

Reinforcement Learning III

Lecture 22

Reminder: Resources

Sutton and Barto, 1998

Reinforcement Learning: An Introduction

Draft of 2018 edition available free online:

<http://www.incompleteideas.net/book/the-book-2nd.html>

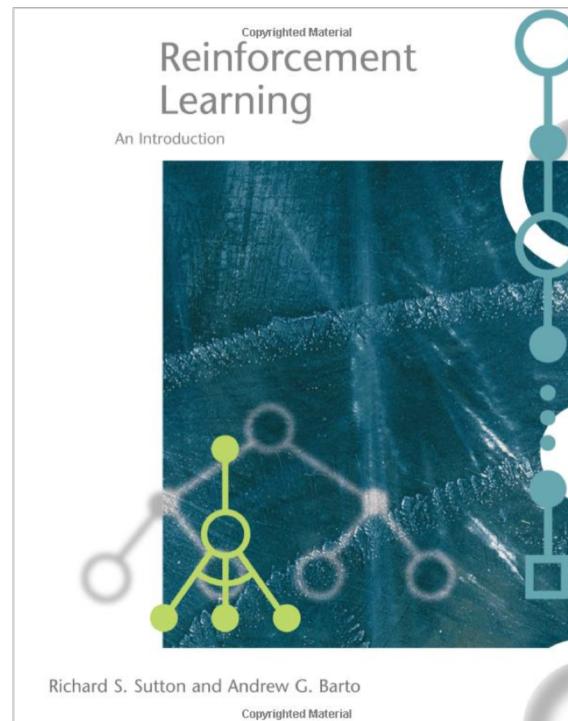


Image from Amazon.com (where the book may be purchased)

This reinforcement learning series draws heavily on these resources

David Silver, 2015

University College London

Advanced Topics 2015 (COMPM050/COMPGI13)

Course website:

<http://www0.cs.ucl.ac.uk/staff/D.Silver/web/Teaching.html>

Video series:

<https://www.youtube.com/watch?v=2pWv7GOvuf0&list=PL7-jPKtc4r78-wCZcQn5lqyuWhBZ8fOxT>

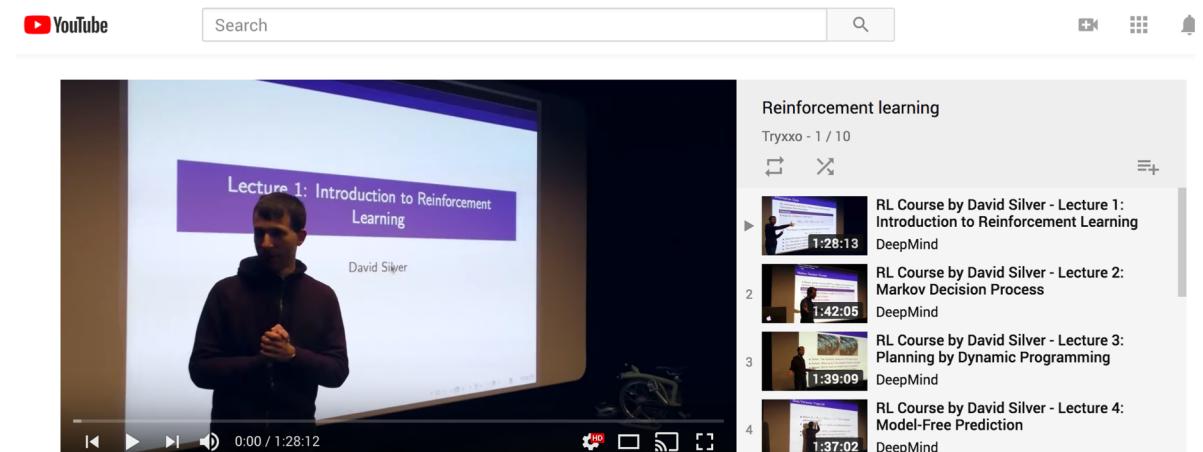


Image from Youtube.com

Markov Models

	States are Fully Observable	States are Partially Observable
Autonomous (no actions; make predictions)	Markov Chain	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, bioinformatics

