

CS6375: Machine Learning

Gautam Kunapuli

Reinforcement Learning



THE UNIVERSITY OF TEXAS AT DALLAS

Erik Jonsson School of Engineering and Computer Science

Reinforcement Learning

Supervised learning: Given **labeled** data $(x_i, y_i), i = 1, \dots, n$, learn a function $f : x \rightarrow y$

- Categorical y : **classification**
- Continuous y : **regression**

Rich feedback from the environment: the learner is told exactly what it should have done

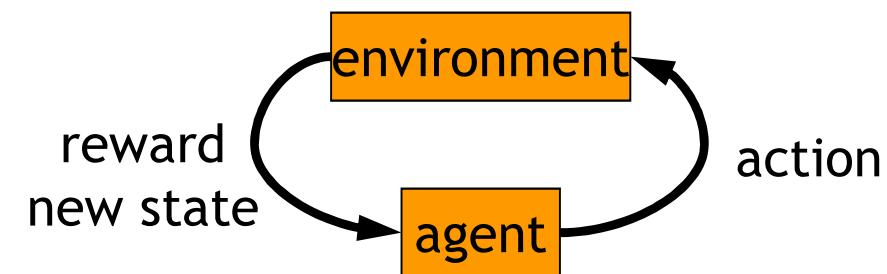
Unsupervised learning: Given **unlabeled** data $x_i, i = 1, \dots, n$, can we infer the underlying structure?

- Clustering
- dimensionality reduction,
- density estimation

No feedback from the environment: the learner receives no labels or any other information

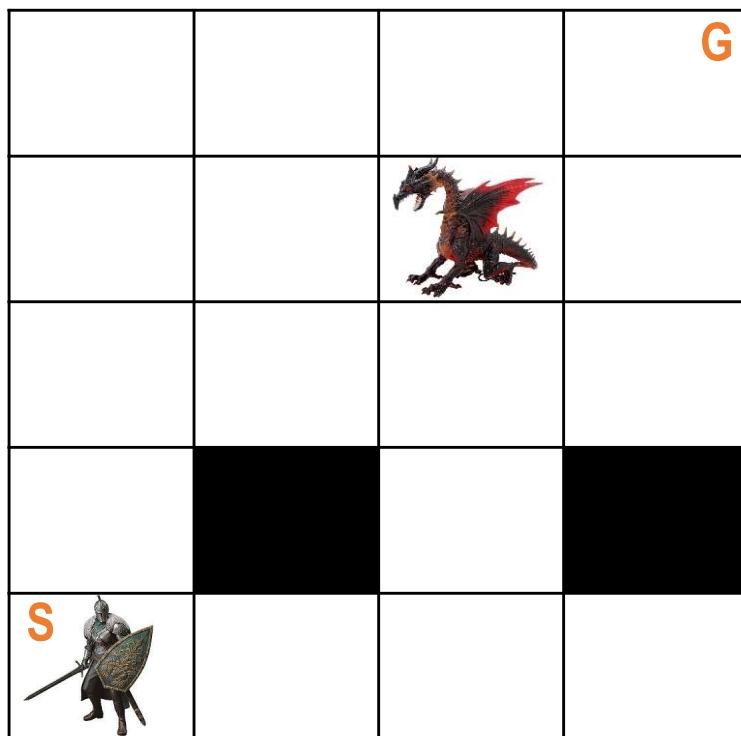
Reinforcement Learning is learning from Interaction: learner (agent) receives feedback about the appropriateness of its actions while interacting with an environment, which provides numeric reward signals

Goal: Learn how to take actions in order to maximize reward



Example: Grid World

Example: Learn to navigate from beginning/start state (S) to goal state (G), while avoiding obstacles



Autonomous “**agent**” interacts with an **environment** through a series of **actions**

- trying to find the way through a maze
- actions include turning and moving through maze
- agent earns rewards from the environment under certain (perhaps unknown) conditions

The agent’s goal is to maximize the reward

- we say that the agent learns if, over time, it improves its performance

actions are what the agent actually wants to do

Actions

→(right)

↑(up)

←(left)

↓(down)

effects are what actually happens after the agent executes the chosen action

Effects

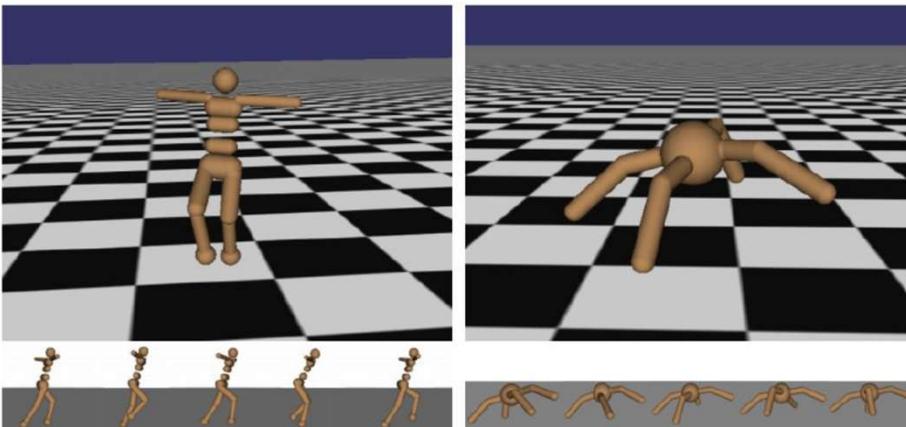
→(60%), ↓(40%)

↑(100%)

←(100%)

↓(70%), ←(30%)

Applications of Reinforcement Learning



Schulman et al (2016)

Robot Locomotion (and other control problems)

Objective: Make the robot move forward

State: Angle and position of the joints

Action: Torques applied on joints

Reward: 1 at each time step upright + forward movement

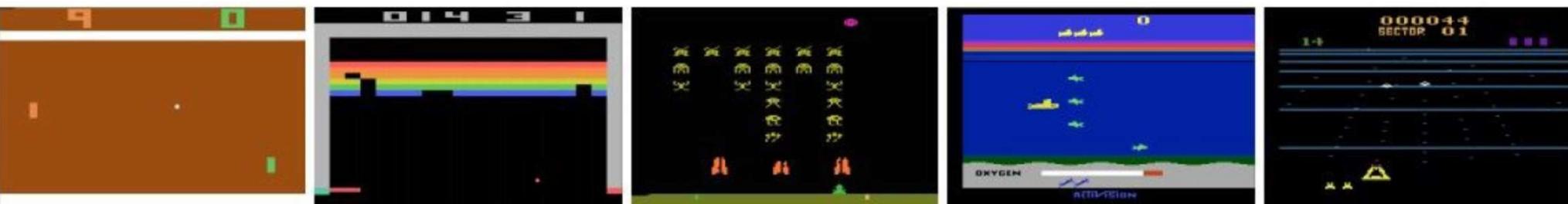
Atari Games

Objective: Complete the game with the highest score

State: Raw pixel inputs of the game state

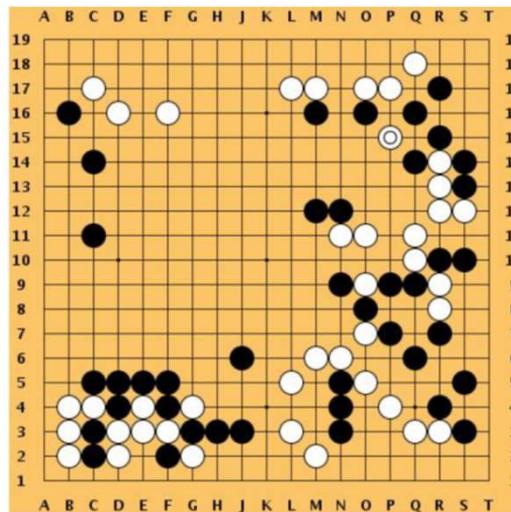
Action: Game controls e.g. Left, Right, Up, Down

Reward: Score increase/decrease at each time step



Adapted from slides by Fei-Fei Li, Justin Johnson and Serena Yeung

Applications of Reinforcement Learning



Go!

Objective: Win the game!

