+

# Machine Learning and Data Mining

## Reinforcement Learning
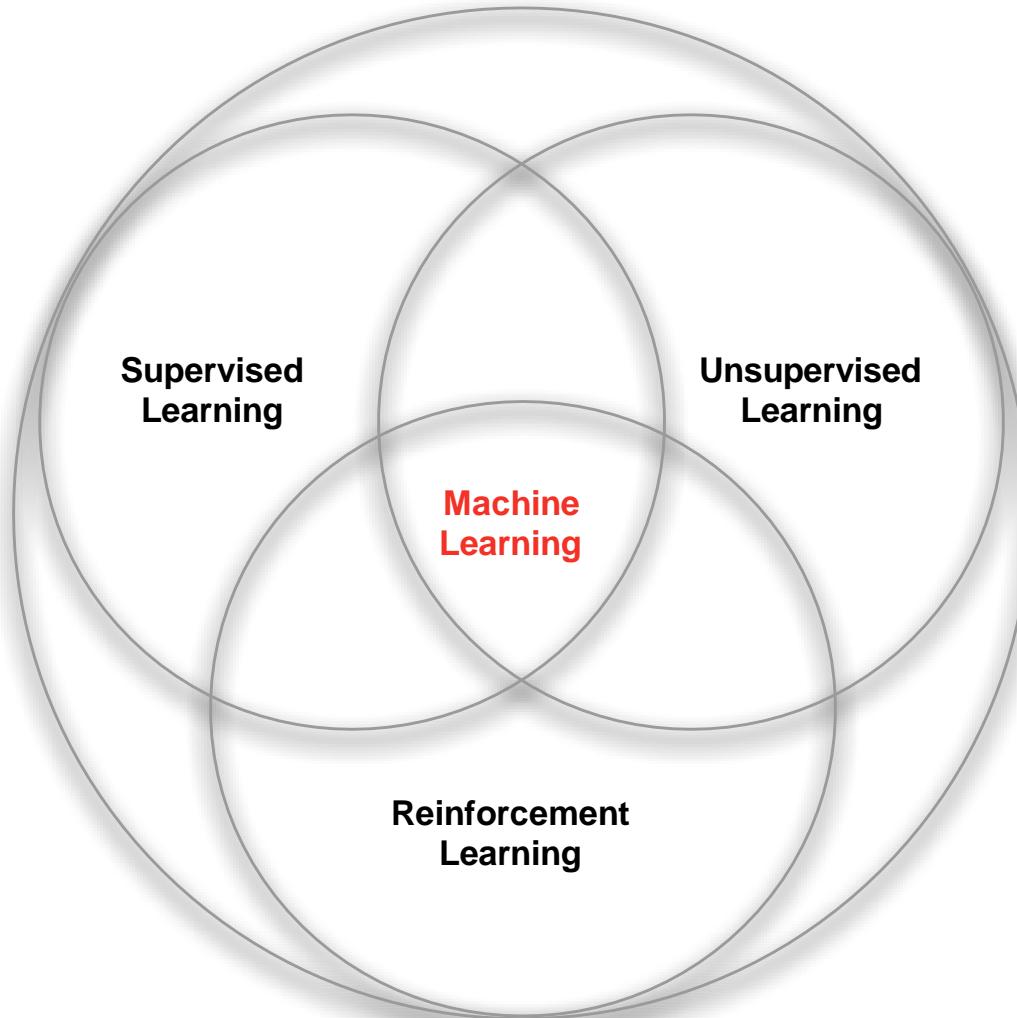## Markov Decision Processes

Kalev Kask

# Overview

- Intro
- Markov Decision Processes
- Reinforcement Learning
  - Sarsa
  - Q-learning
- Exploration vs Exploitation tradeoff

# Resources

- [Book: Reinforcement Learning: An Introduction](#)
  Richard S. Sutton and Andrew G. Barto

- [UCL Course on Reinforcement Learning](#)
  David Silver
  - [https://www.youtube.com/watch?v=2pWv7GOvuf0](https://www.youtube.com/watch?v=2pWv7GOvuf0)
  - [https://www.youtube.com/watch?v=lfHX2hHRMVQ](https://www.youtube.com/watch?v=lfHX2hHRMVQ)
  - [https://www.youtube.com/watch?v=Nd1-UUMVfz4](https://www.youtube.com/watch?v=Nd1-UUMVfz4)
  - [https://www.youtube.com/watch?v=PnHCvfgC_ZA](https://www.youtube.com/watch?v=PnHCvfgC_ZA)
  - [https://www.youtube.com/watch?v=0g4j2k_Ggc4](https://www.youtube.com/watch?v=0g4j2k_Ggc4)
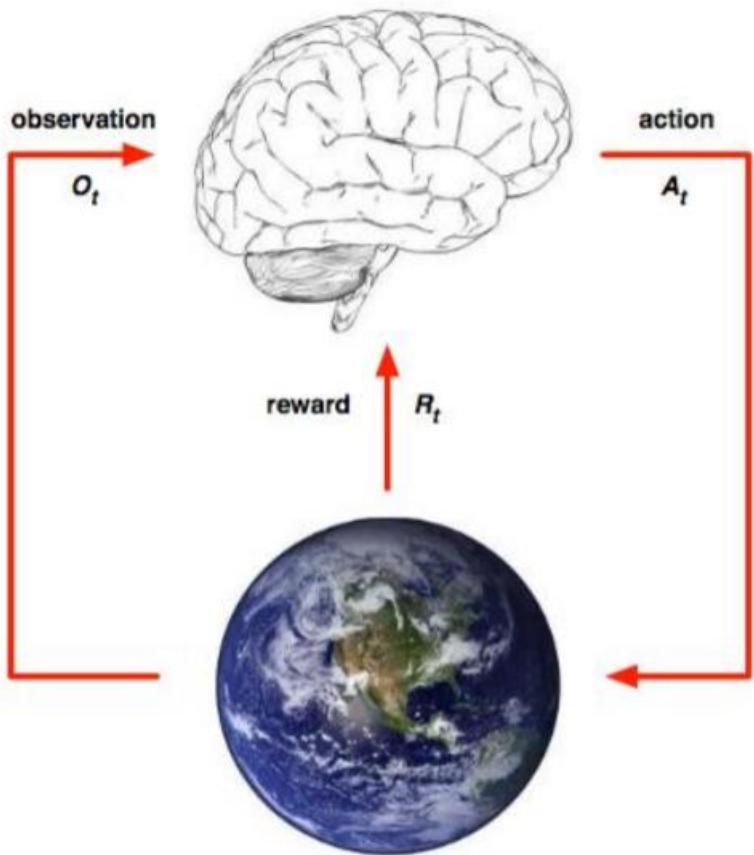  - [https://www.youtube.com/watch?v=UoPei5o4fps](https://www.youtube.com/watch?v=UoPei5o4fps)

# Why is it different

- No target values to predict
- Feedback in the form of rewards
  - May be delayed not instantaneous
- Have a goal : max reward
- Have timeline : actions along arrow of time
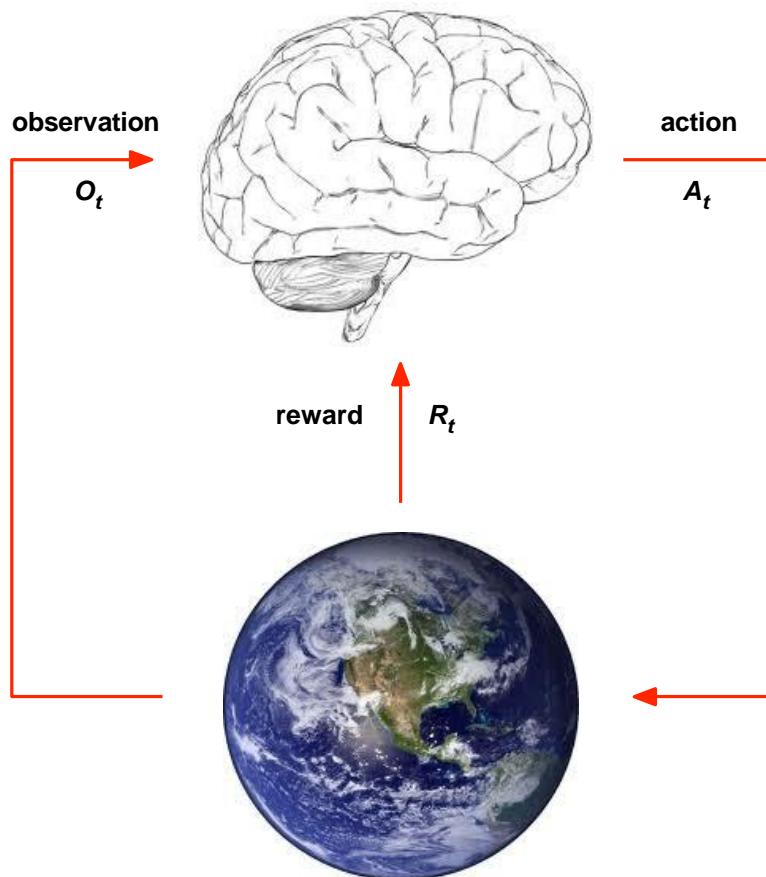- Actions affect what data it will receive

# Agent-Environment



observation $O_t$

action $A_t$

reward $R_t$

## Agent

- decides on an action
- receives next observation
- receives next reward

## Environment

- executes the action
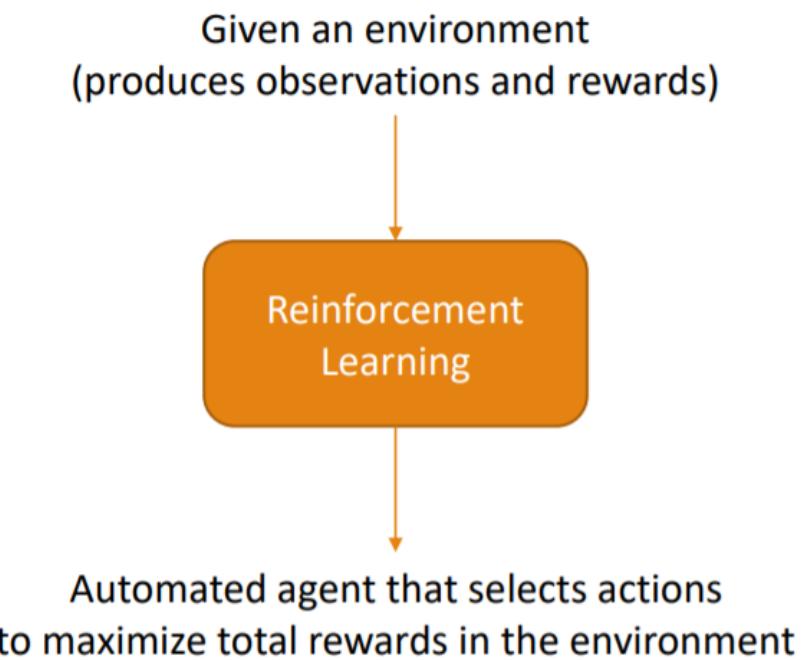- computes next observation
- computes next reward

- observation $O_t$
- action $A_t$
- reward $R_t$

- 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

# Sequential Decision Making

- Actions have long term consequences

- Goal maximize cumulative (long term) reward
  - Rewards may be delayed
  - May need to sacrifice short term reward

- Devise a plan to maximize cumulative reward

Given an environment
(produces observations and rewards)

Reinforcement Learning

Automated agent that selects actions
to maximize total rewards in the environment

Reward

Examples:

- A financial investment (may take months to mature)
- Refuelling a helicopter (might prevent a crash in several hours)
- Blocking opponent moves (might help winning chances many moves from now)

# Reinforcement Learning

Learn a behavior strategy (policy) that maximizes the long term
Sum of rewards **in an unknown and stochastic environment (**Emma Brunskill: )

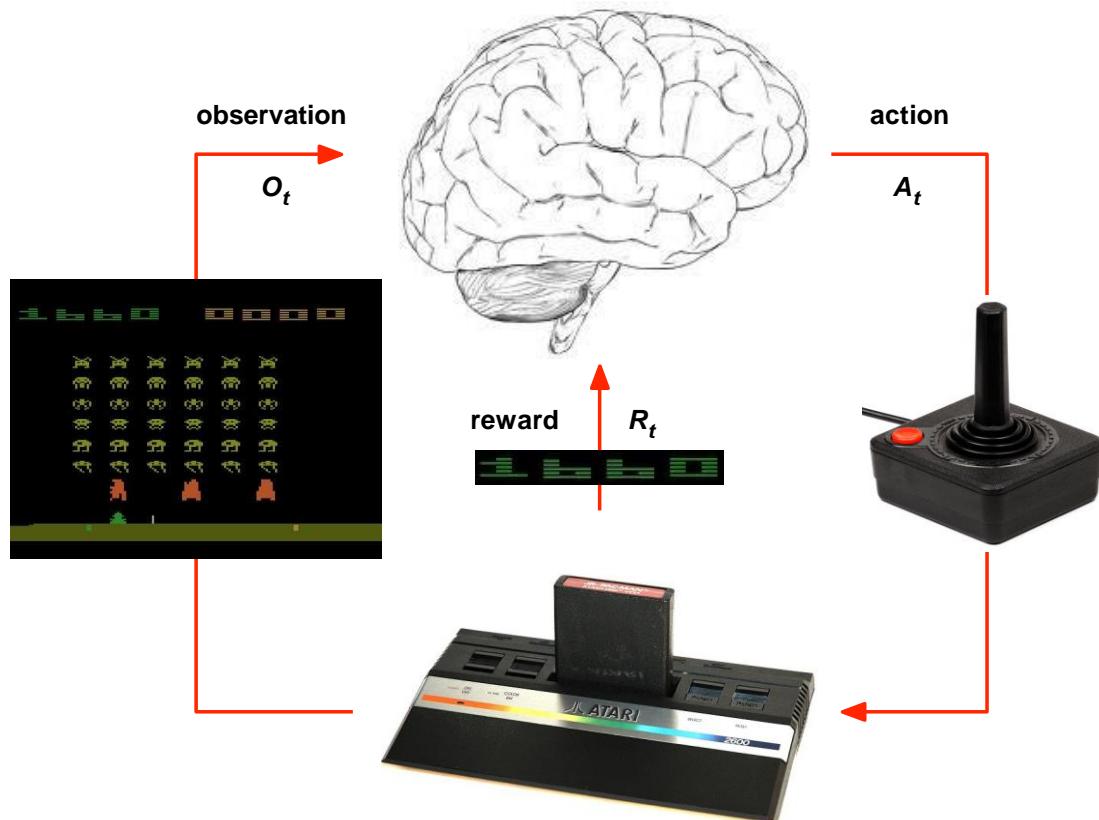## Planning under Uncertainty

Learn a behavior strategy (policy) that maximizes the long term
Sum of rewards **in a known stochastic environment (**Emma Brunskill: )

# Examples: Robotics

observation

$O_t$

reward   $R_t$

action

$A_t$

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

# Demos

Some videos

- [https://www.youtube.com/watch?v=V1eYniJ0Rnk](https://www.youtube.com/watch?v=V1eYniJ0Rnk)
- [https://www.youtube.com/watch?v=CIF2SBVY-J0](https://www.youtube.com/watch?v=CIF2SBVY-J0)
- [https://www.youtube.com/watch?v=I2WFvGI4y8c](https://www.youtube.com/watch?v=I2WFvGI4y8c)

# Markov Property

"The future is independent of the past given the present"

## Definition

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

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

- 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

# State Transition

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

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

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

$$\mathcal{P} = \text{from} \quad \begin{bmatrix} \mathcal{P}_{11} & \cdots & \mathcal{P}_{1n} \\ \vdots & & \\ \mathcal{P}_{n1} & \cdots & \mathcal{P}_{nn} \end{bmatrix}$$

$$\text{to}$$

where each row of the matrix sums to 1.
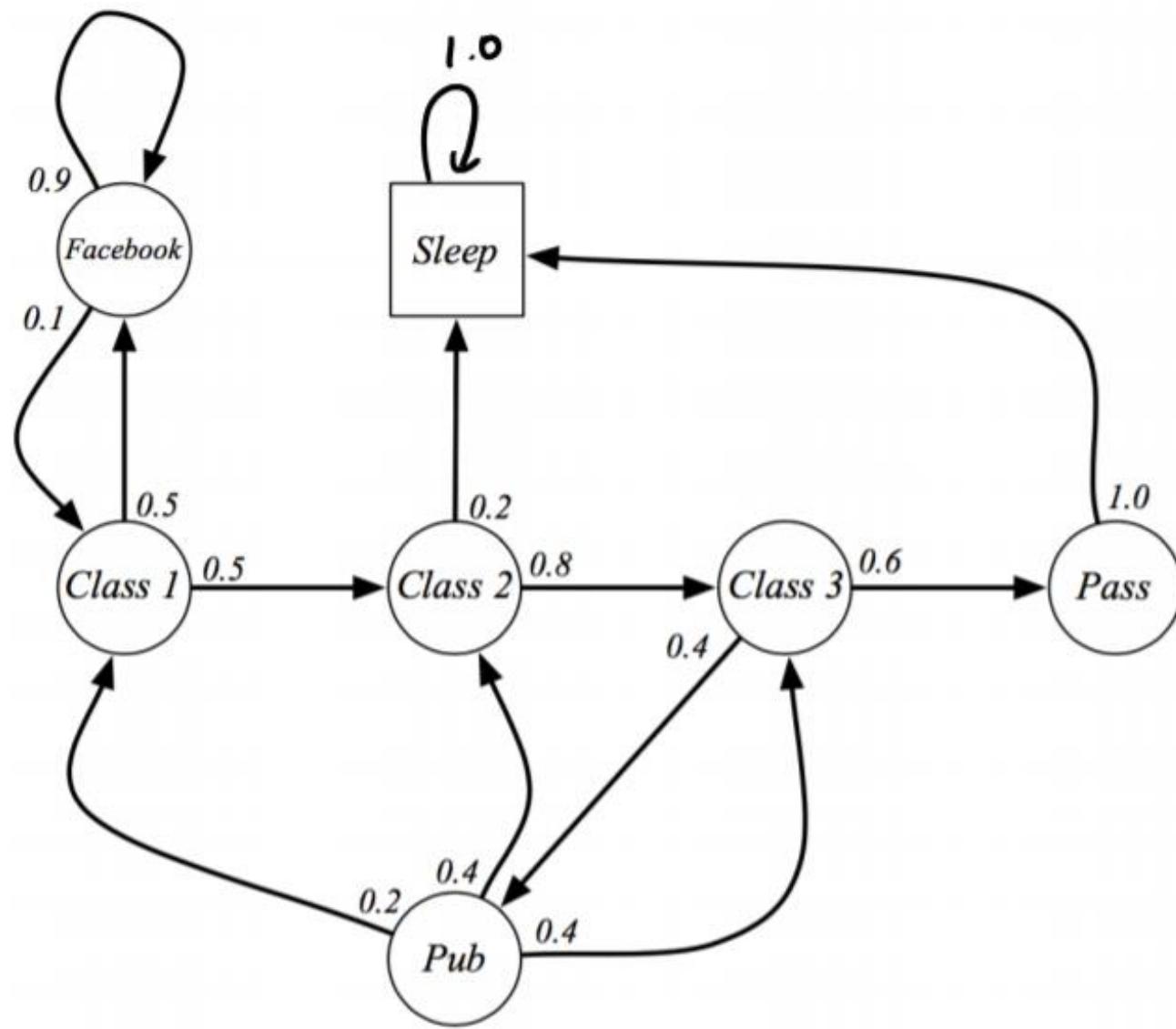
# Markov Process