MDPs: used extensively for reinforcement learning

Building blocks for the full RL problem

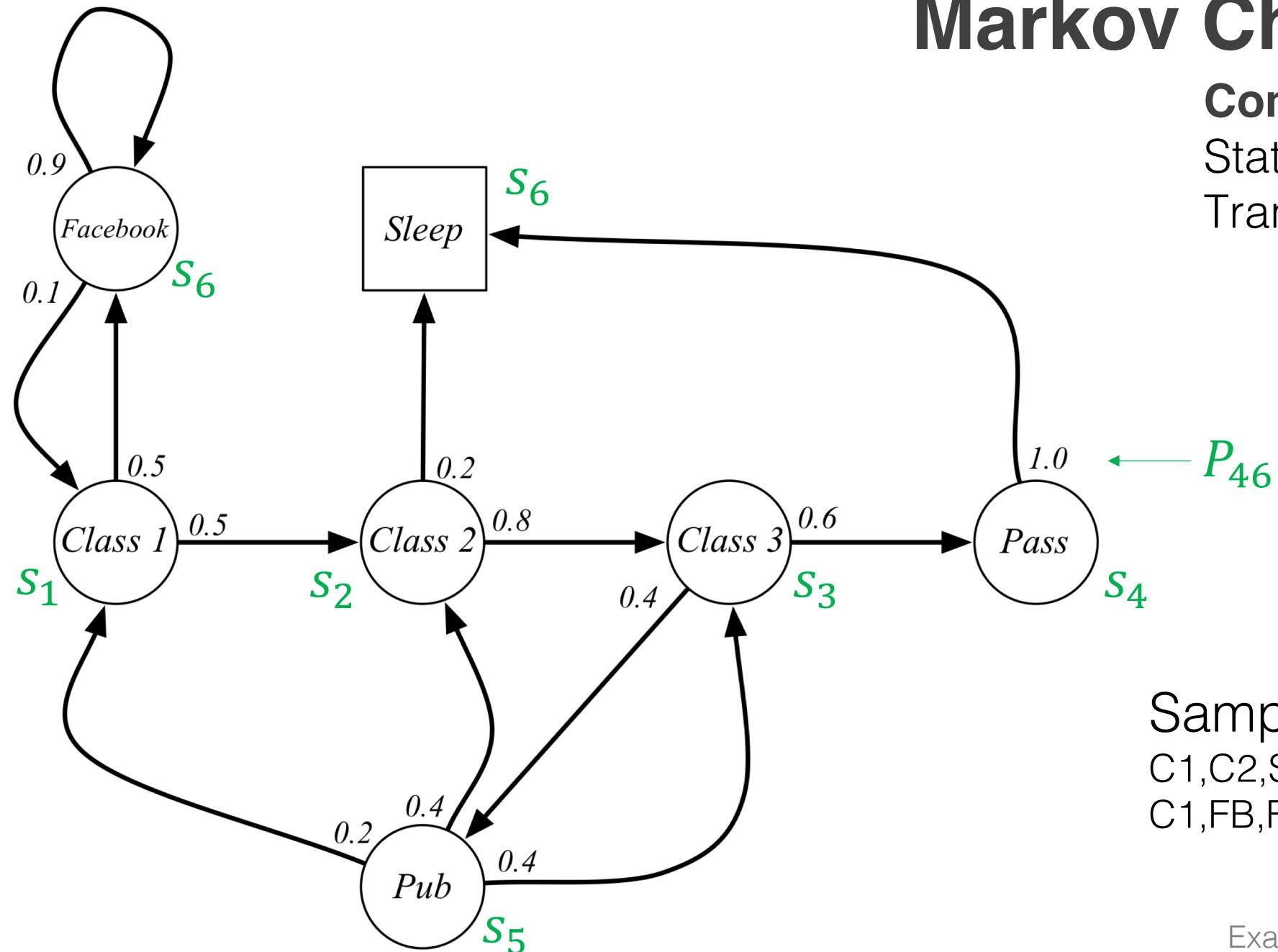
1
2
3

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

MDPs form the basis for most reinforcement learning environments

Adapted from David Silver, 2015

Markov Chain Example



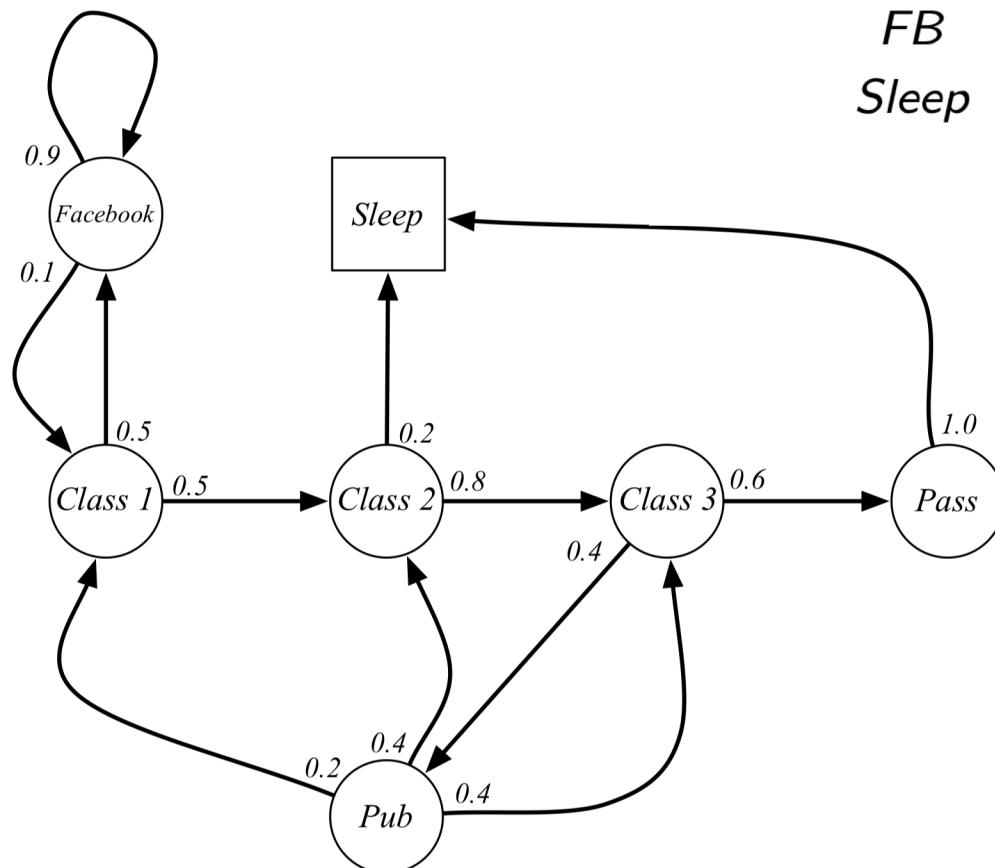
Components:

State space S ,
Transition probabilities P

Sample Episodes:
 $C_1, C_2, Sleep$
 $C_1, FB, FB, FB, C_1, C_2, C_3, Pass, Sleep$

Example from David Silver, UCL, 2015

Markov Chain



$$\mathcal{P} =$$

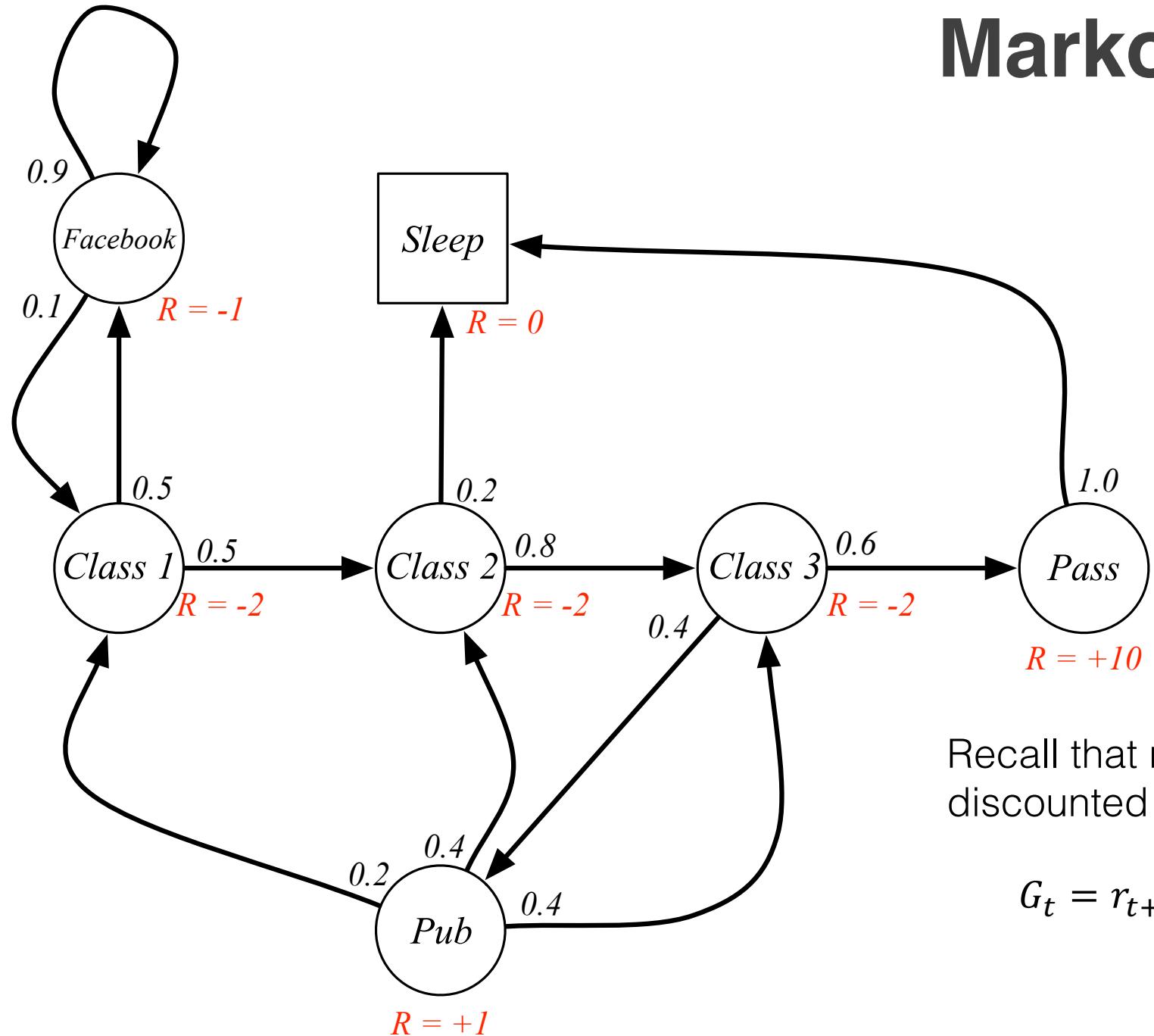
Pass
Pub
FB
Sleep

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

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

Example from David Silver, UCL, 2015

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 :