State: Position of all pieces

Action: Where to put the next piece down

Reward: 1 if win at the end of the game, 0 otherwise

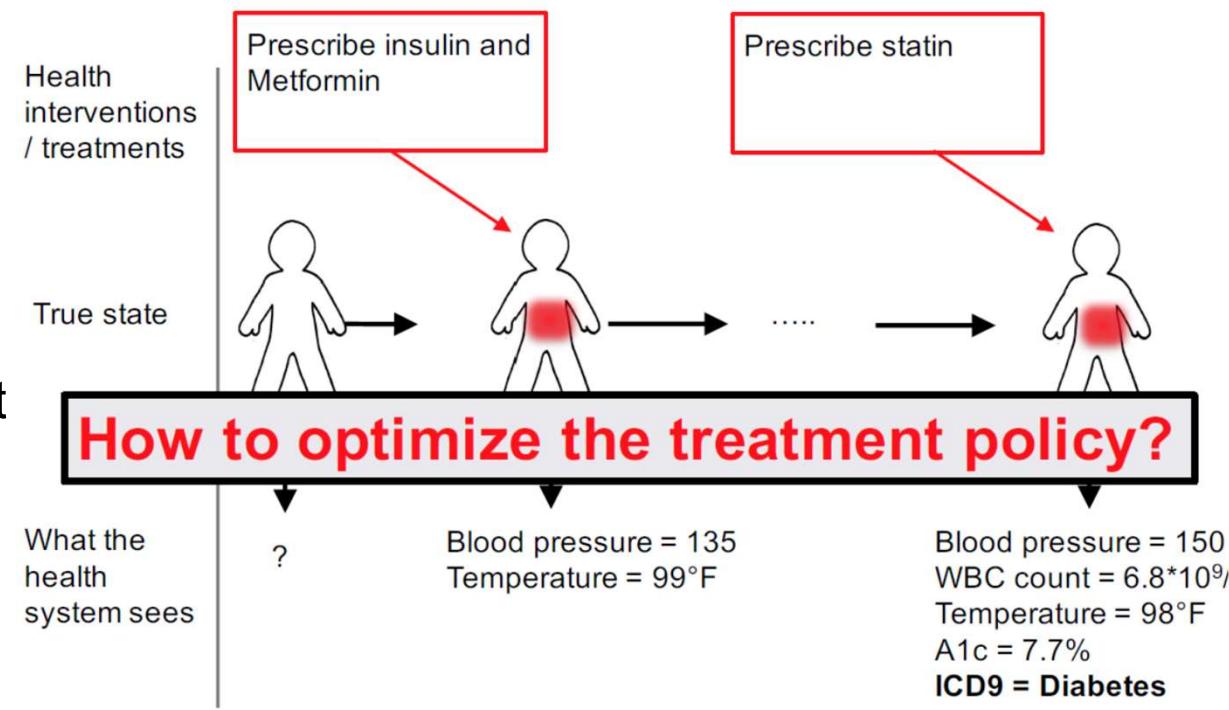
Treatment Planning

Objective: Find the best treatment policy

State: Patient health data every 6 months

Action: Clinical interventions and treatment

Reward: negative rewards for deterioration
positive rewards for improvement



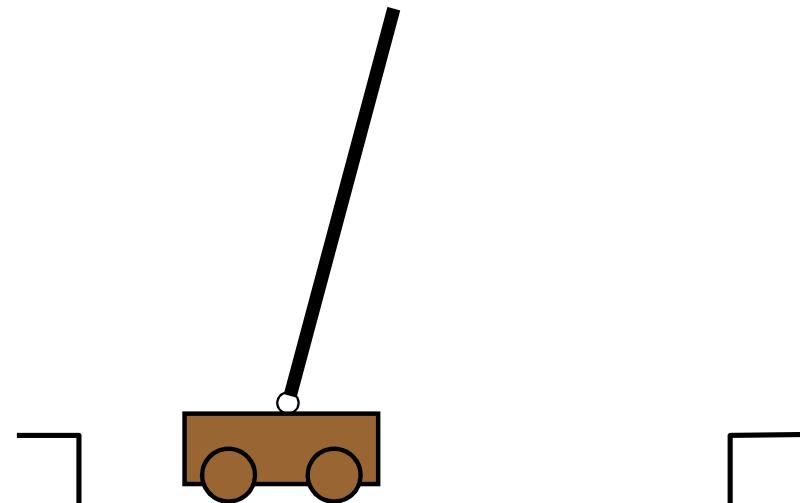
Reinforcement Learning

Other examples

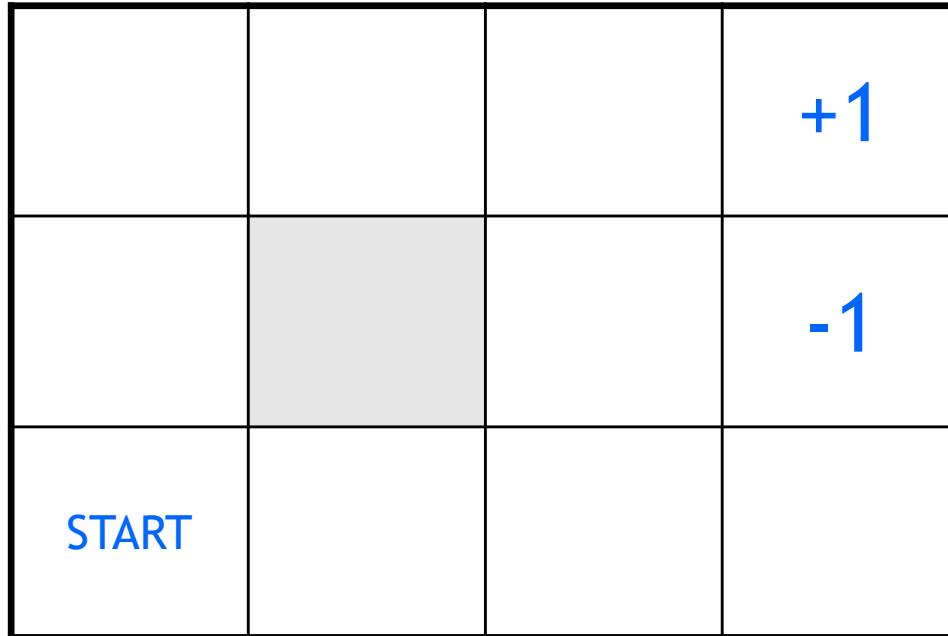
- pole-balancing
- TD-Gammon [Gerry Tesauro]
- helicopter [Andrew Ng]

General challenge: no teacher who would say “good” or “bad”

- is reward “10” good or bad?
- rewards could be delayed
- similar to control theory
 - more general, fewer constraints
- ***explore the environment and learn from experience***
 - not just blind search, try to be smart about it



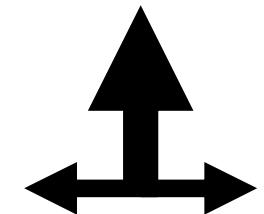
Robot in a room



actions: UP, DOWN, LEFT, RIGHT

UP

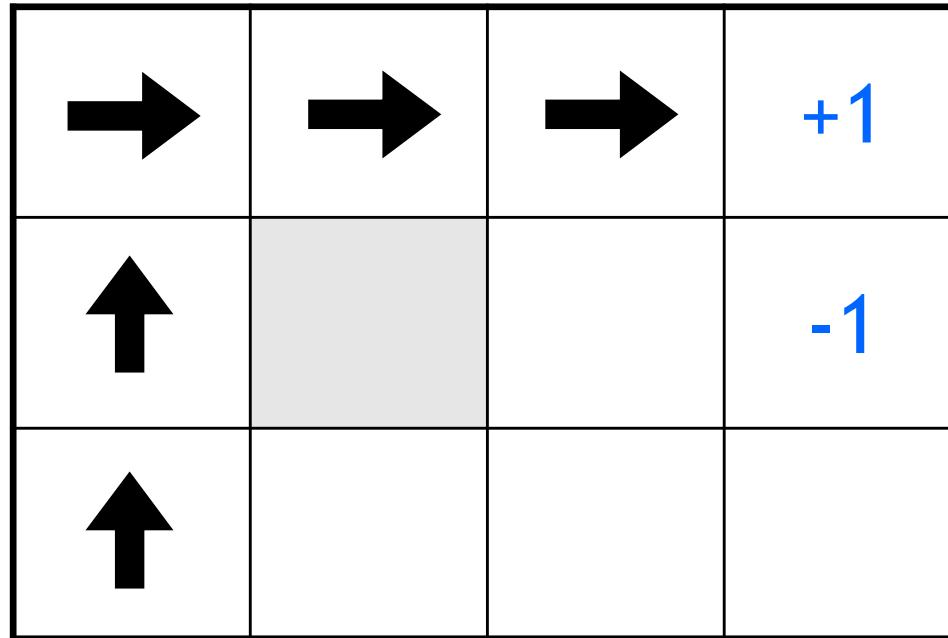
- | | |
|-----|------------|
| 80% | move UP |
| 10% | move LEFT |
| 10% | move RIGHT |



reward +1 at [4,3], -1 at [4,2]
reward -0.04 for each step

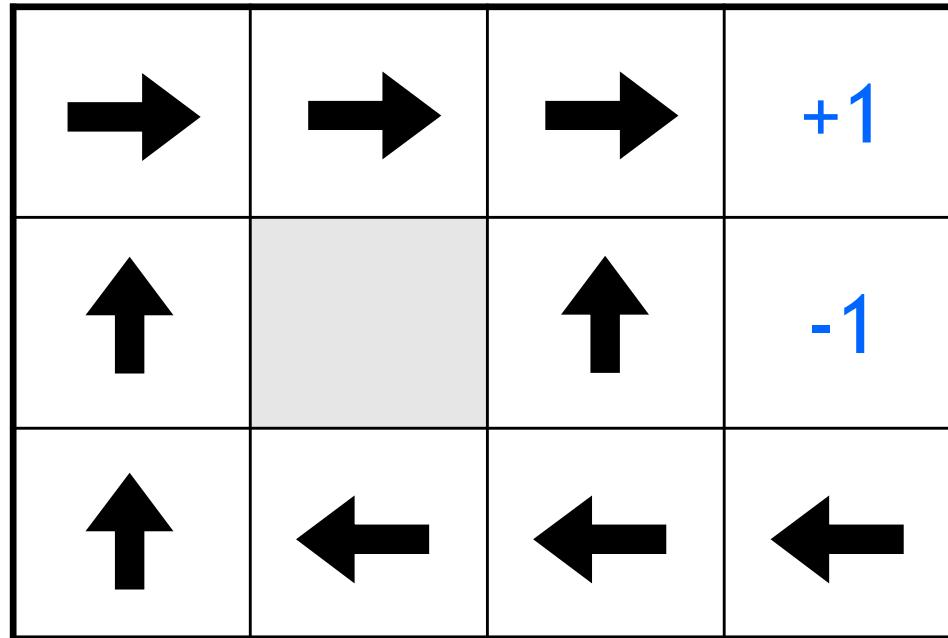
- states
- actions
- rewards
- What is the **solution**? What does the agent **learn**?

