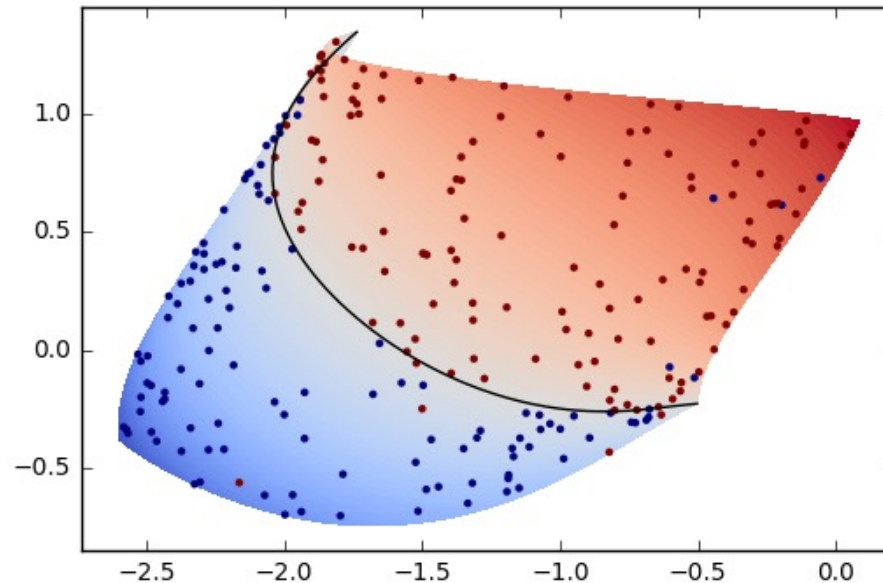# Lecture 26: Reinforcement Learning Part 3



Gavin Kerrigan

Spring 2023

Some slides adapted from Padhraic Smyth, Alex Ihler

# Announcements

- Final course eval
  - evaluations.eee.uci.edu
  - Due 6/11
  - 30/162 students have completed so far

- HW5 due in one week (Friday 6/9)

- Project due in ~1 week (Monday 6/12)

- Next week's "advanced topics" lectures:
  - Monday: Natural Language Processing
  - Weds: Generative Models
  - Friday: Final review

Review

Markov Decision Processes

Policy Evaluation and Improvement

# Markov Reward Process

A Markov reward process is a Markov chain with values.

## Definition

A *Markov Reward Process* is a tuple $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

- $\mathcal{S}$ is a finite set of states
- $\mathcal{P}$ is a state transition probability matrix,
  $\mathcal{P}_{ss'} = \mathbb{P}[S_{t+1} = s' \mid S_t = s]$
- $\mathcal{R}$ is a reward function, $\mathcal{R}_s = \mathbb{E}[R_{t+1} \mid S_t = s]$
- $\gamma$ is a discount factor, $\gamma \in [0, 1]$

Reward can be stochastic or deterministic (here, we often consider deterministic)

$R_s$ is the average reward we receive from being in state s

# Returns

**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 Functions

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

**Definition**

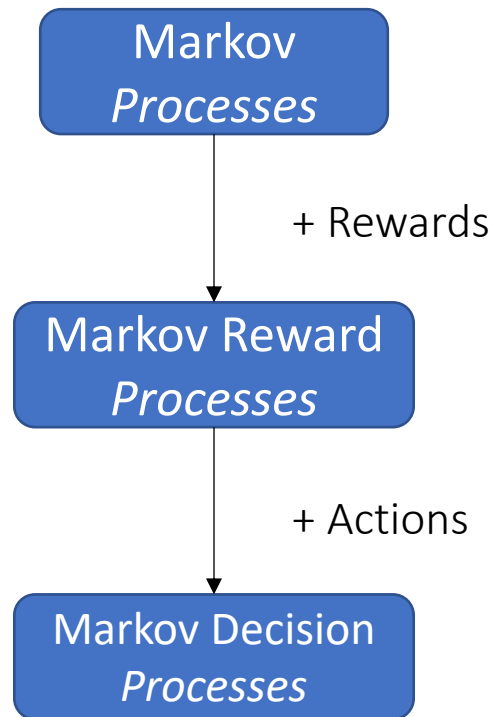The *state value function* $v(s)$ of an MRP is the expected return starting from state $s$

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

Review

Markov Decision Processes

Policy Evaluation and Improvement

# Where We're Headed

Markov *Processes*

+ Rewards

Markov Reward *Processes*

+ Actions

Markov Decision *Processes*

# Markov Decision Processes

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

# Markov Decision Processes

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 $\gamma \in [0, 1]$.

