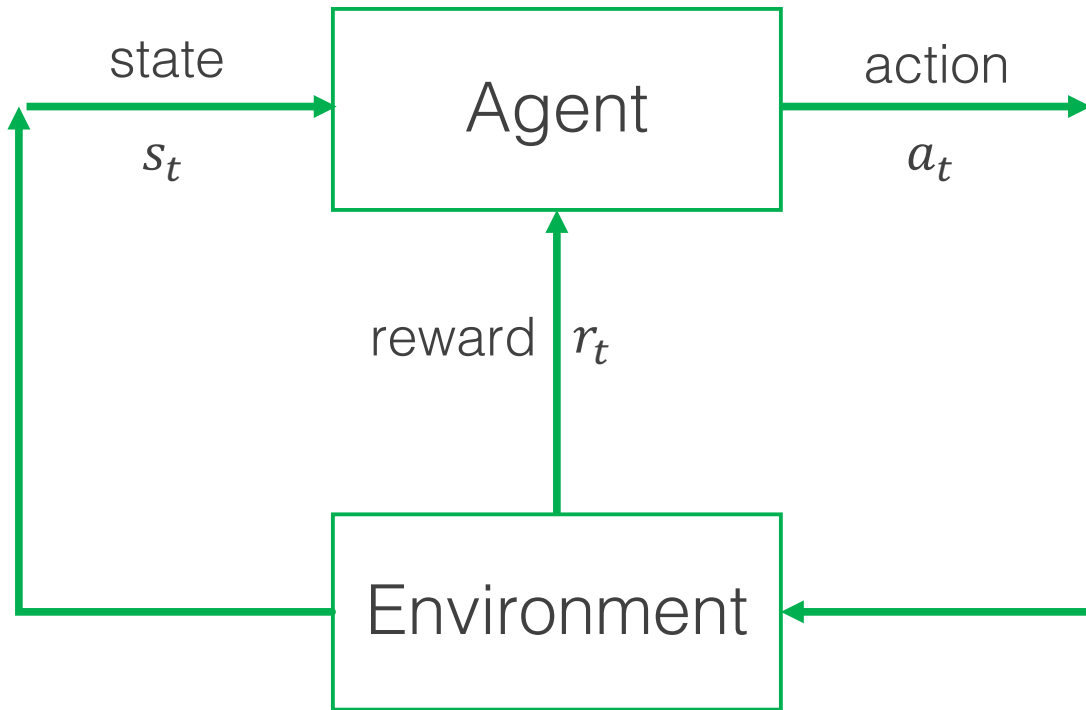
1	Markov Chain	{state space S , transition probabilities P }
2	Markov Reward Process (MRP)	{ S , P , + rewards R , discount rate γ } adds rewards (and values)
3	Markov Decision Process (MDP)	{ S , P , R , γ , + actions A } adds decisions (i.e. the ability to control)

MDPs form the framework for most reinforcement learning environments

Adapted from David Silver, 2015

History

The record of all that has happened in this system



Step 0: s_0, a_0

Step 1: r_1, s_1, a_1

Step 2: r_2, s_2, a_2

\vdots

Step T: r_t, s_t, a_t

History at time t : $H_t = \{s_t, a_t, r_{t-1}, s_{t-1}, a_{t-1}, \dots, r_1, s_1, a_1, s_0, a_0\}$

Markov property

Instead of needing the full history:

$$H_t = \{s_t, a_t, r_{t-1}, s_{t-1}, a_{t-1}, \dots, r_1, s_1, a_1, s_0, a_0\}$$

We can summarize everything in the current state

$$H_t = \{s_t, a_t\}$$

The future is independent of the past given the present

Another way of saying this is:

$$P(s_{t+1}|s_t) = P(s_{t+1}|s_t, s_{t-1}, \dots, s_1, s_0)$$

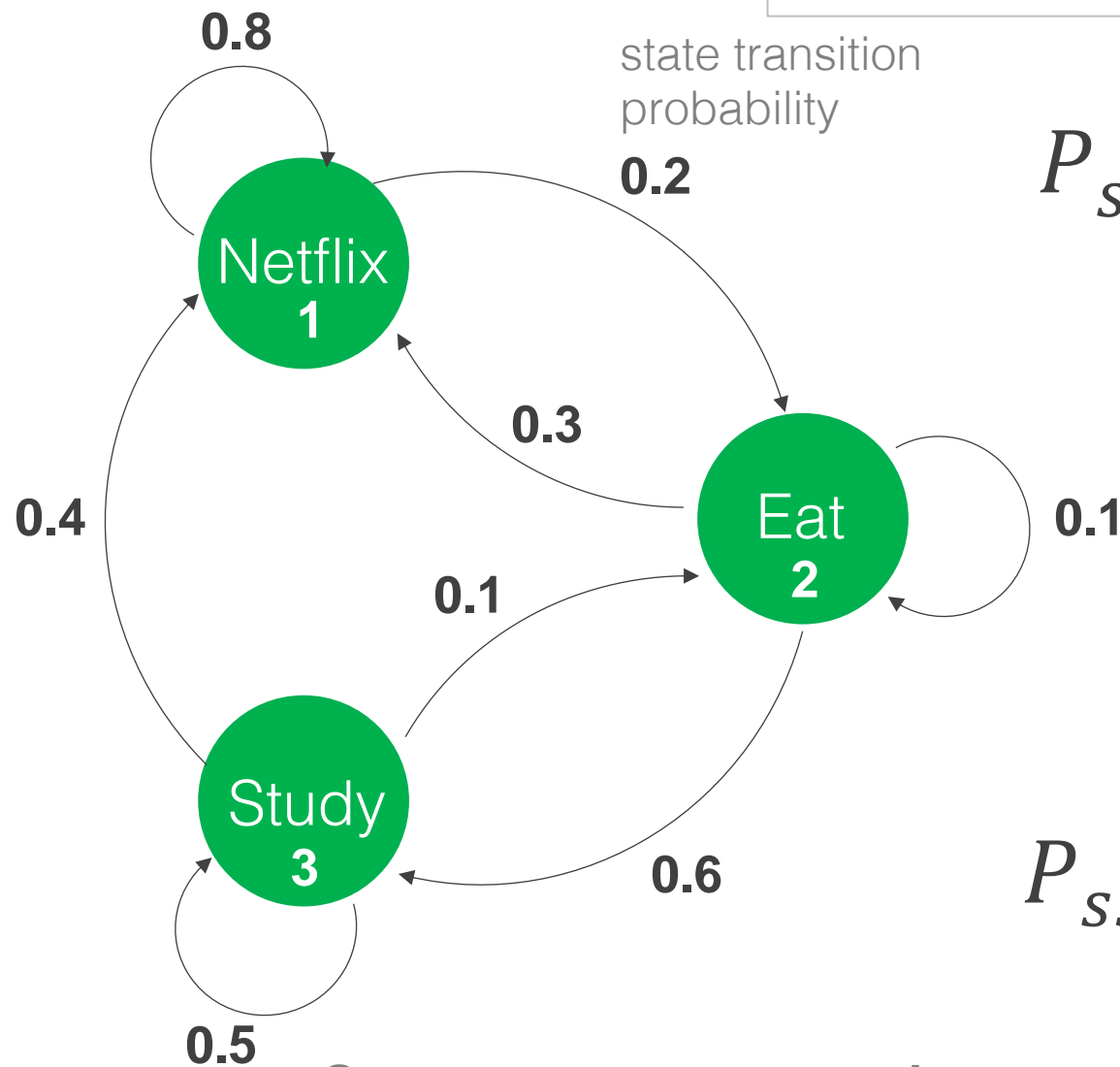
Markov Chains

Example: student life

Two components: $\{S, P\}$

State space, S

Transition matrix, P



State-space representation

$$P_{ss'} =$$

$$P_{ss'} =$$

State transition probabilities

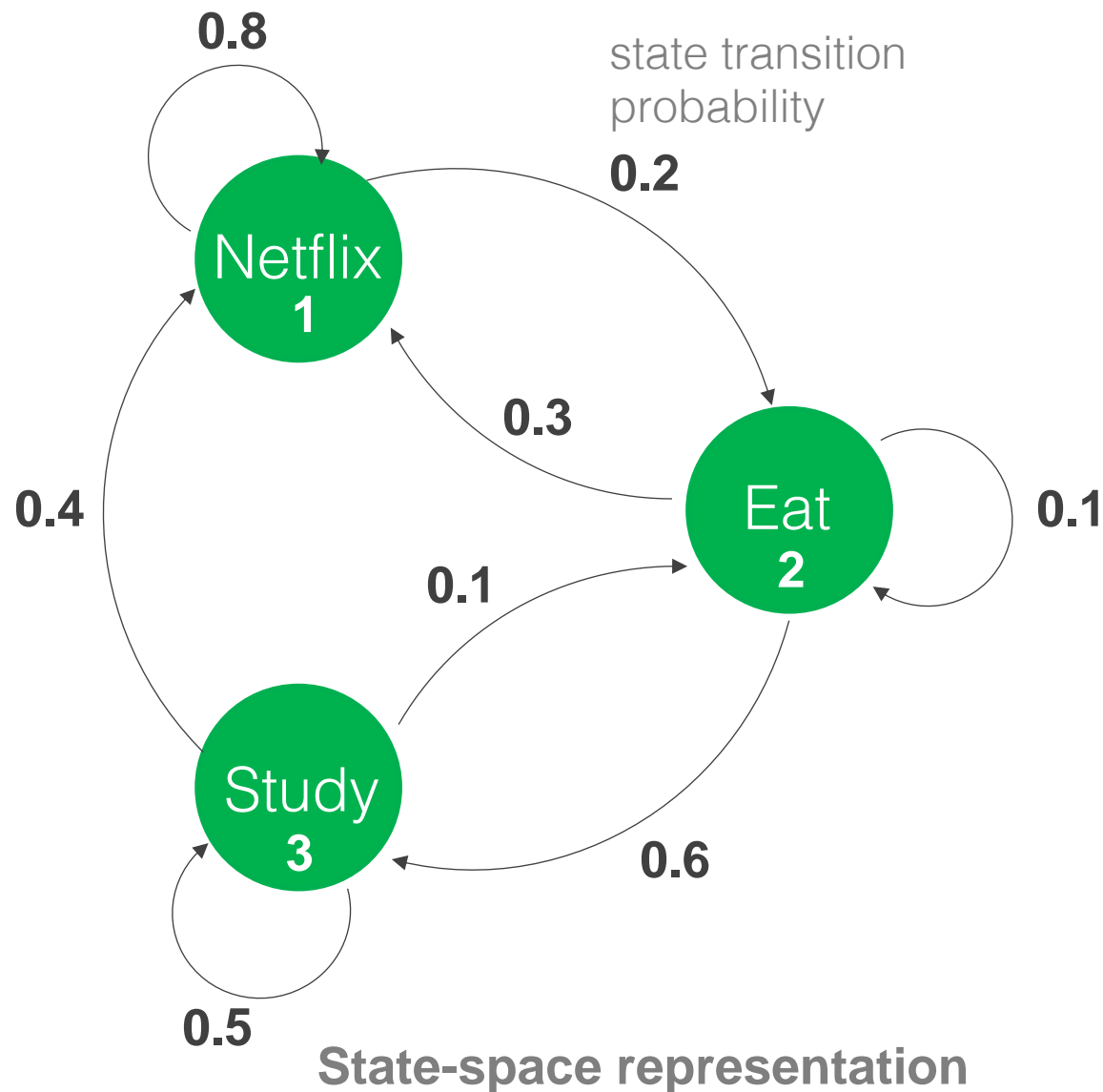
		To state		
		1	2	3
From state	1	p_{11}	p_{12}	p_{13}
	2	p_{21}	p_{22}	p_{23}
	3	p_{31}	p_{32}	p_{33}

Transitions out of each state sum to 1

		To state		
		Netflix	Eat	Study
From state	Netflix	0.8	0.2	0
	Eat	0.3	0.1	0.6
	Study	0.4	0.1	0.5

Markov Chains

Example: student life



If we start in state 1, what's the probability we'll be in each state after one step?