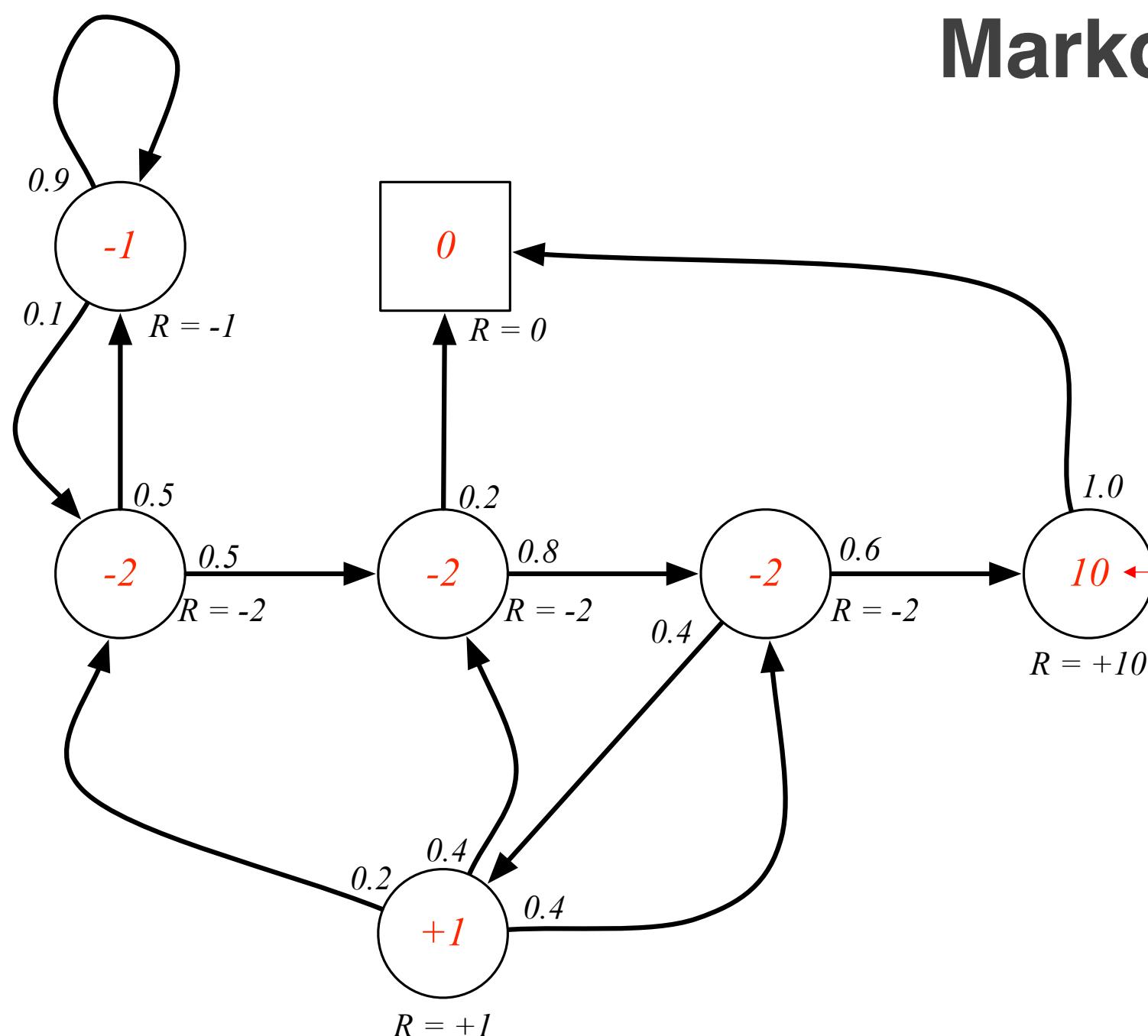
$$G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

Example from David Silver, UCL, 2015

Markov Reward Process

Components:

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



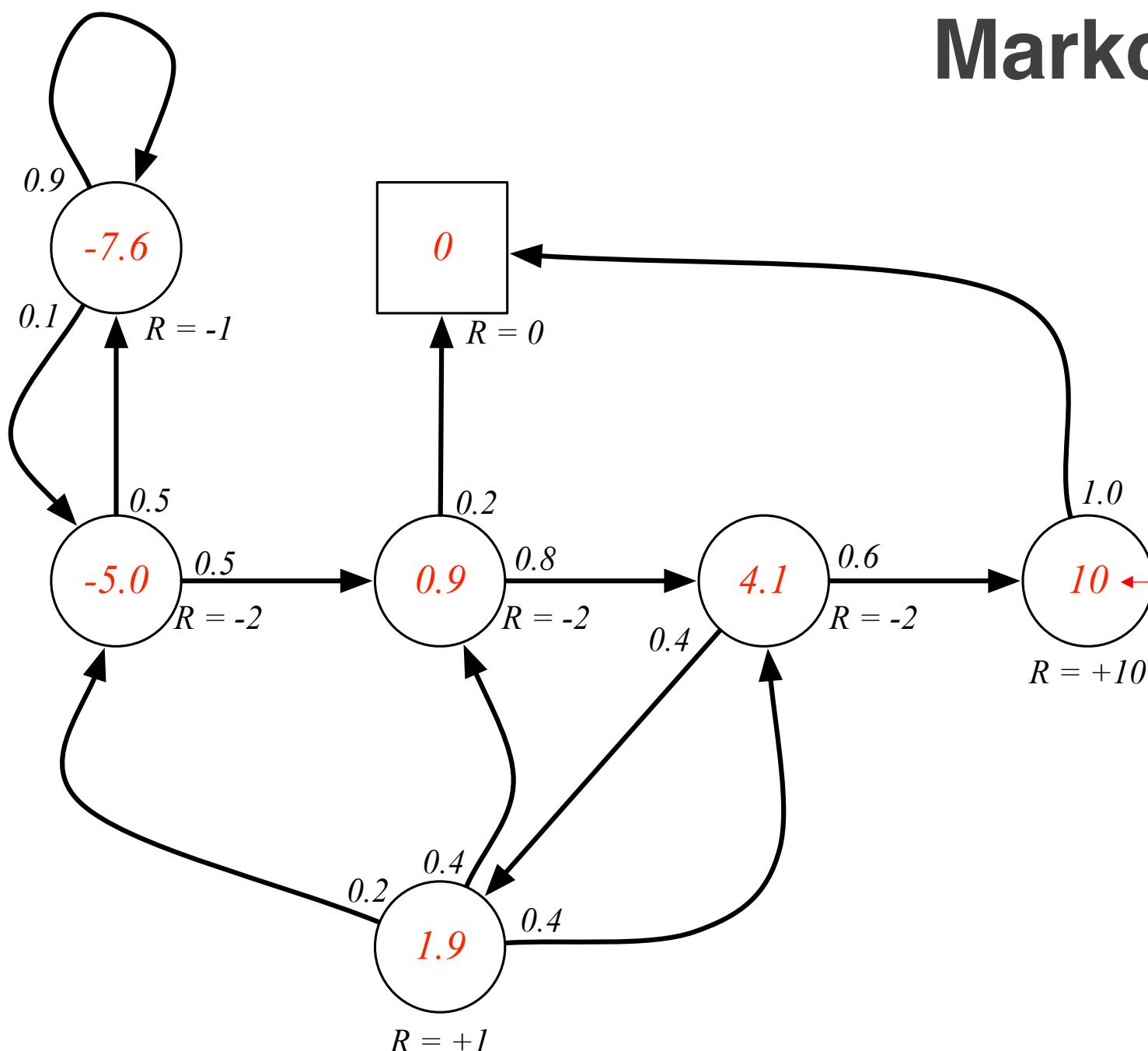
$$v(s) \text{ for } \gamma = 0$$

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

Markov Reward Process



Components:
State space S ,
Transition probabilities, P
Rewards, R
Discount rate, γ

$$v(s) \text{ for } \gamma = 0.9$$

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

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

$$= E[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} \dots | S = s_t]$$

$$= E[r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} \dots) | S = s_t]$$