# Policies

**Definition**

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

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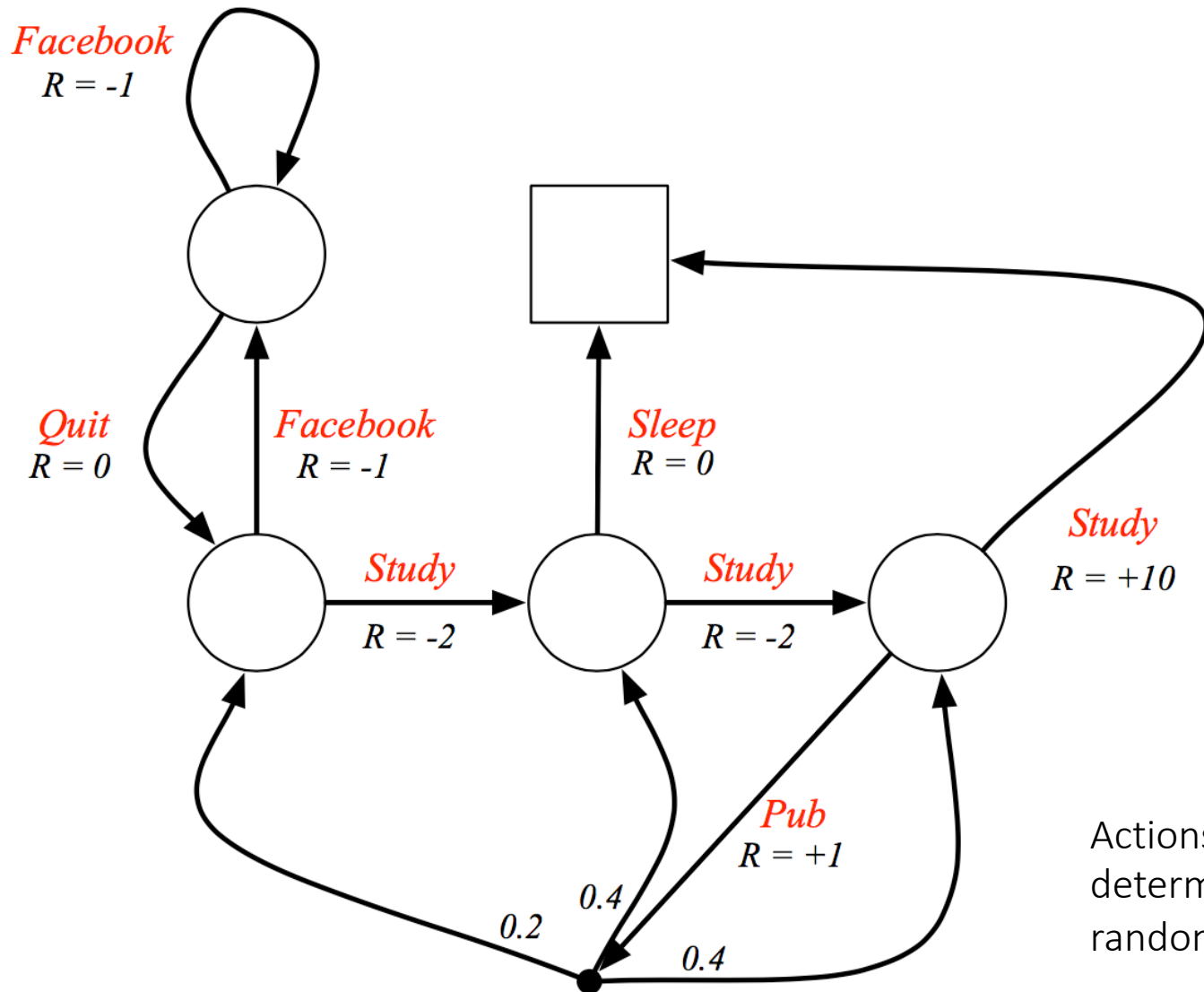
# Policies