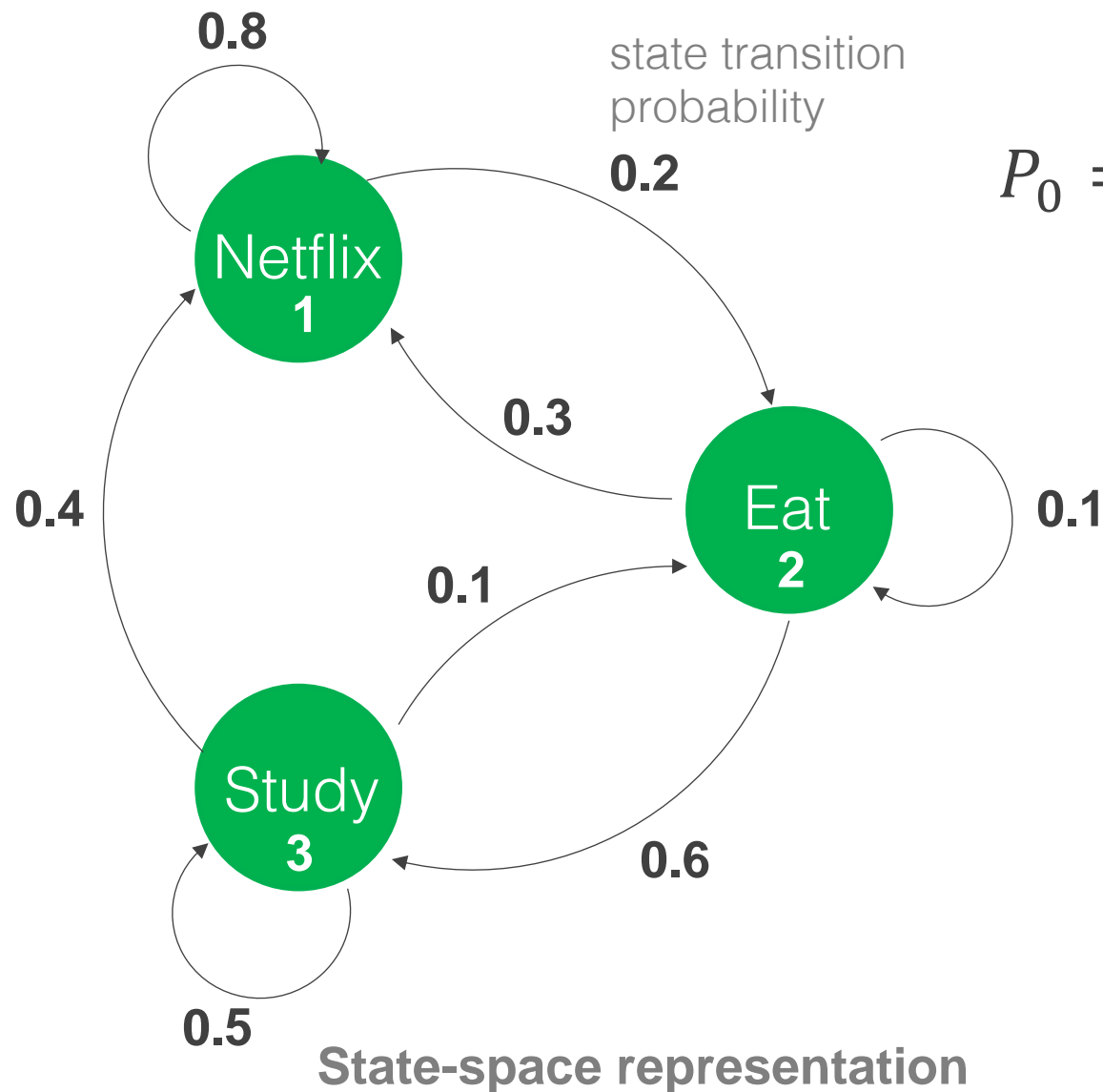
$$P_1 = \begin{bmatrix} 0.8 & 0.2 & 0 \end{bmatrix}$$

This is the first row of the state transition probability matrix

$$P_{ss'} = \begin{matrix} & \text{To state} \\ & \begin{matrix} \text{Netflix} & \text{Eat} & \text{Study} \end{matrix} \\ \begin{matrix} \text{From state} \\ \text{Netflix} \\ \text{Eat} \\ \text{Study} \end{matrix} & \begin{bmatrix} 0.8 & 0.2 & 0 \\ 0.3 & 0.1 & 0.6 \\ 0.4 & 0.1 & 0.5 \end{bmatrix} \end{matrix}$$

Markov Chains

Example: student life



If we started in state 1, we can calculate the probabilities of being in each state at step 1 as:

$$P_0 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}^T$$

$$P_1 = P_0 P_{ss'}$$

$$P_1 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.8 & 0.2 & 0 \\ 0.3 & 0.1 & 0.6 \\ 0.4 & 0.1 & 0.5 \end{bmatrix}$$

$$P_1 = \begin{bmatrix} 0.8 & 0.2 & 0 \end{bmatrix}$$

$$P_{ss'} =$$

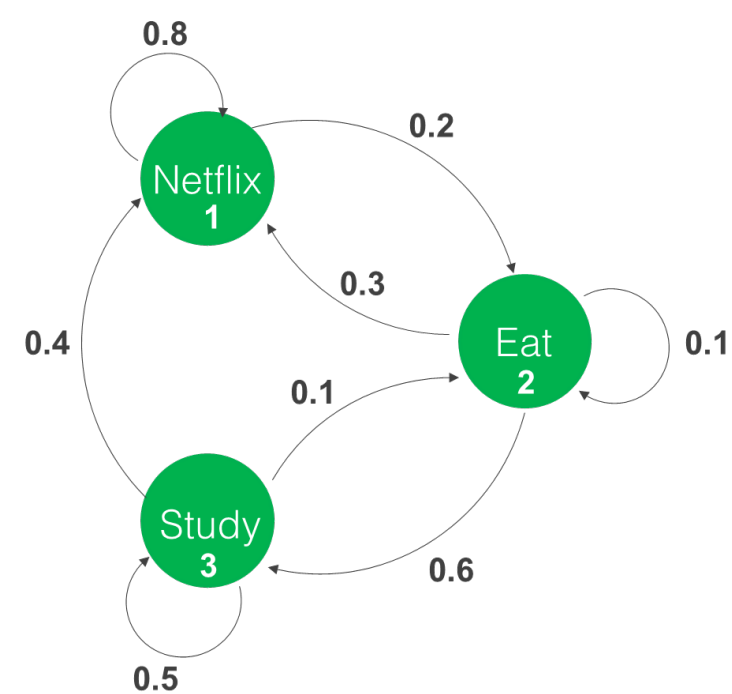
		To state		
		Netflix	Eat	Study
From state	Netflix	0.8	0.2	0
	Eat	0.3	0.1	0.6
	Study	0.4	0.1	0.5

$$\textbf{1} \quad P_1 = P_0 P_{ss'}$$

$$P_1 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.8 & 0.2 & 0 \\ 0.3 & 0.1 & 0.6 \\ 0.4 & 0.1 & 0.5 \end{bmatrix}$$

$$P_1 = [0.8 \quad 0.2 \quad 0]$$

$$P_0 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}^T$$



$$\textbf{2} \quad P_2 = P_1 P_{ss'} = P_0 P_{ss'} P_{ss'} = P_0 P_{ss'}^2$$

$$P_2 = \begin{bmatrix} 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.8 & 0.2 & 0 \\ 0.3 & 0.1 & 0.6 \\ 0.4 & 0.1 & 0.5 \end{bmatrix} \begin{bmatrix} 0.8 & 0.2 & 0 \\ 0.3 & 0.1 & 0.6 \\ 0.4 & 0.1 & 0.5 \end{bmatrix}$$

$$P_2 = [0.7 \quad 0.18 \quad 0.12]$$

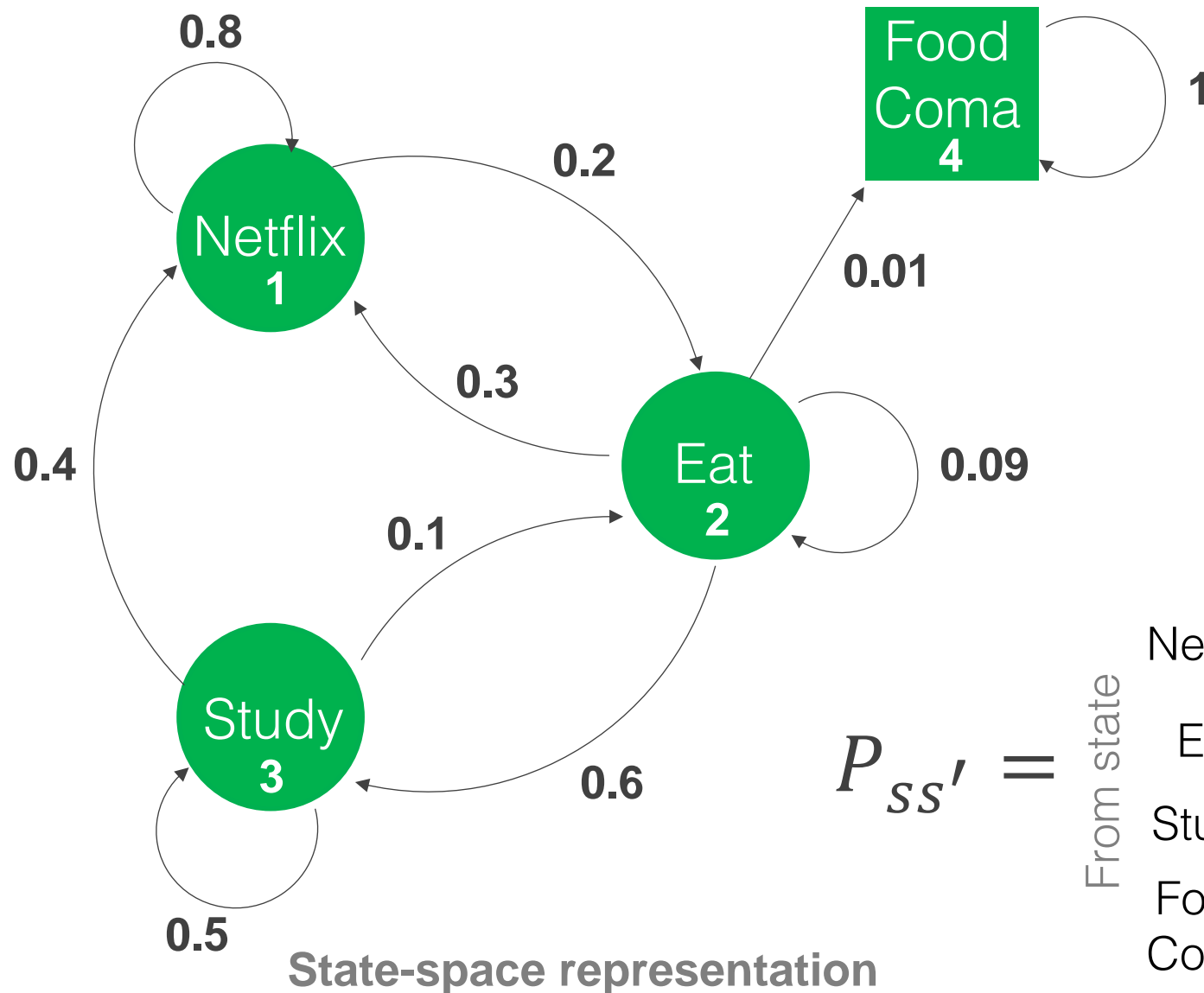
$$\textbf{n} \quad P_n = P_0 P_{ss'}^n$$