A Markov process is a memoryless random process, i.e. a sequence of random states $S_1, S_2, \ldots$ with the Markov property.
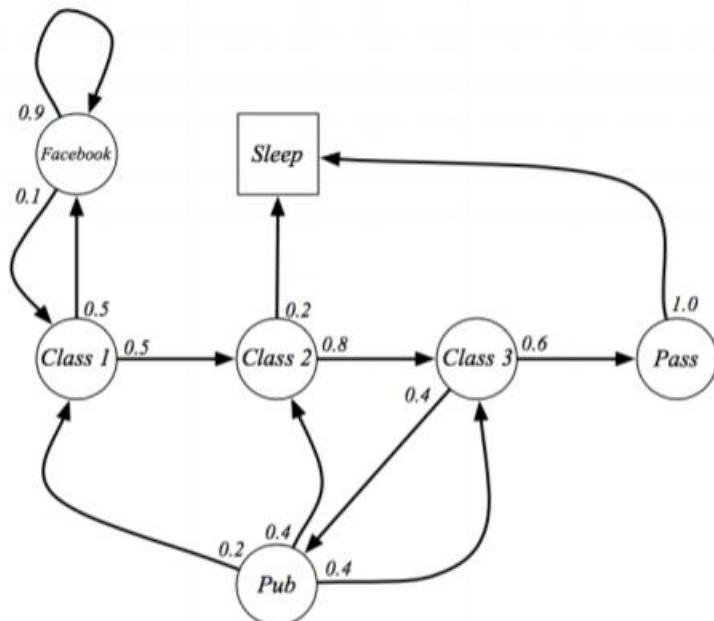
## Definition

A *Markov Process* (or *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,
  $\mathcal{P}_{ss'} = \mathbb{P}\left[S_{t+1} = s' \mid S_t = s\right]$

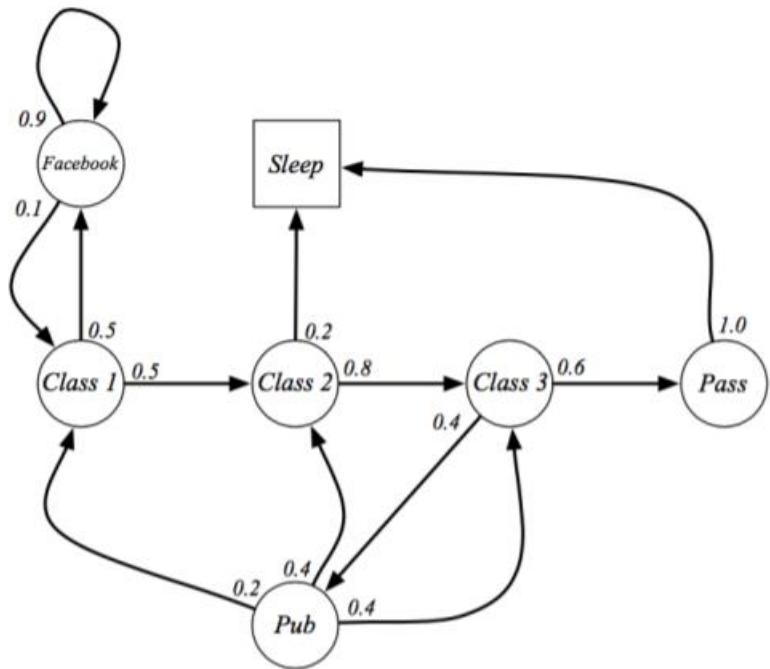# Student Markov Chain

# Student MC : 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

# Student MC : Transition Matrix



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

# Return

## Definition

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

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

- The *discount* $\gamma \in [0, 1]$ 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.
    - $\gamma$ close to 0 leads to "myopic" evaluation
    - $\gamma$ close to 1 leads to "far-sighted" evaluation

# Value

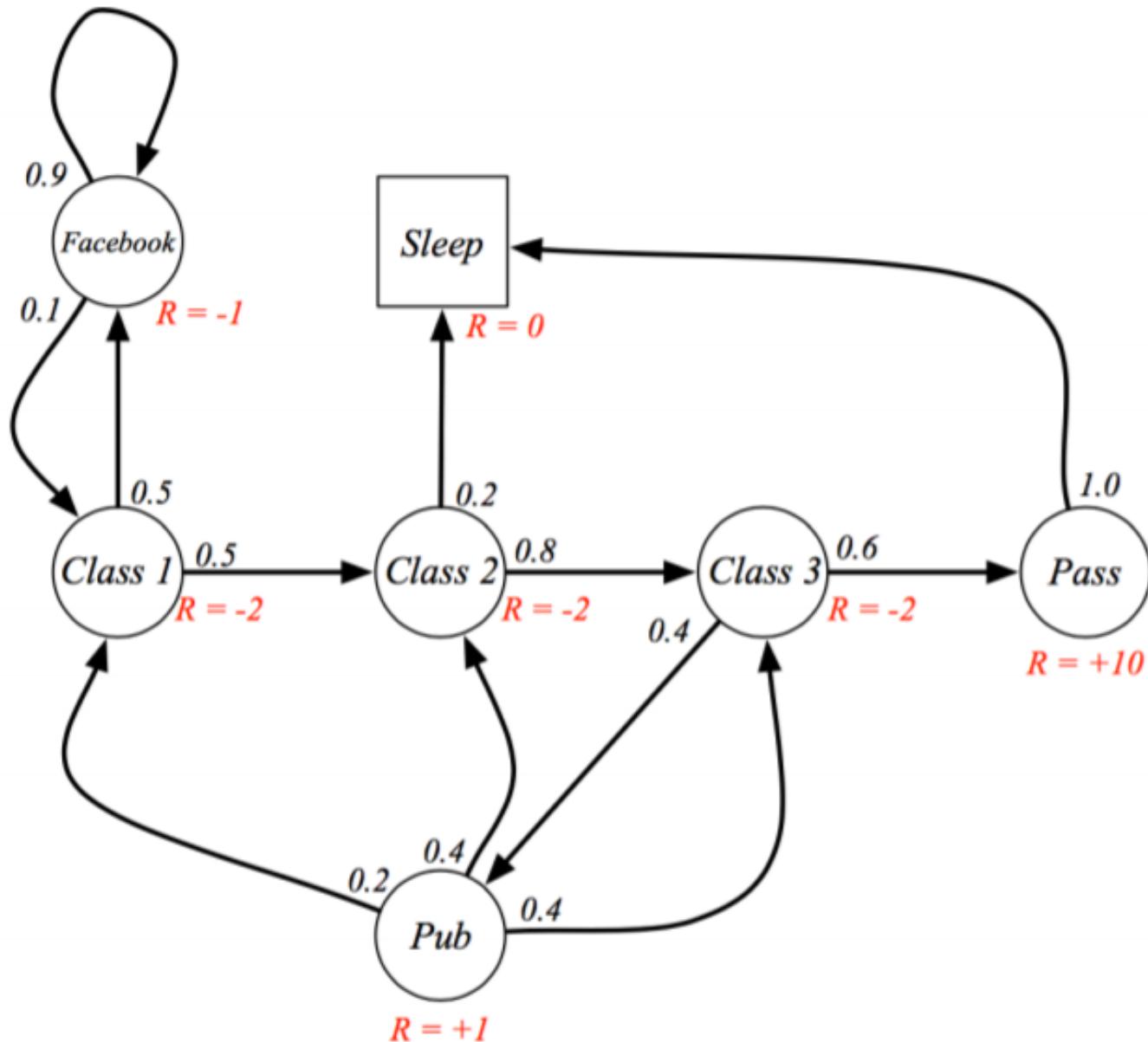The value function $v(s)$ gives the long-term value of state $s$

## Definition

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

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

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t \mid S_t = s] = \mathbb{E}_\pi\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s\right], \quad \text{for all } s \in \mathcal{S},$$

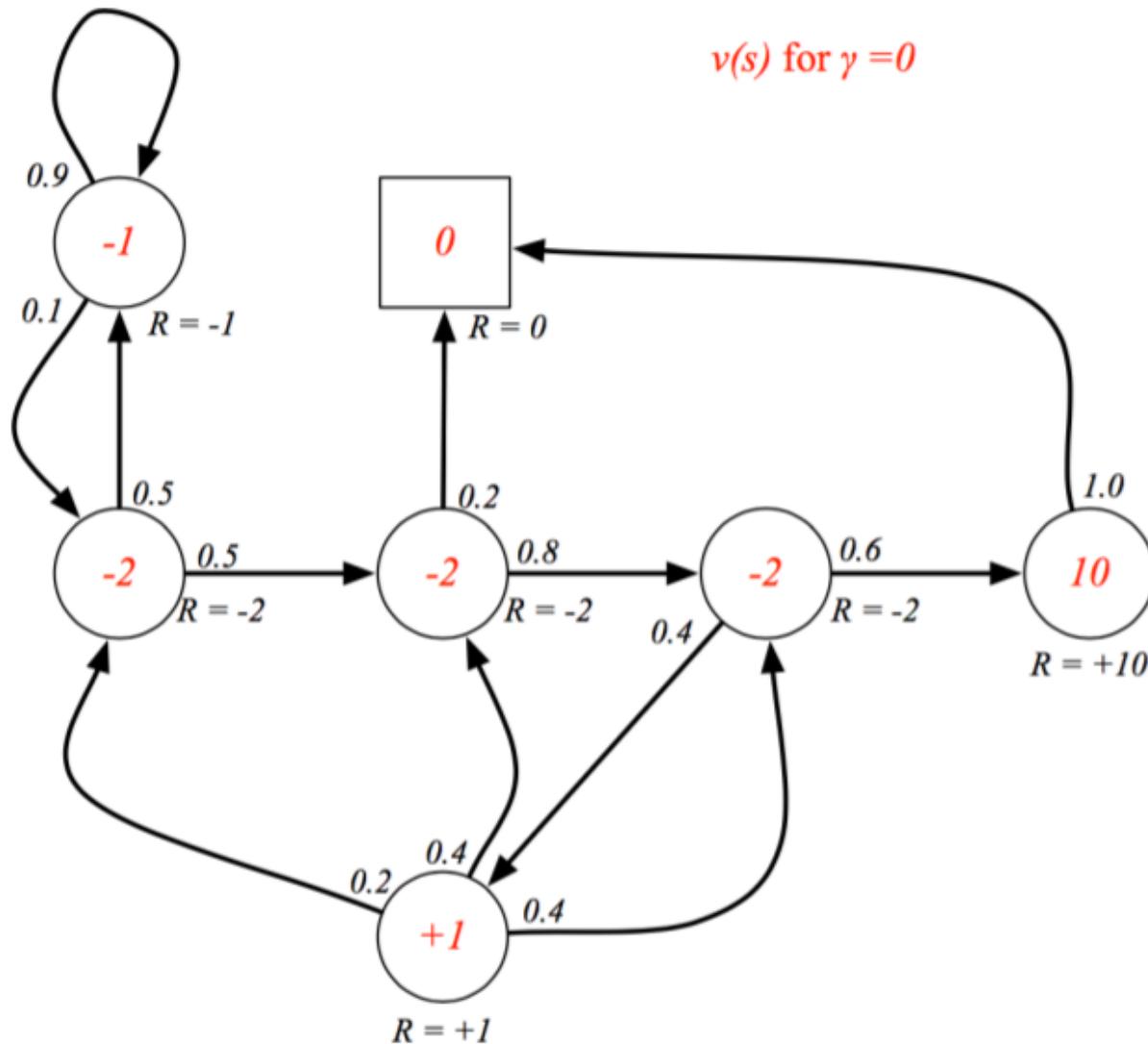# Student MRP

# Student MRP : Returns

Sample returns for Student MRP:
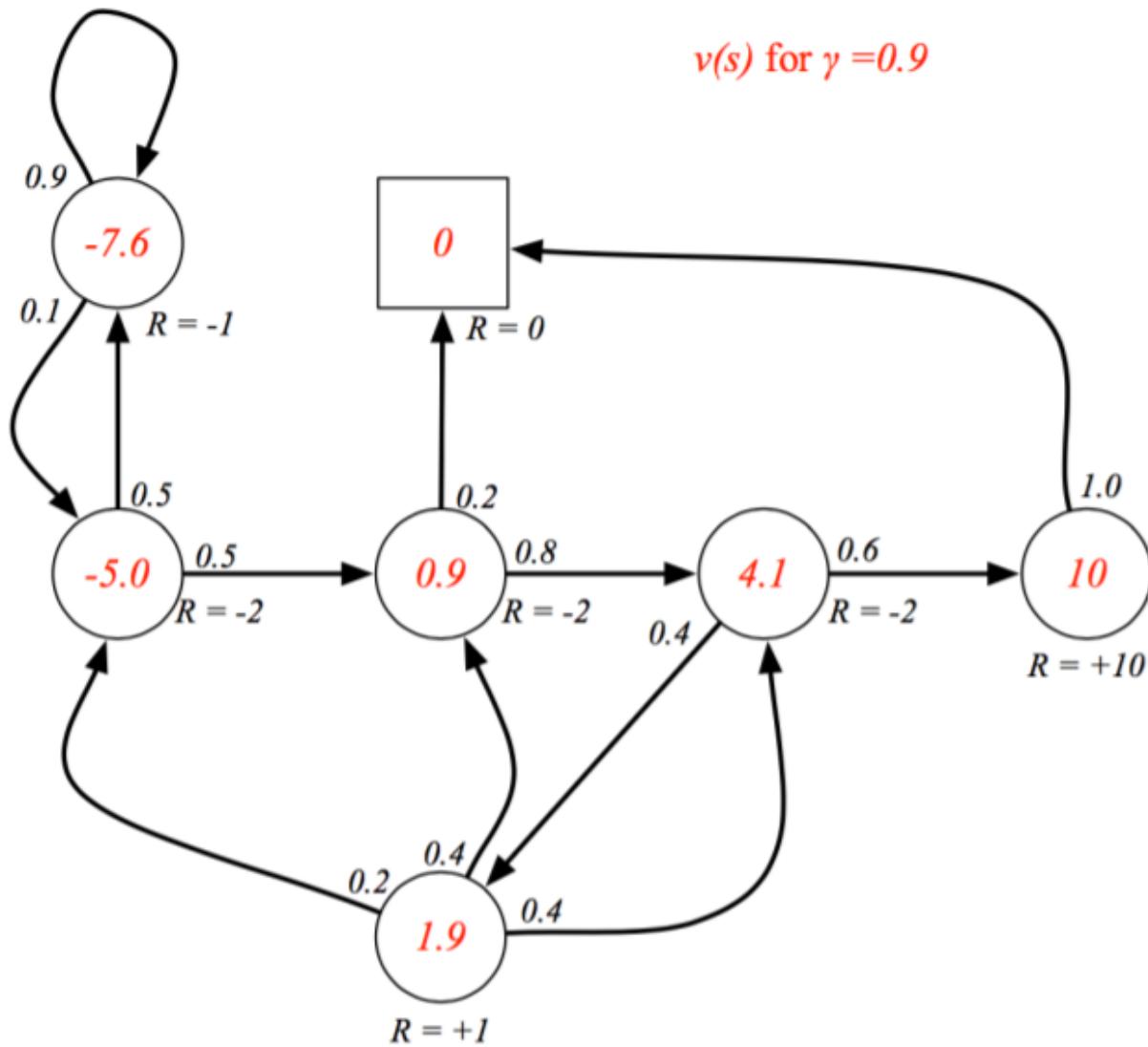Starting from $S_1 = $ C1 with $\gamma = \frac{1}{2}$