Is This A Solution?



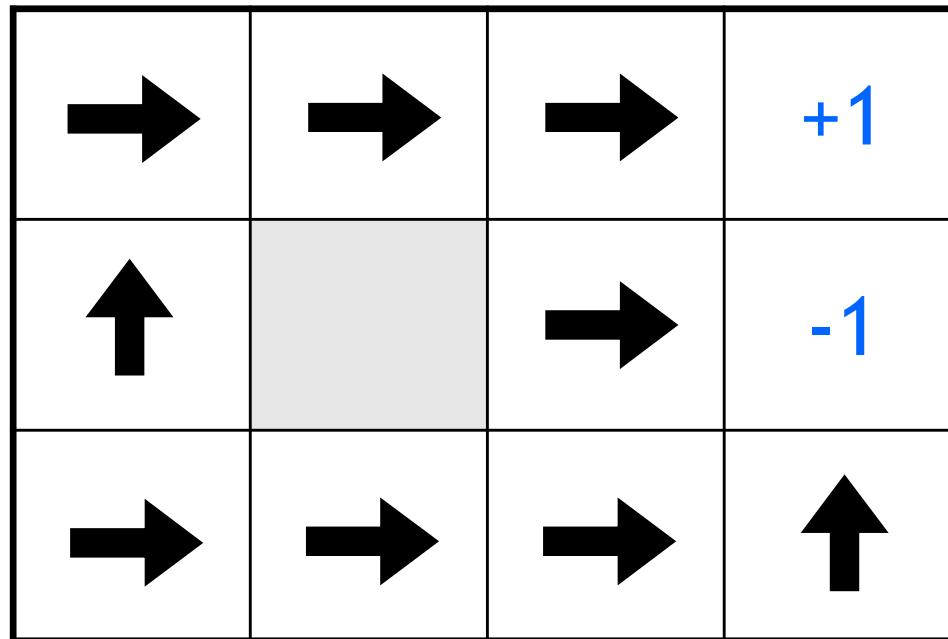
- only if actions **deterministic**
 - this path is a **plan**
 - not guaranteed to work as actions are **stochastic** (actions have probabilistic effects)
- we need a **policy**
 - mapping from each state to an action
 - agent tries to learn an **optimal policy**; but optimal in terms of what?

Optimal Policy

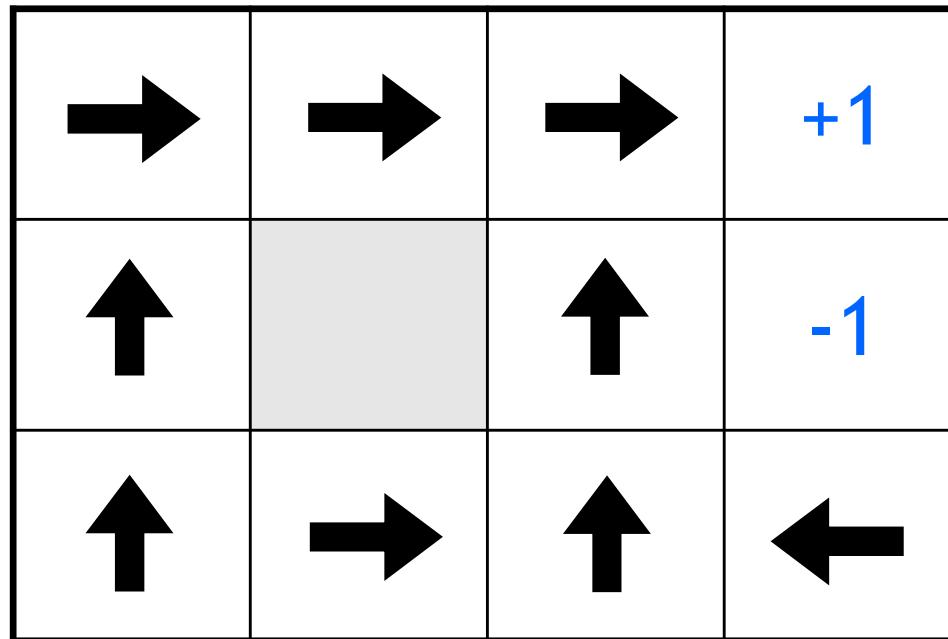


The **optimal policy** will change with the kind of **rewards** the agent receives at each episode!