As $n \rightarrow \infty$, we identify our steady state probabilities

$$P_\infty = [0.64 \quad 0.16 \quad 0.20]$$

Markov Chains with absorbing state

Example: student life

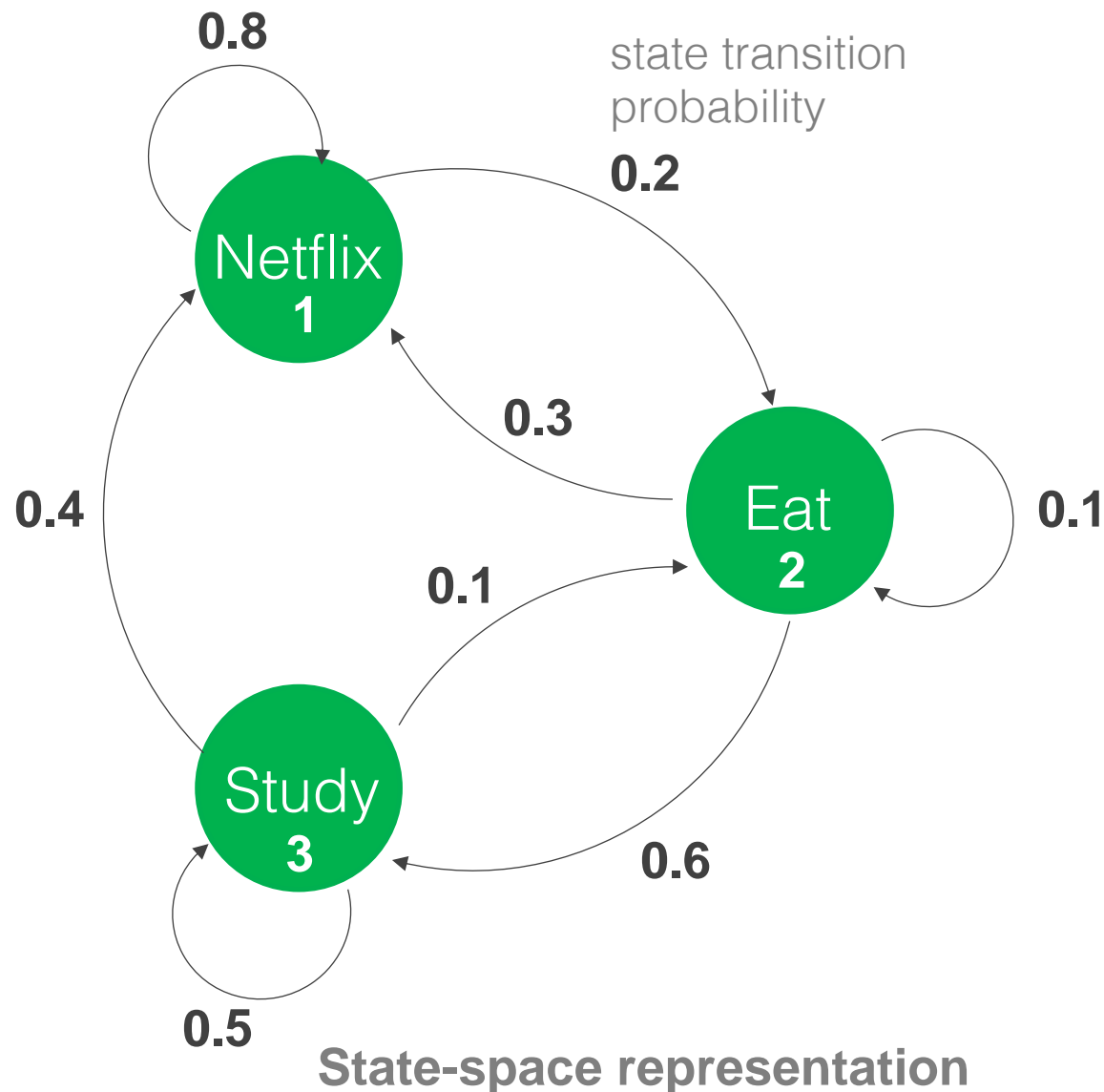


$$P_{ss'} =$$

		To state			
		Netflix	Eat	Study	Food Coma
From state	Netflix	0.8	0.2	0	0
	Eat	0.3	0.09	0.6	0.01
	Study	0.4	0.1	0.5	0
	Food Coma	0	0	0	1

Markov Chains

Example: student life



Markov chains can be used to represent sequential discrete-time data

Can estimate long-term state probabilities

Can simulate state sequences based on the model

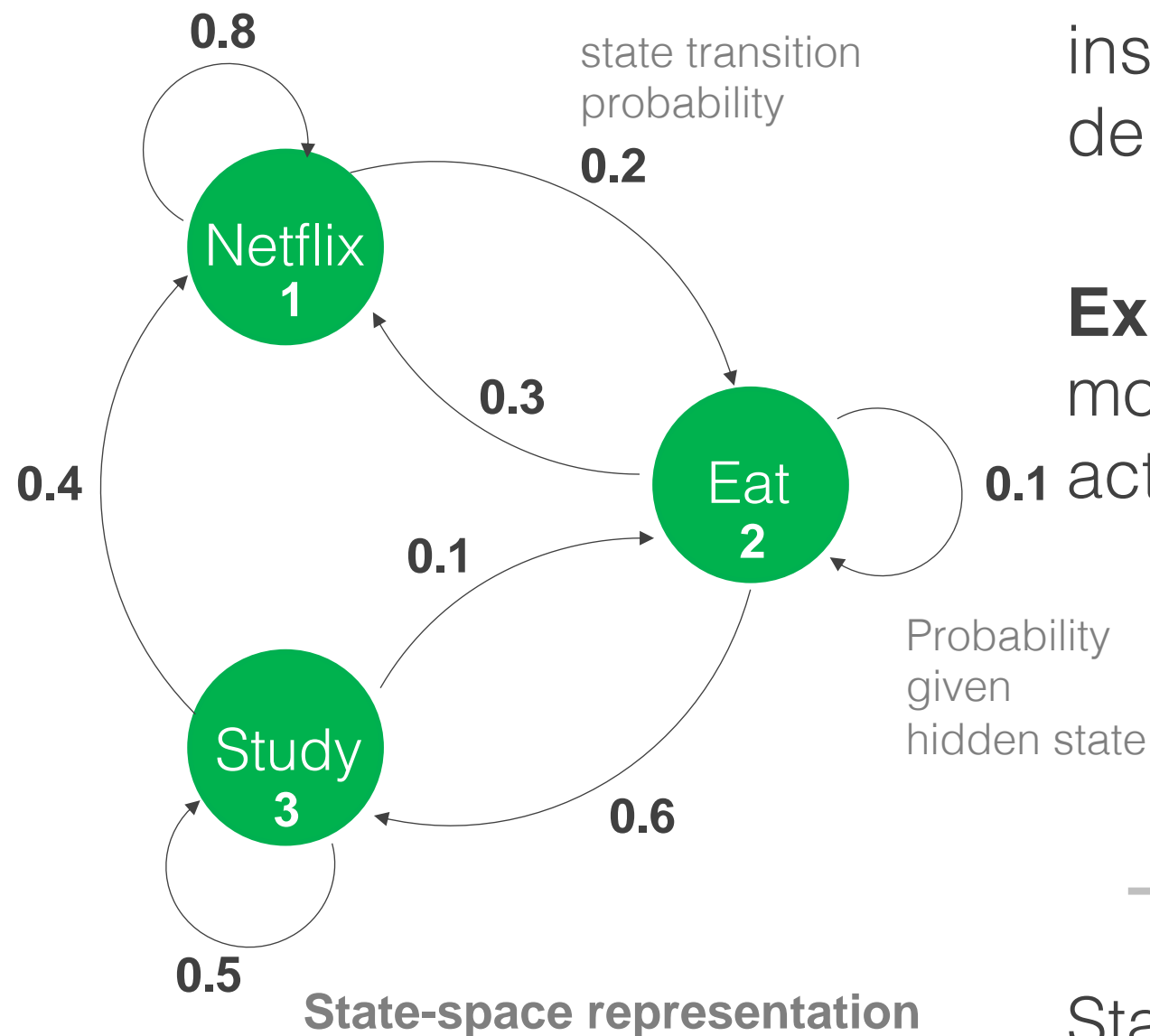
Markov property applies
(current state gives you all the information you need about future states)

$$P(s_{t+1}|s_t) = P(s_{t+1}|s_t, s_{t-1}, \dots, s_1, s_0)$$

Valid if the system is **autonomous** and the states are **fully observable**

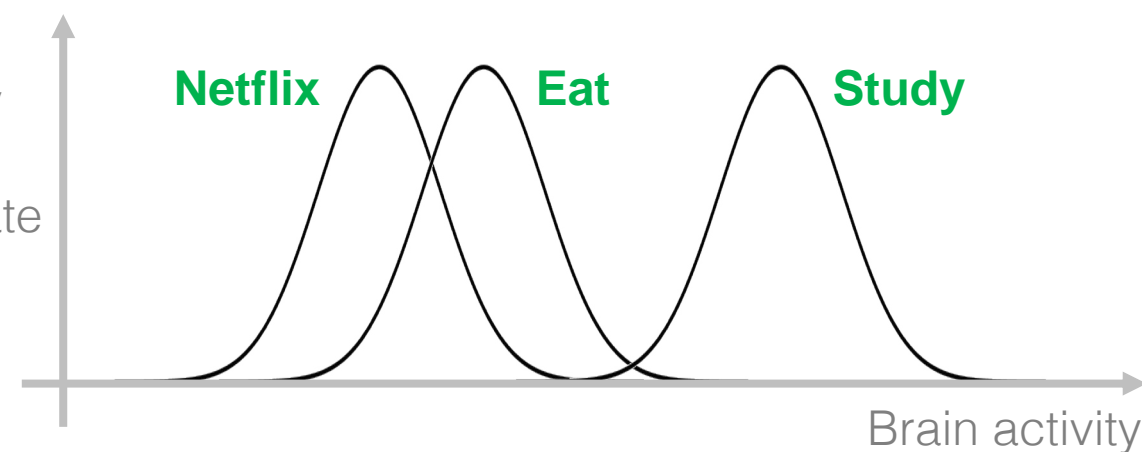
Hidden Markov Models

Example: student life



What if we **don't directly observe what state** the system is in, but instead observe a quantity that depends on the state?

Example: the student wears an EEG monitor, and we see readings of brain activity.



States are **hidden** or **latent variables**

Markov Models

States are
Fully Observable

States are
Partially Observable

Autonomous
(no actions;
make predictions)

Markov Chain,
Markov Reward Process

Hidden Markov Model
(HMM)

Controlled
(can take actions)

Markov Decision
Process (MDP)

Partially Observable
Markov Decision
Process (POMDP)

Applications

HMMs: time series ML, e.g. speech + handwriting recognition

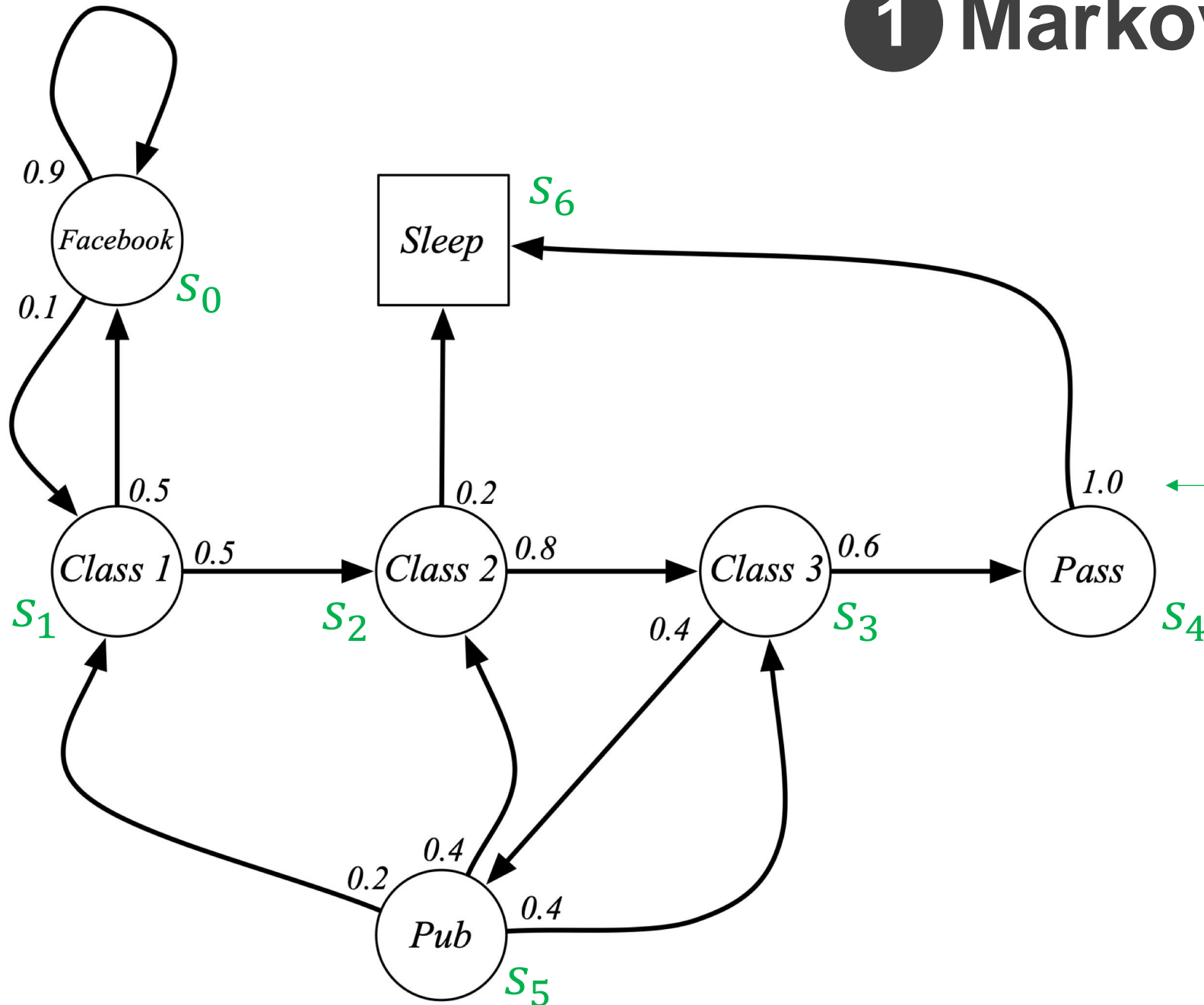
MDPs: framework for reinforcement learning