$$= E[r_{t+1} + \gamma G_{t+1} | S = s_t]$$

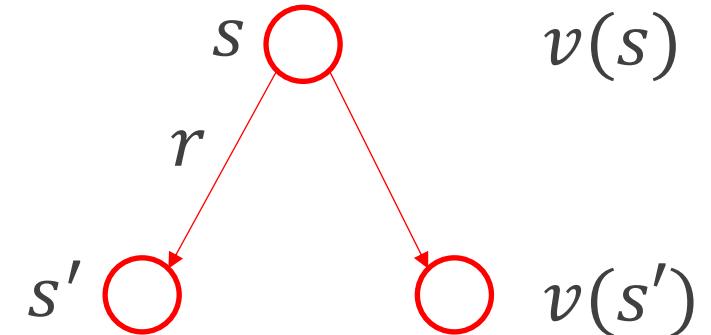
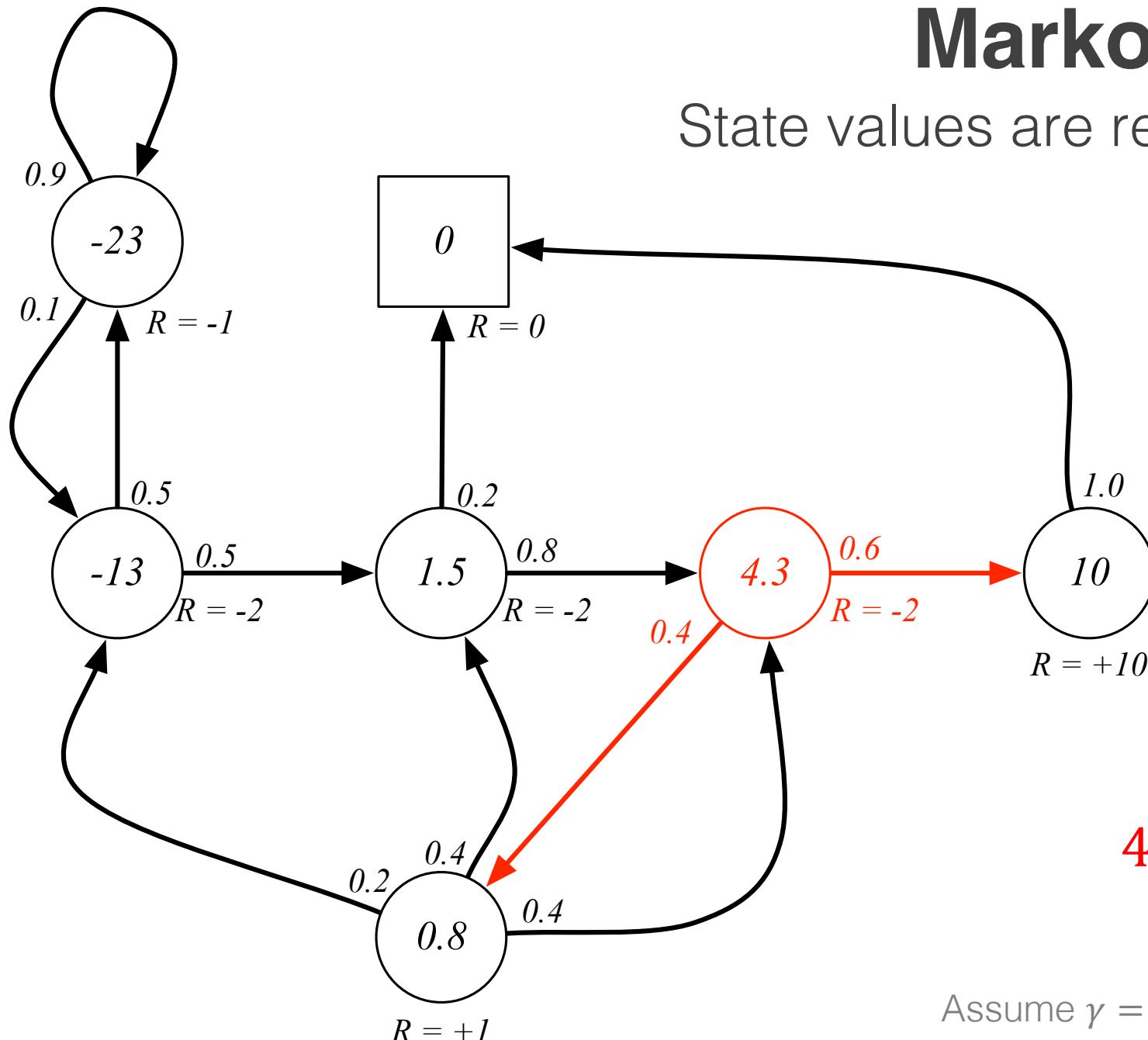
$$= E[r_{t+1} + \gamma v(s_{t+1}) | S = s_t]$$

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

Example from David Silver, UCL, 2015

Markov Reward Process

State values are related to neighboring states



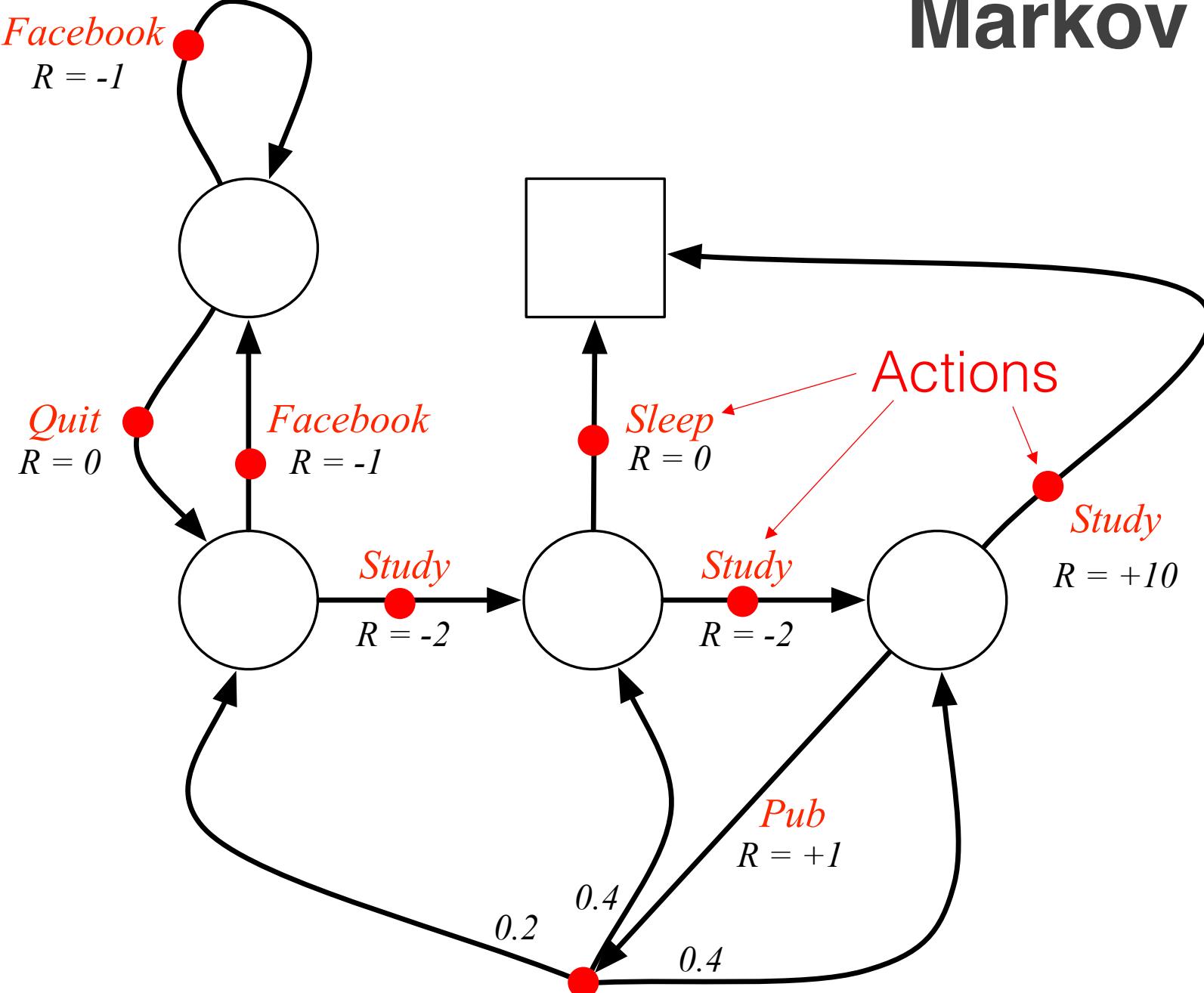
$$4.3 = -2 + 0.6 \times 10 + 0.4 \times 0.8$$

Notation: $s = s_t$ and $s' = s_{t+1}$

$$R_s = E[r_{t+1}|S_t = s]$$

Example from David Silver, UCL, 2015

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

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)