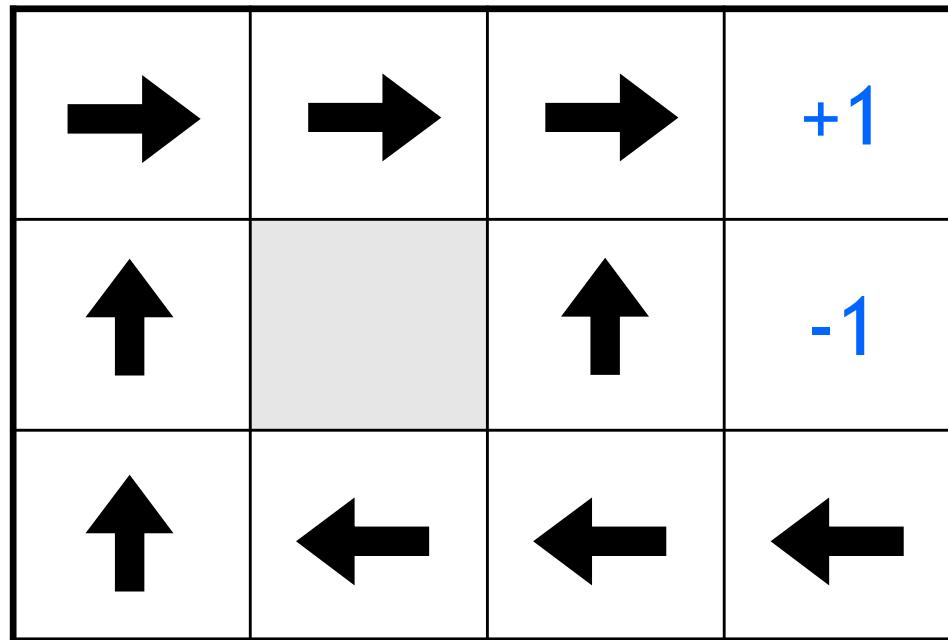
Reward for Each Step: -2.0



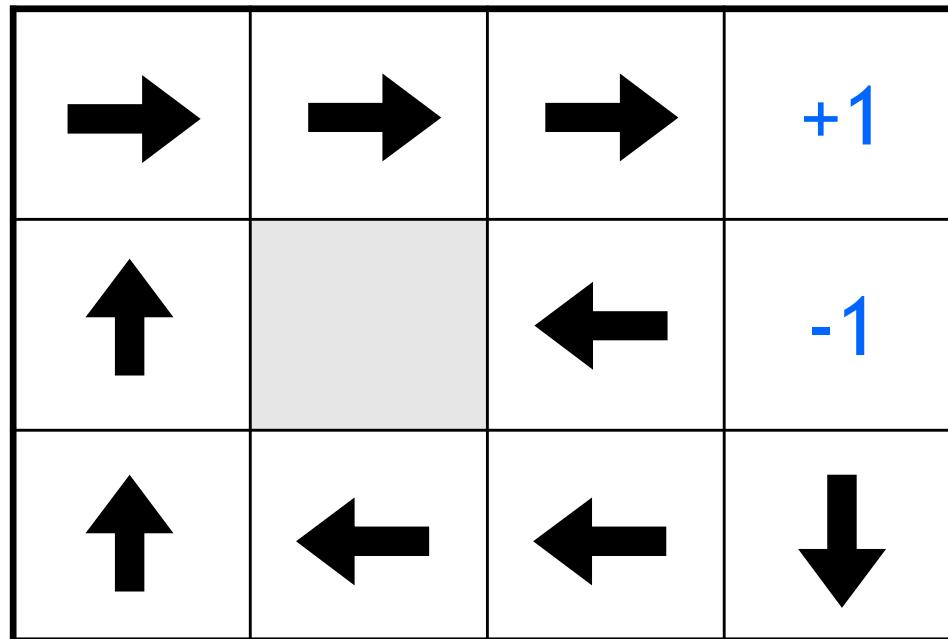
Reward for Each Step: -0.1



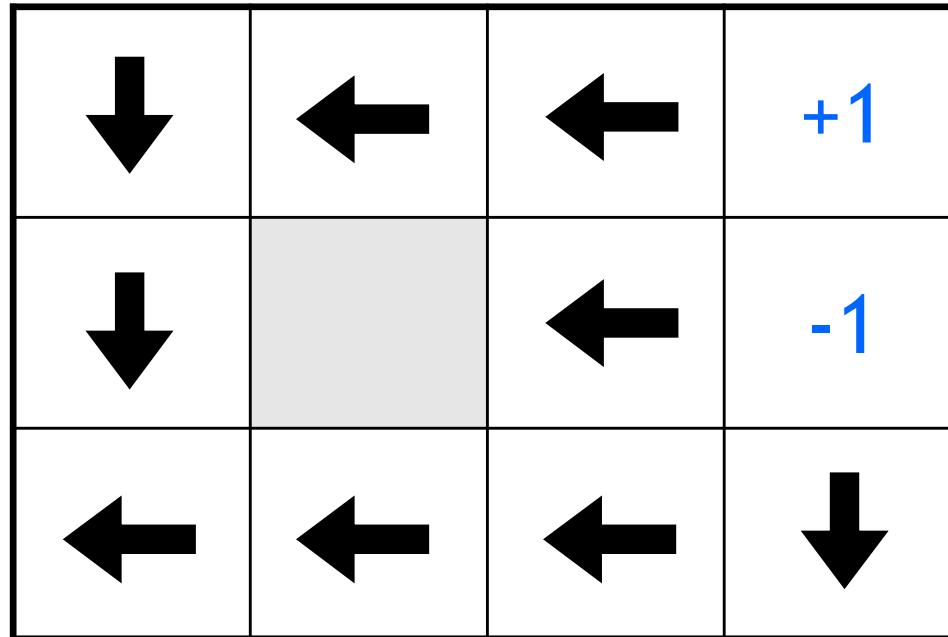
Reward for Each Step: -0.04



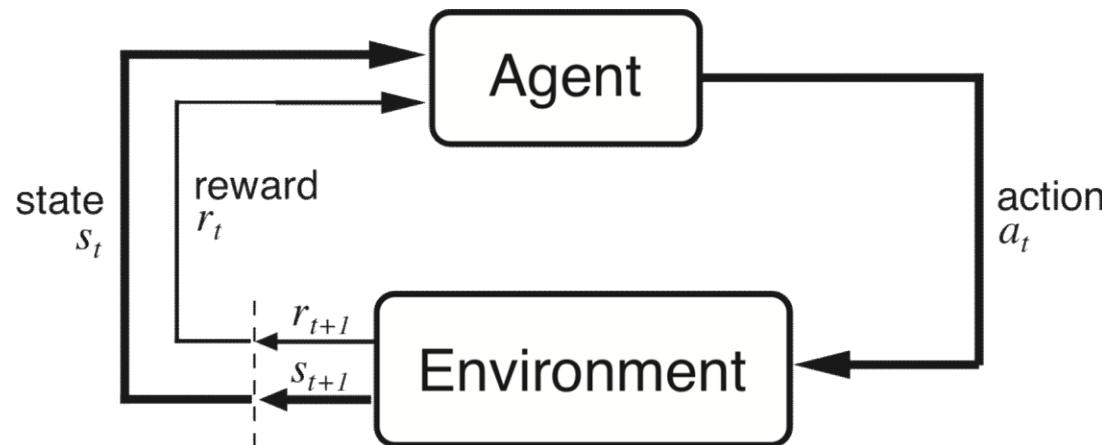
Reward for Each Step: -0.01



Reward for Each Step: +0.01



Formalizing RL: The Agent-Environment Interface



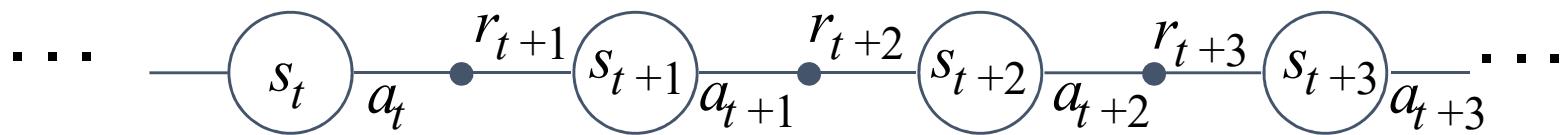
Agent and environment interact at discrete time steps: $t = 0, 1, 2, K$

Agent observes state at step t : $s_t \in S$

produces action at step t : $a_t \in A(s_t)$

gets resulting reward: $r_{t+1} \in \mathcal{R}$

and resulting next state: s_{t+1}

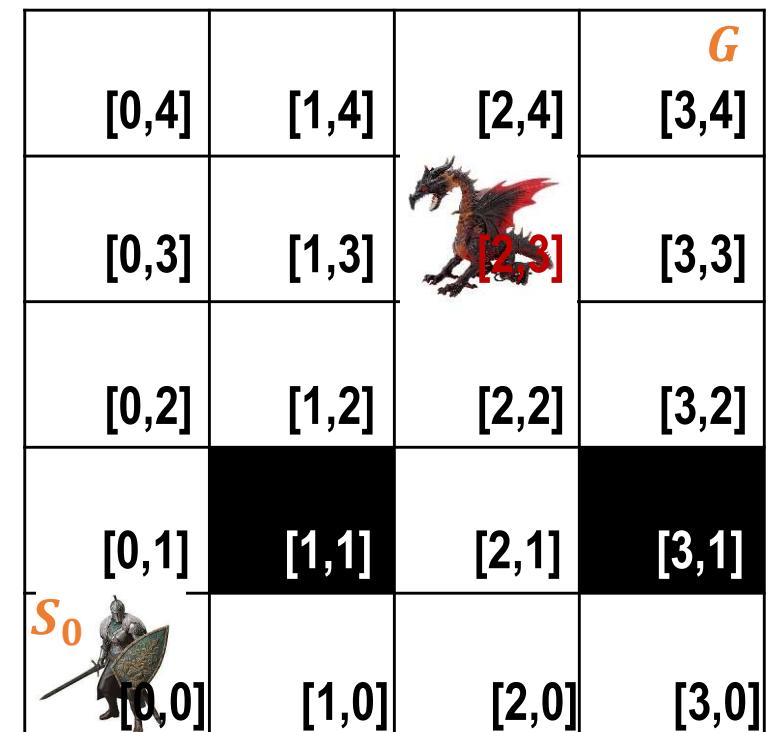


Formalizing RL: Markov Decision Processes

- set of **states** S , set of **actions** A , **initial state** S_0
 - for grid world, can be cell coordinates
- **transition model** $P(s, a, s')$
 - $P([1,1], \uparrow, [1,2]) = 0.8$
- **reward function** $r(s)$
 - $r([3,4]) = +1$
- **goal**: maximize cumulative reward in the long run
- **policy**: mapping from S to A
 - $\pi(s)$ or $\pi(s, a)$ (deterministic vs. stochastic)

Reinforcement Learning

- transitions and rewards usually not available
- how to change the policy based on experience
- how to explore the environment

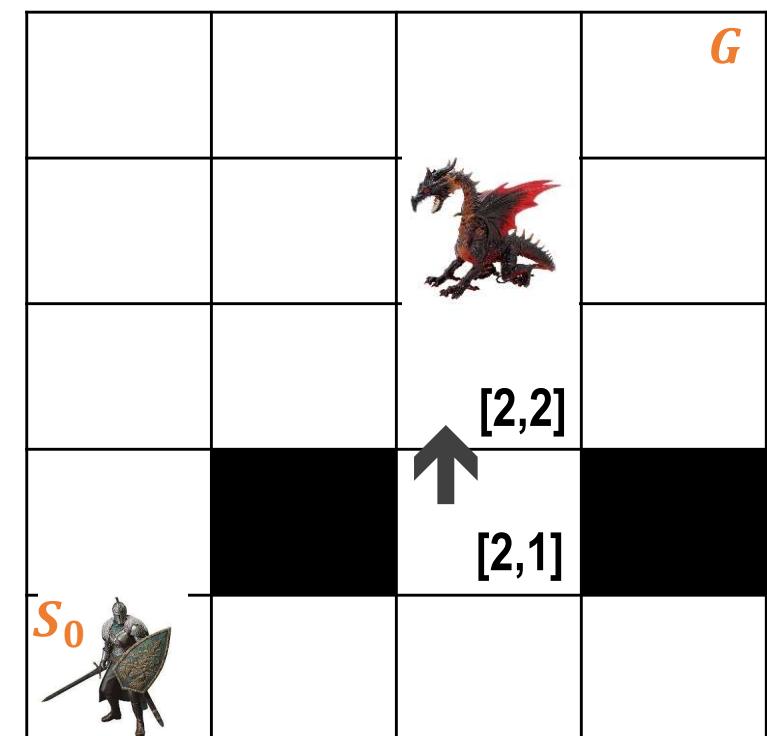


Actions Transition Probabilities

→(right)	→(60%), ↓(40%)
↑(up)	↑(100%)
←(left)	←(100%)
↓(down)	↓(70%), ←(30%)

Formalizing RL: The Markov Property

- “**the state**” at step t , means whatever information is available to the agent at step t about its environment – **snapshot of the world**
- the state can include immediate “**sensations**”, highly processed observations, and structures built up over time from sequences of observations
- ideally, a state should summarize past sensations so as to retain all “**essential**” information, i.e., it should have the **Markov Property**:
 - conditional probability distribution of **future states** depends only upon the **present state**, not on the sequence of events that preceded it*



$$\Pr \left\{ s_{t+1} = s', r_{t+1} = r \mid s_t, a_t, r_t, s_{t-1}, a_{t-1}, \mathcal{H}, r_1, s_0, a_0 \right\} =$$

$$\Pr \left\{ s_{t+1} = s', r_{t+1} = r \mid s_t, a_t \right\}$$

for all s' , r , and histories $s_t, a_t, r_t, s_{t-1}, a_{t-1}, \mathcal{H}, r_1, s_0, a_0$.