1 Markov Chain Example

Components:

State space S ,

Transition probabilities P



$$P_{46} = P_{ss'}$$

Sample Episodes:

C1,C2,Sleep

C1,FB,FB,FB,C1,C2,C3,Pass,Sleep

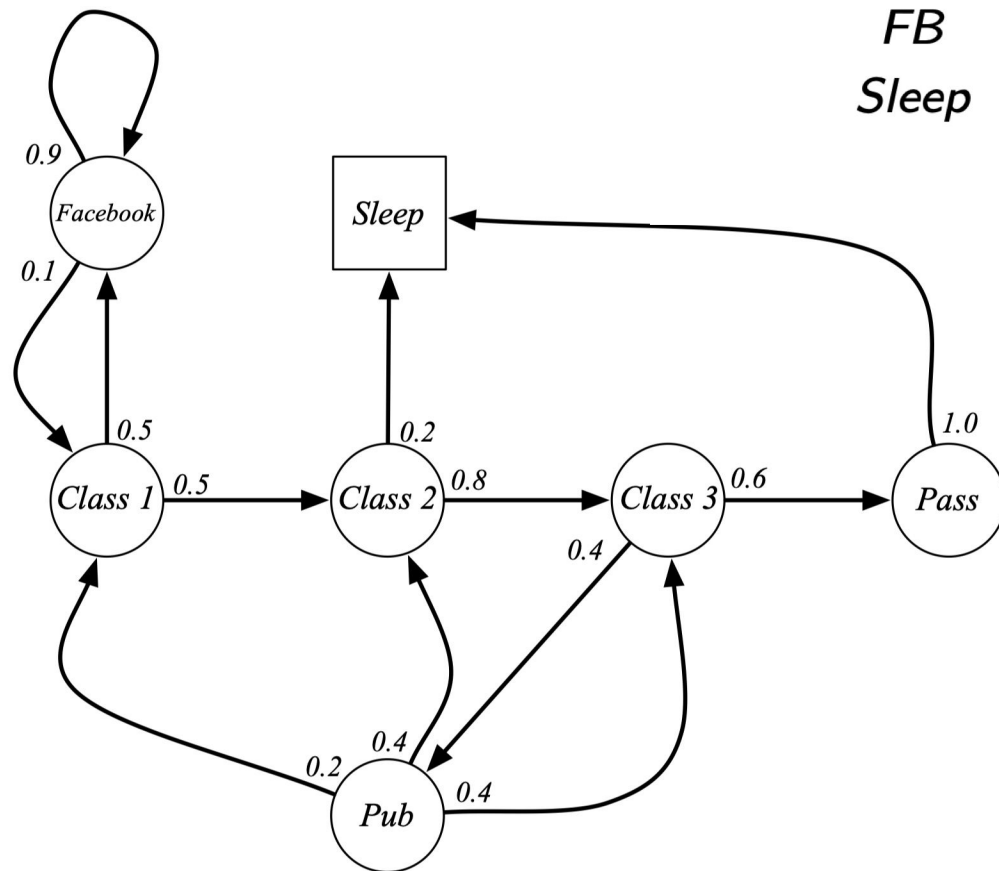
Example from David Silver, UCL, 2015

Markov Chain

$$\mathcal{P} = \begin{matrix} & \begin{matrix} C1 \\ C2 \\ C3 \\ Pass \\ Pub \\ FB \\ Sleep \end{matrix} \end{matrix}$$

	C1	C2	C3	Pass	Pub	FB	Sleep
C1		0.5				0.5	
C2			0.8				0.2
C3				0.6	0.4		
Pass	0.2	0.4	0.4				1.0
Pub	0.1					0.9	
FB							
Sleep							1

State transition probability matrix, $P_{ss'}$



Example from David Silver, UCL, 2015

Building blocks for the full RL problem

1	Markov Chain	{state space S , transition probabilities P }
2	Markov Reward Process (MRP)	{ S , P , + rewards R , discount rate γ } adds rewards (and values)
3	Markov Decision Process (MDP)	{ S , P , R , γ , + actions A } adds decisions (i.e. the ability to control)

MDPs form the framework for most reinforcement learning environments

Adapted from David Silver, 2015

2 Markov Reward Process

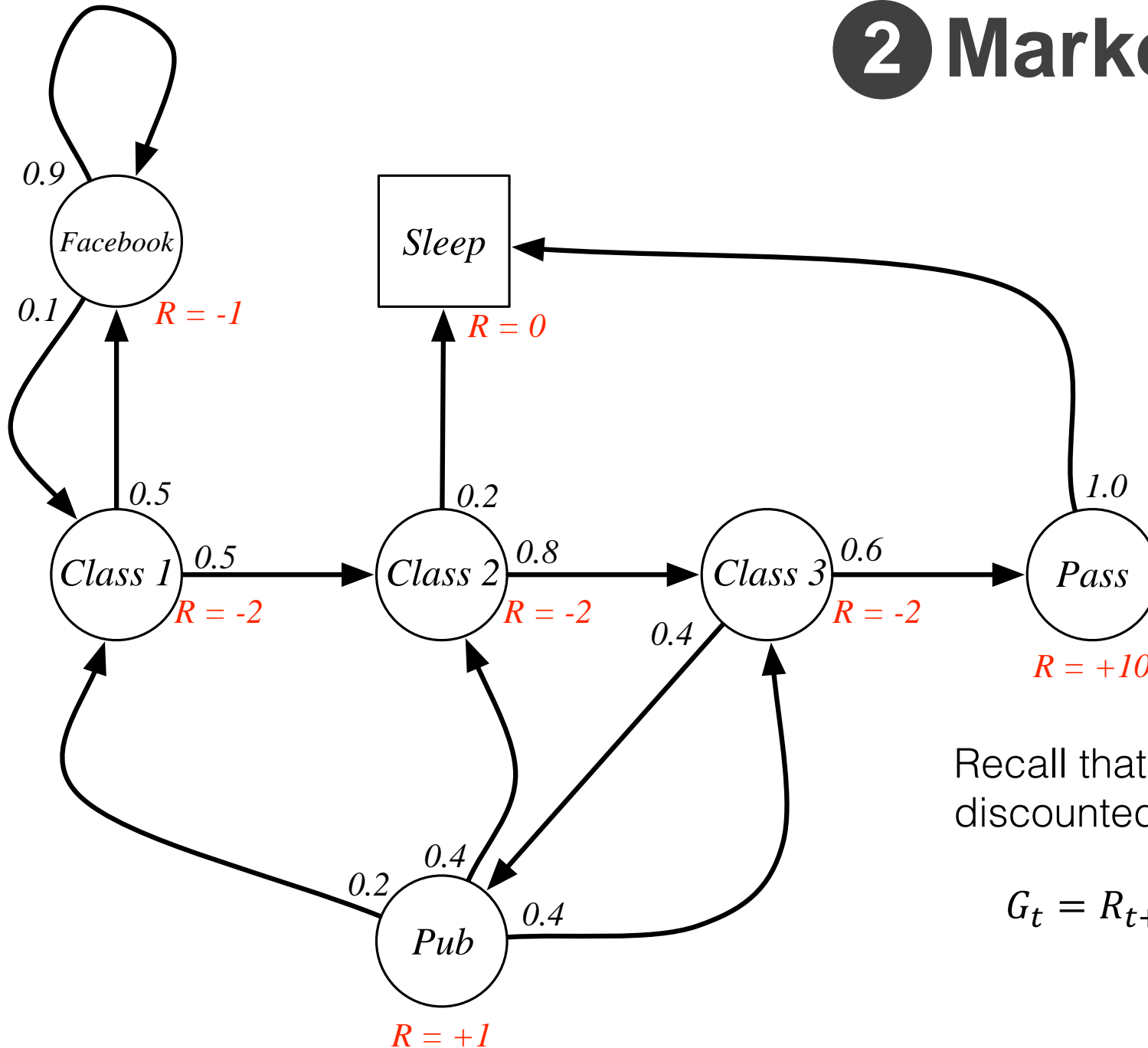
Components:

State space S ,

Transition probabilities, P

Rewards, R

Discount rate, γ



Recall that returns, let's call G_t , are the total discounted rewards from time t :

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Example from David Silver, UCL, 2015

2 Markov Reward Process

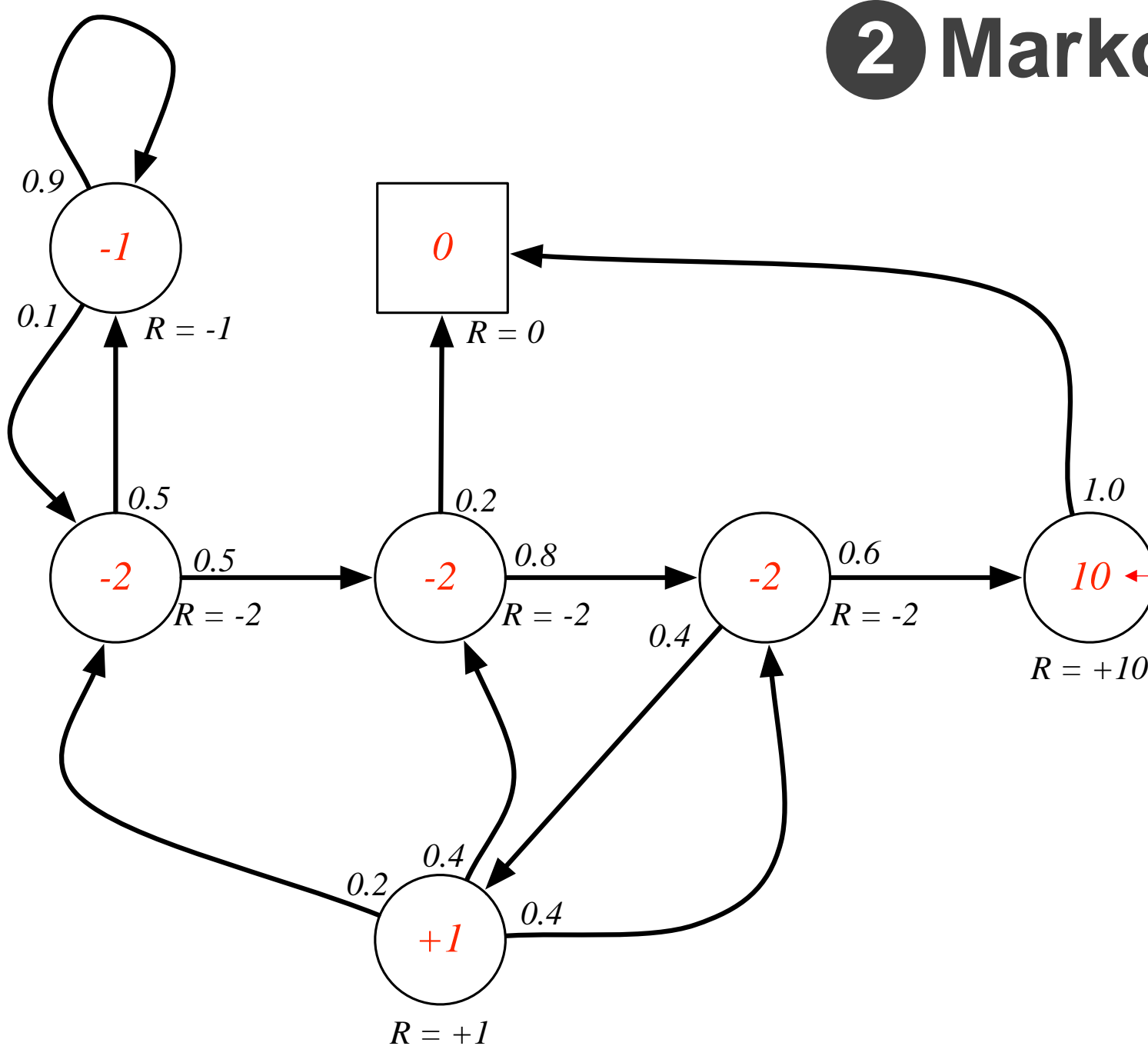
Components:

State space S ,

Transition probabilities, P

Rewards, R

Discount rate, γ



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

State value function $v(s)$ is the expected total reward (into the future)

$$v(s) = E[G_t | S = s_t]$$

Example from David Silver, UCL, 2015

2 Markov Reward Process

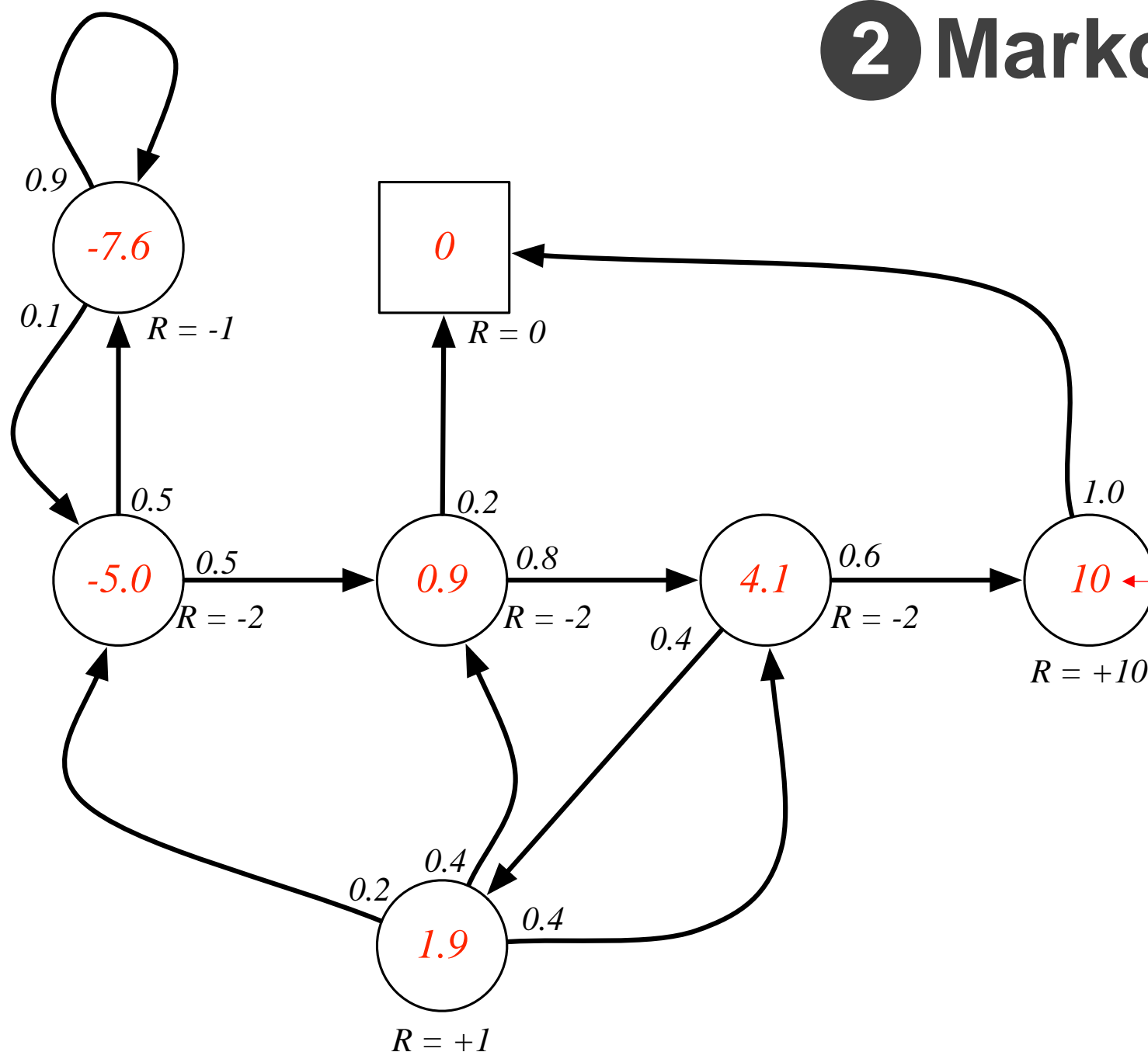
Components:

State space S ,

Transition probabilities, P

Rewards, R

Discount rate, γ



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

State value function $v(s)$ is the expected total reward (into the future)

$$v(s) = E[G_t | S = s_t]$$

Example from David Silver, UCL, 2015

“Backup” property of state value functions

$$v(s) \triangleq E[G_t | S_t = s] \quad \text{where } G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots$$

$$= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots | S_t = s]$$

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

$$= E[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

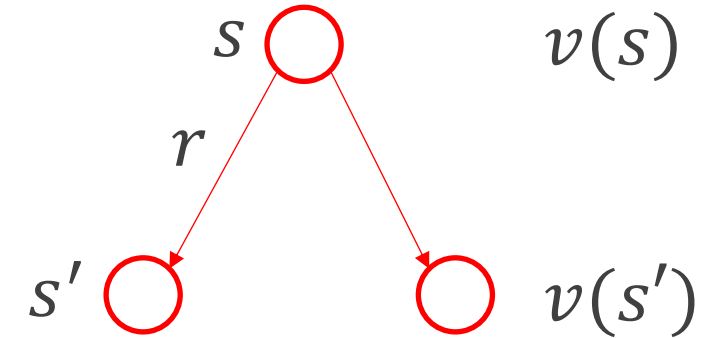
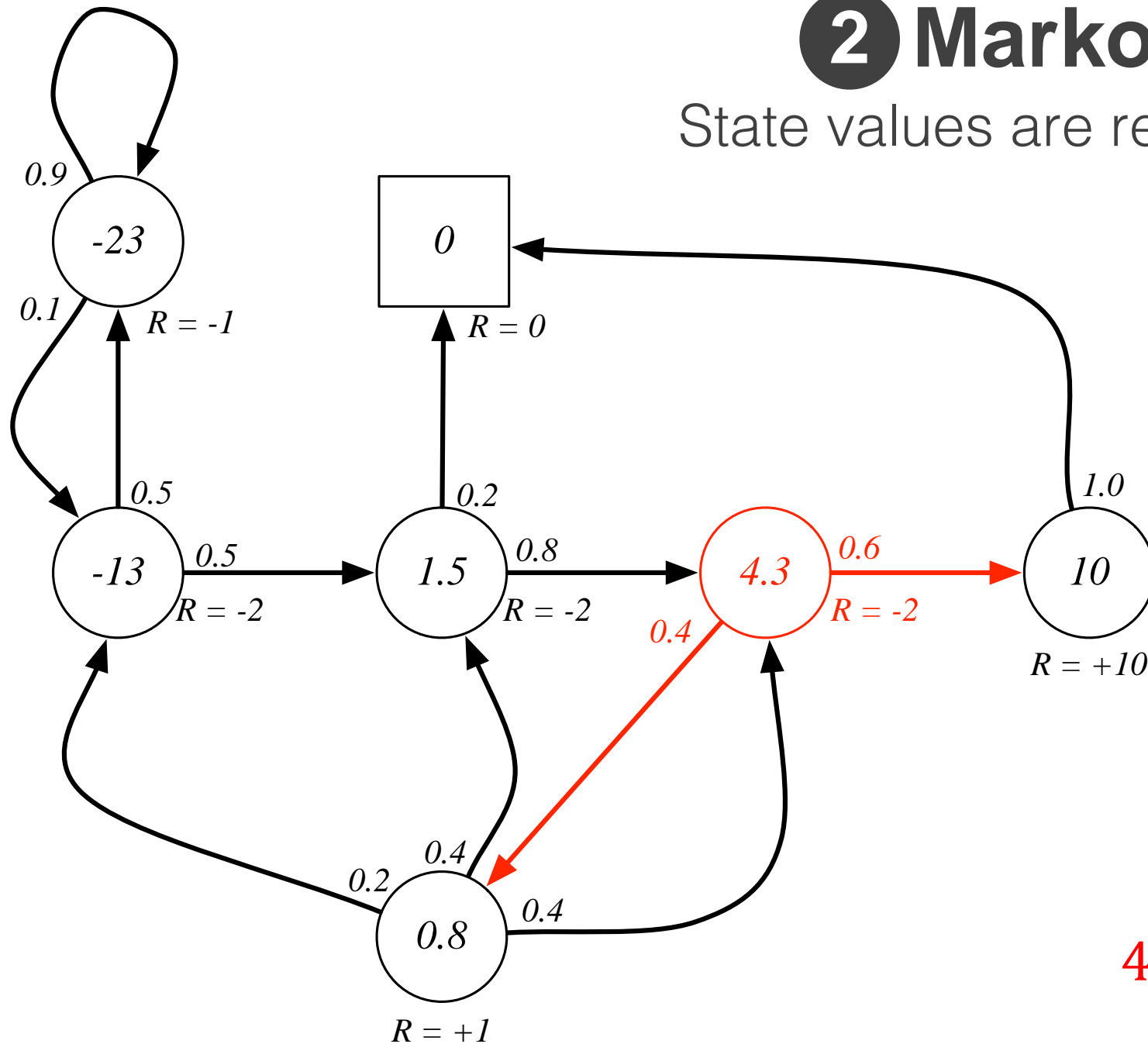
$$= E[R_{t+1} + \gamma v(S_{t+1}) | S_t = s]$$

This recursive relationship is a version of the **Bellman Equation**

Example from David Silver, UCL, 2015

2 Markov Reward Process

State values are related to neighboring states



possible states we could transition to from s

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

$$v(s) = r_s + \gamma \sum_{s'} P_{ss'} v(s')$$

$$r_s = E[R_{t+1} | S_t = s]$$

Assume $\gamma = 1$

$$4.3 = -2 + 0.6 \times 10 + 0.4 \times 0.8$$

Example from David Silver, UCL, 2015

Building blocks for the full RL problem

1	Markov Chain	{state space S , transition probabilities P }
2	Markov Reward Process (MRP)	{ S , P , + rewards R , discount rate γ } adds rewards (and values)
3	Markov Decision Process (MDP)	{ S , P , R , γ , + actions A } adds decisions (i.e. the ability to control)