**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),

# Student MDP



Facebook
R = -1

Quit
R = 0

Facebook
R = -1

Sleep
R = 0

Study
R = +10

Study
R = -2

Study
R = -2

Pub
R = +1

0.4

0.2

0.4

Actions can result in deterministic or random next states

# State-Value Functions

**Definition**

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

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

- Policy determines which states we will enter
- Hence value of a state depends on the policy

# State-Value Functions

**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 $\pi$

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

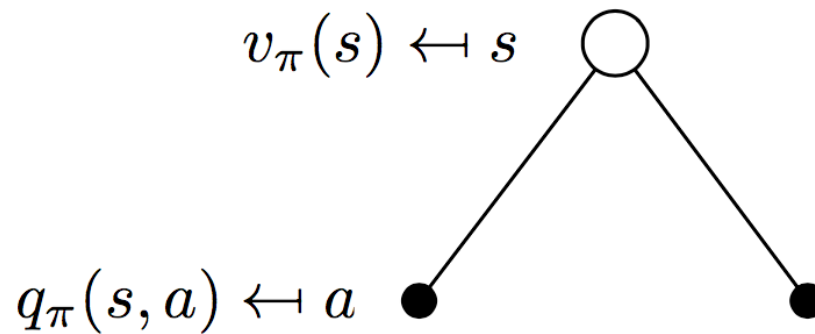# Questions?

# Bellman Equations

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

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

# Bellman Equations

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,

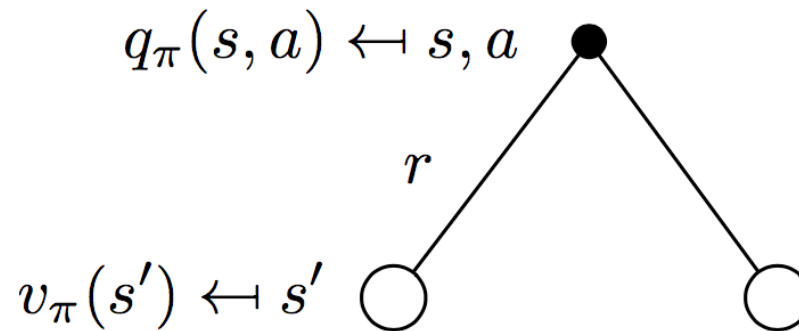$$q_\pi(s, a) = \mathbb{E}_\pi \left[ R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a \right]$$

# Backup Diagrams for V



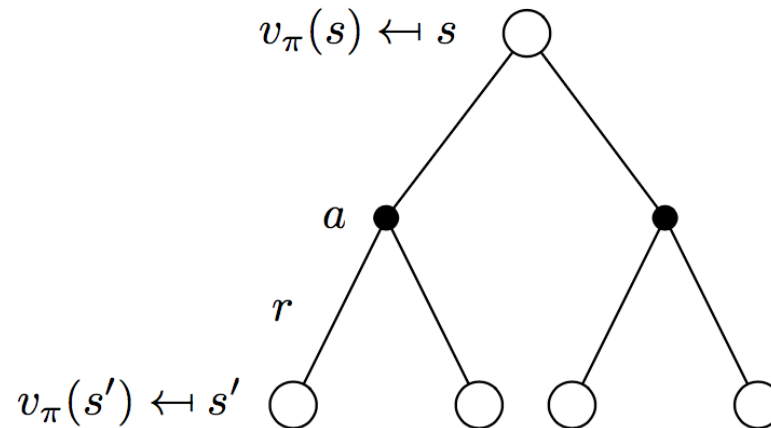$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) q_\pi(s, a)$$