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

$$v_\pi(s) = E[R_s^a + \gamma v_\pi(s')|s]$$

Action value function

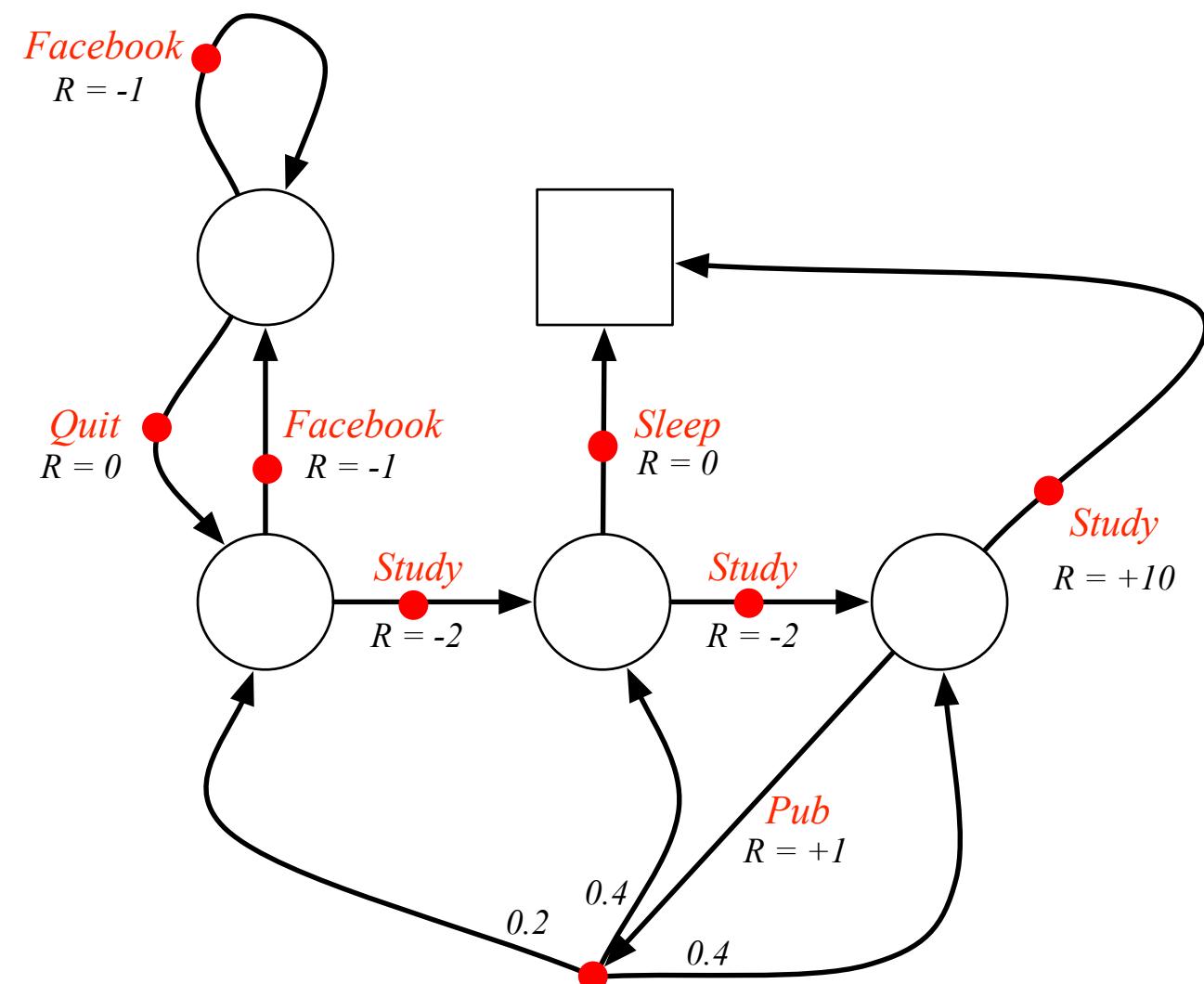
(expected return from state s , taking action a , and following policy π)

$$q_\pi(s, a) = E[G_t|s, a]$$

$$q_\pi(s, a) = E[R_s^a + \gamma q_\pi(s', a')|s, a]$$

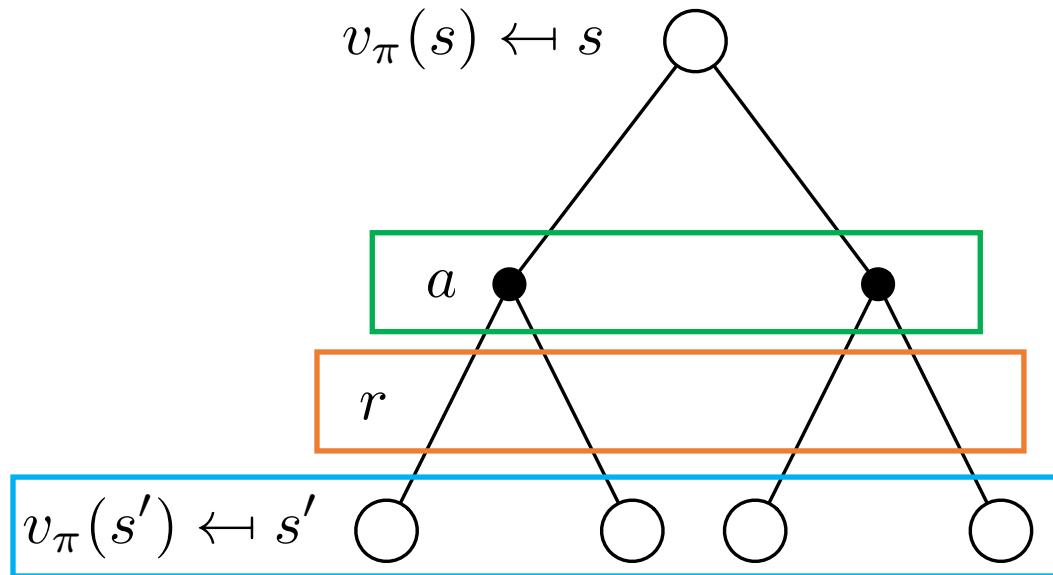
$$R_s^a = E[r_{t+1}|S_t = s, A_t = a]$$

Example from David Silver, UCL, 2015



Bellman Expectation Equations for the state value function

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



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

$$v_\pi(s) = E[R_s^a + \gamma v_\pi(s')|s]$$

$$R_s^a = E[r_{t+1}|S_t = s, A_t = a]$$

Expectation over the possible actions

Expectation over the rewards

(based on state and choice of action)

Expectation over the next possible states

$$v_\pi(s) = \sum_a \pi(a|s) \left(R_s^a + \gamma \sum_{s'} P_{ss'}^a v_\pi(s') \right)$$

— — —

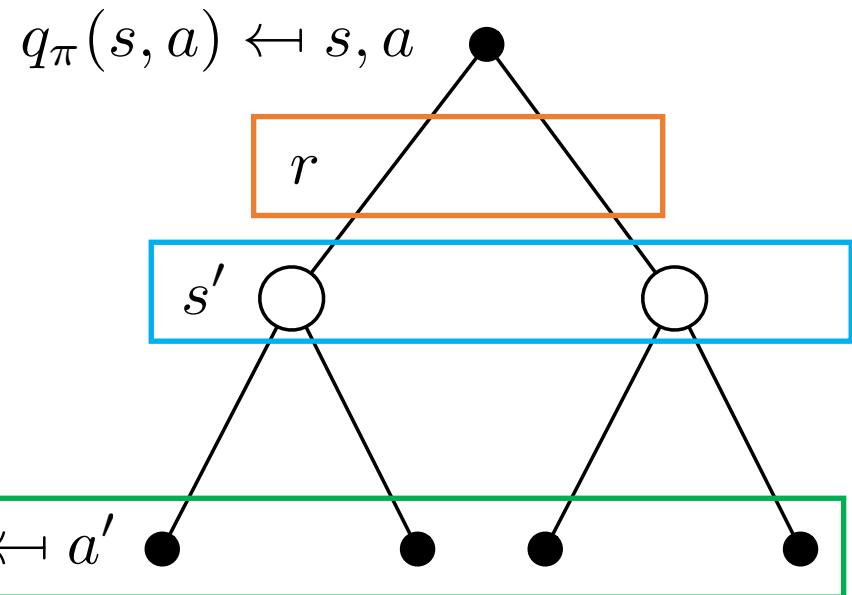
Bellman Expectation Equations for the action value function