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

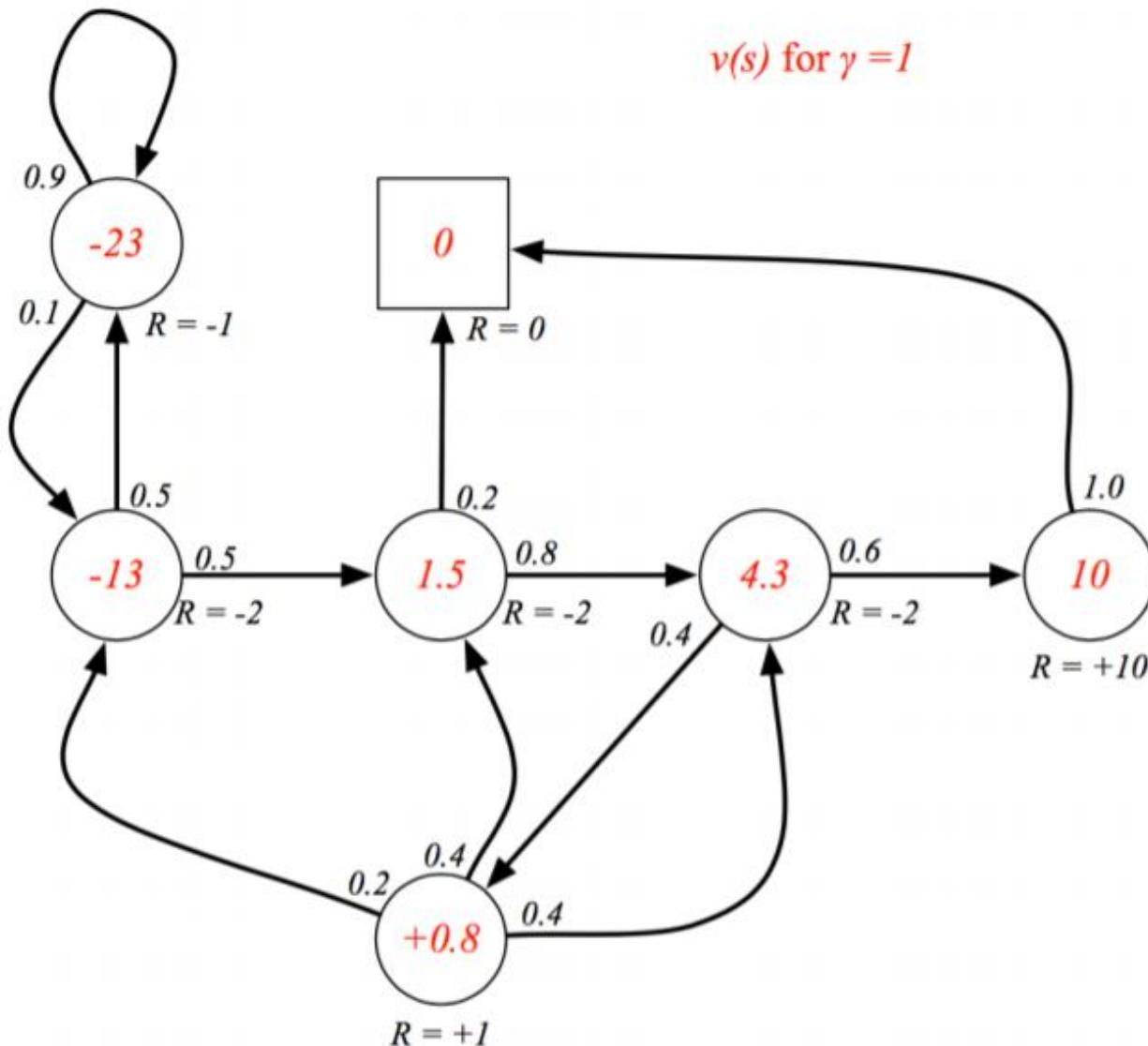| | | |
|---|---|---|
| C1 C2 C3 Pass Sleep | $v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 10 * \frac{1}{8}$ | $= \quad -2.25$ |
| C1 FB FB C1 C2 Sleep | $v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16}$ | $= \quad -3.125$ |
| C1 C2 C3 Pub C2 C3 Pass Sleep | $v_1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

# Student MRP : Value Function



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

# Student MRP : Value Function



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

# Student MRP : Value Function
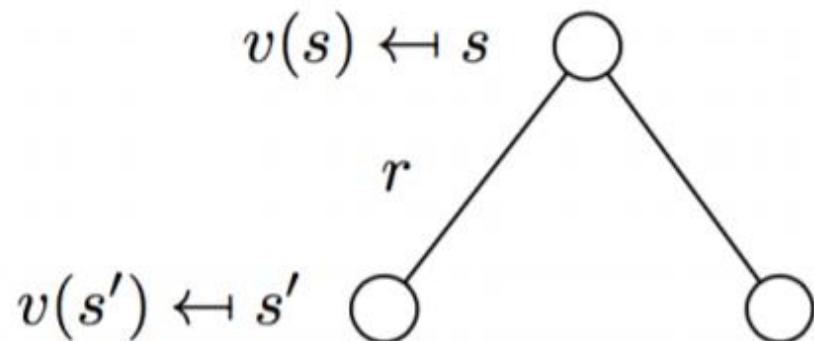
# Bellman Equation for MRP

The value function can be decomposed into two parts:

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

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

# Backup Diagrams for MRP

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

$$v(s) \hookleftarrow s$$

$$v(s') \hookleftarrow s'$$

$$r$$

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

# Bellman Eq: Student MRP



$$4.3 = -2 + 0.6*10 + 0.4*0.8$$

# Bellman Eq: Student MRP

The Bellman equation can be expressed concisely using matrices,

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

where $v$ is a column vector with one entry per state

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

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

$$v = R + \gamma P v$$
$$(I - \gamma P)\, v = R$$
$$v = (I - \gamma P)^{-1}\, R$$

- 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
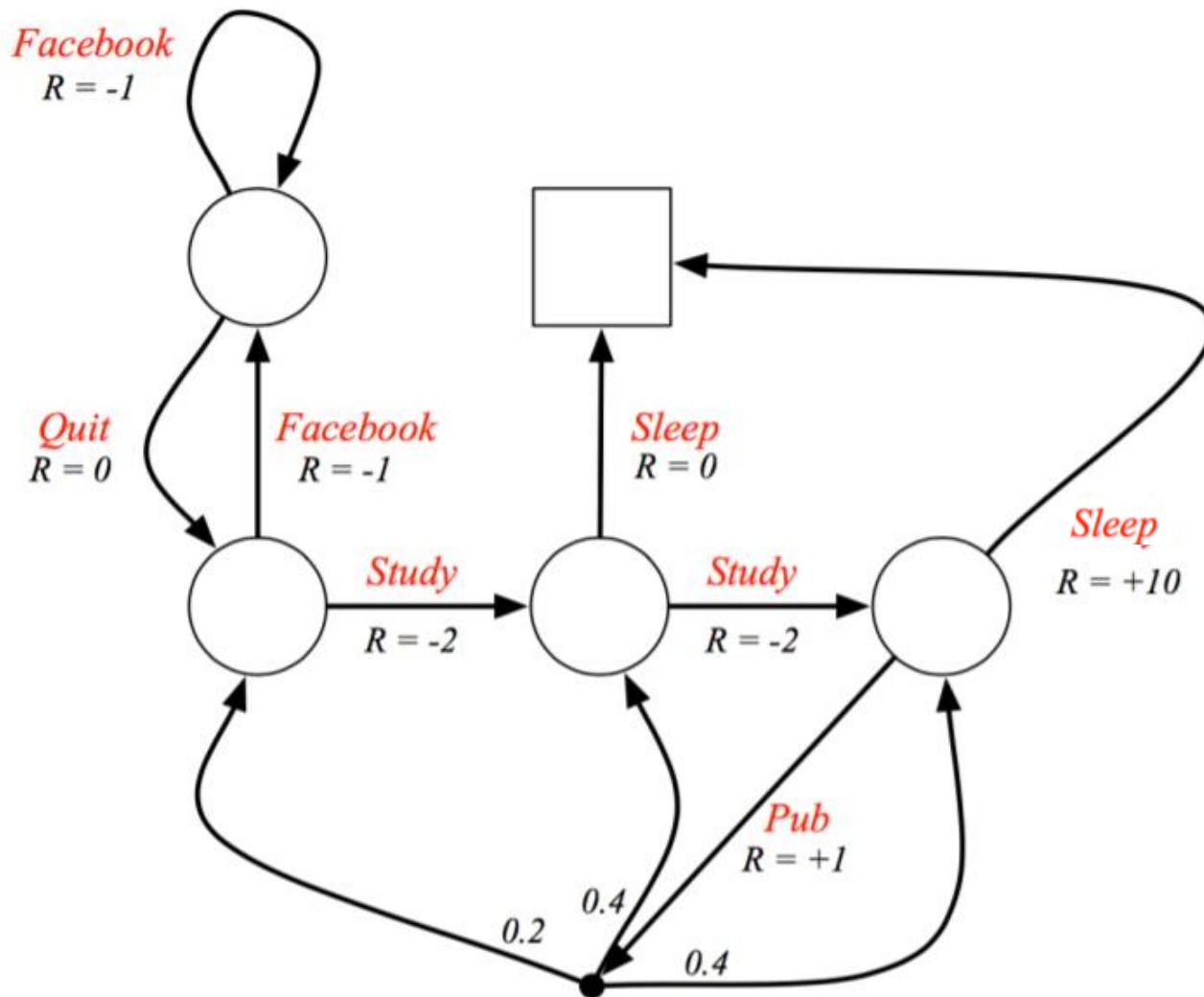
31

# Markov Decision Process

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

## 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



Facebook
R = -1

Quit
R = 0

Facebook
R = -1

Sleep
R = 0

Sleep
R = +10

Study
R = -2

Study
R = -2

Pub
R = +1

0.4

0.2

0.4

# Policies

## Definition

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

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

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

# MP $\rightarrow$ MRP $\rightarrow$ MDP

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

**Definition**

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

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

**Definition**

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

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

# Bellman Eq for MDP

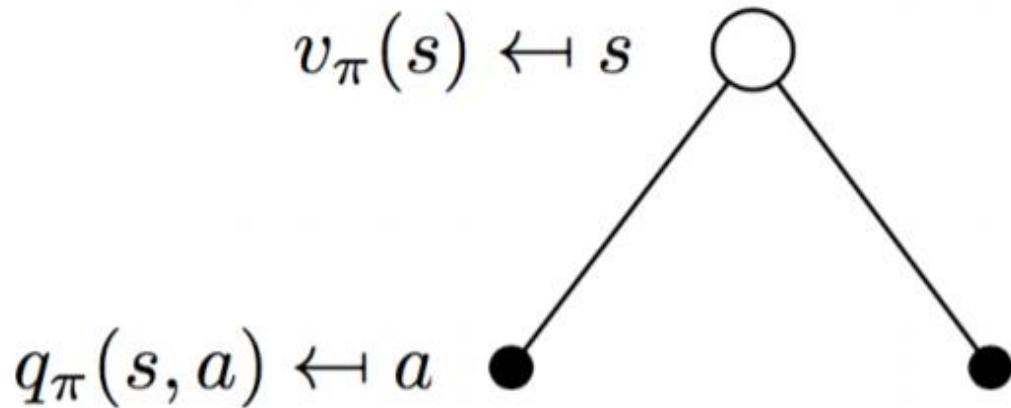The state-value function can again be decomposed into immediate reward plus discounted value of successor state,

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

The action-value function can similarly be decomposed,

$$q_\pi(s, a) = \mathbb{E}_\pi \left[ R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a \right]$$
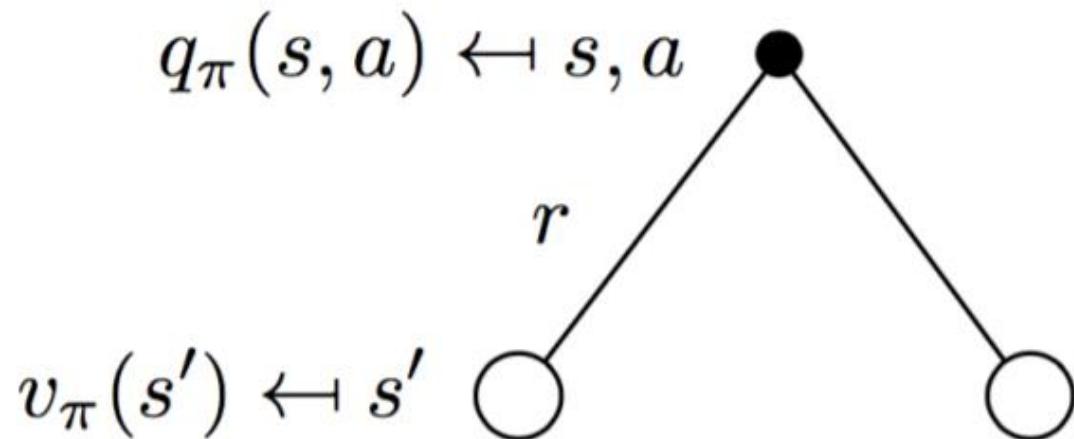
Evaluating Bellman equation translates into 1-step lookahead
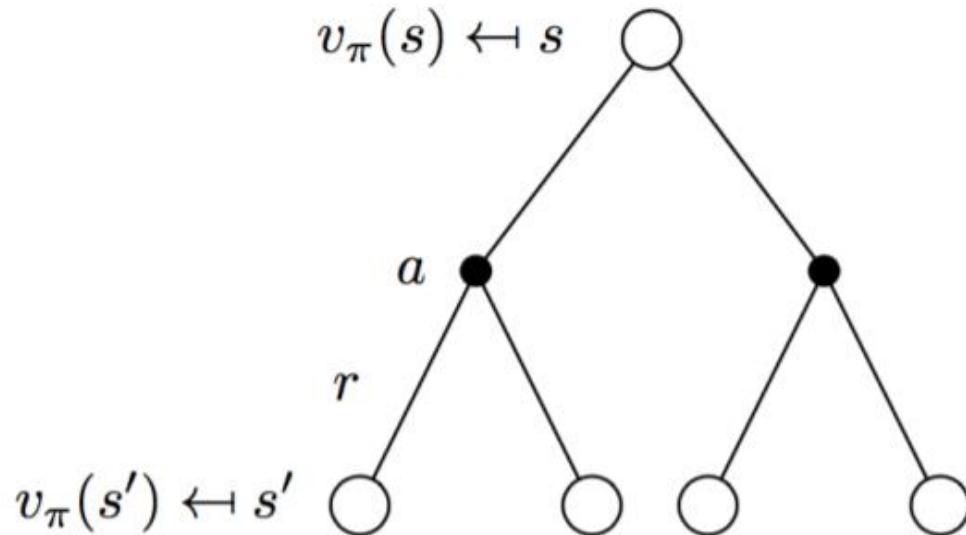
# Bellman Eq, V