# Backup Diagrams for Q



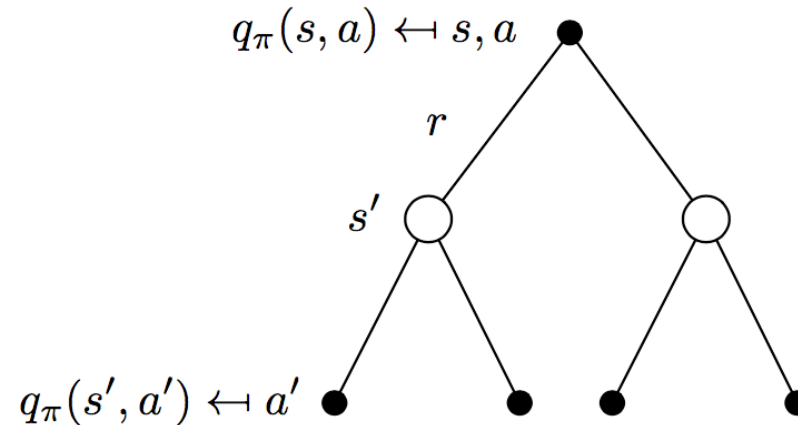$$q_\pi(s, a) \leftarrow s, a$$

$$r$$

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

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

# Bellman Equations for V



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

# Bellman Equations for Q



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

# Optimal Value Functions

**Definition**

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

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

# Optimal Value Functions

## Definition

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

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

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

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

# Optimal Value Functions

## 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.
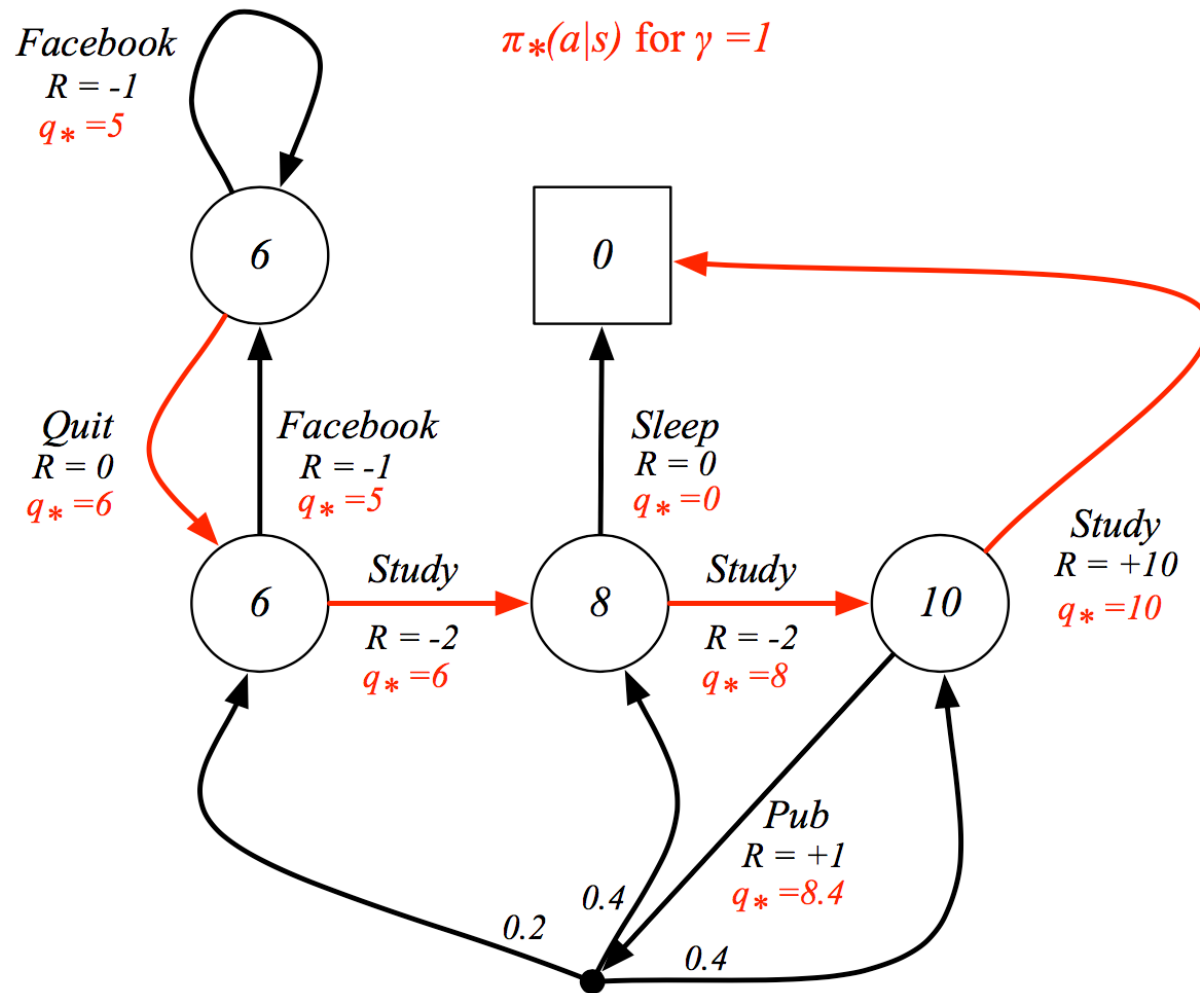- An MDP is "solved" when we know the optimal value fn.