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

$$q_\pi(s, a) = E[G_t | s, a]$$

$$q_\pi(s, a) = E[R_s^a + \gamma q_\pi(s', a') | s, a]$$

$$R_s^a = E[r_{t+1} | S_t = s, A_t = a]$$



$$q_\pi(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \sum_{a'} \pi(a'|s') q_\pi(s', a')$$

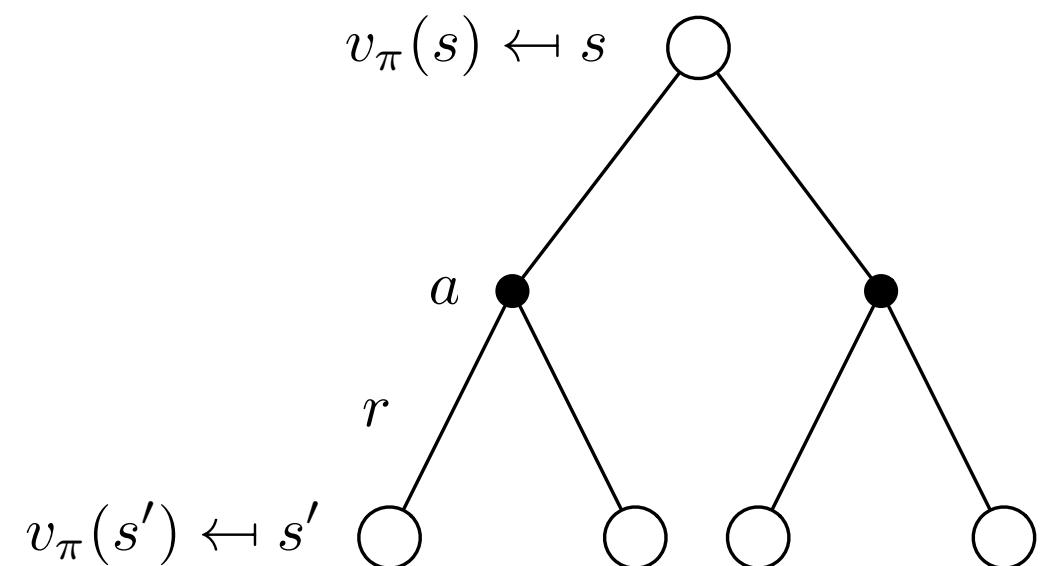
Bellman Expectation Equations

State value function

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

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

$$v_\pi(s) = E[R_s^a + \gamma v_\pi(s')|s]$$



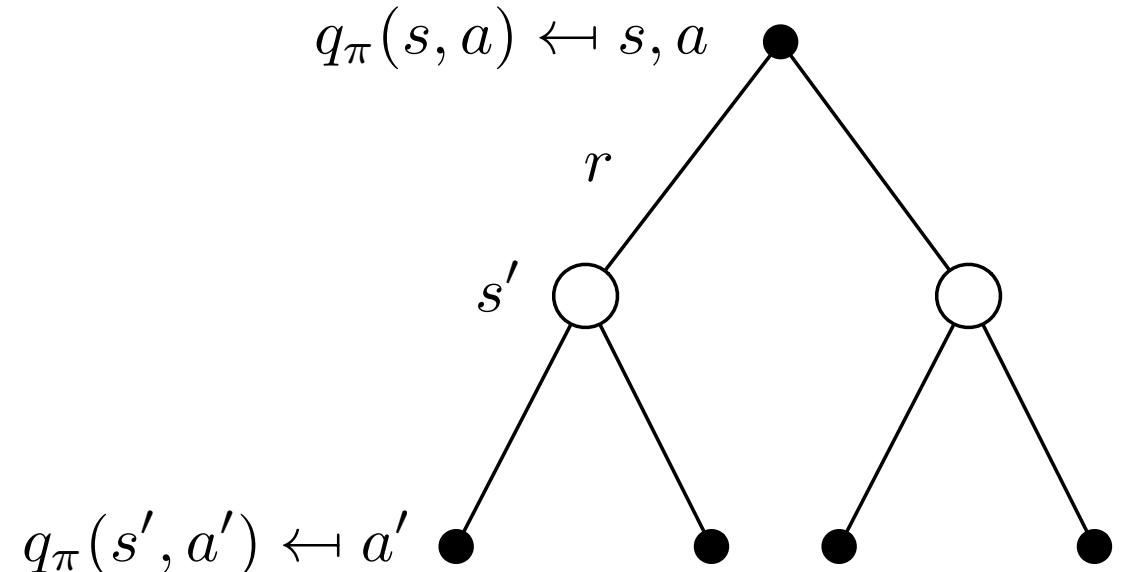
$$v_\pi(s) = \sum_a \pi(a|s) \left(R_s^a + \gamma \sum_{s'} P_{ss'}^a v_\pi(s') \right)$$

Action value function

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

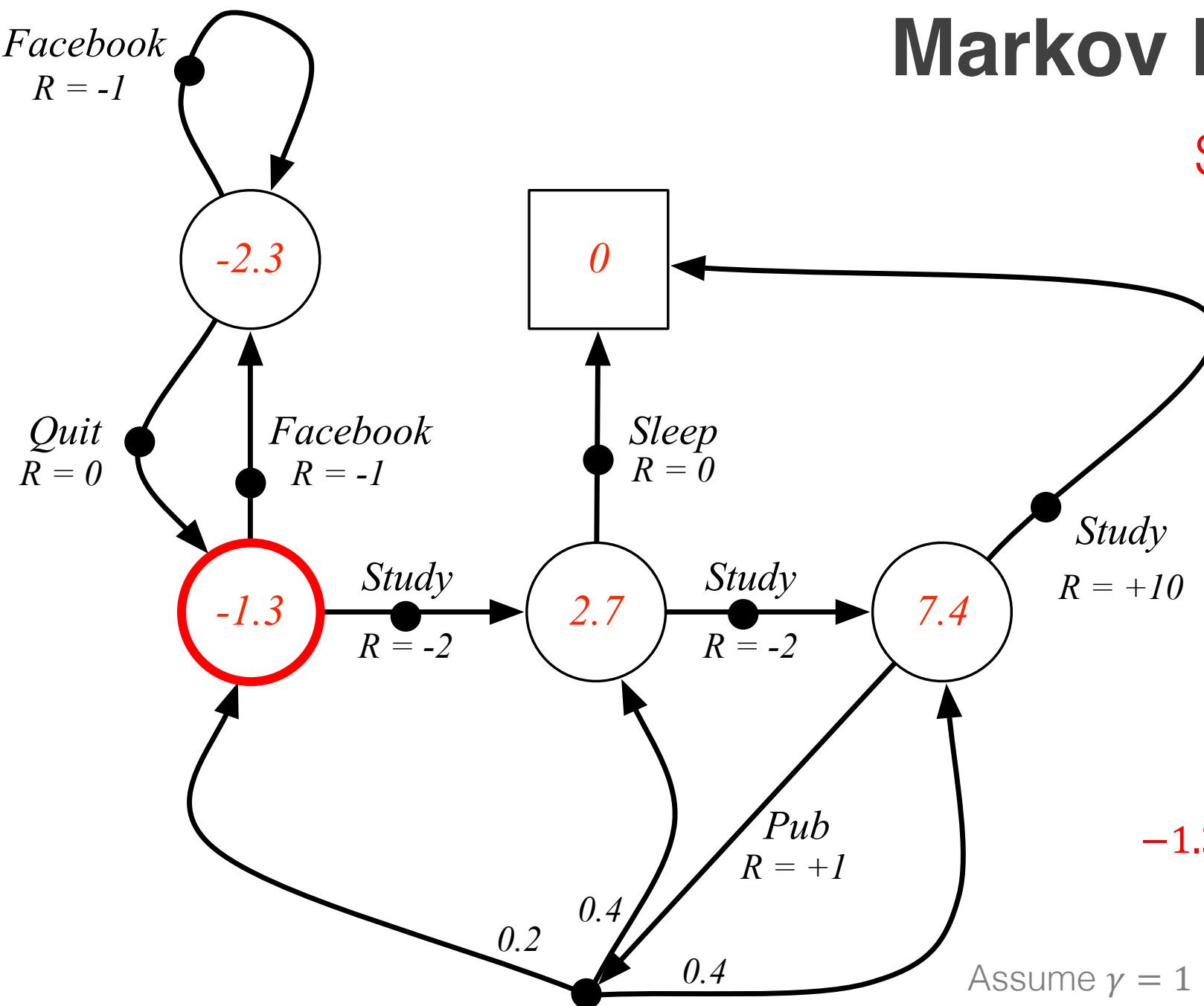
$$q_\pi(s, a) = E[G_t|s, a]$$

$$q_\pi(s, a) = E[R_s^a + \gamma q_\pi(s', a')|s, a]$$



$$q_\pi(s, a) = R_s^a + \gamma \sum_{s'} \sum_{a'} P_{ss'}^a \pi(a'|s') q_\pi(s', a')$$