$$v_\pi(s) \hookleftarrow s$$

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

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

# Bellman Eq, q
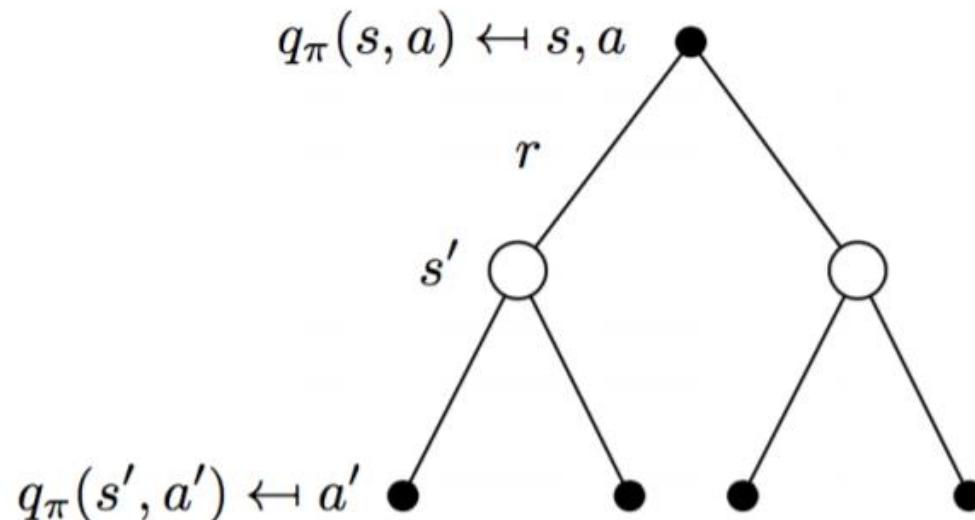
$$q_\pi(s, a) \leftarrowtail s, a$$

$$r$$

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

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

# Bellman Eq, V

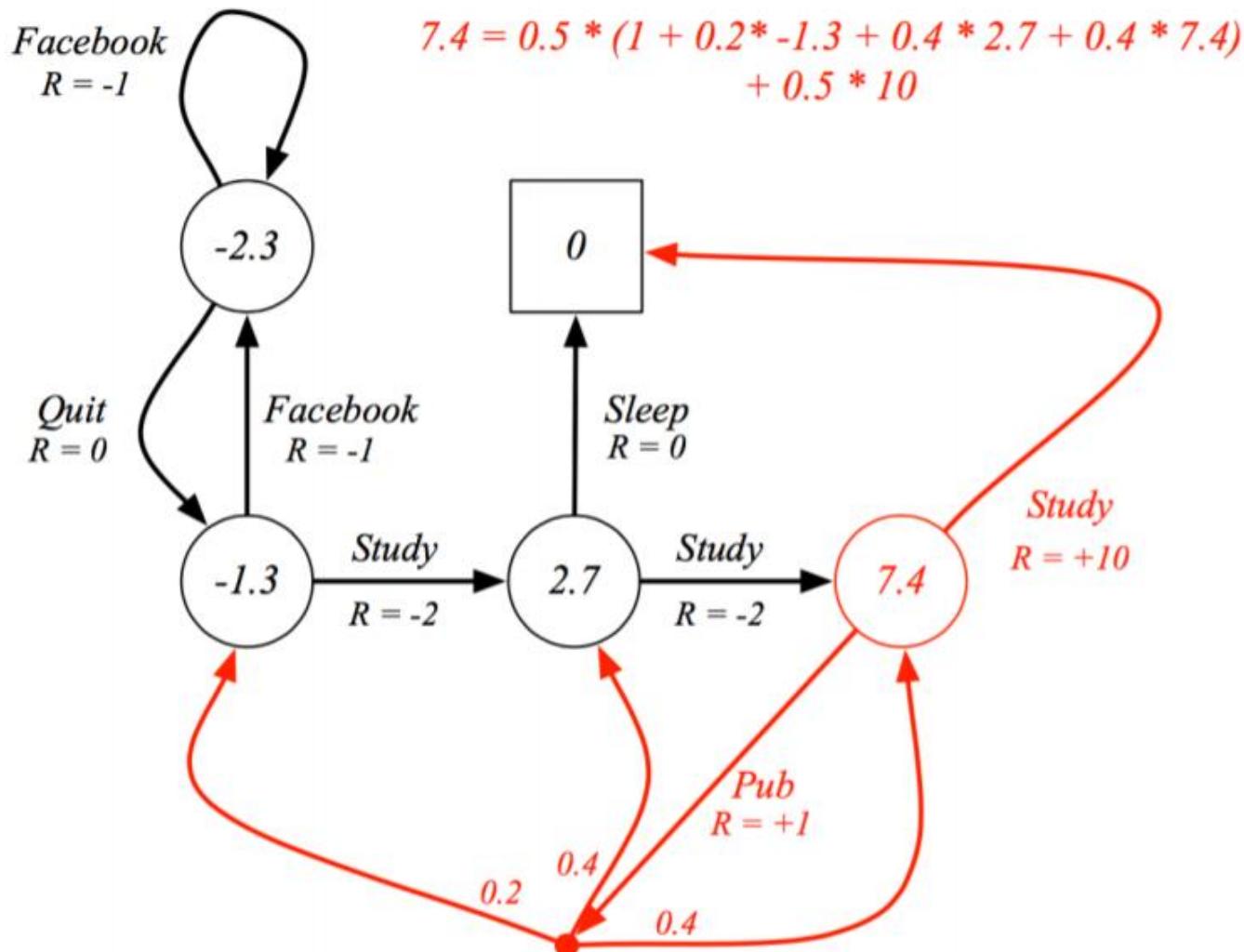$$v_\pi(s) \leftarrow s$$

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

$$a$$

$$r$$

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

40

# Bellman Eq, q

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

$$r$$

$$s'$$

$$q_\pi(s',a') \hookleftarrow a'$$

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

# Student MDP : Bellman Eq



Facebook
R = -1

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

-2.3

0

Quit
R = 0

Facebook
R = -1

Sleep
R = 0

Study
R = +10

-1.3

Study
R = -2

2.7

Study
R = -2

7.4

Pub
R = +1

0.4

0.2

0.4

# Bellman Eq : Matrix Form

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

$$v_\pi = \mathcal{R}^\pi + \gamma \mathcal{P}^\pi v_\pi$$

with direct solution

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

# Optimal Value Function

## Definition

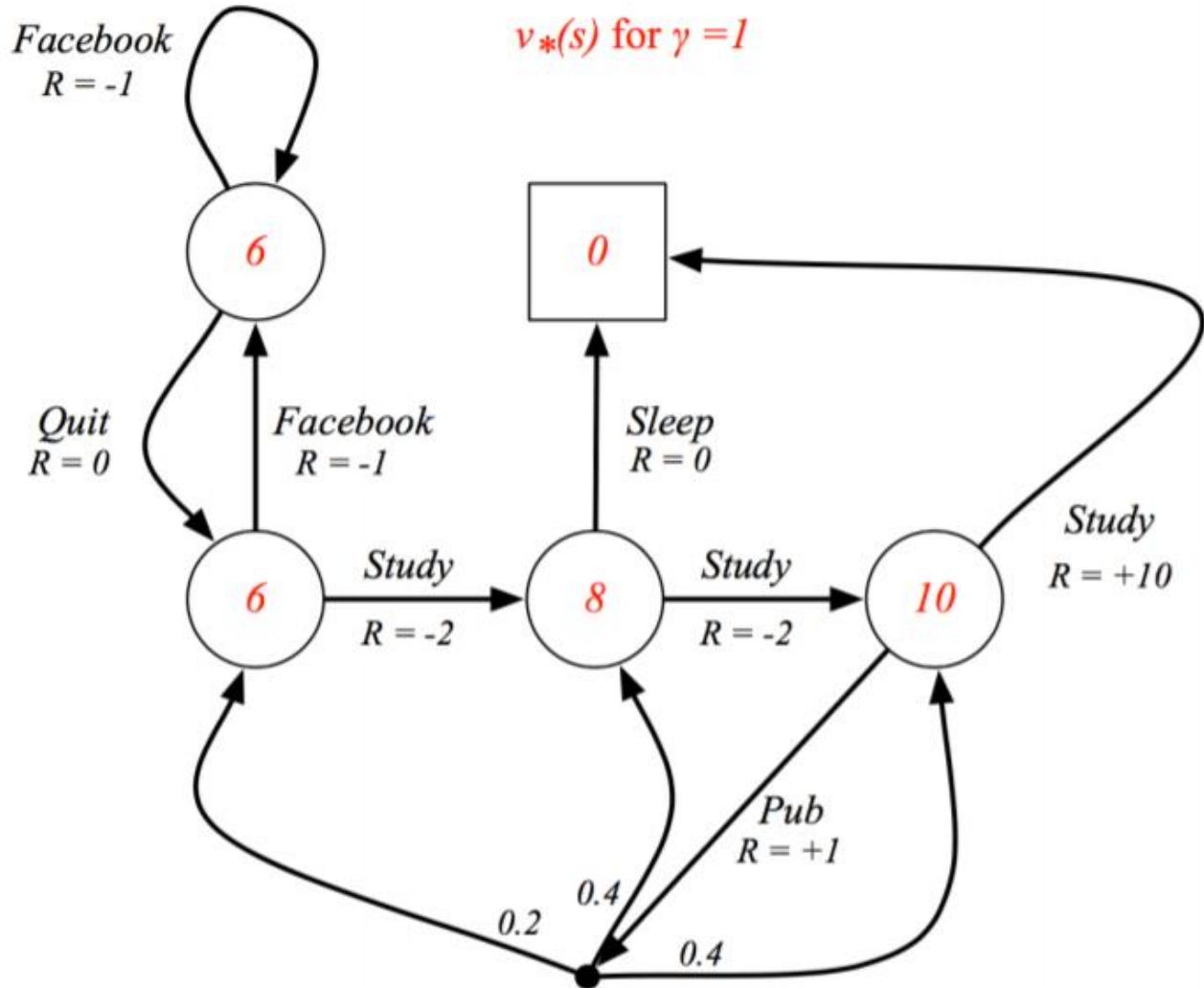The *optimal state-value function* $v_*(s)$ is the maximum value function over all policies

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

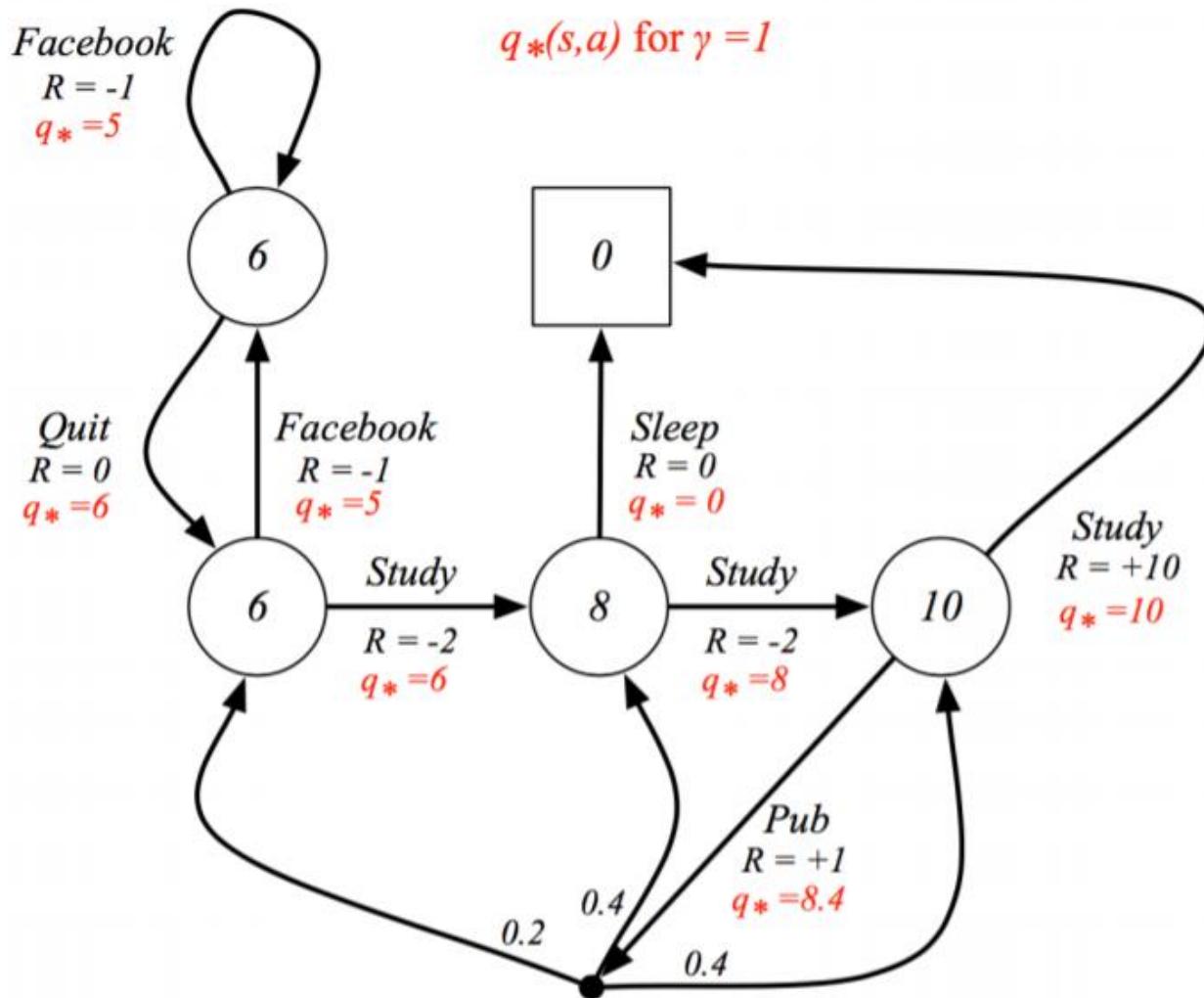The *optimal action-value function* $q_*(s, a)$ is the maximum action-value function over all policies

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

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

# Student MDP : Optimal V

# Student MDP : Optimal Q



Facebook
R = -1
$q_* = 5$

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

6

0

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

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, $v_{\pi_*}(s) = v_*(s)$*

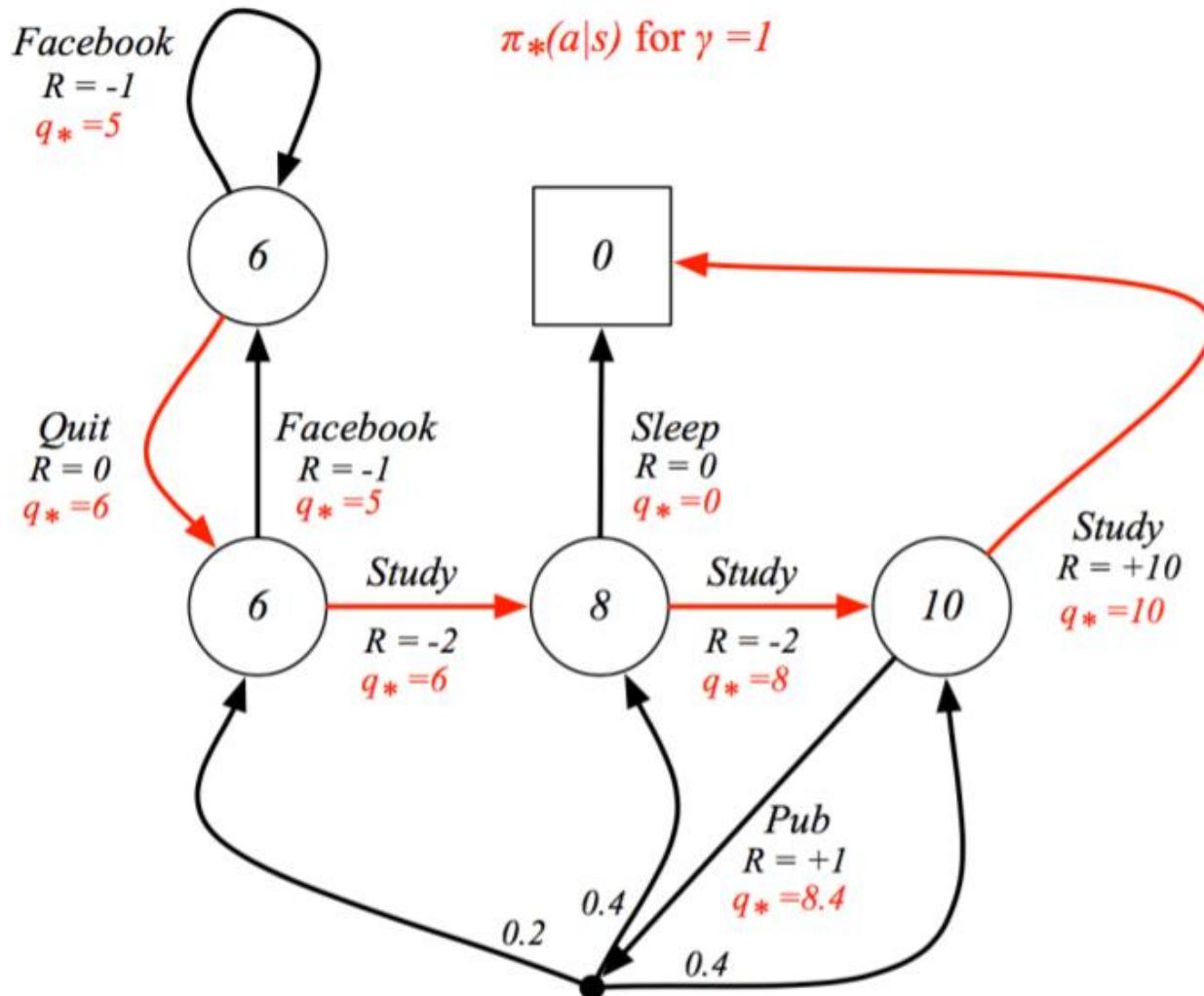- *All optimal policies achieve the optimal action-value function, $q_{\pi_*}(s, a) = q_*(s, a)$*

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

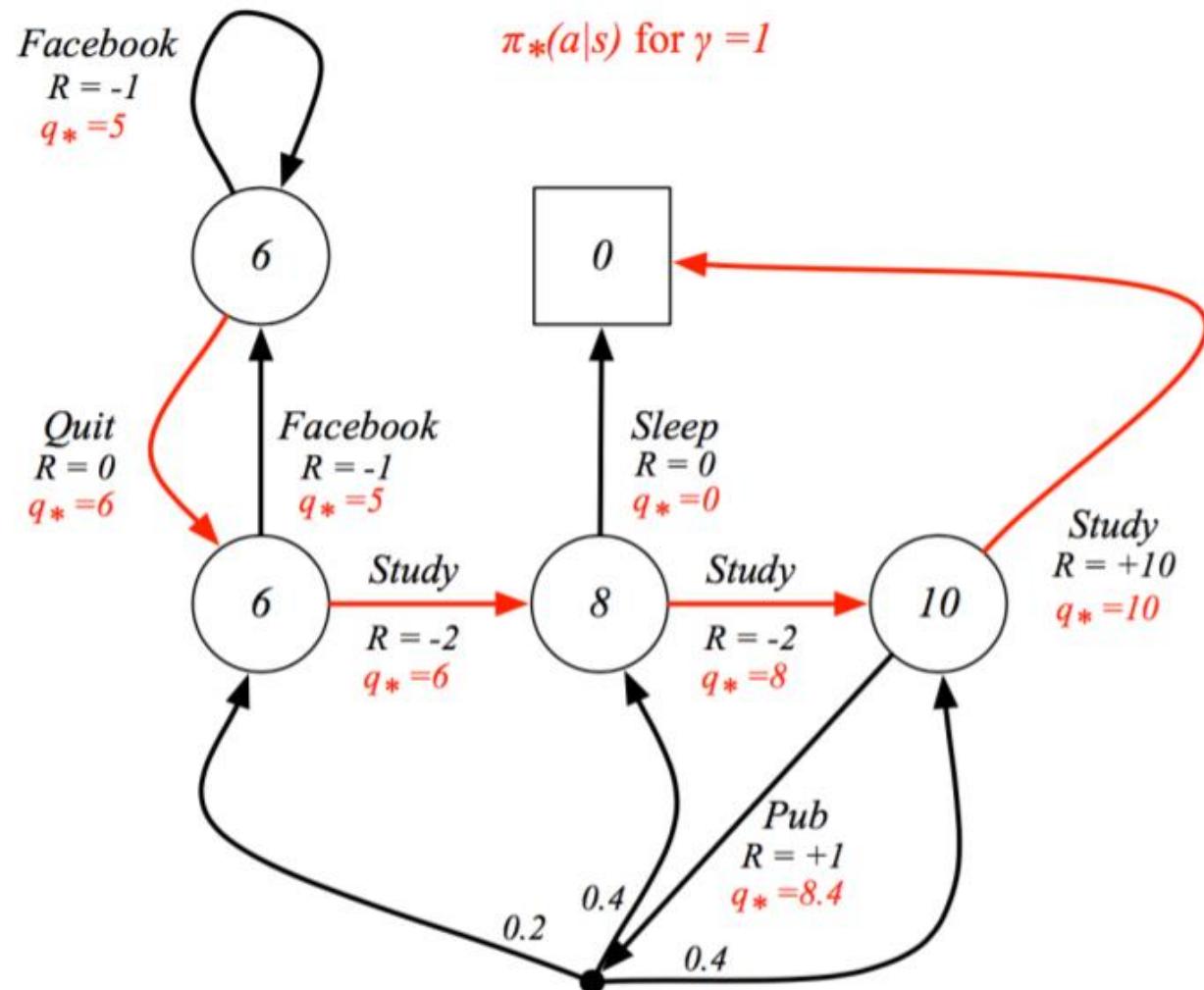$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \underset{a \in \mathcal{A}}{\text{argmax }} q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