Formalizing RL: Rewards and Policy

- **episodic tasks:** interaction breaks naturally into episodes, e.g., plays of a game, trips through a maze
- **non-episodic tasks:** no episodes; infinite game. e.g. a self-driving car
- **additive rewards**
 - $R = r(s_0) + r(s_1) + r(s_2) + \dots$
 - infinite value for continuing tasks
- **discounted rewards**
 - rewards are discounted by a discount factor $\gamma \in [0, 1)$
 - $R = r(s_0) + \gamma * r(s_1) + \gamma^2 * r(s_2) + \dots$

Learning Problem: Find a policy that maximizes the **total expected reward**,
 $E_{\pi} [\sum_{t=1}^{\infty} \gamma^t r_t]$

A policy $\pi(s)$ or $\pi(s, a)$ is the prescription by which the agent selects an action to perform

- **Deterministic:** the agent observes the state of the system and chooses an action
- **Stochastic:** the agent observes the state of the system and then selects an action, at random, from some probability distribution over possible actions

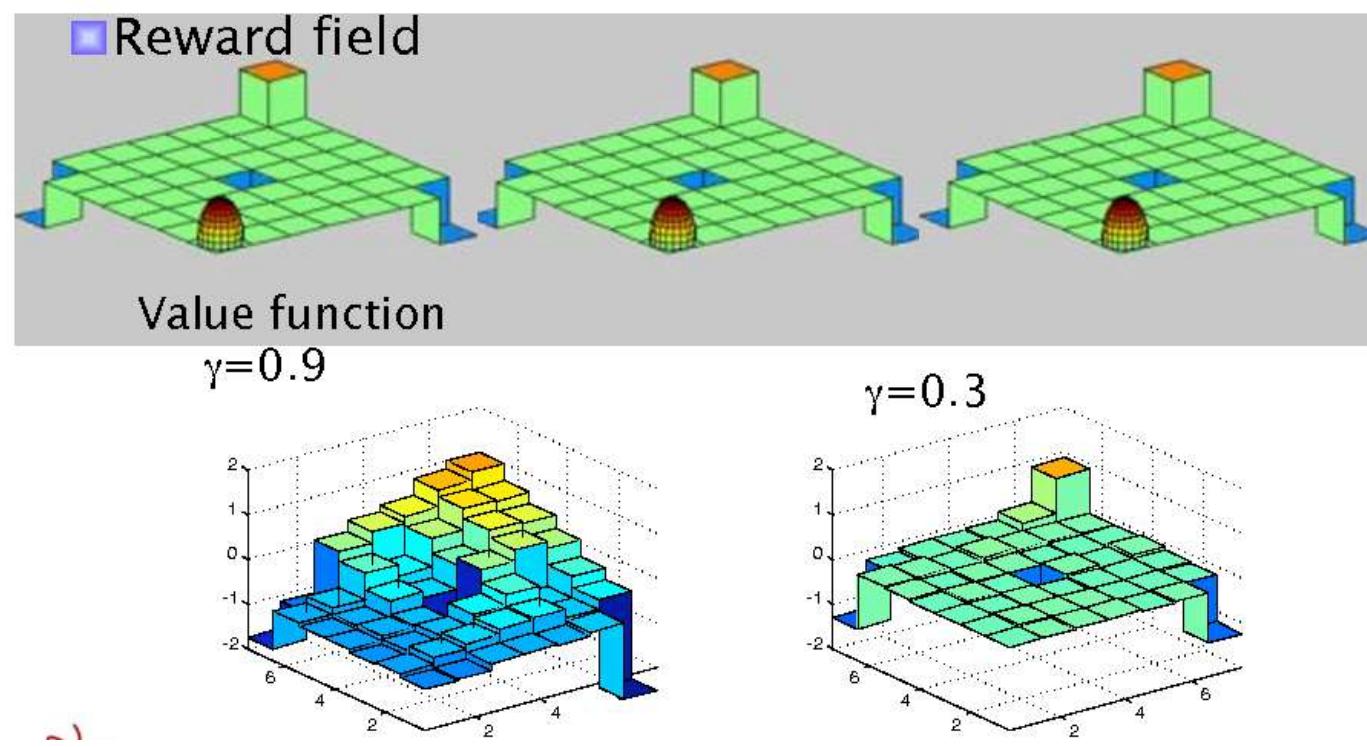
The (State) Value Function

A **value function** assigns a real number to **each state** called its **value**. The **value of a state (s)** is the expected reward *starting from that state (s) and then following the policy π* .



$$v^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_{t+1}, a_{t+1}) \mid s_0 = s \right]$$

The value function helps us evaluate the quality of a policy π . Informally, the value of a state indicates how much better it is to be in that state than other states, when following π .

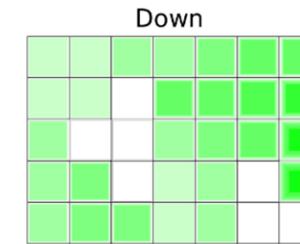
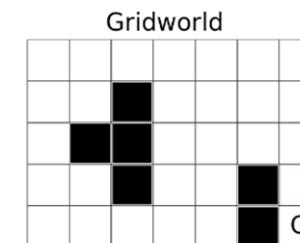
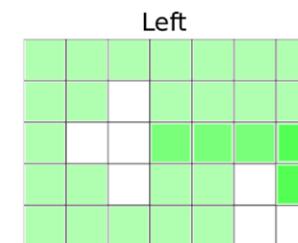
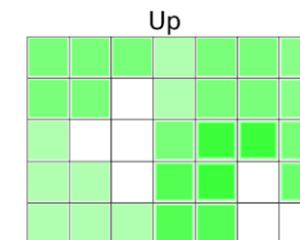


The (State-Action) Value Function

The **state-action value function** assigns a real number to **each state-action pair** called its **q-value**. The **q-value of a state** is the expected reward *starting from that state (s), executing that action (a) and then following the policy π* .

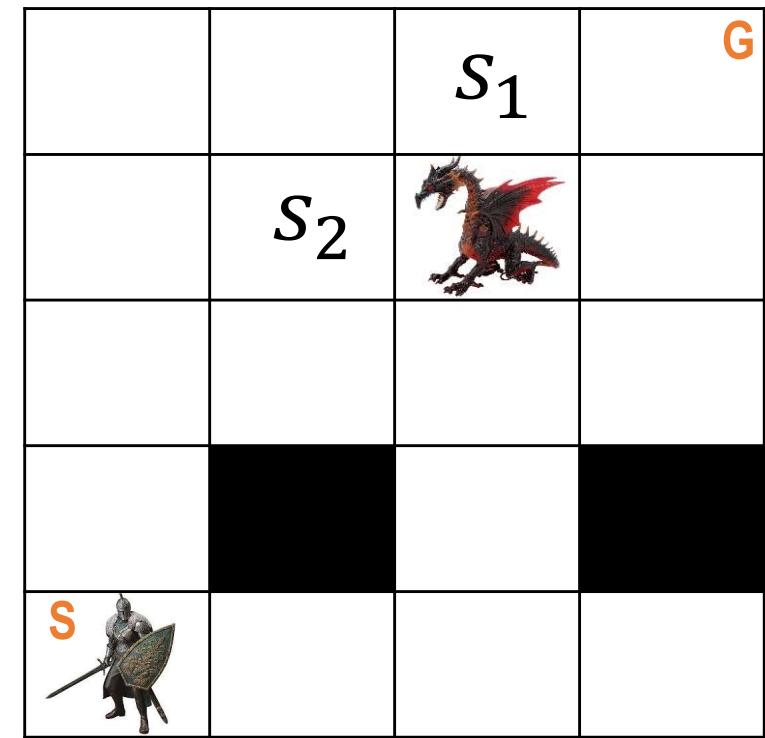


$$q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_{t+1}, a_{t+1}) \mid s_0 = s, a_0 = a \right]$$



Value Functions

- Which states should have a higher value? s_1 or s_2 ?
- Which action should have a higher value in s_1 ? \rightarrow or \downarrow ?
- Which action should have a higher value in s_2 ? \rightarrow or \uparrow ?