MDPs form the basis for most reinforcement learning environments

Adapted from David Silver, 2015

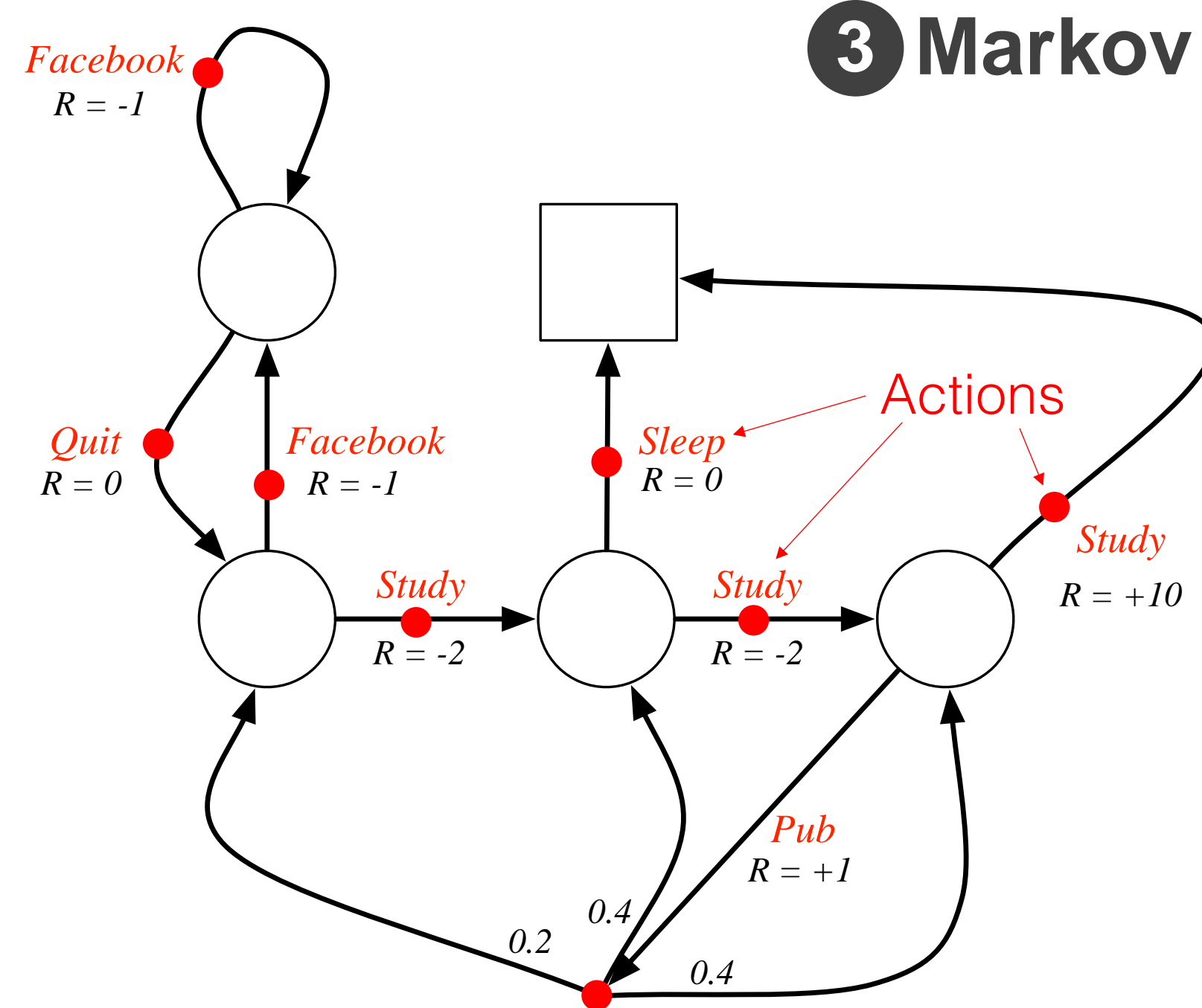
3 Markov Decision Process

Components:

State space S ,
Transition probabilities, P
Rewards, R
Discount rate, γ
Actions, A

Adds interaction with the environment

An agent in a state chooses an action, the environment (the MDP) provides a reward and the next state



Example from David Silver, UCL, 2015

3 Markov Decision Process

Policy (how we choose actions)

(can be stochastic or deterministic)

$$\pi(a|s) = P(a|s)$$

State value function

(expected return from state s , and following policy π)

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

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

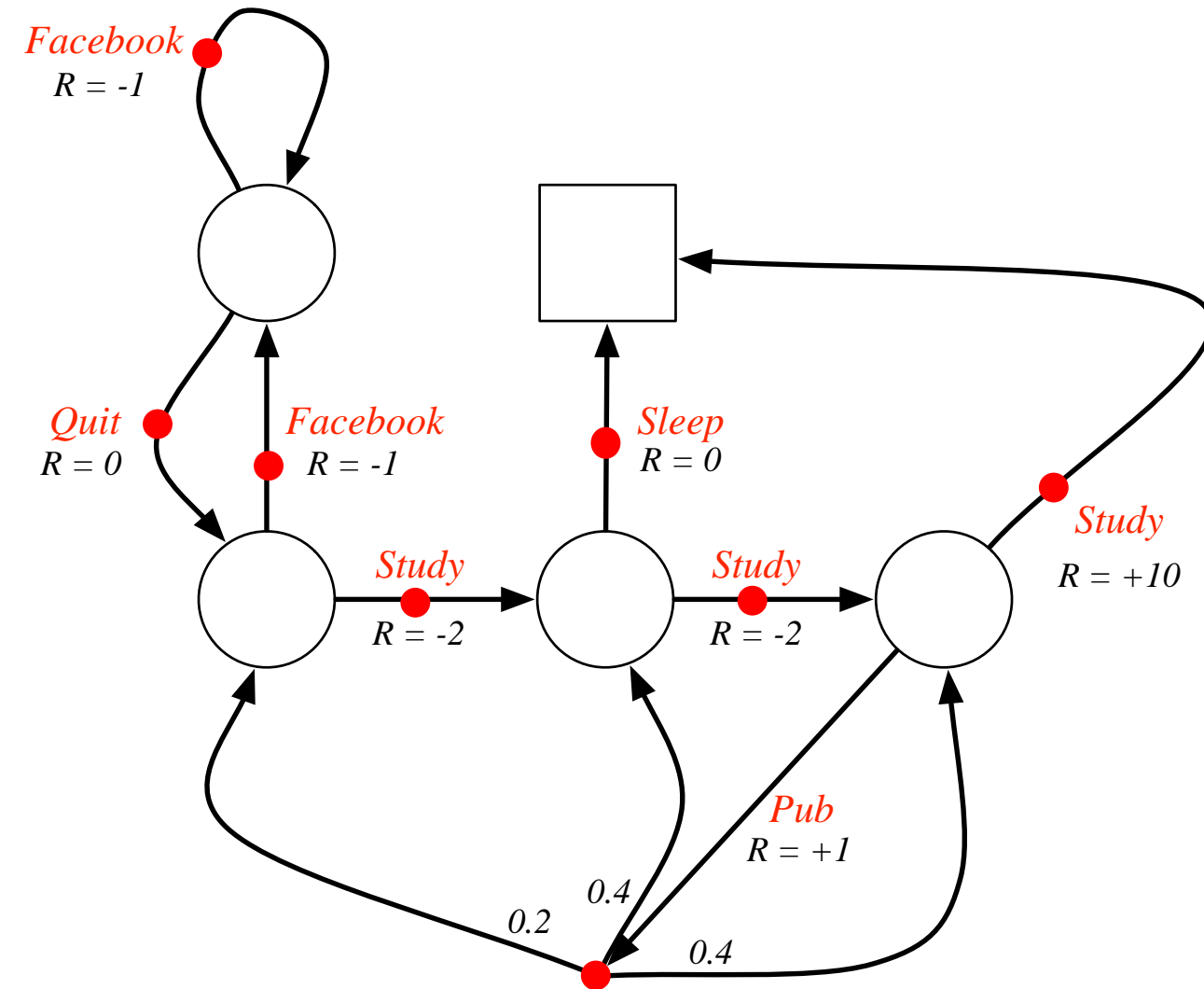
Action value function

(expected return from state s , taking action a , and following policy π)

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

$$q_{\pi}(s, a) = E_{\pi}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1})|s, a]$$

Example from David Silver, UCL, 2015



“Backup” property of state value functions

$$\begin{aligned} v_{\pi}(s) &\triangleq E_{\pi}[G_t | S_t = s] \quad \text{where } G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots \\ &= E_{\pi}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \dots | S_t = s] \\ &= E_{\pi}[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} \dots) | S_t = s] \\ &= E_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= E_{\pi}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s] \end{aligned}$$

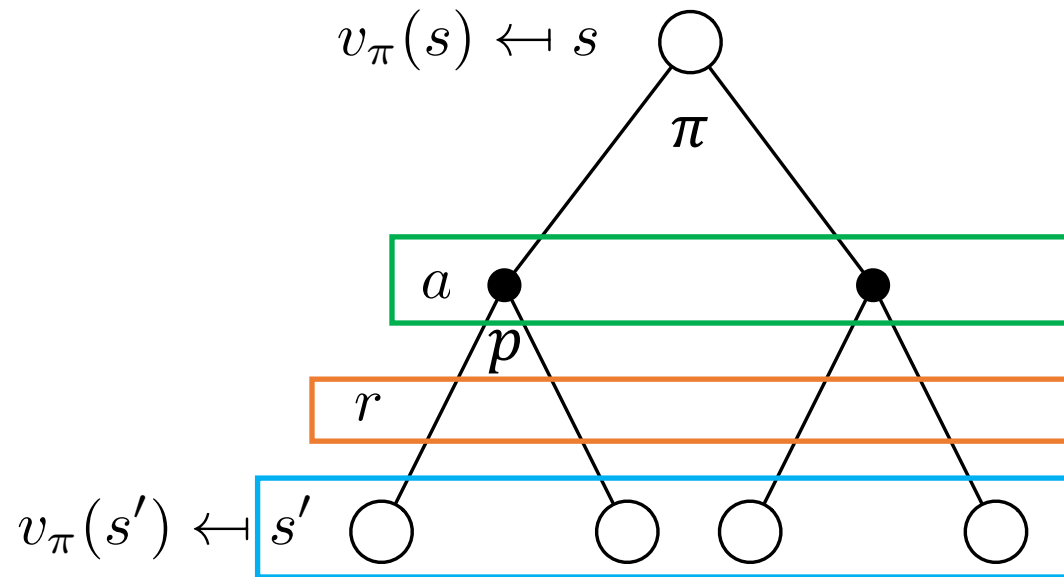
This recursive relationship is a version of the **Bellman Equation**

Example from David Silver, UCL, 2015

Bellman Expectation Equations for the **state value** function

(expected return from state s , and following policy π)

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



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

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

Expectation over the possible actions

Expectation over the rewards

(based on state and choice of action)

Expectation over the next possible states

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

$p(s', r | s, a)$ is the joint distribution of state transitions and rewards given you start in state s and take action a

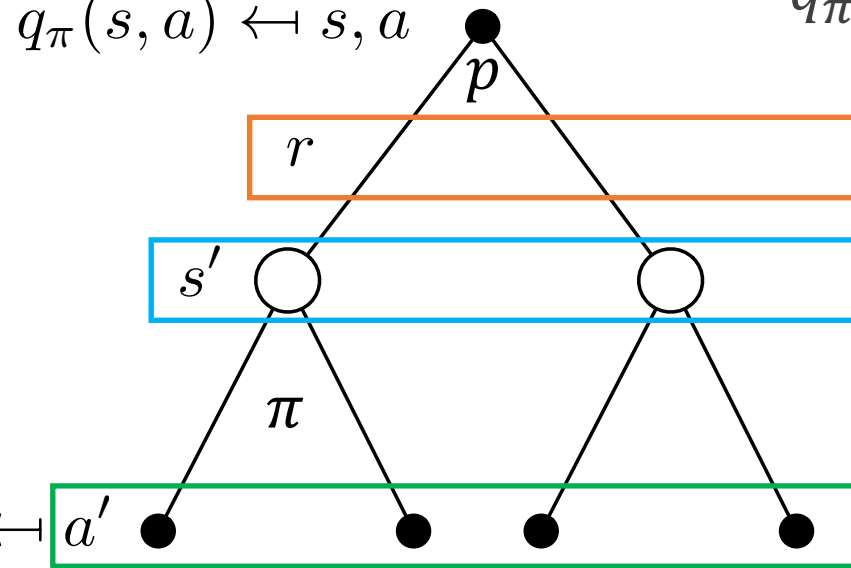
Bellman Expectation Equations for the **action value** function

(expected return from state s , taking action a , then following policy π)

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

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

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



Expectation over the rewards

(based on state and choice of action)

Expectation over the next possible states

Expectation over the possible actions

$$q_{\pi}(s, a) = \sum_{\underline{s'}} \sum_{\underline{r}} p(s', r | s, a) \left[r + \gamma \sum_{\underline{a'}} \pi(a' | s') q_{\pi}(s', a') \right]$$

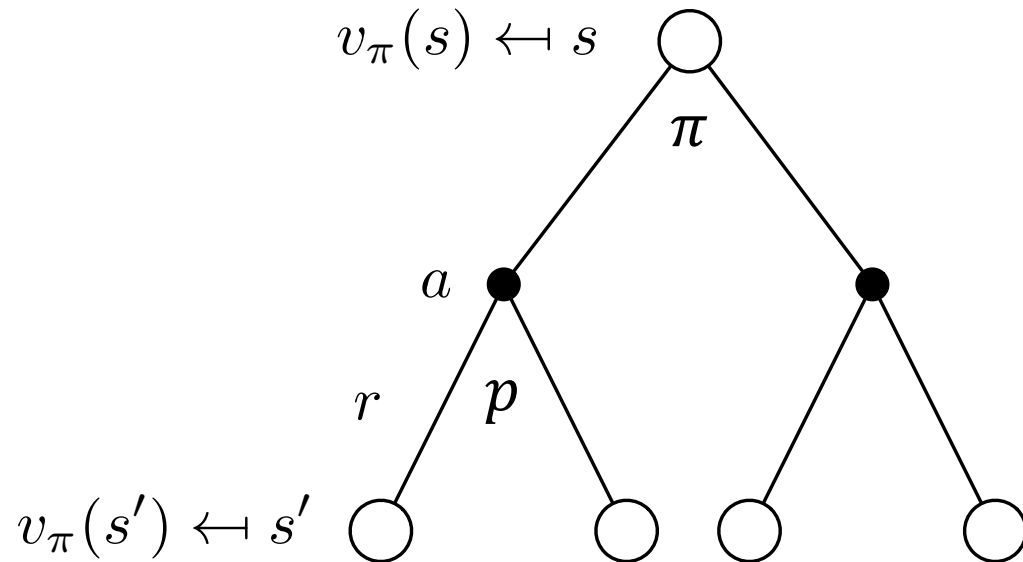
Bellman Expectation Equations

State value function

(expected return from state s , and following policy π)