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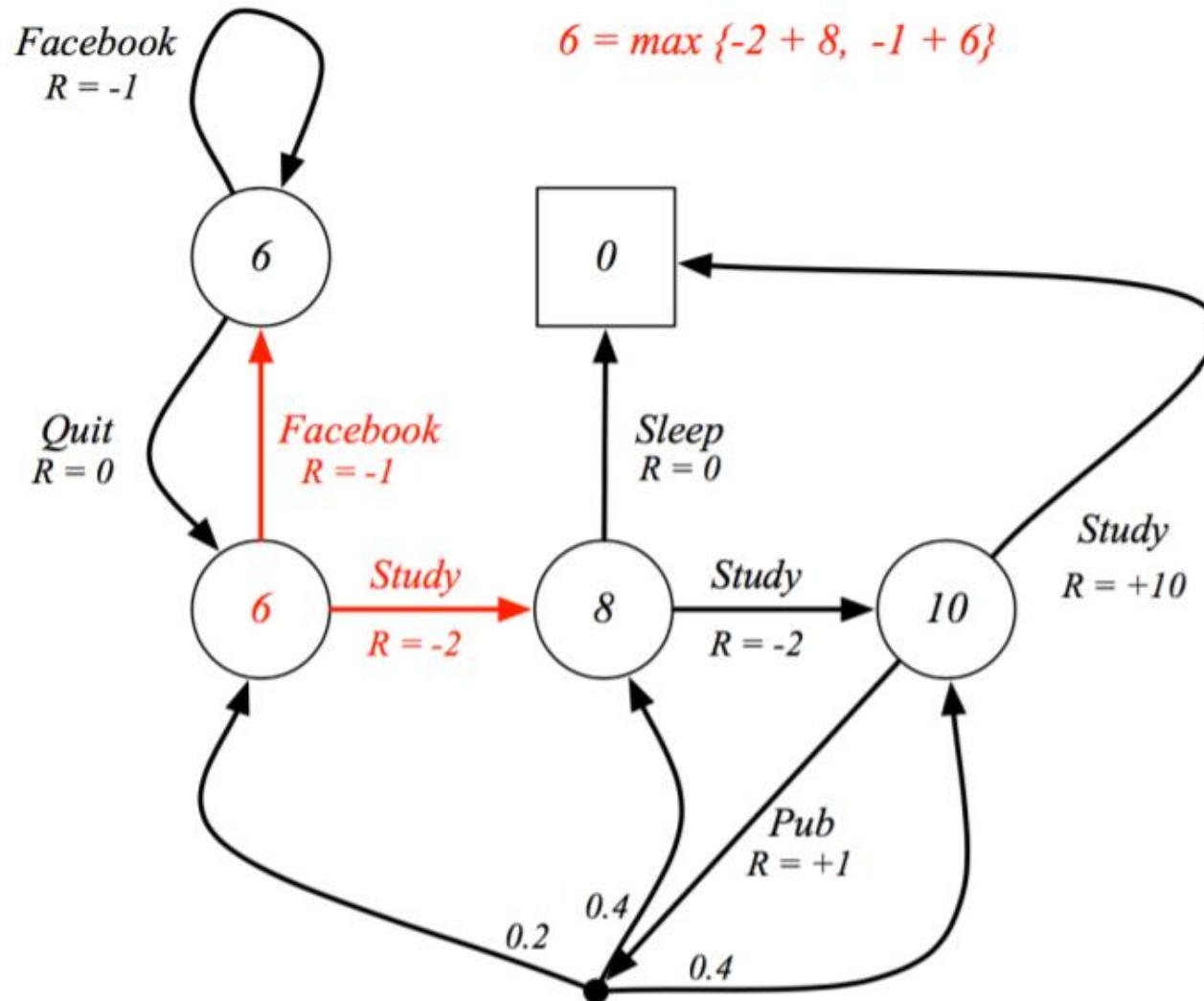- If we know $q_*(s, a)$, we immediately have the optimal policy

# Student MDP : Optimal Policy



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

Facebook
R = -1
$q_* = 5$

Quit
R = 0
$q_* = 6$

Facebook
R = -1
$q_* = 5$

Sleep
R = 0
$q_* = 0$

Study
R = -2
$q_* = 6$

Study
R = -2
$q_* = 8$

Study
R = +10
$q_* = 10$

Pub
R = +1
$q_* = 8.4$

6

0

6

8

10

0.4

0.2

0.4

# Bellman Optimality Eq, V



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

Facebook
R = -1
$q_* = 5$

6

0

Quit
R = 0
$q_* = 6$

Facebook
R = -1
$q_* = 5$

Sleep
R = 0
$q_* = 0$

6

Study
R = -2
$q_* = 6$

8

Study
R = -2
$q_* = 8$

10

Study
R = +10
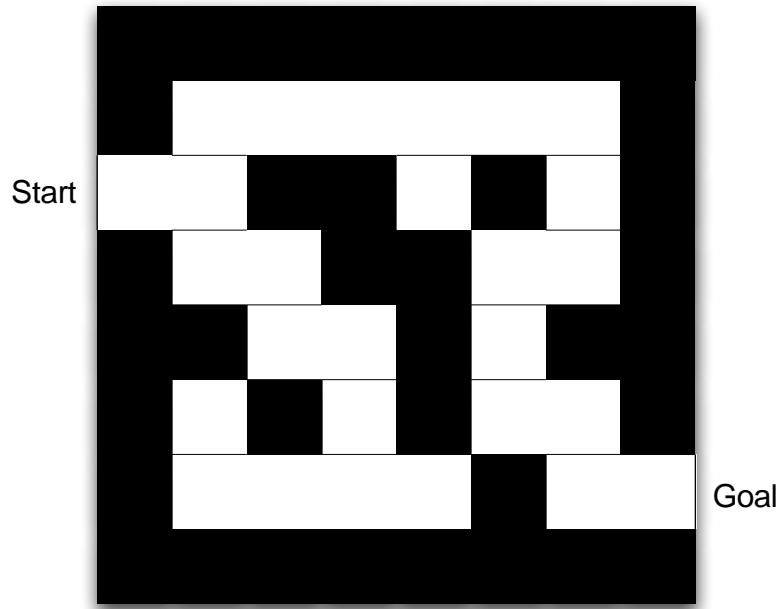$q_* = 10$

Pub
R = +1
$q_* = 8.4$

0.4

0.2

0.4

# Student MDP : Bellman Optimality

# Solving the Bellman Optimality Equation
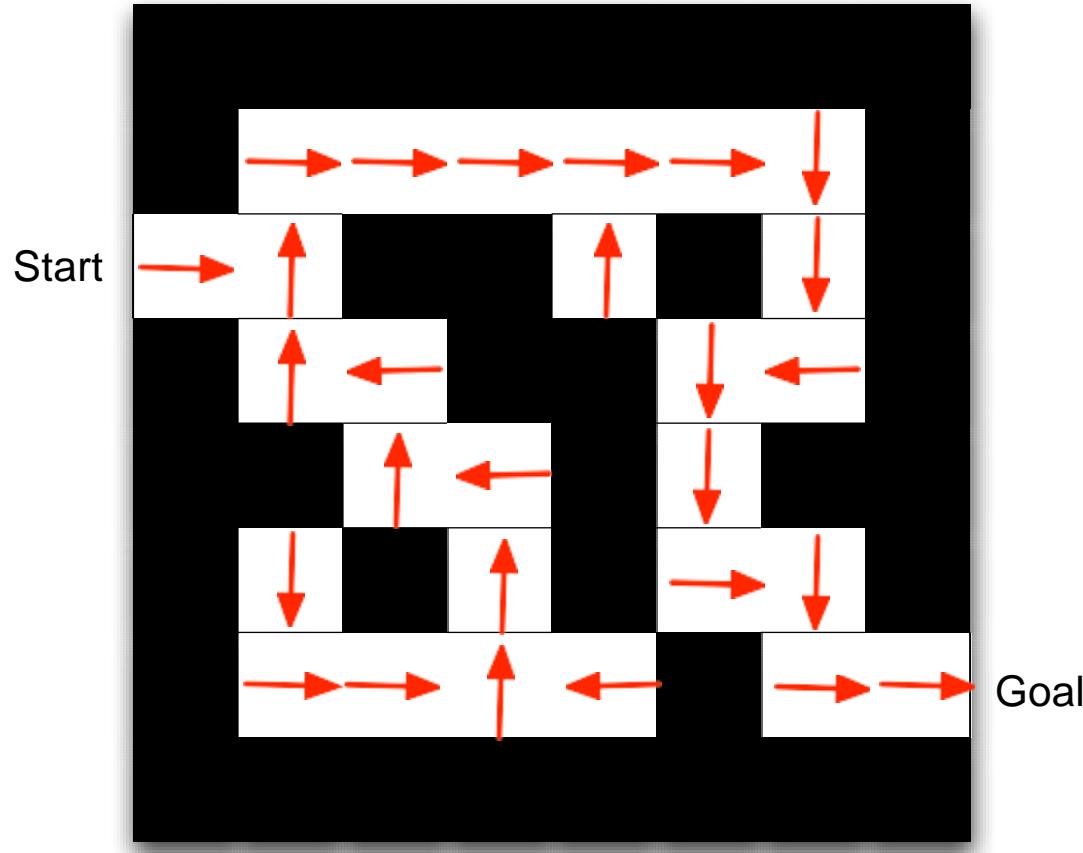
Not easy

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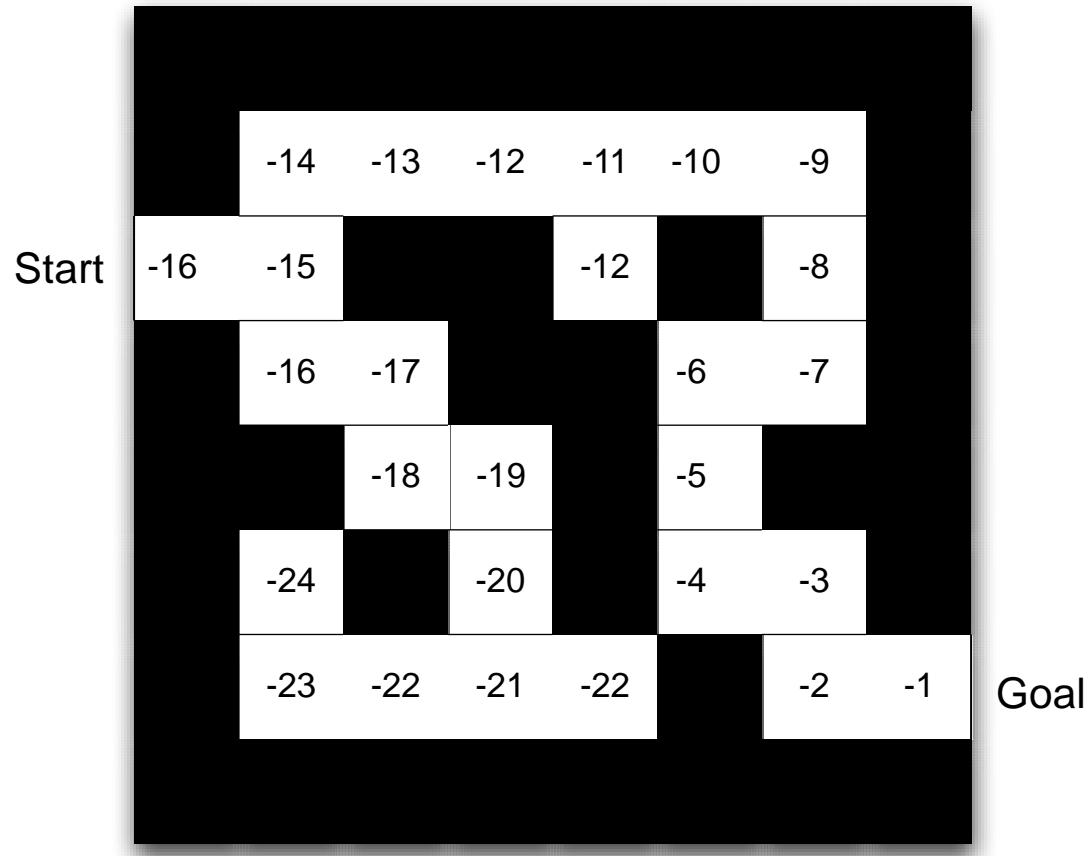  - Q-learning
  - Sarsa

# Maze Example

Start

Goal

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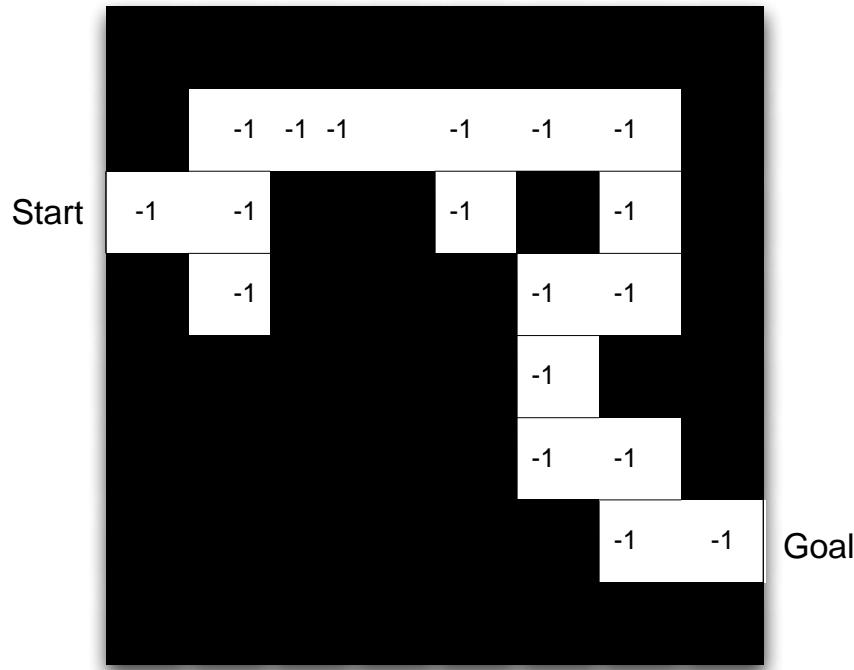- Actions: N, E, S, W
- States: Agent's location

Start

Goal

■ Arrows represent policy π(*s*) for each state *s*

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

Start

Goal

- 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 $P_{ss'}^a$
- Numbers represent immediate reward $R_s^a$ from each state $s$ (same for all $a$)

# Algorithms for MDPs

**MDPs**      States, Transitions, Actions, Rewards

**Prediction**      Given Policy $\pi$, Estimate State Value Functions, Action Value Functions

**Control**      Estimate Optimal Value Functions, Optimal Policy

**Does the agent know the MDP?**

**Yes!**    It's "planning"          **No!**    It's "Model-free RL"
Agent knows everything                  Agent observes everything as it goes

# 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.

$$\mathcal{P}_{ss'}^{a} = \mathbb{P}[S_{t+1} = s' \mid S_t = s, A_t = a]$$
$$\mathcal{R}_s^a = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$$

# Algorithms cont.

|  | Prediction | Control |
|---|---|---|
|  | Evaluate Policy, π | Find Best Policy, π* |
| **Planning** MDP Known | Policy Evaluation | Policy/Value Iteration |
| MDP Unknown | MC and TD Learning | Q-Learning |

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

# Major Components of an RL Agent

- An RL agent may include one or more of these components:
    - Policy: agent's behaviour function
    - Value function: how good is each state and/or action
    - Model: agent's representation of the environment

# Dynamic Programming

**Dynamic** sequential or temporal component to the problem

**Programming** optimising a "program", i.e. a policy

- c.f. linear programming

- A method for solving complex problems
- By breaking them down into subproblems
    - Solve the subproblems
    - Combine solutions to subproblems

# Requirements for DP

Dynamic Programming is a very general solution method for problems which have two properties:

- Optimal substructure
    - *Principle of optimality* applies
    - Optimal solution can be decomposed into subproblems
- Overlapping subproblems
    - Subproblems recur many times
    - Solutions can be cached and reused
- Markov decision processes satisfy both properties
    - Bellman equation gives recursive decomposition
    - Value function stores and reuses solutions

# Applications for DPs

Dynamic programming is used to solve many other problems, e.g.

- Scheduling algorithms
- String algorithms (e.g. sequence alignment)
- Graph algorithms (e.g. shortest path algorithms)
- Graphical models (e.g. Viterbi algorithm)
- Bioinformatics (e.g. lattice models)

- Dynamic programming assumes full knowledge of the MDP

- It is used for *planning* in an MDP

- For prediction:
    - Input: MDP $(S, A, P, R, \gamma)$ and policy $\pi$
    - or: MRP $(S, P^{\pi}, R^{\pi}, \gamma)$
    - Output: value function $v_{\pi}$

- Or for control:
    - Input: MDP $(S, A, P, R, \gamma)$
    - Output: optimal value function $v_*$
    - and: optimal policy $\pi_*$

# Policy Evaluation (Prediction)

- Problem: evaluate a given policy $\pi$

- Solution: iterative application of Bellman expectation backup

- $v_1 \rightarrow v_2 \rightarrow \ldots \rightarrow v_\pi$

- Using *synchronous* backups,
  - At each iteration $k + 1$
  - For all states $s \in S$
  - Update $v_{k+1}(s)$ from $v_k(s')$
  - where $s'$ is a successor state of $s$

- We will discuss *asynchronous* backups later

- Convergence to $v_\pi$ can be proven

66

# Iterative policy Evaluation

**Iterative policy evaluation**

Input $\pi$, the policy to be evaluated
Initialize an array $V(s) = 0$, for all $s \in \mathcal{S}^+$
Repeat
   $\Delta \leftarrow 0$
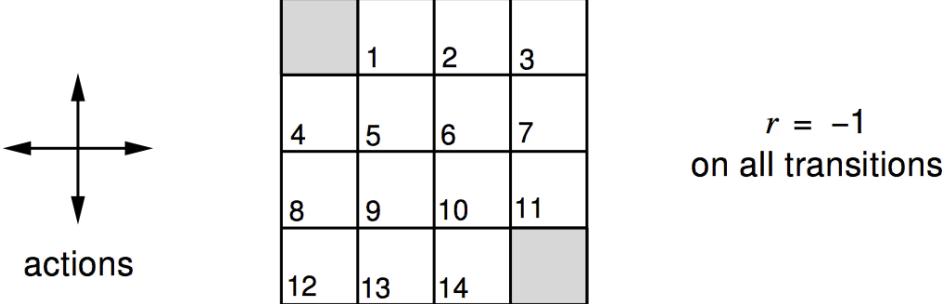   For each $s \in \mathcal{S}$:
      $v \leftarrow V(s)$
      $V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$
      $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
until $\Delta < \theta$ (a small positive number)
Output $V \approx v_\pi$

- Undiscounted episodic MDP $(\gamma = 1)$
- Nonterminal states $1, \ldots, 14$
- One terminal state (shown twice as shaded squares)
- Actions leading out of the grid leave state unchanged
- Reward is $-1$ until the terminal state is reached
- Agent follows uniform random policy

$$\pi(n|\cdot) = \pi(e|\cdot) = \pi(s|\cdot) = \pi(w|\cdot) = 0.25$$

# Policy Evaluation : Grid World

$v_k$ for the
Random Policy

$k = 0$

| 0.0 | 0.0 | 0.0 | 0.0 |
|-----|-----|-----|-----|
| 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 |
| 0.0 | 0.0 | 0.0 | 0.0 |

Time 0 : do nothing, stop; no cost.