Value Functions

Unroll the discounted reward:

$$\begin{aligned} R_t &= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots \\ &= r_t + \gamma(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots) \\ &= r_t + \gamma R_{t+1} \text{ (recurrence)} \end{aligned}$$

Recall the definition of the value function:

$$\begin{aligned} V_\pi(s) &= E_\pi[\mathbb{I}[R_t | s_t = s]] = E_\pi[R_t + \gamma V_\pi(s_{t+1}) | s_t = s] \\ &\text{(another recurrence relation)} \end{aligned}$$

Unrolling the expectation using transition probabilities:

$$V_\pi(s) = R(s) + \gamma \sum_{s'} P(s'|s, a) V_\pi(s')$$

immediate reward

value of a state is the expected sum of discounted rewards when starting from that state

expected sum of discounted rewards after the first step from s taking into account all possible next states s' from all possible next actions $a \in A(s)$

This is one of the **Bellman equations**.

This equation can also be written as:

$$V_\pi(s) = R(s) + \gamma \sum_{a \in A(s)} \pi(s, a) \sum_{s'} P(s'|s, a) V_\pi(s')$$

agent: each action is chosen with (policy) probability $\pi(s, a)$

a

r

environment: next state is obtained with (transition) probability $P(s'|s, a)$

s'

s

a

r

s'

s

a

Value Functions

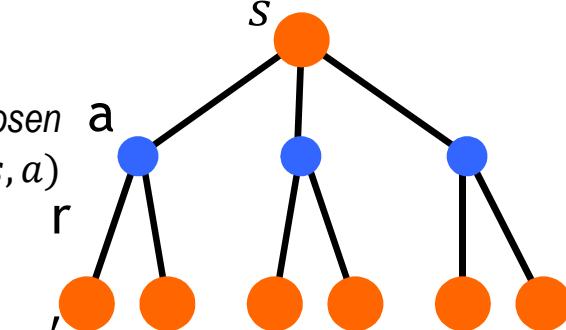
Recall the definition of the Q-value function:

$$\begin{aligned} Q_{\pi}(s, a) &= E_{\pi}[\mathbb{E}[R_t | s_t = s, a_t = a]] \\ &= E_{\pi}[R_t + \gamma V_{\pi}(s_{t+1}) | s_t = s, a_t = a] \end{aligned}$$

agent: each action is chosen with (policy) probability $\pi(s, a)$

a

environment: next state is obtained with (transition) probability $P(s' | s, a)$



Unrolling the expectation using transition probabilities:

$$Q_{\pi}(s, a) = R(s) + \gamma \sum_{s'} P(s' | s, a) V_{\pi}(s')$$

immediate reward
 ↓
 value of a state-action pair is the expected sum of discounted rewards when **starting from that state and executing that action**
 ↓
 expected sum of discounted rewards **after the first step from s taking into account all possible next states s' from executing action a**

This is another of the **Bellman equations**.

Compare with the state-value function, which considers all actions $a \in A(s)$:

$$V_{\pi}(s) = R(s) + \gamma \sum_{a \in A(s)} \pi(s, a) \sum_{s'} P(s' | s, a) V_{\pi}(s')$$

Optimal Value Functions

V^π defines a **partial ordering on policies**, that is, **value functions** are useful for finding the **optimal policy**.

Learning problem: find a policy $\pi^*: S \rightarrow A$ such that

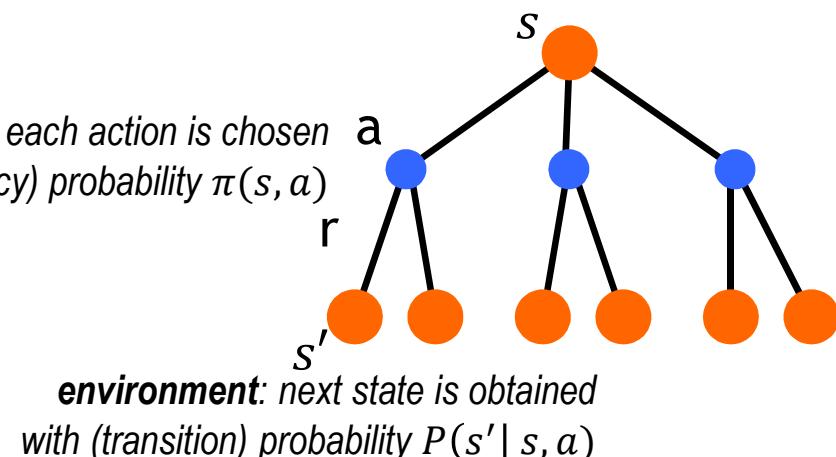
$$V^{\pi^*}(s) \geq V^\pi(s)$$

for all $s \in S$ and all policies π .

- any policy satisfying this condition is called an **optimal policy** (*may not be unique*)
- there **always** exists an optimal policy
- optimal policies share the same optimal value function

$$V^*(s) = \max_{\pi} V^\pi(s)$$

agent: each action is chosen with (policy) probability $\pi(s, a)$



environment: next state is obtained with (transition) probability $P(s' | s, a)$

Bellman Optimality Equations

$$V^*(s) = \max_{a \in A(s)} R(s) + \gamma \sum_{s'} P(s'|s, a) V^*(s')$$

$$Q^*(s, a) = \max_{a' \in A(s)} R(s) + \gamma \sum_{s'} P(s'|s, a) Q^*(s', a')$$

- system of n ($=$ number of states) non-linear equations describing a recurrence relation between current and next states
- for a finite-state MDP, we obtain a system of linear equations
- solve for $V^*(s)$
- easy to extract the optimal policy

The **optimal policy** π^* that satisfies these equations is **the optimal policy for all states s** . That is, it does not matter if we start in a state s or a different state s' , that is, we can use the same policy π^* no matter the initial state of our MDP.

Solving the MDP: Policy Iteration

Policy iteration: iteratively perform **policy evaluation** + **policy improvement**, which are repeated iteratively until policy converges

1. Initialization

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

► initialize to a random policy

2. Policy Evaluation

Repeat

$$\Delta \leftarrow 0$$

For each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s',r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

► find the value function corresponding to this policy using iterative policy evaluation

► can also be done by solving a system of equations

$$V^\pi(s) = R(s) + P_{s\pi(s)}V(s)$$

until $\Delta < \theta$ (a small positive number)

► improve the policy based on the new values

3. Policy Improvement

policy-stable \leftarrow true

For each $s \in \mathcal{S}$:

$$\text{old-action} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$$

If $\text{old-action} \neq \pi(s)$, then *policy-stable* \leftarrow false

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Solving the MDP: Value Iteration

Value iteration: directly find optimal value function and extract the optimal policy from it

Initialize array V arbitrarily (e.g., $V(s) = 0$ for all $s \in \mathcal{S}^+$)
 ▶ initialize values to zero

Repeat

$$\Delta \leftarrow 0$$

For each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$ (a small positive number)

▶ effectively combines, in each of its sweeps, one sweep of policy evaluation and one sweep of policy improvement

Output a deterministic policy, $\pi \approx \pi_*$, such that

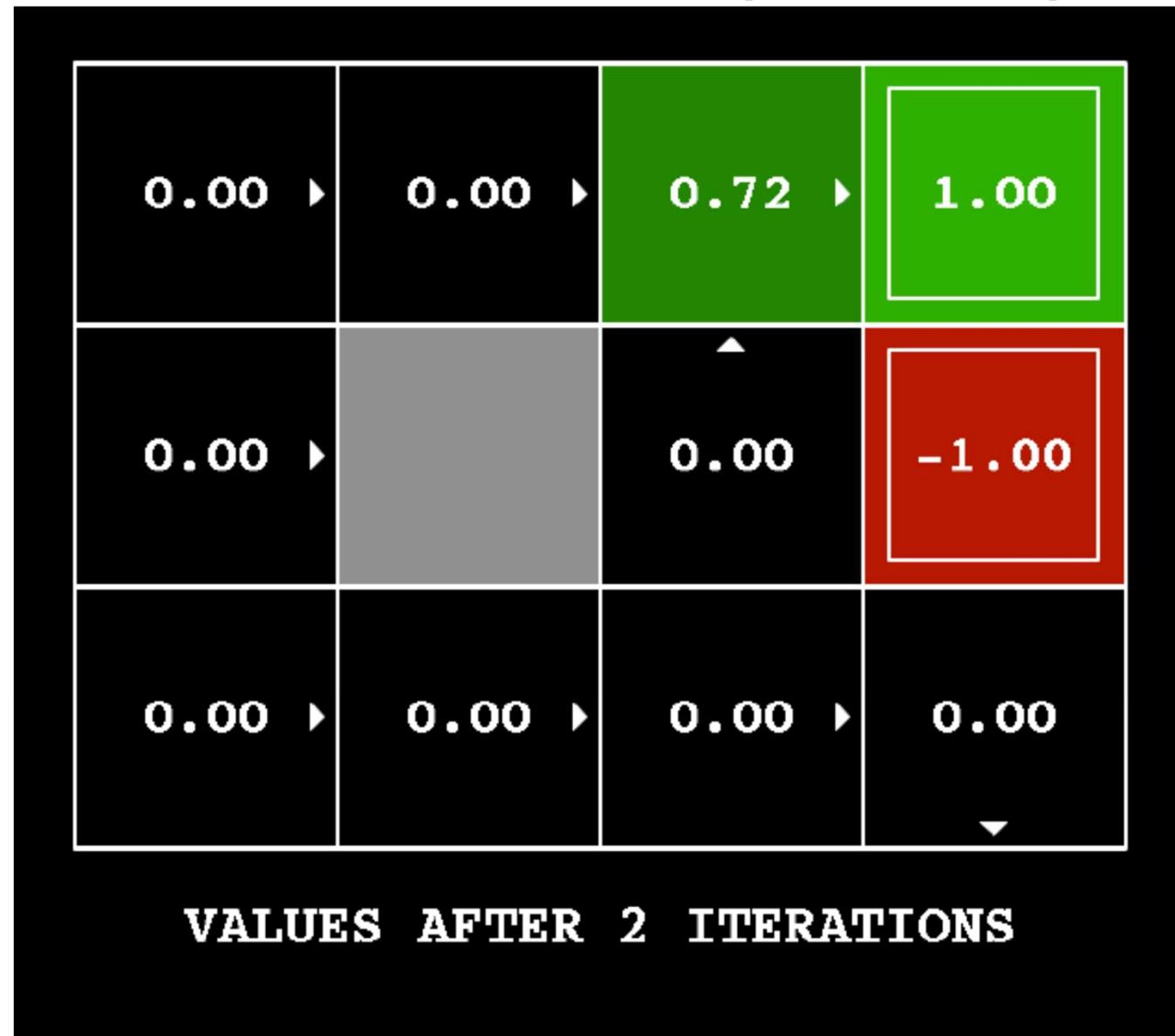
$$\pi(s) = \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

