# Approximate Dynamic Programming
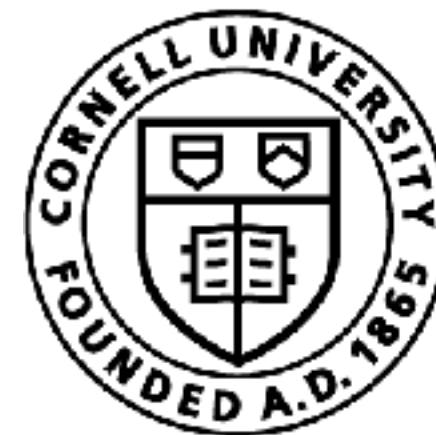
Sanjiban Choudhury
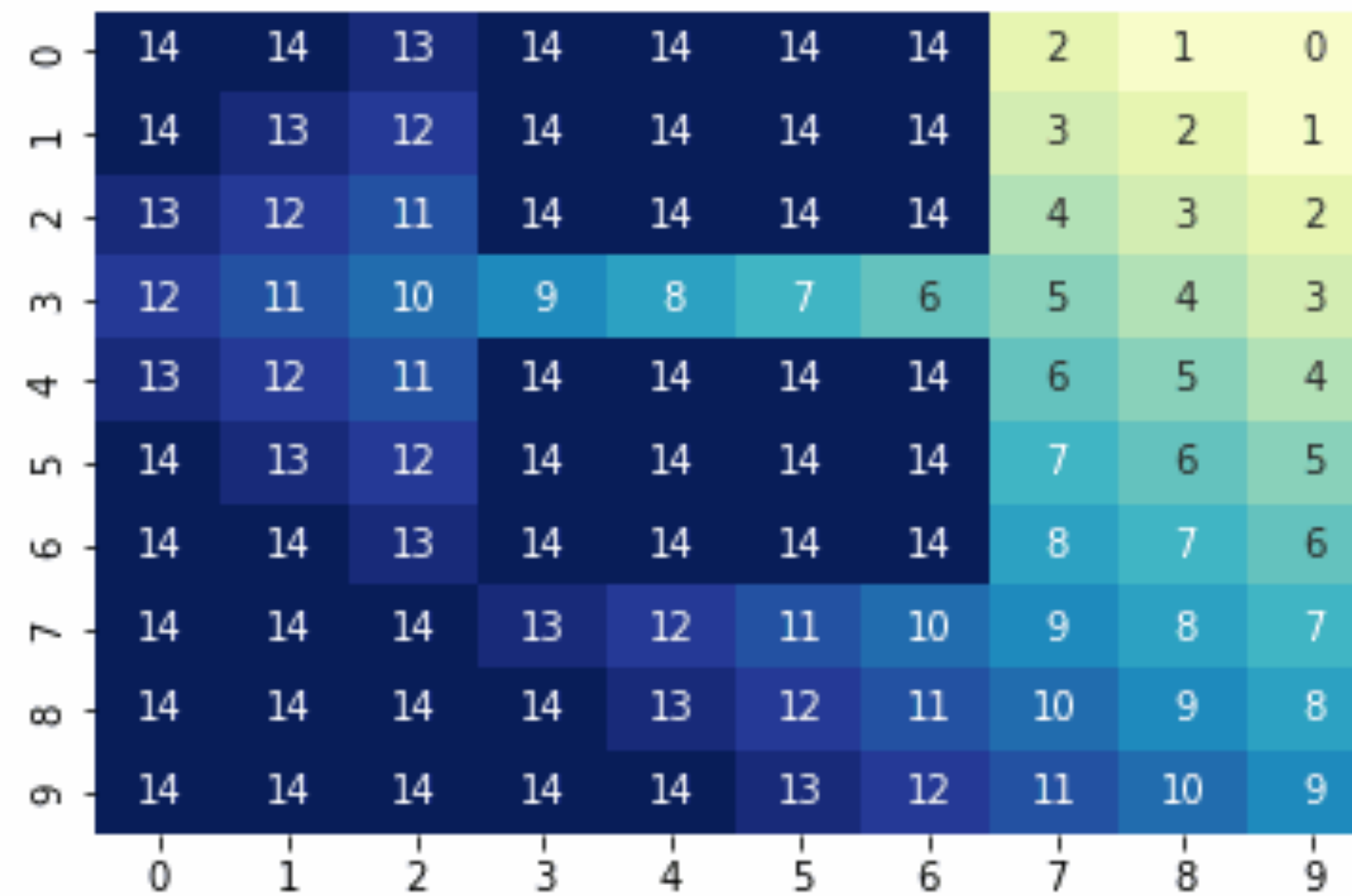
# When the MDP is known:

# Two Fundamental Ways to Solve for Optimal Policy

# Value Iteration



$$V^*(s) = \min_a [c(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s,a)} V^*(s')]$$

# Policy Iteration

# Which one converges faster: value/policy?



Values

Policy

Policy converges faster
than the value

Can we iterate over policies?