$k = 1$

| 0.0 | -1.0 | -1.0 | -1.0 |
|-----|------|------|------|
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | -1.0 |
| -1.0 | -1.0 | -1.0 | 0.0 |

Time 1 : move (reward -1); then k=0
   Unless in goal: reward 0

$k = 2$

| 0.0 | -1.7 | -2.0 | -2.0 |
|-----|------|------|------|
| -1.7 | -2.0 | -2.0 | -2.0 |
| -2.0 | -2.0 | -2.0 | -1.7 |
| -2.0 | -2.0 | -1.7 | 0.0 |

Time 2 : move (reward -1); then k=1
   Most: move (-1) + [v1 = -1]  = -2
   Some: move (-1) + ¾ [v1 = -1] + ¼ [v1=0] = 1.75

# Policy Evaluation : Grid World

$v_k$ for the
Random Policy

$k = 3$

| 0.0 | -2.4 | -2.9 | -3.0 |
|------|------|------|------|
| -2.4 | -2.9 | -3.0 | -2.9 |
| -2.9 | -3.0 | -2.9 | -2.4 |
| -3.0 | -2.9 | -2.4 | 0.0 |

$k = 10$

| 0.0 | -6.1 | -8.4 | -9.0 |
|------|------|------|------|
| -6.1 | -7.7 | -8.4 | -8.4 |
| -8.4 | -8.4 | -7.7 | -6.1 |
| -9.0 | -8.4 | -6.1 | 0.0 |

$k = \infty$

| 0.0 | -14. | -20. | -22. |
|------|------|------|------|
| -14. | -18. | -20. | -20. |
| -20. | -20. | -18. | -14. |
| -22. | -20. | -14. | 0.0 |

# Policy Evaluation : Grid World

# Policy Evaluation : Grid World

$v_k$ for the
Random Policy

Greedy Policy
w.r.t. $v_k$

$k = 3$

| 0.0 | -2.4 | -2.9 | -3.0 |
|---|---|---|---|
| -2.4 | -2.9 | -3.0 | -2.9 |
| -2.9 | -3.0 | -2.9 | -2.4 |
| -3.0 | -2.9 | -2.4 | 0.0 |

$k = 10$

| 0.0 | -6.1 | -8.4 | -9.0 |
|---|---|---|---|
| -6.1 | -7.7 | -8.4 | -8.4 |
| -8.4 | -8.4 | -7.7 | -6.1 |
| -9.0 | -8.4 | -6.1 | 0.0 |

optimal
policy

$k = \infty$

| 0.0 | -14. | -20. | -22. |
|---|---|---|---|
| -14. | -18. | -20. | -20. |
| -20. | -20. | -18. | -14. |
| -22. | -20. | -14. | 0.0 |

In general:
best policy & value for
"one step, then
follow random policy"
(always better policy than random!)

Most of the story in a nutshell:

# Will Value Iteration Converge?

- Yes, if discount factor is < 1 or end up in a terminal state with probability 1

- Bellman equation is a contraction
- If apply it to two different value functions, distance between value functions shrinks after apply Bellman equation to each

# Finding Best Policy

|  | Evaluate Policy, $\pi$ | Find Best Policy, $\pi^*$ |
|---|---|---|
| MDP Known | Policy Evaluation | Policy/Value Iteration |
| MDP Unknown | MC and TD Learning | Sarsa + Q-Learning |

- Given a policy $\pi$
  - Evaluate the policy $\pi$

    $$v_\pi(s) = E\left[R_{t+1} + \gamma R_{t+2} + ... | S_t = s\right]$$

  - Improve the policy by acting greedily with respect to $v_\pi$

    $$\pi' = greedy(v_\pi)$$

- In Small Gridworld improved policy was optimal, $\pi' = \pi*$
- In general, need more iterations of improvement / evaluation
- But this process of policy iteration always converges to $\pi*$

# Policy Iteration

$$\pi_0 \xrightarrow{\text{E}} v_{\pi_0} \xrightarrow{\text{I}} \pi_1 \xrightarrow{\text{E}} v_{\pi_1} \xrightarrow{\text{I}} \pi_2 \xrightarrow{\text{E}} \cdots \xrightarrow{\text{I}} \pi_* \xrightarrow{\text{E}} v_*,$$

where $\xrightarrow{\text{E}}$ denotes a policy *evaluation* and $\xrightarrow{\text{I}}$ denotes a policy *improvement*. Each policy is guaranteed to be a strict improvement over the previous one (unless it is already optimal). Because a finite MDP has only a finite number of policies, this process must converge to an optimal policy and optimal value function in a finite number of iterations.

---

**Policy iteration (using iterative policy evaluation)**

1. Initialization
   $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation
   Repeat
   $\quad \Delta \leftarrow 0$
   $\quad$ For each $s \in \mathcal{S}$:
   $\quad\quad v \leftarrow V(s)$
   $\quad\quad V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s))\big[r + \gamma V(s')\big]$
   $\quad\quad \Delta \leftarrow \max(\Delta, |v - V(s)|)$
   until $\Delta < \theta$ (a small positive number)

3. Policy Improvement
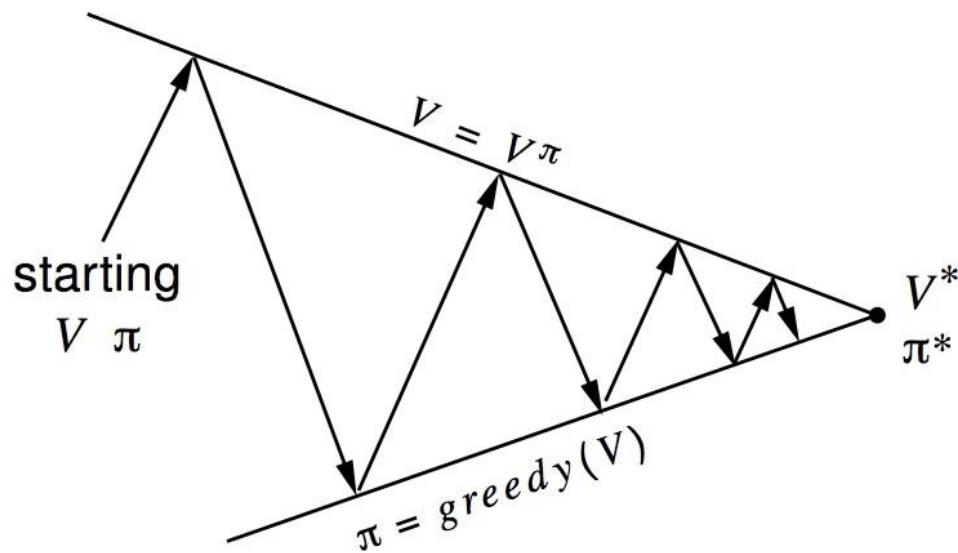   *policy-stable* $\leftarrow$ *true*
   For each $s \in \mathcal{S}$:
   $\quad$ *old-action* $\leftarrow \pi(s)$
   $\quad \pi(s) \leftarrow \text{argmax}_a \sum_{s',r} p(s',r|s,a)\big[r + \gamma V(s')\big]$
   $\quad$ If *old-action* $\neq \pi(s)$, then *policy-stable* $\leftarrow$ *false*
   If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2
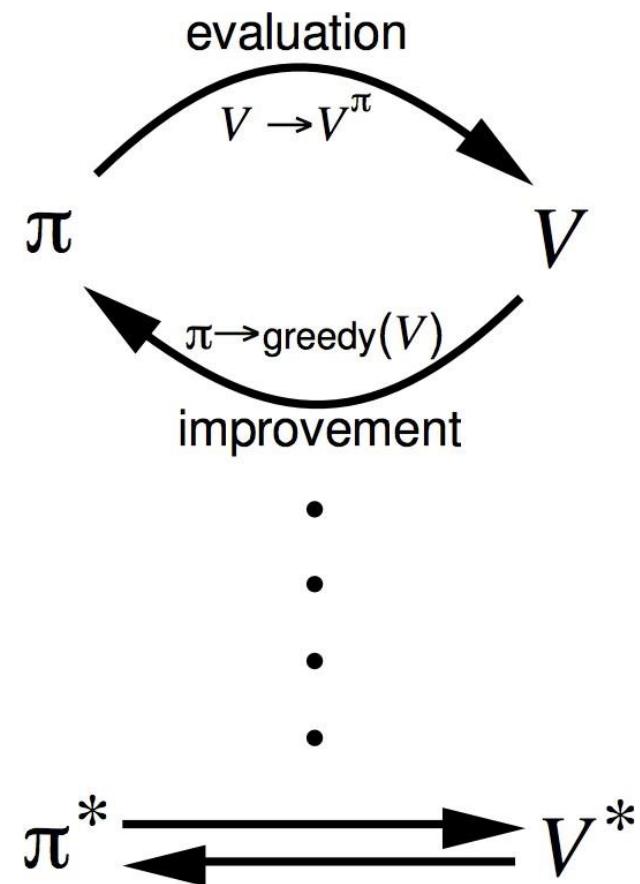
# Policy Iteration



Policy evaluation Estimate $v_\pi$
  Iterative policy evaluation

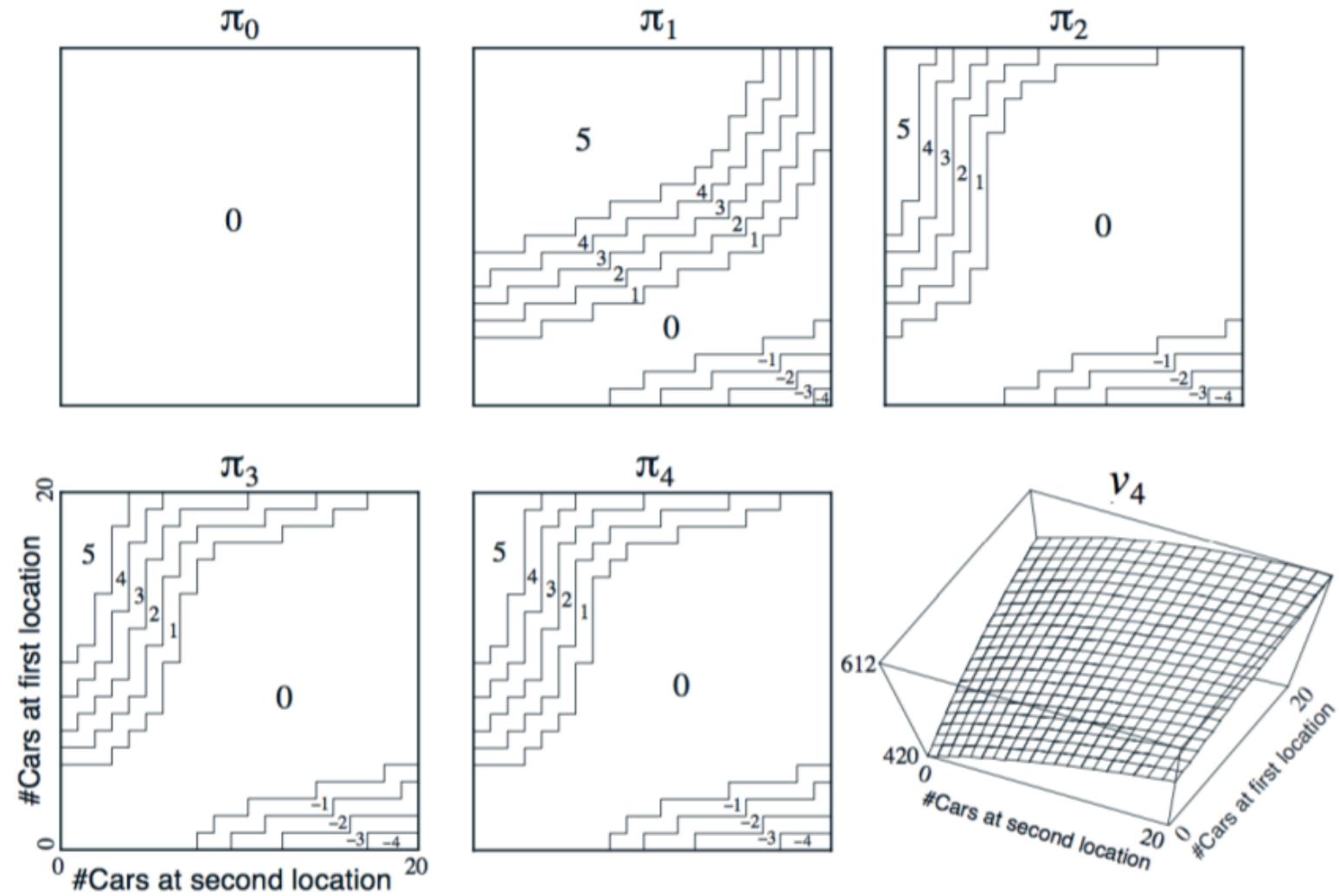Policy improvement Generate $\pi^I \geq \pi$
  Greedy policy improvement

# Jack's Car Rental



- States: Two locations, maximum of 20 cars at each
- Actions: Move up to 5 cars between locations overnight
- Reward: $10 for each car rented (must be available)
- Transitions: Cars returned and requested randomly
    - Poisson distribution, $n$ returns/requests with prob $\frac{\lambda^n}{n!}e^{-\lambda}$
    - 1st location: average requests $= 3$, average returns $= 3$
    - 2nd location: average requests $= 4$, average returns $= 2$

# Policy Iteration in Car Rental

- Consider a deterministic policy, $a = \pi(s)$
- We can *improve* the policy by acting greedily

$$\pi'(s) = \underset{a \in \mathcal{A}}{\text{argmax}} \; q_\pi(s, a)$$

- This improves the value from any state $s$ over one step,

$$q_\pi(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_\pi(s, a) \geq q_\pi(s, \pi(s)) = v_\pi(s)$$

- It therefore improves the value function, $v_{\pi'}(s) \geq v_\pi(s)$

$$
\begin{aligned}
v_\pi(s) &\leq q_\pi(s, \pi'(s)) = \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s\right] \\
&\leq \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma q_\pi(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s\right] \\
&\leq \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 q_\pi(S_{t+2}, \pi'(S_{t+2})) \mid S_t = s\right] \\
&\leq \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma R_{t+2} + \ldots \mid S_t = s\right] = v_{\pi'}(s)
\end{aligned}
$$

- If improvements stop,

$$q_\pi(s, \pi'(s)) = \max_{a \in A} q_\pi(s, a) = q_\pi(s, \pi(s)) = v_\pi(s)$$

- Then the Bellman optimality equation has been satisfied

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

- Therefore $v_\pi(s) = v_*(s)$ for all $s \in S$
- so $\pi$ is an optimal policy

- How do we know that value iteration converges to $v_*$?
- Or that iterative policy evaluation converges to $v_\pi$?
- And therefore that policy iteration converges to $v_*$?
- Is the solution unique?