# Finding an Optimal Policy

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
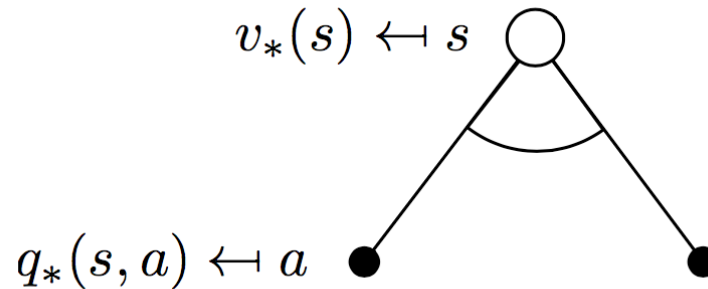
# Student MDP: Optimal Policy

# Bellman Optimality for V

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

$$v_*(s) \leftarrowtail s \; \bigcirc$$

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

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

# Bellman Optimality for V



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

Non-Linear equation: can't be solved easily
- Most of RL is focused on solving this problem!

# Questions?

Review

Markov Decision Processes

Policy Evaluation and Improvement

# An Overview of RL Problems

MDPs in general model states, transitions, actions rewards

**Prediction**: Given policy $\pi$ : Estimate State/Action value functions

**Control**: Estimate *optimal* value functions, *optimal* policy

Is the MDP known?

- Yes: Agent is then "planning"; everything is known about environment

- No: "Model-Free RL"; agent observes as it goes

# An Overview of RL Problems

|  | Evaluate Policy, π (Prediction) | Find Best Policy, π* (Control) |
|---|---|---|
| MDP Known | Policy Evaluation | Policy Iteration |
| MDP Unknown (Model-free) | Monte Carlo and Temporal Difference Learning | Q-Learning |

# Iterative Policy Evaluation

- Problem: evaluate a given policy $\pi$
- Solution: iterative application of Bellman expectation backup
- $v_1 \rightarrow v_2 \rightarrow ... \rightarrow v_\pi$
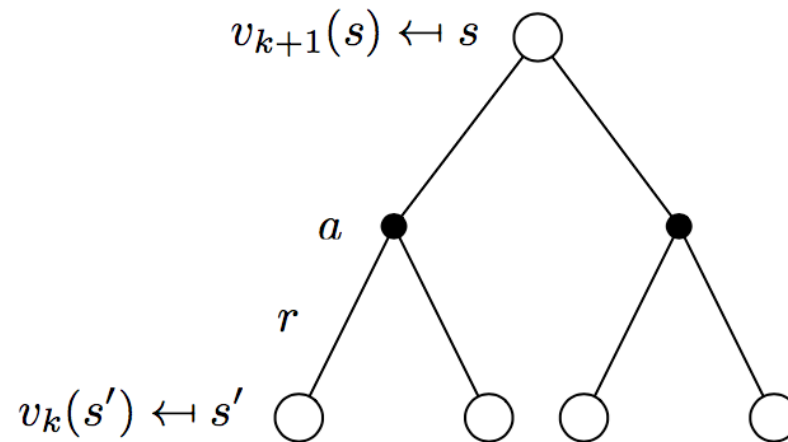
# Iterative Policy Evaluation

- 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 $s \in \mathcal{S}$
    - Update $v_{k+1}(s)$ from $v_k(s')$
    - where $s'$ is a successor state of $s$