▶ one-time policy extraction

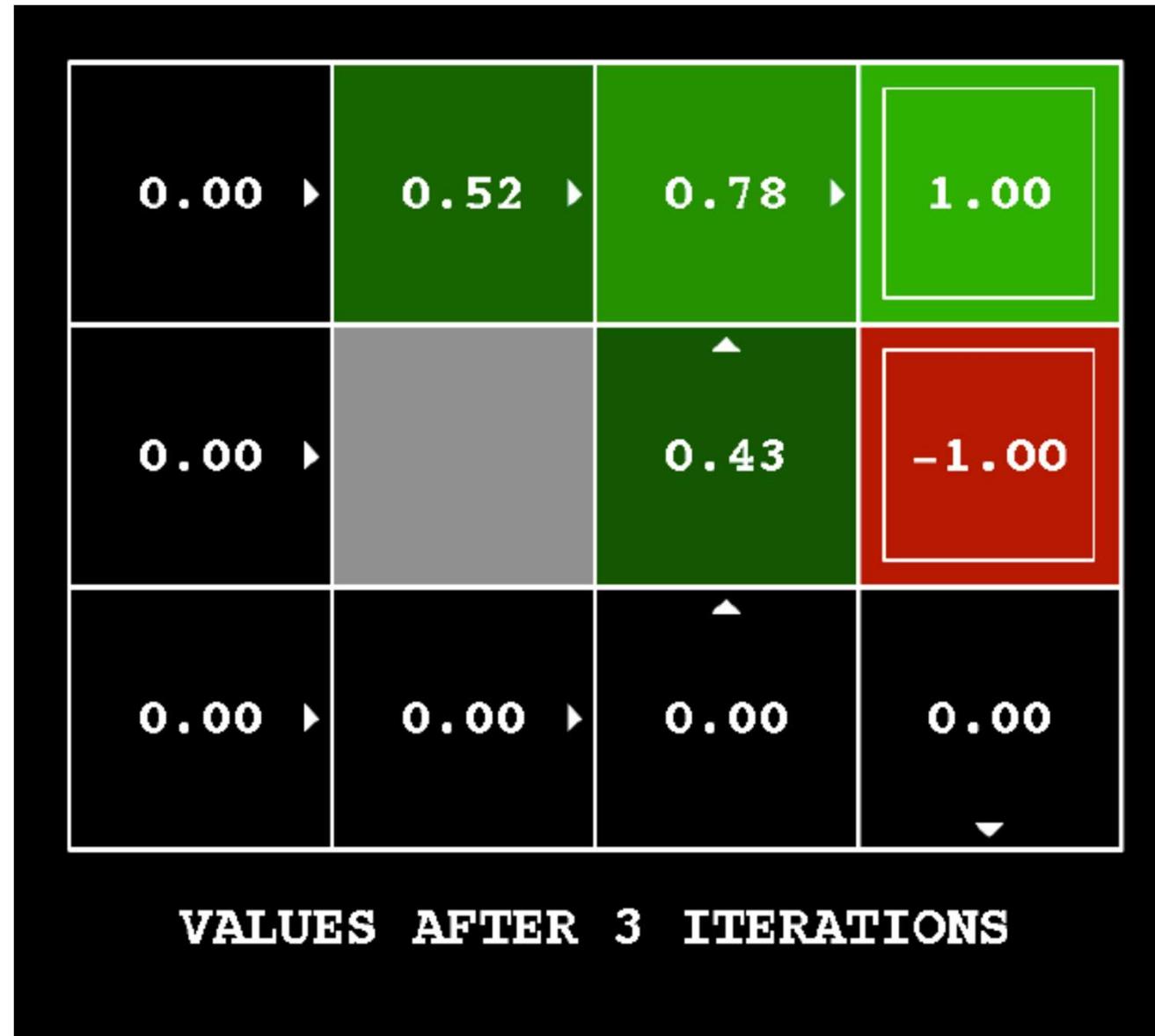
Value Iteration in GridWorld ($\gamma = 0.9$)



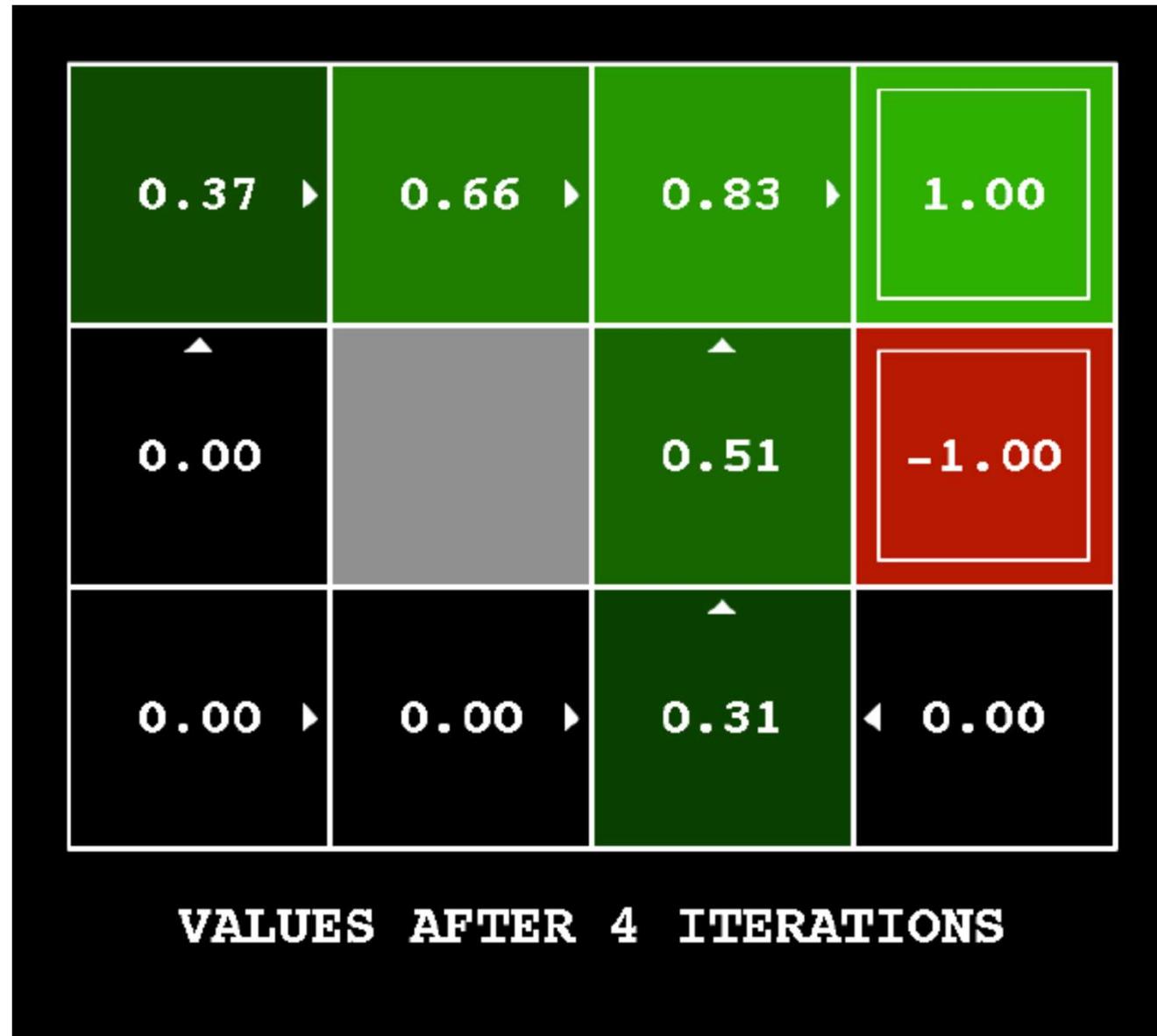
Value Iteration in GridWorld ($\gamma = 0.9$)