- How fast do these algorithms converge?
- These questions are resolved by *contraction mapping theorem*

- Consider the vector space $V$ over value functions
- There are $|S|$ dimensions
- Each point in this space fully specifies a value function $v(s)$
- What does a Bellman backup do to points in this space?
- We will show that it brings value functions *closer*
- And therefore the backups must converge on a unique solution

- We will measure distance between state-value functions $u$ and $v$ by the $\infty$-norm

- i.e. the largest difference between state values,

$$||u - v||_\infty = \max_{s \in S} |u(s) - v(s)|$$

- Approximate the value function
- Using a *function approximator* $\hat{v}(s, \mathbf{w})$
- Apply dynamic programming to $\hat{v}(\cdot, \mathbf{w})$
- e.g. Fitted Value Iteration repeats at each iteration $k$,
  - Sample states $\tilde{\mathcal{S}} \subseteq \mathcal{S}$
  - For each state $s \in \tilde{\mathcal{S}}$, estimate target value using Bellman optimality equation,

$$\tilde{v}_k(s) = \max_{a \in \mathcal{A}} \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \hat{v}(s', \mathbf{w_k}) \right)$$

  - Train next value function $\hat{v}(\cdot, \mathbf{w_{k+1}})$ using targets $\{\langle s, \tilde{v}_k(s) \rangle\}$

# Contraction Mapping Theorem

## Theorem (Contraction Mapping Theorem)

*For any metric space $\mathbb{V}$ that is complete (i.e. closed) under an operator $T$ $(v)$, where $T$ is a $\gamma$-contraction,*

- *$T$ converges to a unique fixed point*
- *At a linear convergence rate of $\gamma$*

- The Bellman expectation operator $T^{\pi}$ has a unique fixed point
- $v_{\pi}$ is a fixed point of $T^{\pi}$ (by Bellman expectation equation)
- By contraction mapping theorem
- Iterative policy evaluation converges on $v_{\pi}$
- Policy iteration converges on $v_*$

- Define the *Bellman optimality backup operator $T*$*,

$$T*(v) = \max_{a \in A} \mathrm{R}^a + \gamma \mathrm{P}^a v$$

- This operator is a $\gamma$-contraction, i.e. it makes value functions closer by at least $\gamma$ (similar to previous proof)

$$||T*(u) - T*(v)||_\infty \leq \gamma ||u - v||_\infty$$

- The Bellman optimality operator $T$ *has a unique fixed point
- $v_*$ is a fixed point of $T$ *(by Bellman optimality equation)  By
- contraction mapping theorem
- Value iteration converges on $v_*$

Most of the story in a nutshell:

# Value Iteration Converges

- If discount factor < 1
- Bellman is a contraction
- Value iteration converges to unique solution which is optimal value function

Most of the story in a nutshell:

# Properties of Contraction

- Only has 1 fixed point
  - If had two, then would not get closer when apply contraction function, violating definition of contraction
- When apply contraction function to any argument, value must get closer to fixed point
  - Fixed point doesn't move
  - Repeated function applications yield fixed point

Most of the story in a nutshell:

# Bellman Operator is a Contraction

$\| V - V' \|$ = Infinity norm
(find max diff
Over all states)

$$\| BV - BV' \| = \left\| \begin{array}{c} \max_a \left[ R(s,a) + \gamma \sum_{s_j \in S} p(s_j \mid s_i, a) V(s_j) \right] \\ - \max_{a'} \left[ R(s,a') - \gamma \sum_{s_j \in S} p(s_j \mid s_i, a') V'(s_j) \right] \end{array} \right\|$$

$$\leq \left\| \max_a \left[ R(s,a) + \gamma \sum_{s_j \in S} p(s_j \mid s_i, a) V(s_j) - R(s,a) + \gamma \sum_{s_j \in S} p(s_j \mid s_i, a) V'(s_j) \right] \right\|$$

$$\leq \gamma \left\| \max_a \left[ \sum_{s_j \in S} p(s_j \mid s_i, a) V(s_j) - \sum_{s_j \in S} p(s_j \mid s_i, a) V'(s_j) \right] \right\|$$

$$= \gamma \max_a \left\| \left[ \sum_{s_j \in S} p(s_j \mid s_i, a) (V(s_j) - V'(s_j)) \right] \right\|$$

$$\leq \gamma \max_{a, s_i} \sum_{s_j \in S} p(s_j \mid s_i, a) \left| V(s_j) - V'(s_j) \right|$$

$$\leq \gamma \max_{a, s_i} \sum_{s_j \in S} p(s_j \mid s_i, a) \| V - V' \|$$

$$= \gamma \| V - V' \|$$

- Does policy evaluation need to converge to $v_\pi$?
- Or should we introduce a stopping condition
    - e.g. $\varepsilon$-convergence of value function
- Or simply stop after $k$ iterations of iterative policy evaluation?
- For example, in the small gridworld $k = 3$ was sufficient to achieve optimal policy
- Why not update policy every iteration? i.e. stop after $k = 1$
    - This is equivalent to *value iteration* (next section)
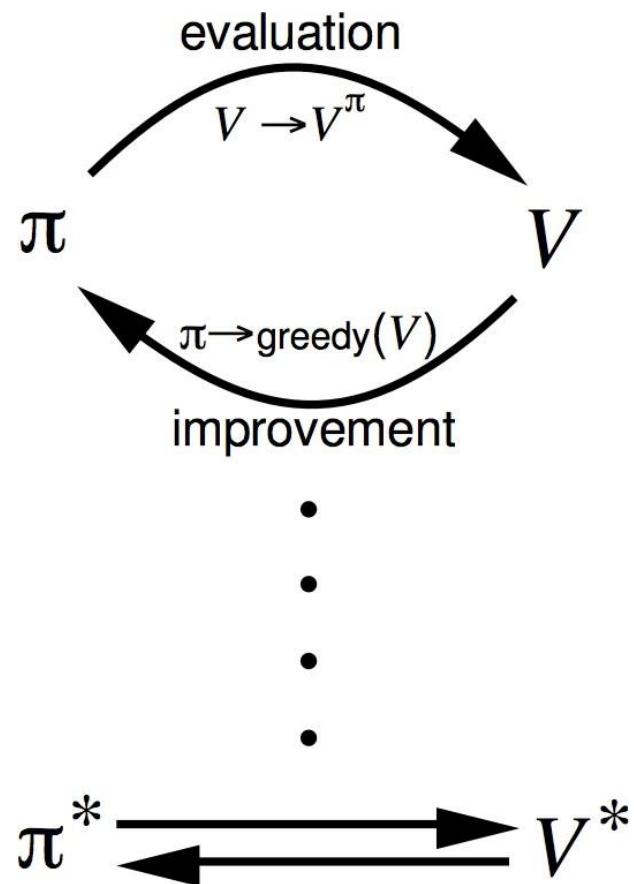
Extensions to Policy Iteration



Policy evaluation Estimate $v_\pi$
  Any policy evaluation algorithm
Policy improvement Generate $\pi' \geq \pi$
  Any policy improvement algorithm

# Value Iteration

- Problem: find optimal policy $\pi$

- Solution: iterative application of Bellman optimality backup

- $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_*$

- Using synchronous backups
    - At each iteration $k + 1$
    - For all states $s \in S$
    - Update $v_{k+1}(s)$ from $v_k(s')$

- Convergence to $v_*$ will be proven later

- Unlike policy iteration, there is no explicit policy

- Intermediate value functions may not correspond to any policy

$$v_{k+1}(s) = \max_{a \in \mathcal{A}} \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_k(s') \right)$$

$$\mathbf{v}_{k+1} = \max_{a \in \mathcal{A}} \mathcal{R}^a + \gamma \mathcal{P}^a \mathbf{v}_k$$

# Asynchronous Dynamic Programming

- DP methods described so far used *synchronous* backups
- i.e. all states are backed up in parallel
- *Asynchronous DP* backs up states individually, in any order
- For each selected state, apply the appropriate backup
- Can significantly reduce computation
- Guaranteed to converge if all states continue to be selected

99

# Asynchronous Dynamic Programming

Three simple ideas for asynchronous dynamic programming:

- *In-place* dynamic programming
- *Prioritised sweeping*
- *Real-time* dynamic programming

# In-Place Dynamic Programming

- Synchronous value iteration stores two copies of value function

  for all $s$ in $\mathcal{S}$

$$v_{new}(s) \leftarrow \max_{a \in \mathcal{A}} \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{old}(s') \right)$$

  $v_{old} \leftarrow v_{new}$

- In-place value iteration only stores one copy of value function

  for all $s$ in $\mathcal{S}$

$$v(s) \leftarrow \max_{a \in \mathcal{A}} \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v(s') \right)$$
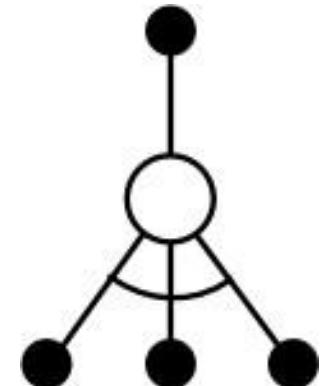
└─ Asynchronous Dynamic Programming

- Use magnitude of Bellman error to guide state selection, e.g.

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

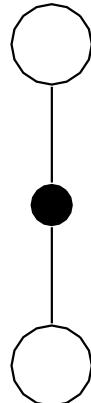- Backup the state with the largest remaining Bellman error
- Update Bellman error of affected states after each backup
- Requires knowledge of reverse dynamics (predecessor states)
- Can be implemented efficiently by maintaining a priority queue

# Real-Time Dynamic Programming

- Idea: only states that are relevant to agent
- Use agent's experience to guide the selection of states
- After each time-step $S_t, A_t, R_{t+1}$
- Backup the state $S_t$

$$v(S_t) \leftarrow \max_{a \in \mathcal{A}} \left( \mathcal{R}_{S_t}^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{S_t s'}^a v(s') \right)$$

# Full-Width Backups

- DP uses *full-width* backups
- For each backup (sync or async)
  - Every successor state and action is considered
  - Using knowledge of the MDP transitions and reward function
- DP is effective for medium-sized problems (millions of states)
- For large problems DP suffers Bellman's *curse of dimensionality*
  - Number of states $n = |S|$ grows exponentially with number of state variables
- Even one backup can be too expensive

$$v_{k+1}(s) \leftarrow s$$

$$a$$

$$r$$

$$v_k(s') \leftarrow s'$$

- In subsequent lectures we will consider *sample backups*

- Using sample rewards and sample transitions
  $(S, A, R, S')$

- Instead of reward function $R$ and transition dynamics $P$

- Advantages:
  - Model-free: no advance knowledge of MDP required
  - Breaks the curse of dimensionality through sampling
  - Cost of backup is constant, independent of $n = |S|$

- Approximate the value function
- Using a *function approximator* $\hat{v}(s, \mathbf{w})$
- Apply dynamic programming to $\hat{v}(\cdot, \mathbf{w})$
- e.g. Fitted Value Iteration repeats at each iteration $k$,
  - Sample states $\tilde{\mathcal{S}} \subseteq \mathcal{S}$
  - For each state $s \in \tilde{\mathcal{S}}$, estimate target value using Bellman optimality equation,

$$\tilde{v}_k(s) = \max_{a \in \mathcal{A}} \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \hat{v}(s', \mathbf{w_k}) \right)$$

  - Train next value function $\hat{v}(\cdot, \mathbf{w_{k+1}})$ using targets $\{\langle s, \tilde{v}_k(s) \rangle\}$

# Monte Carlo Learning

|  | Evaluate Policy, π | Find Best Policy, π* |
|---|---|---|
| MDP Known | Policy Evaluation | Policy/Value Iteration |
| MDP Unknown | MC and TD Learning | Sarsa + Q-Learning |

# Monte-Carlo Reinforcement Learning

MC methods can solve the RL problem by averaging sample returns

- MC methods learn directly from episodes of experience
- MC is *model-free*: no knowledge of MDP transitions / rewards
- MC learns from *complete* episodes: no bootstrapping
- MC uses the simplest possible idea: value = mean return
- Caveat: can only apply MC to *episodic* MDPs
  - All episodes must terminate

MC is incremental episode by episode but not step by step

Approach: adapting general policy iteration to sample returns

First policy evaluation, then policy improvement, then control

# Monte-Carlo Policy Evaluation

- Goal: learn $v_\pi$ from episodes of experience under policy $\pi$

$$S_1, A_1, R_2, ..., S_k \sim \pi$$

- Recall that the *return* is the total discounted reward:

$$G_t = R_{t+1} + \gamma R_{t+2} + ... + \gamma^{T-1} R_T$$

- Recall that the value function is the expected return:

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

- Monte-Carlo policy evaluation uses *empirical mean* return instead of *expected* return, because we do not have the model

# Every Visit MC Policy Evaluation

- To evaluate state $s$
- **Every** time-step $t$ that state $s$ is visited in an episode,
- Increment counter $N(s) \leftarrow N(s) + 1$
- Increment total return $S(s) \leftarrow S(s) + G_t$
- Value is estimated by mean return $V(s) = S(s)/N(s)$
- Again, $V(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$

Equivalent, "incremental tracking" form:
$$V(s) \leftarrow V(s) + \frac{1}{N(s)}(G_t - V(s))$$
Looks like SGD to minimize MSE from the mean value...

# Blackjack Example

- States (200 of them):

    - Current sum (12-21)
    - Dealer's showing card (ace-10)
    - Do I have a "useable" ace? (yes-no)

- Action stand  Stop receiving cards (and terminate)
- Action hit  : Take another card (no replacement)

- Reward for stand

    - $+1$ if sum of cards $>$ sum of dealer cards
    - 0 if sum of cards $=$ sum of dealer cards
    - -1 if sum of cards $<$ sum of dealer cards

- Reward for hit  :

    - -1 if sum of cards $>$ 21 (and terminate)
    - 0 otherwise

- Transitions: automatically hit  if sum of cards $<$ 12

# Blackjack Value Function



After 10,000 episodes      After 500,000 episodes

Usable ace

No usable ace

+1

−1

Dealer showing   A ... 10   12   Player sum   21

Policy: stand if sum of cards $\geq$ 20, otherwise hit

# Temporal Difference Learning

- TD methods learn directly from episodes of experience
- TD is *model-free*: no knowledge of MDP transitions / rewards
- TD learns from *incomplete* episodes, by *bootstrapping*
- TD updates a guess towards a guess

# MC and TD

- Goal: learn $v_\pi$ online from experience under policy $\pi$
- Incremental every-visit Monte-Carlo
    - Update value $V(S_t)$ toward *actual* return $G_t$

$$V(S_t) \leftarrow V(S_t) + \alpha\,(G_t - V(S_t))$$

- Simplest temporal-difference learning algorithm: TD(0)
    - Update value $V(S_t)$ toward *estimated* return $R_{t+1} + \gamma V(S_{t+1})$

$$V(S_t) \leftarrow V(S_t) + \alpha\,(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

    - $R_{t+1} + \gamma V(S_{t+1})$ is called the *TD target*
    - $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$ is called the *TD error*

# Driving Home Example

| State | Elapsed Time (minutes) | Predicted Time to Go | Predicted Total Time |
|---|---|---|---|
| leaving office | 0 | 30 | 30 |
| reach car, raining | 5 | 35 | 40 |
| exit highway | 20 | 15 | 35 |
| behind truck | 30 | 10 | 40 |
| home street | 40 | 3 | 43 |
| arrive home | 43 | 0 | 43 |

# Driving Home: MC vs TD



Changes recommended by Monte Carlo methods ($\alpha=1$)

Changes recommended by TD methods ($\alpha=1$)

# Finite Episodes: AB Example

Two states $A$, $B$; no discounting; 8 episodes of experience

$A, 0, B, 0$
$B, 1$
$B, 1$
$B, 1$
$B, 1$
$B, 1$
$B, 1$
$B, 0$

What is $V(A)$, $V(B)$?

MC & TD can give different answers on fixed data:

V(B) = 6 / 8

V(A) = 0 ?        (Direct MC estimate)

V(A) = 6 / 8?   (TD estimate)

# MC vs TD

**Monte Carlo**

- Wait till end of episode to learn
  - Only for *terminating* worlds

- High-variance, low bias
  - Not sensitive to initial value
  - Good convergence properties

- Doesn't exploit Markov property

- Minimizes squared error

**Temporal Difference**

- Learn online after every step
  - Non-*terminating* worlds ok

- Low variance, high bias
  - Sensitive to initial value
  - Much more efficient

- Exploits Markov Property

- Maximizes log-likelihood

# Unified View: Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t - V(S_t) \right)$$

# Unified View: TD Learning

$$V(S_t) \leftarrow V(S_t) + \alpha \left( R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right)$$

# Unified View: Dynamic Prog.

$$V(S_t) \leftarrow \mathbb{E}_\pi \left[ R_{t+1} + \gamma V(S_{t+1}) \right]$$
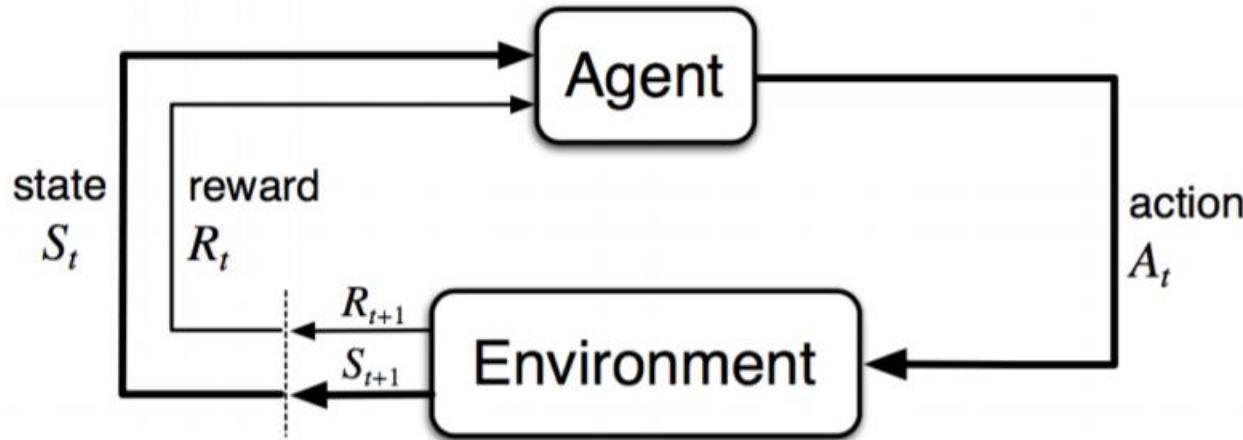
# Unified View of RL (Prediction)
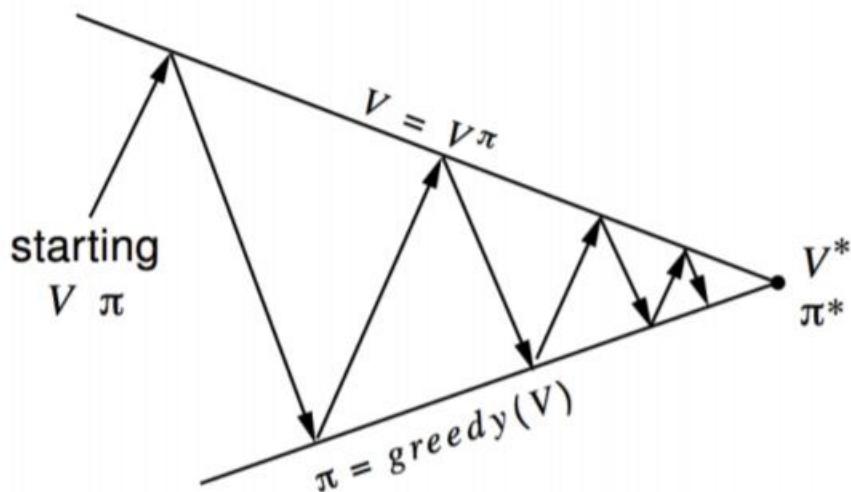
# Overview

|  | Evaluate Policy, $\pi$ | Find Best Policy, $\pi^*$ |
|---|---|---|
| **MDP Known** | Policy Evaluation | Policy/Value Iteration |
| **MDP Unknown** | MC and TD Learning | Sarsa + Q-Learning |

# Which Policy Evaluation?

- Temporal-difference (TD) learning has several advantages over Monte-Carlo (MC)
    - Lower variance
    - Online
    - Incomplete sequences
- Natural idea: use TD
    - Apply TD to $Q(S, A)$
    - Use $\epsilon$-greedy policy improvement
    - Update every time-step

# Model-free Control



state $S_t$

reward $R_t$

$R_{t+1}$

$S_{t+1}$

Agent

Environment

action $A_t$

Learn a policy $\pi$ to maximize rewards in the environment

# Generalized Policy Iteration



Policy evaluation  Estimate $v_\pi$
  e.g. Iterative policy evaluation

Policy improvement  Generate $\pi' \geq \pi$
  e.g. Greedy policy improvement

# Gen Policy Improvement?



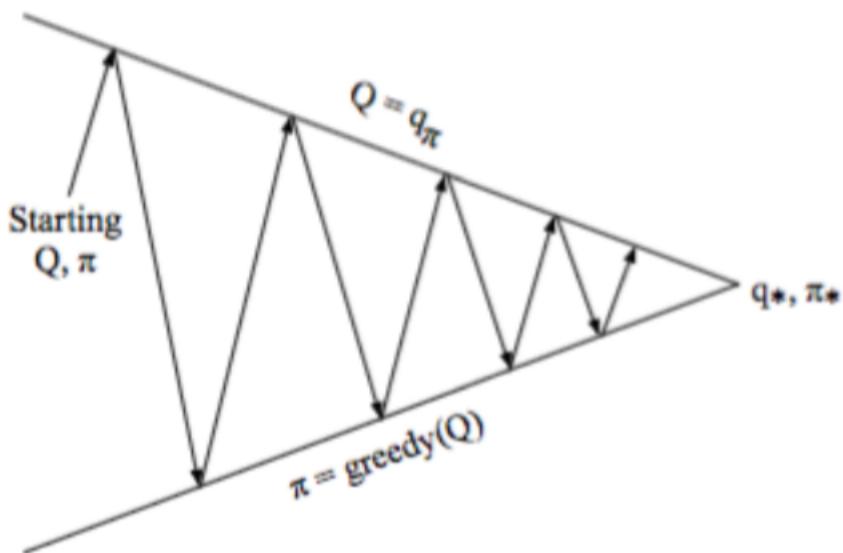Policy evaluation Monte-Carlo policy evaluation, $V = v_\pi$?