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

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



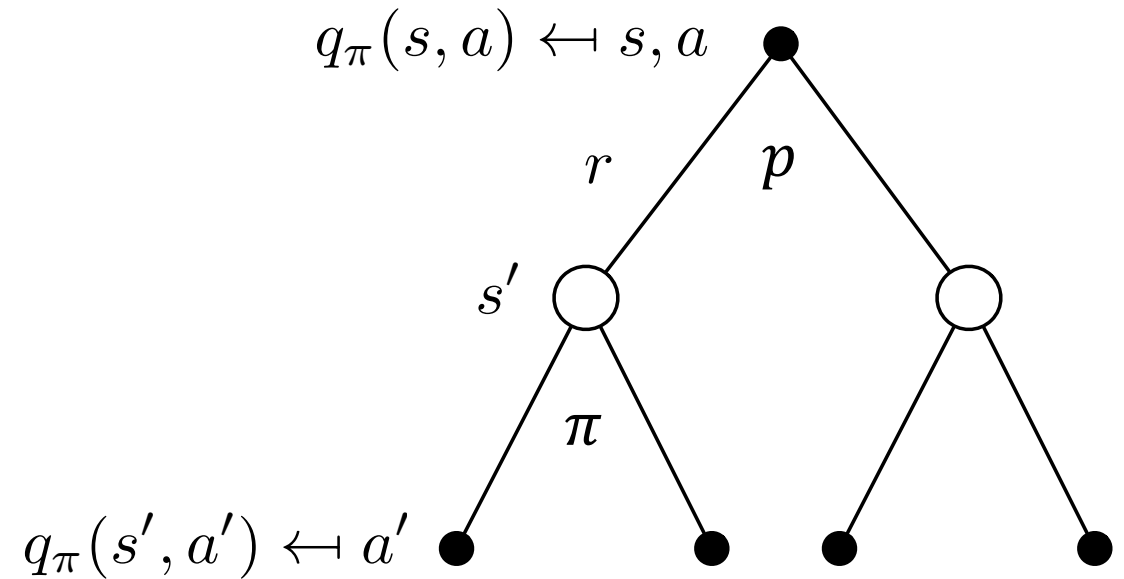
$$v_{\pi}(s) = \sum_a \pi(a|s) \underbrace{\sum_{s'} \sum_r p(s', r | s, a) [r + \gamma v_{\pi}(s')]}_{q_{\pi}(s, a)}$$

Action value function

(expected return from state s , taking action a , then following policy π)

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

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

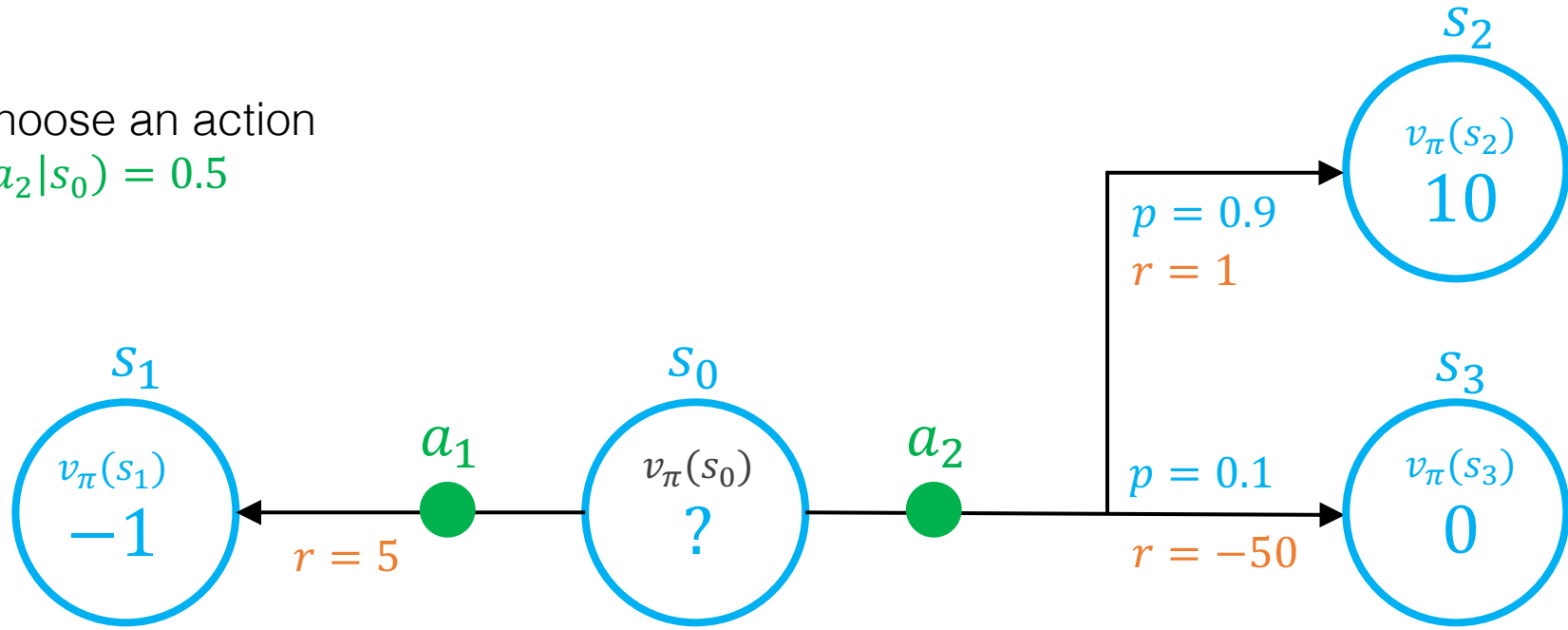


$$q_{\pi}(s, a) = \sum_{s'} \sum_r p(s', r | s, a) \left[r + \gamma \underbrace{\sum_{a'} \pi(a' | s') q_{\pi}(s', a')}_{v_{\pi}(s')} \right]$$

Example

Policy: randomly choose an action

$$\pi(a_1|s_0) = \pi(a_2|s_0) = 0.5$$



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

$$v_{\pi}(s_0) = \underbrace{\pi(a_1|s_0)}_{0.5} \underbrace{(\underbrace{5}_{r} - \underbrace{1}_{v_{\pi}(s_1)})}_{q_{\pi}(s_0, a_1)} + \underbrace{\pi(a_2|s_0)}_{0.5} \underbrace{[(\underbrace{0.9}_{p})(\underbrace{1}_{r}) + (\underbrace{0.1}_{p})(\underbrace{-50}_{r}) + (\underbrace{0.9}_{p})(\underbrace{10}_{v_{\pi}(s_2)}) + (\underbrace{0.1}_{p})(\underbrace{0}_{v_{\pi}(s_3)})]}_{q_{\pi}(s_0, a_2)}] \quad \gamma = 1$$

3 Markov Decision Process

State value function, $v_{\pi}(s)$

Assumptions:

$$\pi(a|s) = \frac{1}{\text{\# of reachable states}}$$

(randomly select an action)

$$\gamma = 1$$

(no discounting)