Markov Decision Process



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

Assumptions:

$$\pi(a|s) = 0.5$$

(randomly select an action)

$$\gamma = 1$$

(no discounting)

Bellman Expectation Equation

$$v(s) = E[R_s^a + \gamma v(s')|s]$$

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

Assume $\gamma = 1$

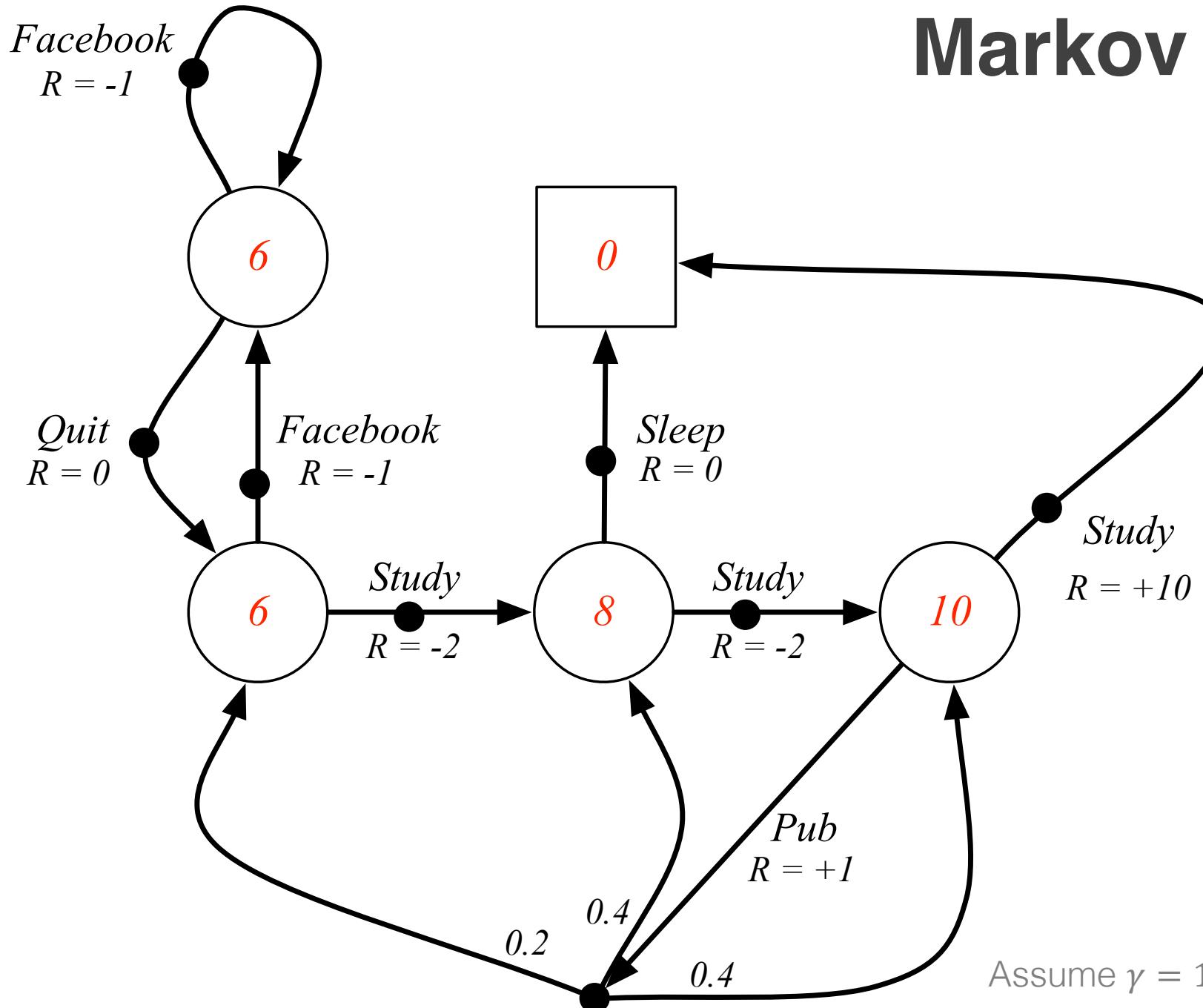
Example from David Silver, UCL, 2015

Markov Decision Process

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

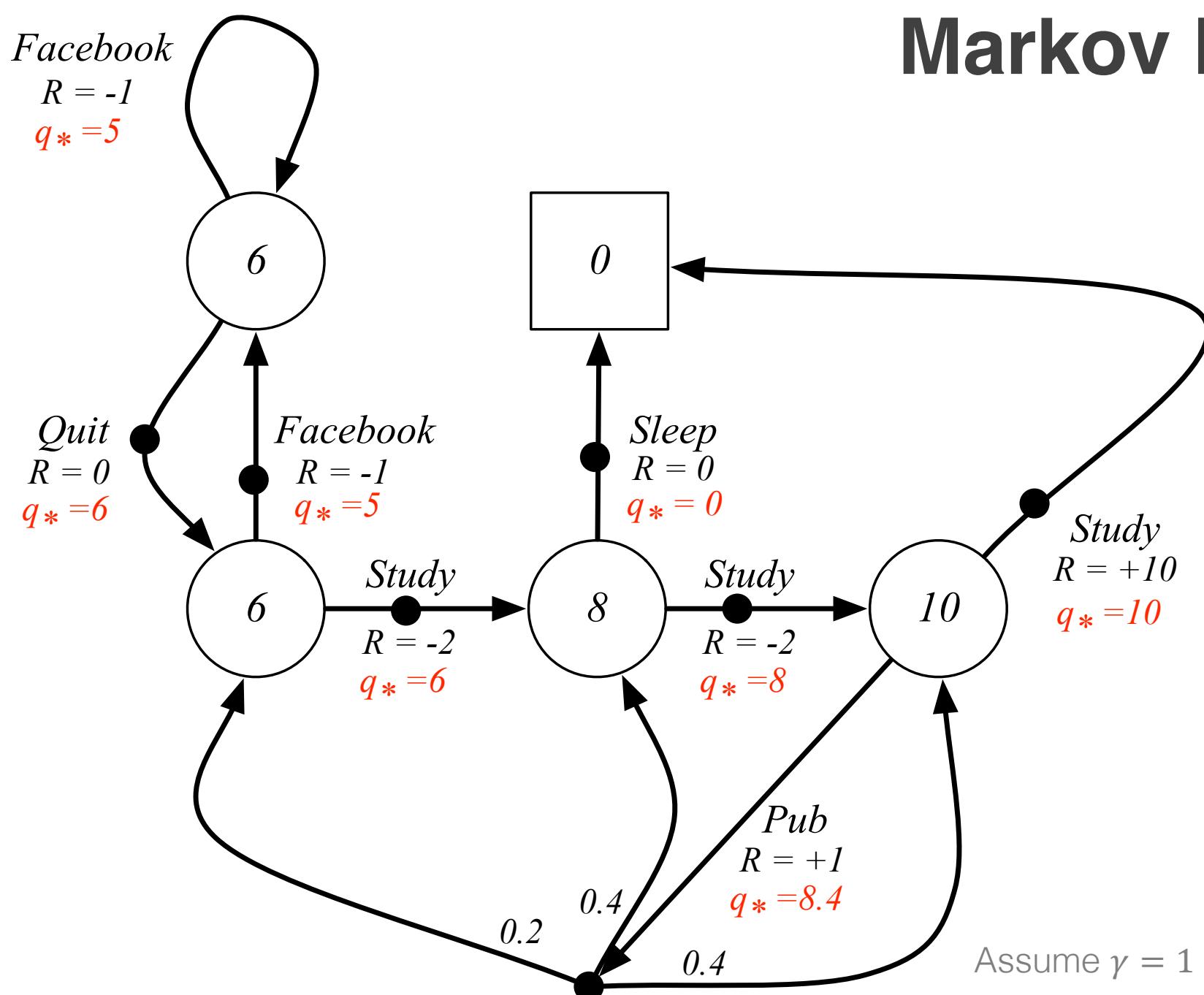
Maximum value function over all policies