# Iterative Policy Evaluation

$$v_{k+1}(s) \leftarrowtail s \quad \bigcirc$$

$$a$$

$$r$$

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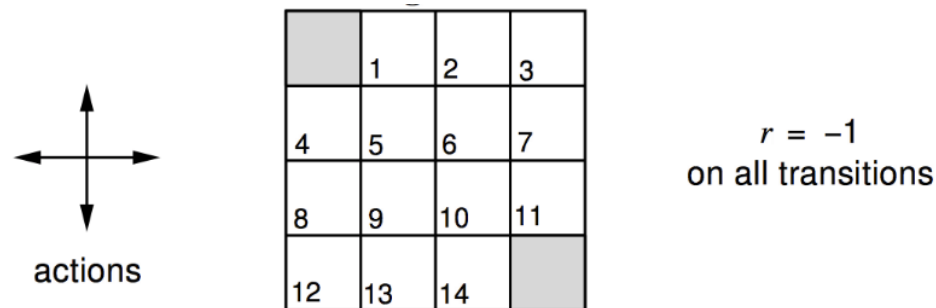
$$v_{k+1}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left( \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_k(s') \right)$$

$$\mathbf{v}^{k+1} = \boldsymbol{\mathcal{R}}^{\boldsymbol{\pi}} + \gamma \boldsymbol{\mathcal{P}}^{\boldsymbol{\pi}} \mathbf{v}^k$$
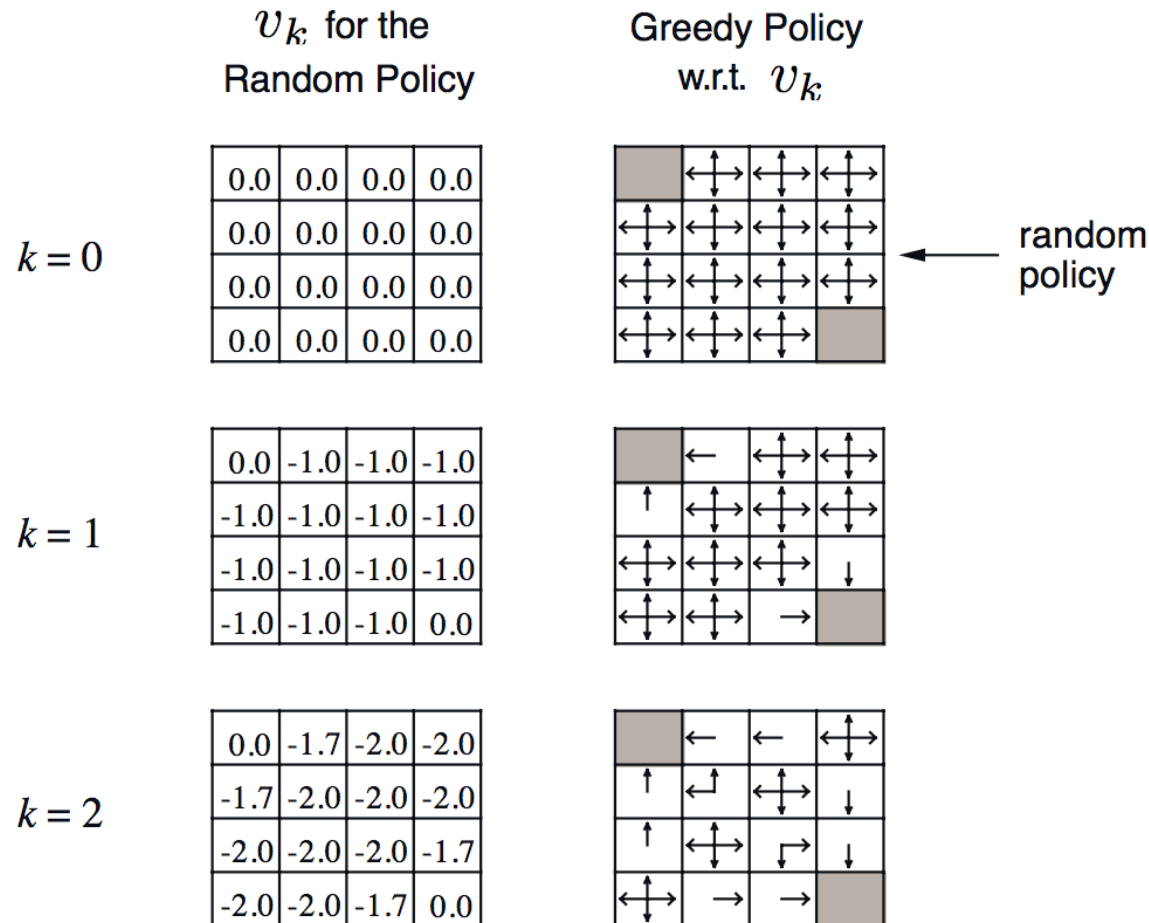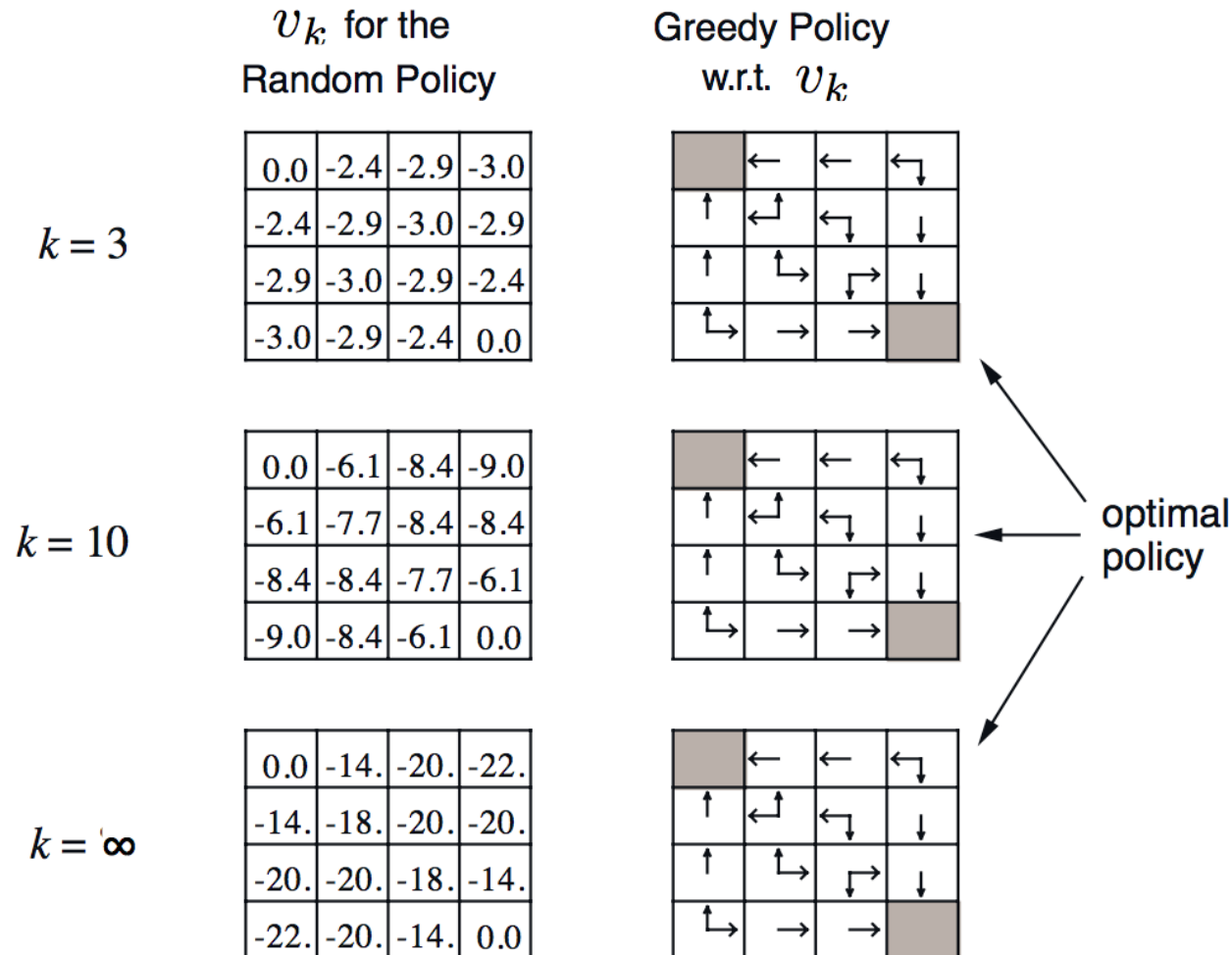
# Grid World