Policy improvement Greedy policy improvement?

# Not quite!

- Greedy policy improvement over $V(s)$ requires model of MDP

$$\pi'(s) = \underset{a \in \mathcal{A}}{\mathrm{argmax}} \; \mathcal{R}_s^a + \mathcal{P}_{ss'}^a V(s')$$

- Greedy policy improvement over $Q(s, a)$ is model-free
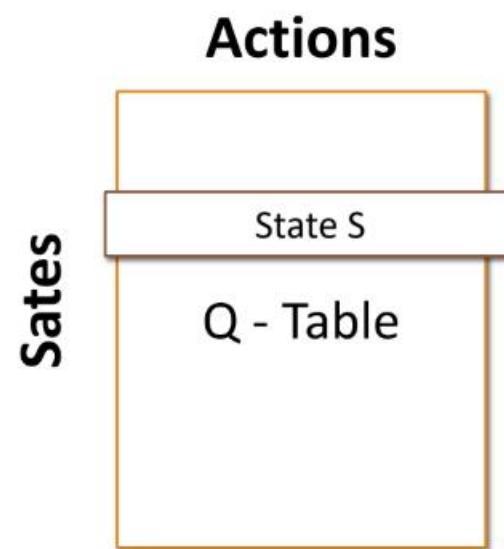
$$\pi'(s) = \underset{a \in \mathcal{A}}{\mathrm{argmax}} \; Q(s, a)$$

# Learn Q function directly...



Policy evaluation Monte-Carlo policy evaluation, $Q = q_\pi$
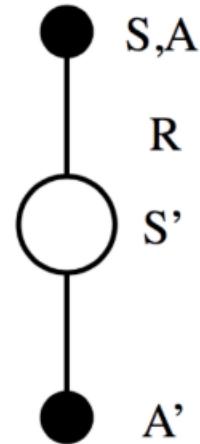
Policy improvement Greedy policy improvement?

# Q-Learning

**Actions**

**Sates**

State S

Q - Table

# On and Off Policy Learning

- **On-policy** learning
  - "Learn on the job"
  - Learn about policy $\pi$ from experience sampled from $\pi$
- **Off-policy** learning
  - "Look over someone's shoulder"
  - Learn about policy $\pi$ from experience sampled from $\mu$

# Sarsa: TD for Policy Evaluation



$$Q(S, A) \leftarrow Q(S, A) + \alpha \left( R + \gamma Q(S', A') - Q(S, A) \right)$$

# SARSA

**Sarsa (on-policy TD control) for estimating $Q \approx q_*$**

Initialize $Q(s, a)$, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(terminal\text{-}state, \cdot) = 0$
Repeat (for each episode):
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
    Repeat (for each step of episode):
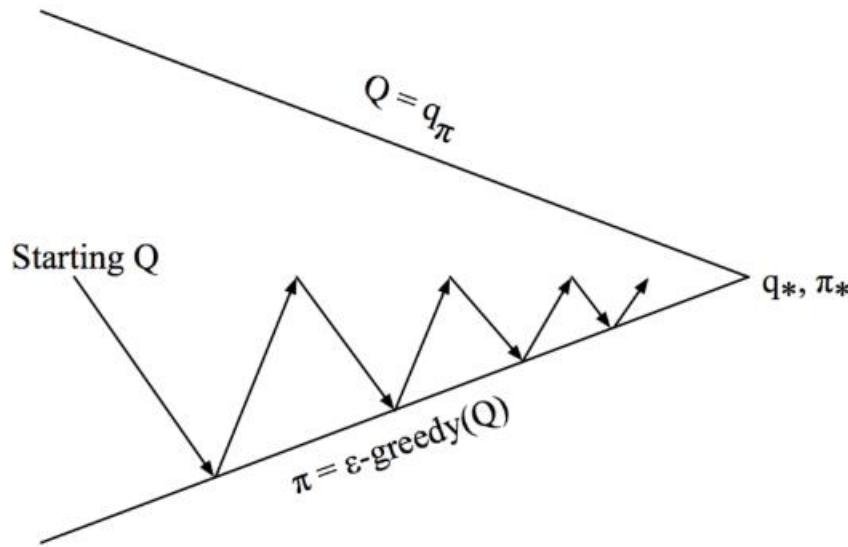        Take action $A$, observe $R, S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

# On-Policy Control w/ Sarsa



Every time-step:

Policy evaluation Sarsa, $Q \approx q_\pi$

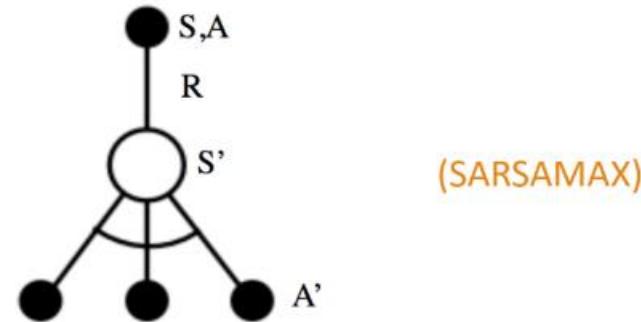Policy improvement $\epsilon$-greedy policy improvement

# Q-Learning

Learning
Rate

Future Reward

$$Q(S, a) \leftarrow Q(S, a) + \alpha \left( R + \gamma \max_{a'} Q(S', a') - Q(S, a) \right)$$

Current
Value

Reward

Current
Value Offset

# Q-Learning Control Algorithm



(SARSAMAX)

$$Q(S, A) \leftarrow Q(S, A) + \alpha \left( R + \gamma \max_{a'} Q(S', a') - Q(S, A) \right)$$

# Q-Learning

**Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$**

Initialize $Q(s, a)$, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\textit{terminal-state}, \cdot) = 0$
Repeat (for each episode):
  Initialize $S$
  Repeat (for each step of episode):
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\epsilon$-greedy)
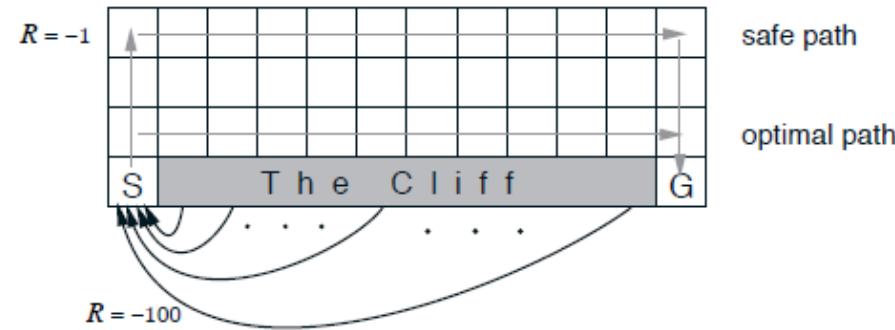    Take action $A$, observe $R$, $S'$
    $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
    $S \leftarrow S'$
  until $S$ is terminal

# Q-Learning vs. Sarsa

# Greedy Action Selection?

- There are two doors in front of you.
- You open the left door and get reward 0
  $V(left) = 0$
- You open the right door and get reward $+1$
  $V(right) = +1$
- You open the right door and get reward $+3$
  $V(right) = +2$
- You open the right door and get reward $+2$
  $V(right) = +2$

$$\vdots$$

- Are you sure you've chosen the best door?

# $\epsilon$-Greedy Exploration

- Simplest idea for ensuring continual exploration
- All $m$ actions are tried with non-zero probability
- With probability $1 - \epsilon$ choose the greedy action
- With probability $\epsilon$ choose an action at random
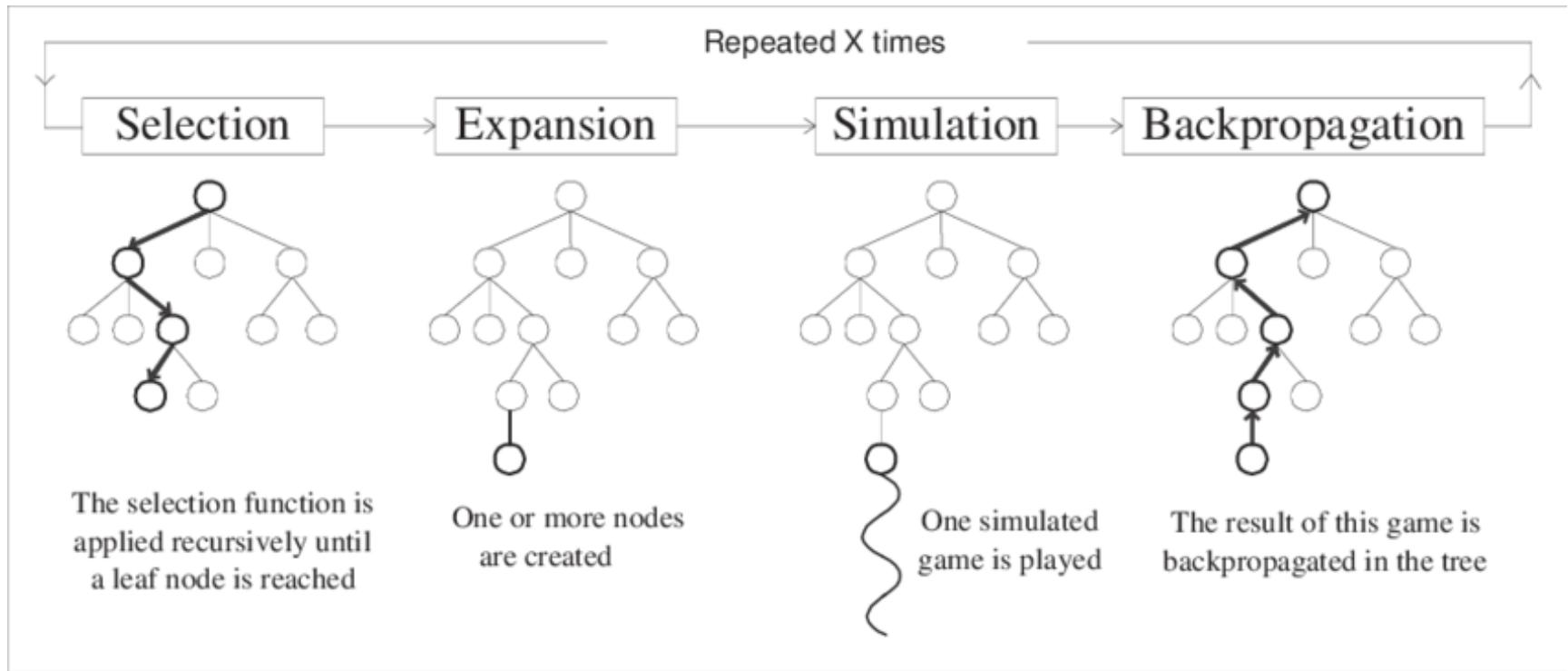
# Relation between DP and TD



|  | Full Backup (DP) | Sample Backup (TD) |
|---|---|---|
| Bellman Expectation Equation for $v_\pi(s)$ | Iterative Policy Evaluation | TD Learning |
| Bellman Expectation Equation for $q_\pi(s, a)$ | Q-Policy Iteration | Sarsa |
| Bellman Optimality Equation for $q_*(s, a)$ | Q-Value Iteration | Q-Learning |

# Update Eqns for DP and TD

| Full Backup (DP) | Sample Backup (TD) |
|---|---|
| Iterative Policy Evaluation | TD Learning |
| $V(s) \leftarrow \mathbb{E}\left[R + \gamma V(S') \mid s\right]$ | $V(S) \overset{\alpha}{\leftarrow} R + \gamma V(S')$ |
| Q-Policy Iteration | Sarsa |
| $Q(s, a) \leftarrow \mathbb{E}\left[R + \gamma Q(S', A') \mid s, a\right]$ | $Q(S, A) \overset{\alpha}{\leftarrow} R + \gamma Q(S', A')$ |
| Q-Value Iteration | Q-Learning |
| $Q(s, a) \leftarrow \mathbb{E}\left[R + \gamma \max\limits_{a' \in \mathcal{A}} Q(S', a') \mid s, a\right]$ | $Q(S, A) \overset{\alpha}{\leftarrow} R + \gamma \max\limits_{a' \in \mathcal{A}} Q(S', a')$ |

where $x \overset{\alpha}{\leftarrow} y \equiv x \leftarrow x + \alpha(y - x)$

# Monte Carlo Tree Search

# Large-Scale RL

Reinforcement learning can be used to solve *large* problems, e.g.

- Backgammon: $10^{20}$ states
- Computer Go: $10^{170}$ states
- Helicopter: continuous state space

How can we scale up the model-free methods for *prediction* and *control* from the last two lectures?

# Value Function Approximation

- So far we have represented value function by a *lookup table*
    - Every state $s$ has an entry $V(s)$
    - Or every state-action pair $s, a$ has an entry $Q(s, a)$
- Problem with large MDPs:
    - There are too many states and/or actions to store in memory
    - It is too slow to learn the value of each state individually
- Solution for large MDPs:
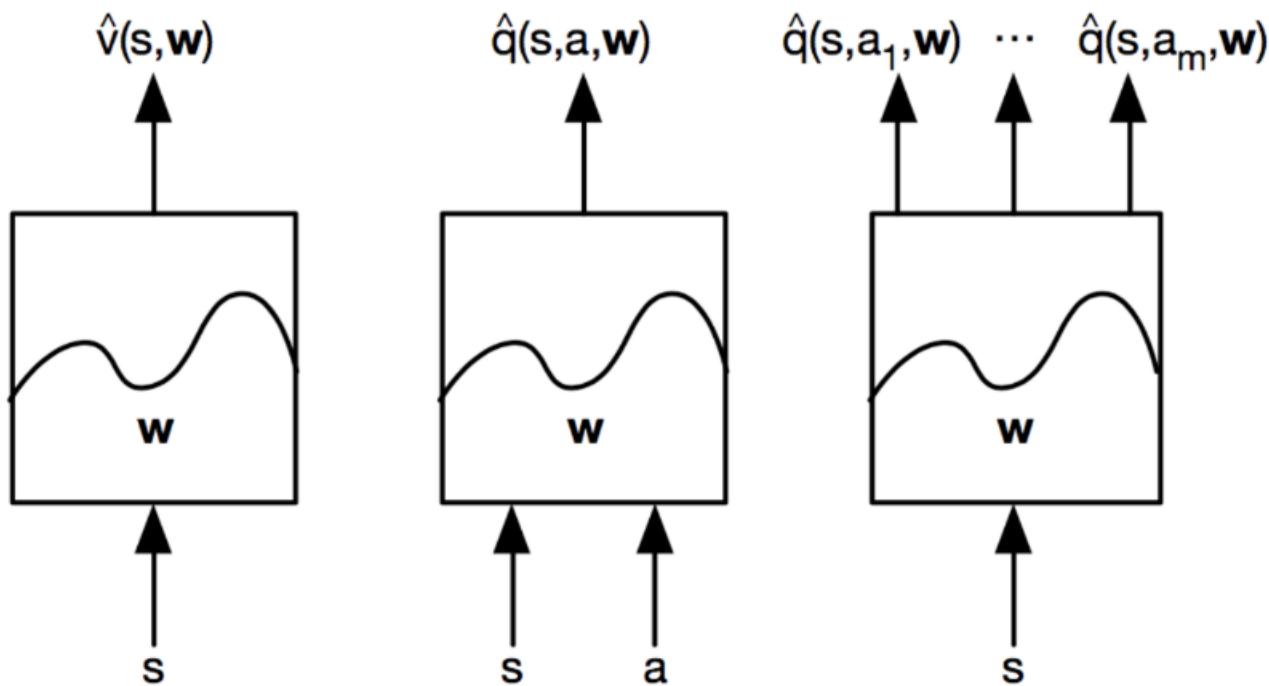    - Estimate value function with *function approximation*

$$\hat{v}(s, \mathbf{w}) \approx v_\pi(s)$$
$$\text{or } \hat{q}(s, a, \mathbf{w}) \approx q_\pi(s, a)$$

- *Generalise* from seen states to unseen states
- *Update* parameter $\mathbf{w}$ using MC or TD learning

# Types of Function Approx.

# Which Approximator?

There are many function approximators, e.g.

- Linear combinations of features
- Neural network
- Decision tree
- Nearest neighbour
- Fourier / wavelet bases
- ...

# Deep-Q learning

Use deep neural network architectures for Q(s,a)

Ex: Atari game playing (DeepMind)
◦ Input: pixel images of current state
◦ Output: joystick actions