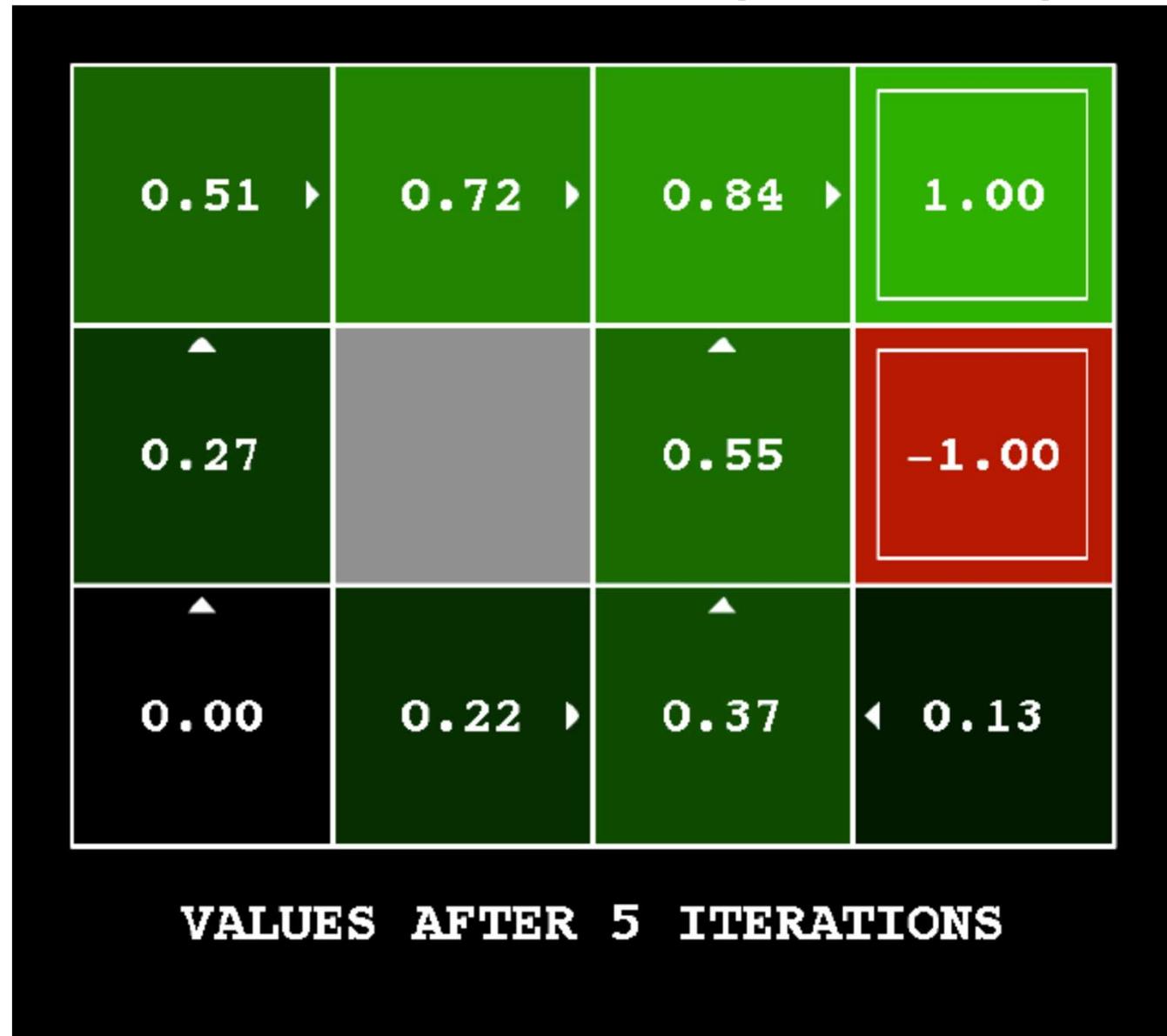
Value Iteration in GridWorld ($\gamma = 0.9$)



Value Iteration in GridWorld ($\gamma = 0.9$)



Value Iteration in GridWorld ($\gamma = 0.9$)



Value Iteration in GridWorld ($\gamma = 0.9$)



Q-Learning

Full reinforcement learning

- You don't know the **transitions** $P(s, a, s')$
- You don't know the **rewards** $R(s, a, s')$
- You can choose any actions you like
- **Goal:** learn the optimal policy / values
 - Learn the MDP first, then use value/policy iteration (requires learning the MDP: transition and reward functions)
 - **Learn only the values (don't learn the MDP or explicitly model it)**
 - Learner makes choices: **exploration vs. exploitation**
 - This is not offline planning; you **take actions in the world** and find out what happens!

Value iteration: find successive approximate optimal values

$$V_{i+1}(s) = \max_a \sum_{s'} P(s, a, s') [R(s) + \gamma V_i(s)]$$

Q-values are more useful!

$$Q_{i+1}(s, a) = \sum_{s'} P(s, a, s') \left[R(s) + \gamma \max_{a'} Q_i(s', a') \right]$$

Q-Learning

How should we pick an action to take based on Q ?

- Shouldn't always be greedy
 - we won't explore much of the state space this way
- Shouldn't always be random
 - will take a long time to generate a good Q
- **ϵ -greedy strategy**: with some small probability choose a random action (exploration), otherwise select the greedy action (exploitation)

Initialize:

- Choose an initial state-value function $Q(s, a)$
- Let s be the initial state of the environment

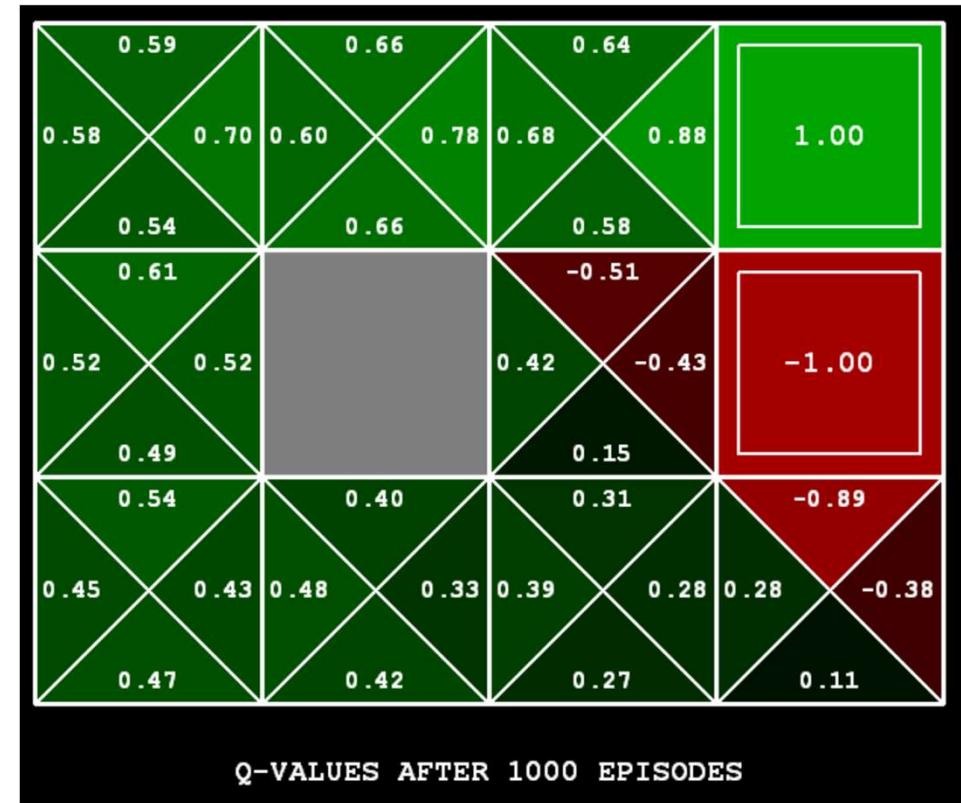
Repeat until convergence:

- Choose an action a_t for the current state s_t based on Q
- Take action a_t and observe the reward r_t and the new state s_{t+1}
- Update Q

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{(1 - \alpha) \cdot Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \overbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} \right)}^{\text{learned value}}$$


Q-Learning

Q-learning produces tables of q-values



Initialize:

- Choose an initial state-value function $Q(s, a)$
- Let s be the initial state of the environment

Repeat until convergence:

- Choose an action a_t for the current state s_t based on Q
- Take action a_t and observe the reward r_t and the new state s_{t+1}
- Update Q

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{(1 - \alpha) \cdot Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\substack{\text{learned value} \\ \text{estimate of optimal future value}}} \right)}_{\text{new value}}$$

Learning Rate (α) and Discount Factor (γ)

Explore vs exploit:

- learning rate (α) determines to what extent newly acquired information overrides old information
- $\alpha = 0$ makes the agent learn nothing (exclusively exploiting prior knowledge)
- $\alpha = 1$ makes the agent consider only the most recent information (ignoring prior knowledge to explore possibilities).
- in practice, often a constant learning rate is used

Discount factor:

- discount factor (γ) determines the importance of future rewards
- $\gamma = 0$ will make the agent "myopic" (or short-sighted) by only considering current rewards
- $\gamma = 1$ will make the agent strive for a long-term high reward
- if $\gamma > 1$, action values may diverge