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



Example from David Silver, UCL, 2015

Markov Decision Process



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

Maximum value function over all policies

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

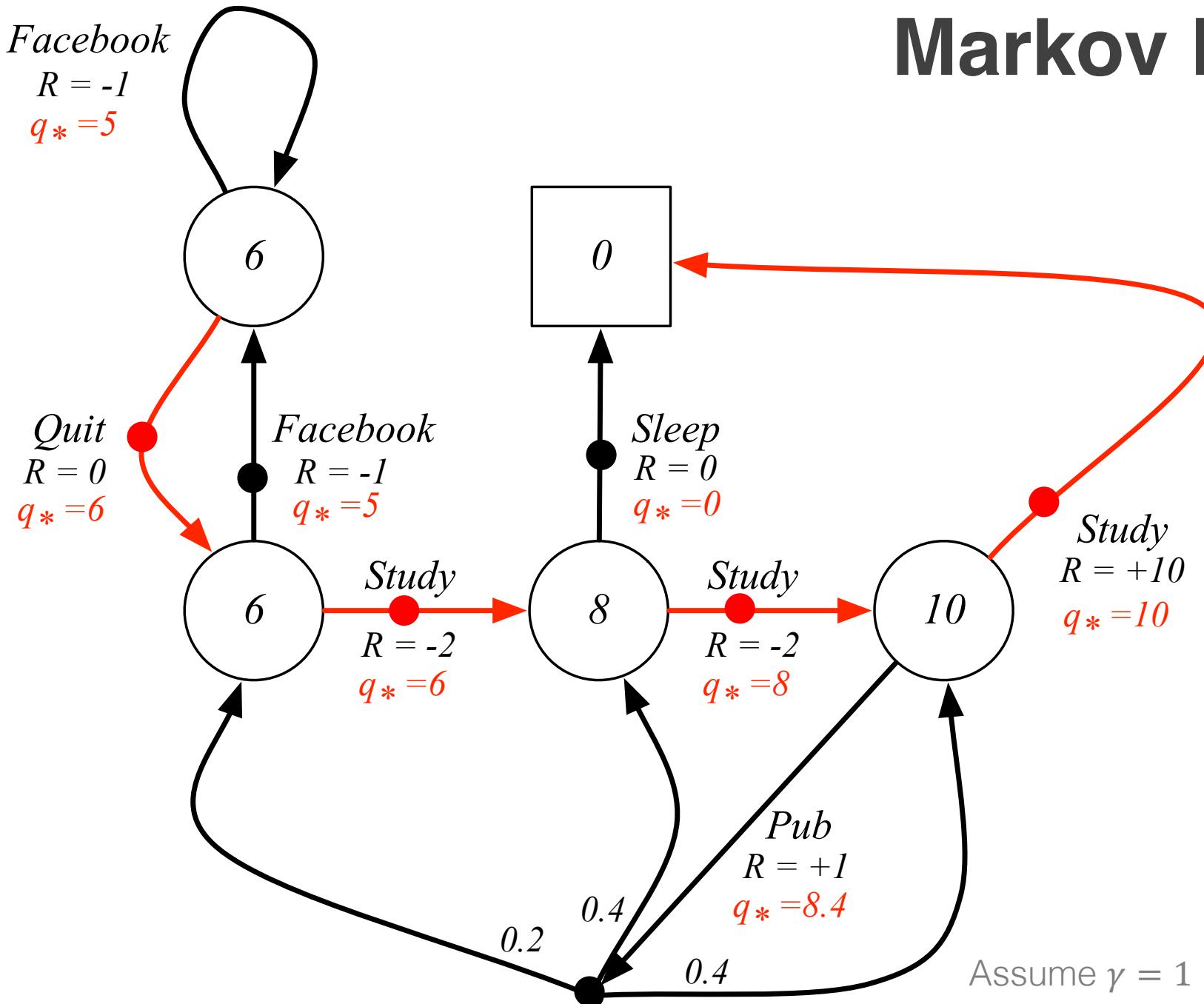
Optimal **action-value** function, $q_*(s, a)$

Maximum value function over all policies

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

Example from David Silver, UCL, 2015

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 our value functions, we've “solved” the MDP!

Example from David Silver, UCL, 2015

Learning strategy

Knowledge of **Environment**

No knowledge
Must learn from experience

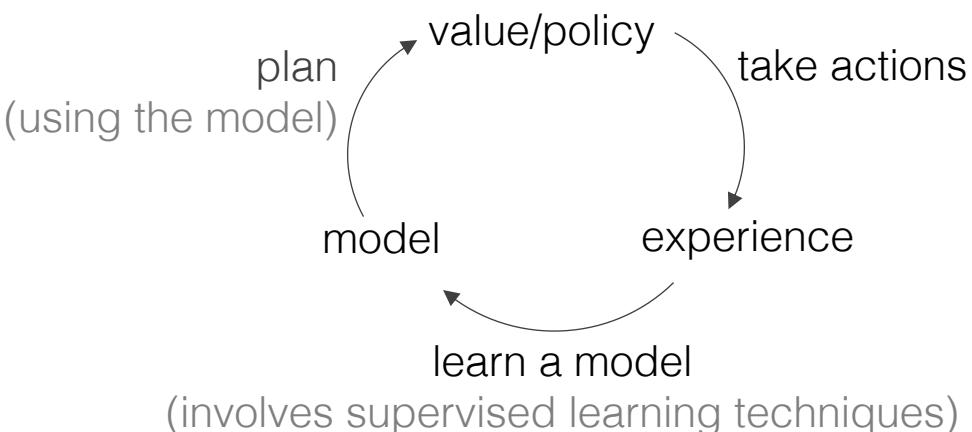
Perfect knowledge
Known MDP

Model-based
(planning)

Model-free

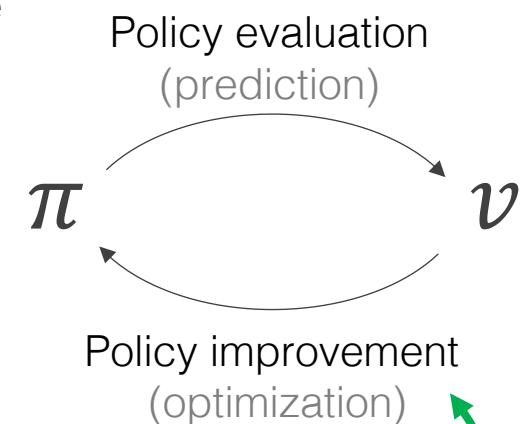
Reinforcement Learning

Simulation-based search



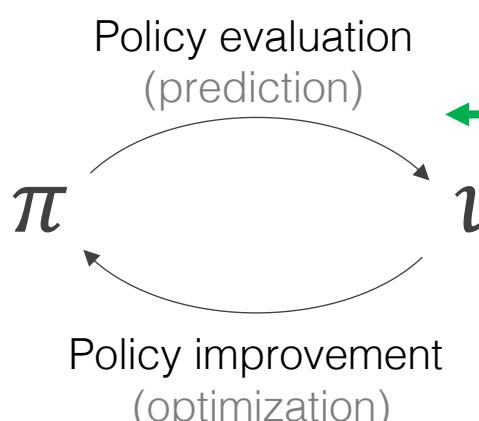
Monte Carlo Learning

Temporal Difference
Learning



Dynamic Programming

Policy iteration
Value iteration



Next class:

1. How to **compute optimal policies** for known MDPs?
2. How to extend this to the case **without full knowledge** of MDPs (model-free learning)