- 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

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

$v_k$ for the Random Policy — Greedy Policy w.r.t. $v_k$

$k = 0$

| 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 |

random policy

$k = 1$

| 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 |

$k = 2$

| 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 |

$v_k$ 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 |

$k = \infty$

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

optimal policy

# Questions?

# An Overview of RL Problems

|  | Evaluate Policy, π (Prediction) | Find Best Policy, π* (Control) |
|---|---|---|
| MDP Known | Policy Evaluation | Policy Iteration |
| MDP Unknown (Model-free) | Monte Carlo and Temporal Difference Learning | Q-Learning |

# Improving a Policy

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

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

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

    $$\pi' = \text{greedy}(v_\pi)$$

# Improving a Policy

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

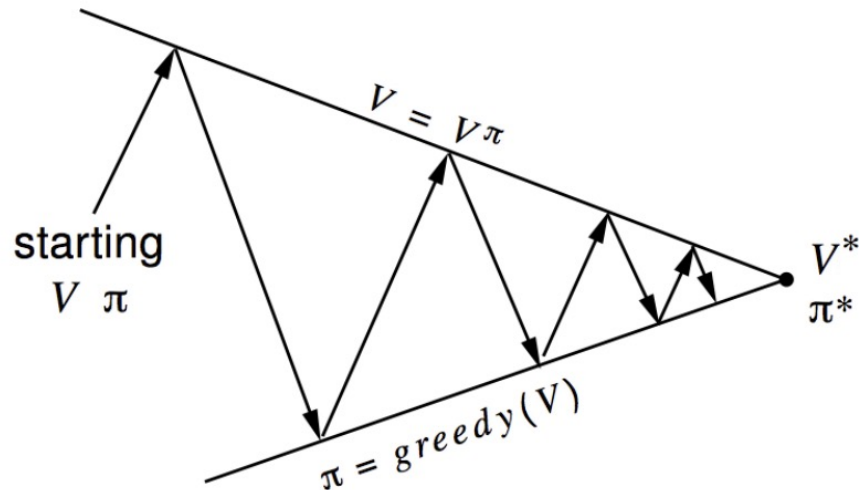$$v_\pi(s) = \mathbb{E}\left[R_{t+1} + \gamma R_{t+2} + \ldots | S_t = s\right]$$

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

$$\pi' = \text{greedy}(v_\pi)$$

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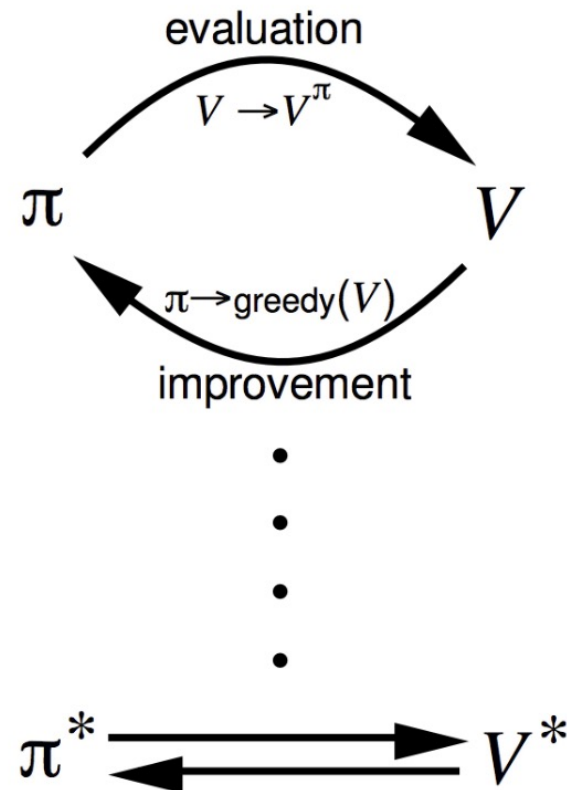- But this process of policy iteration always converges to $\pi*$

# Policy Iteration



Policy evaluation Estimate $v_\pi$
  Iterative policy evaluation

Policy improvement Generate $\pi' \geq \pi$
  Greedy policy improvement

# Questions?

# Wrapup

### Reinforcement Learning is built on MDPs

- States, Actions, Rewards, Transitions, Discount

### Is the underlying MDP known?

- Yes: Agent needs to find optimal policy ("planning")
- No: Agent must also discover the MDP ("model-free RL")

### If the MDP is known: learn optimal policy iteratively

- Evaluate policy
- Improve policy by behaving greedily

### RL is a huge field: we have barely covered the basics