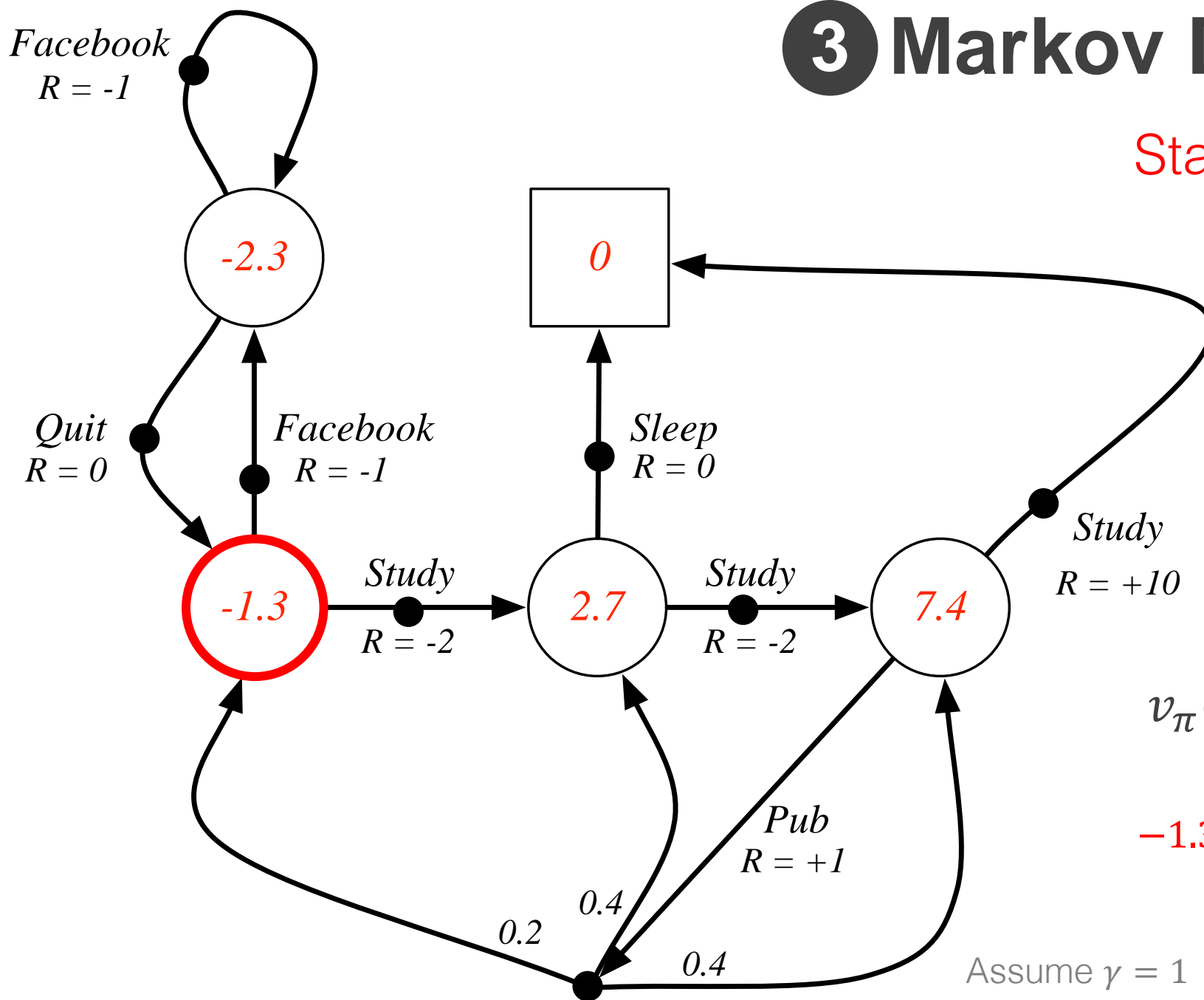
Bellman Expectation Equation

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

$$-1.3 \approx 0.5 \times (-1 - 2.3) + 0.5 \times (-2 + 2.7)$$

Assume $\gamma = 1$

Example from David Silver, UCL, 2015

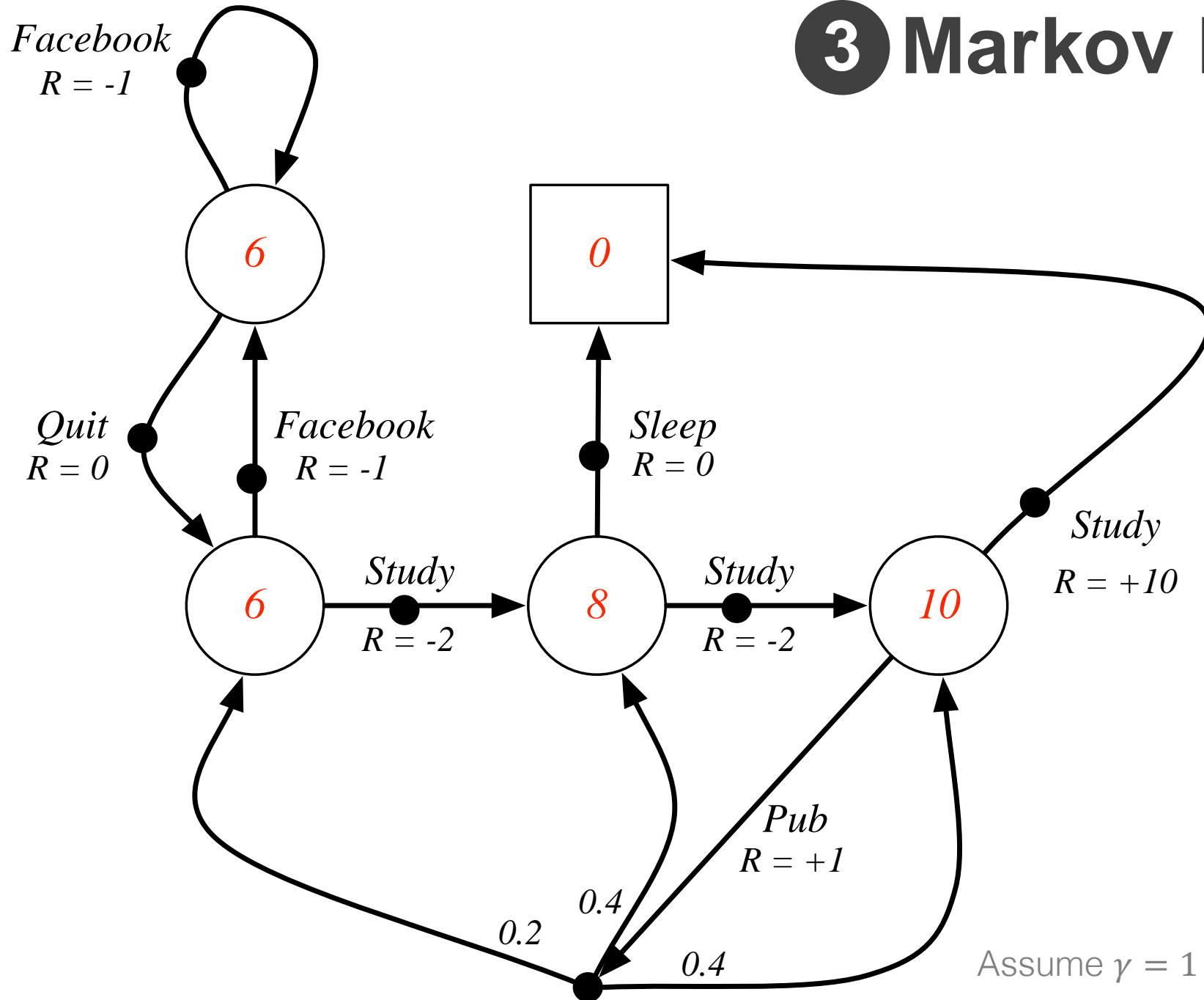


3 Markov Decision Process

Optimal **state-value** function, $v_*(s)$

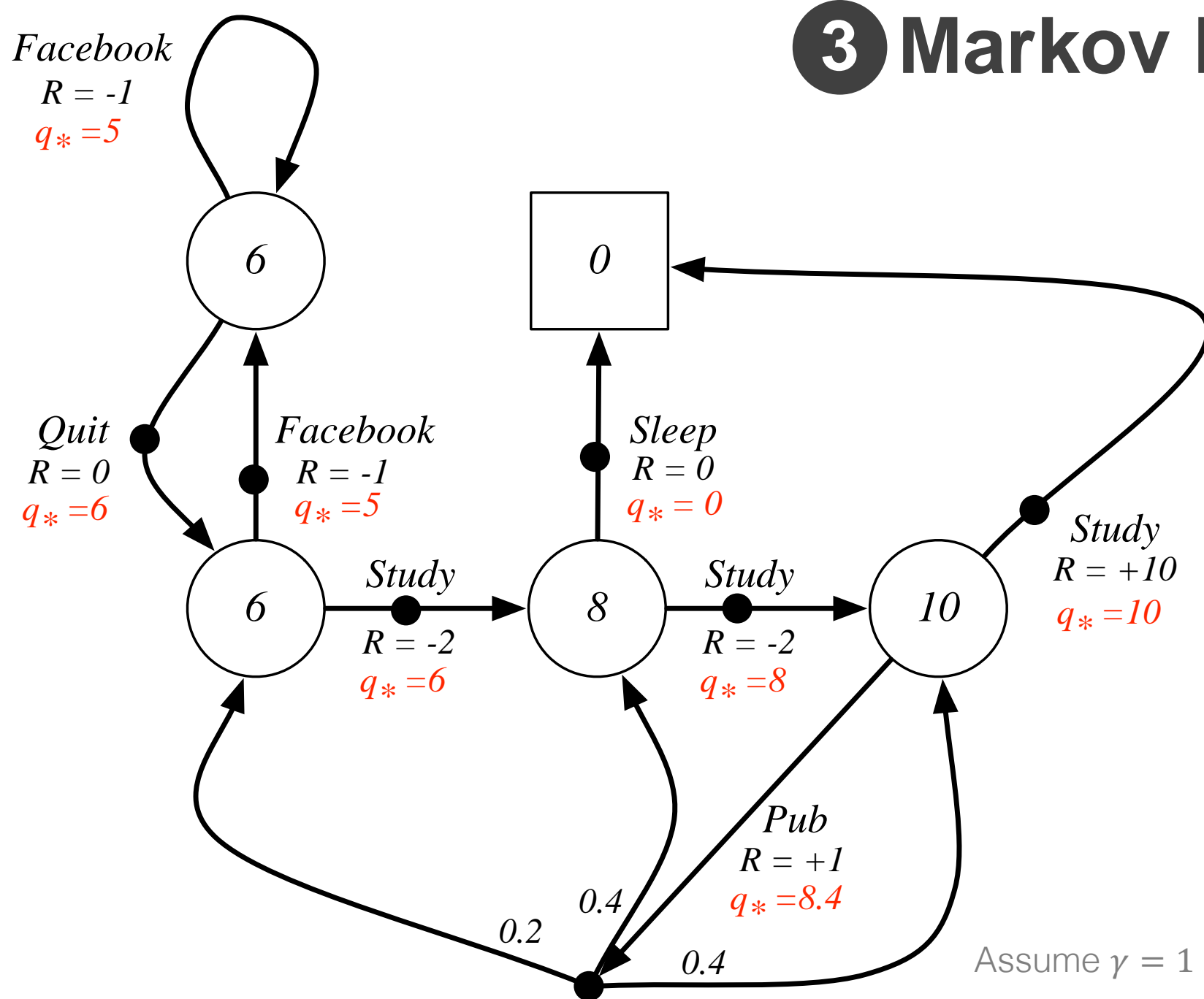
Maximum value function over all policies

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



Example from David Silver, UCL, 2015

3 Markov Decision Process



Optimal **state-value** function, $v_*(s)$

Maximum value function over all policies

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

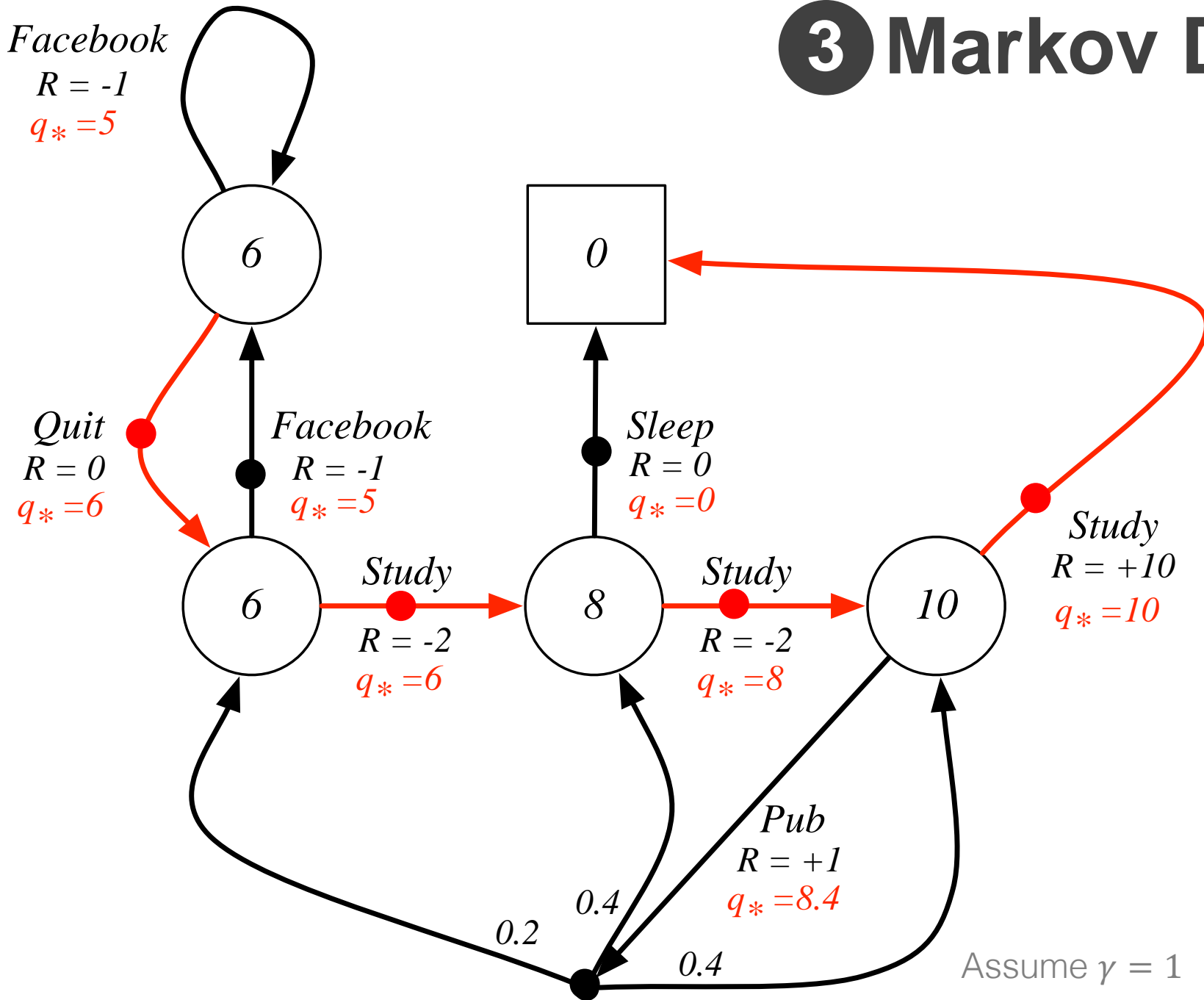
Optimal **action-value** function, $q_*(s, a)$

Maximum value function over all policies

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

Example from David Silver, UCL, 2015

3 Markov Decision Process



Optimal **policy**, $\pi_*(s)$
Which action to take at each moment

$$\pi_*(s) = \arg \max_a q_*(s, a)$$

Once we have the optimal value functions, we've "solved" the MDP!

Example from David Silver, UCL, 2015

Building blocks for the full RL problem

1	Markov Chain	{state space S , transition probabilities P }
2	Markov Reward Process (MRP)	{ S , P , + rewards R , discount rate γ } adds rewards (and values)
3	Markov Decision Process (MDP)	{ S , P , R , γ , + actions A } adds decisions (i.e. the ability to control)

- RL methods **do NOT ASSUME** knowledge of P or R (while dynamic programming does)
- RL learns/approximates that knowledge

Adapted from David Silver, 2015

Reinforcement Learning Roadmap

Knowledge of **Environment**

Perfect knowledge

Known Markov
Decision Process



No knowledge

Must learn from
experience

1

Core concepts in reinforcement learning

Actions, Rewards, Value, Environments, and Policies

2

Markov decision processes

...and Markov chains and Markov reward processes

3

Dynamic Programming

How do we find optimal policies?
(Bellman equations)

4

Monte Carlo Control

How do we estimate our value functions?
How do we use the value functions to choose actions?
How do we learn optimal policies from experience?