# Policy Iteration

Init with some policy $\pi$

Repeat forever

    Evaluate policy

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim \mathscr{T}(s,a)} V^\pi(s')]$$

    Improve policy

$$\pi^+(s) = \arg\min_a c(s, a) + \gamma \mathbb{E}_{s' \sim \mathscr{T}(s,a)} V^\pi(s')]$$

# Init with some policy $\pi$

# Iteration 1: Compute the value of the policy



Iter: 1

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s,a)} V^\pi(s')]$$

$$\pi^+(s) = \arg \min_a c(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s,a)} V^\pi(s')]$$

# Policy Iteration
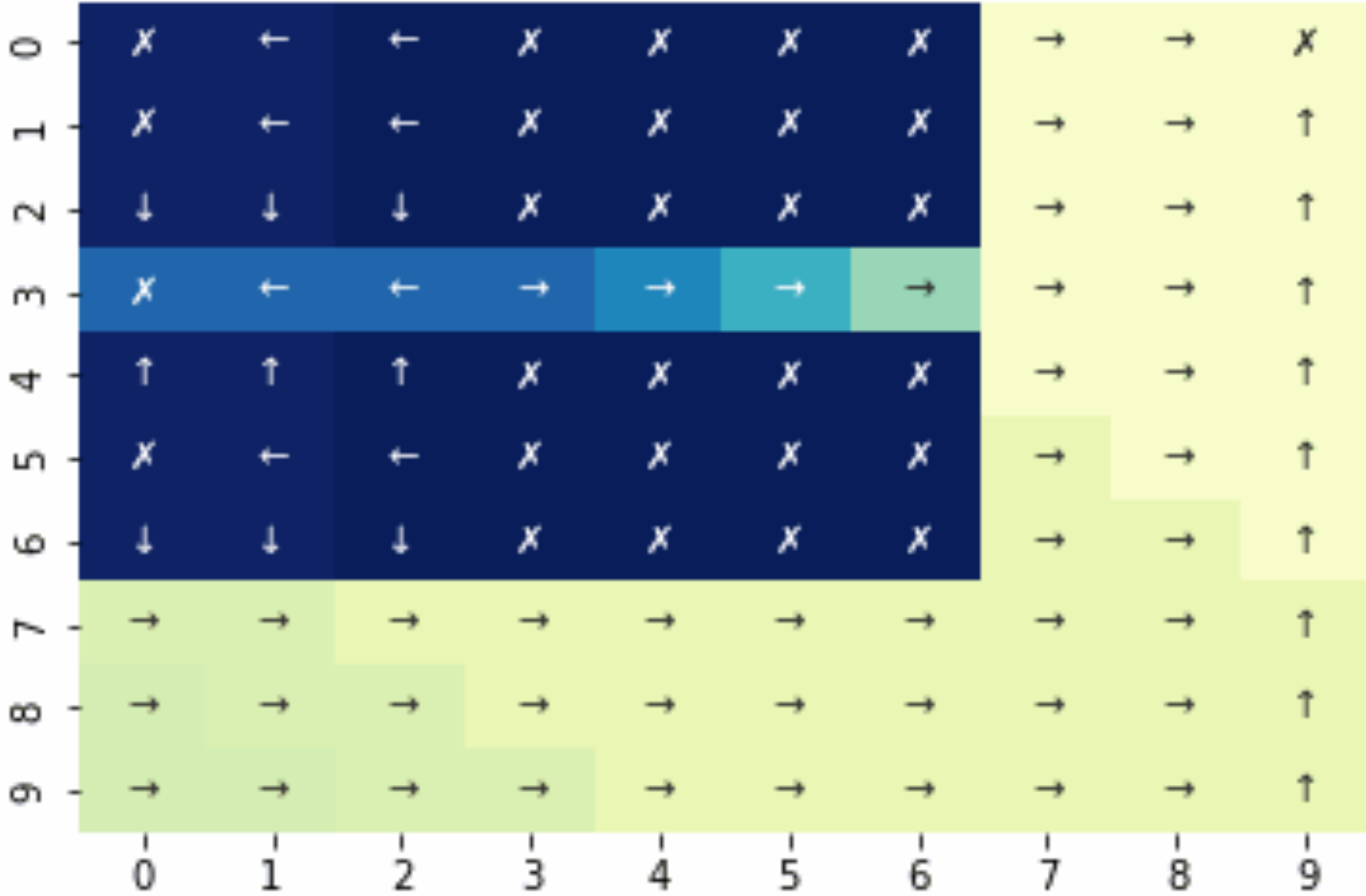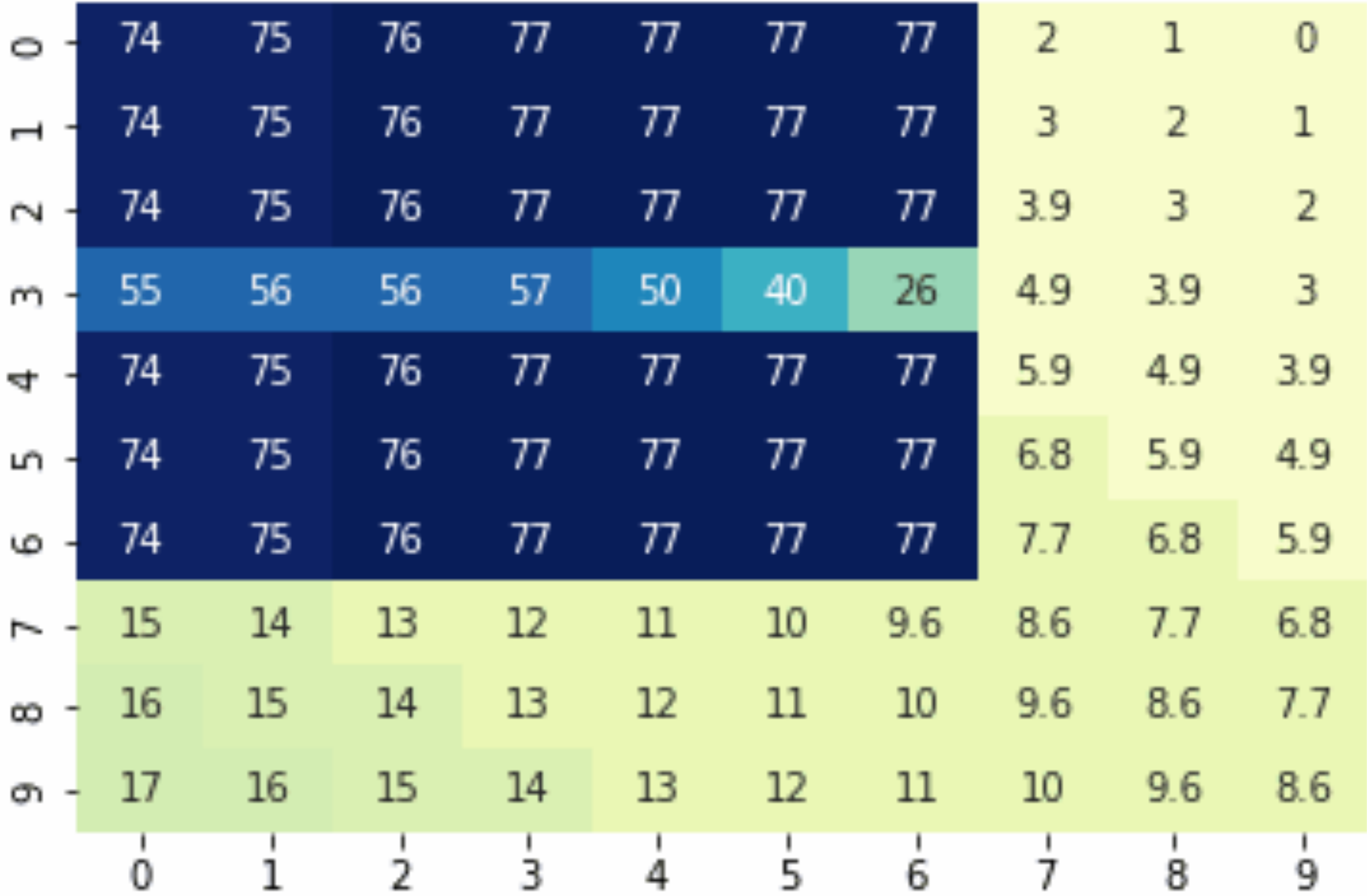


Iter: 0

$$V^\pi(s) = c(s, \pi(s)) + \gamma\mathbb{E}_{s' \sim \mathscr{T}(s,a)} V^\pi(s')]$$
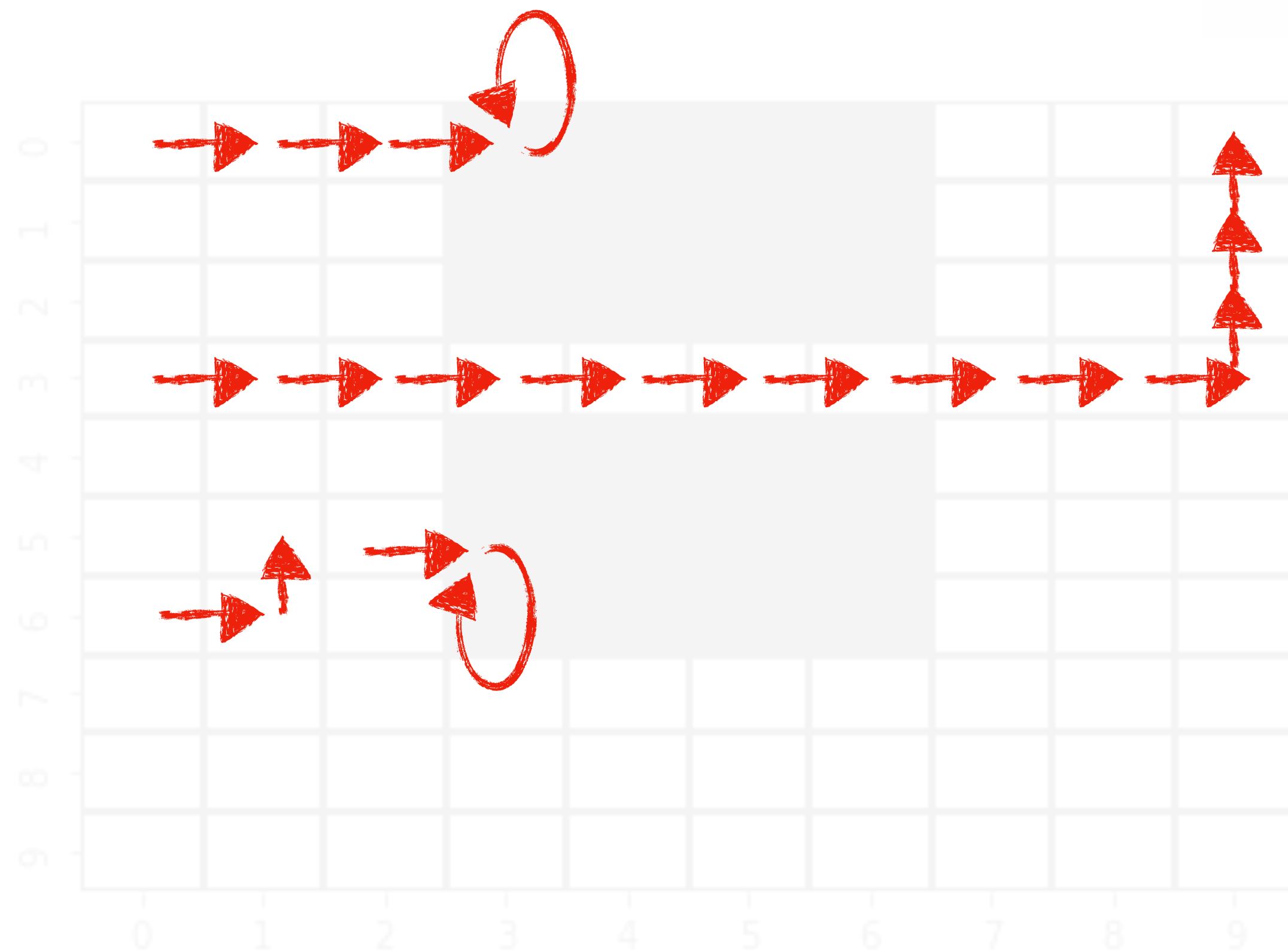
$$\pi^+(s) = \arg\min_a c(s, a) + \gamma\mathbb{E}_{s' \sim \mathscr{T}(s,a)} V^\pi(s')]$$

# When the MDP is ~~known~~ <span style="color:red">unknown:</span>
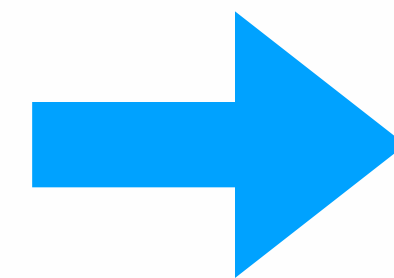
Restricted access to the transition function

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s,a)} V^\pi(s')]$$

# *Estimate the value* of policy from sample rollouts



Roll outs

Value $V^\pi(s)$

# Monte-Carlo

$$V(s) \leftarrow V(s) + \alpha(G_t - V(s))$$

## Zero Bias

## High Variance

## Always convergence
(Just have to wait till heat death of the universe)

# Temporal Difference

$$V(s) \leftarrow V(s) + \alpha(c + \gamma V(s') - V(s))$$

## Can have bias

## Low Variance

## May *not* converge if using function approximation
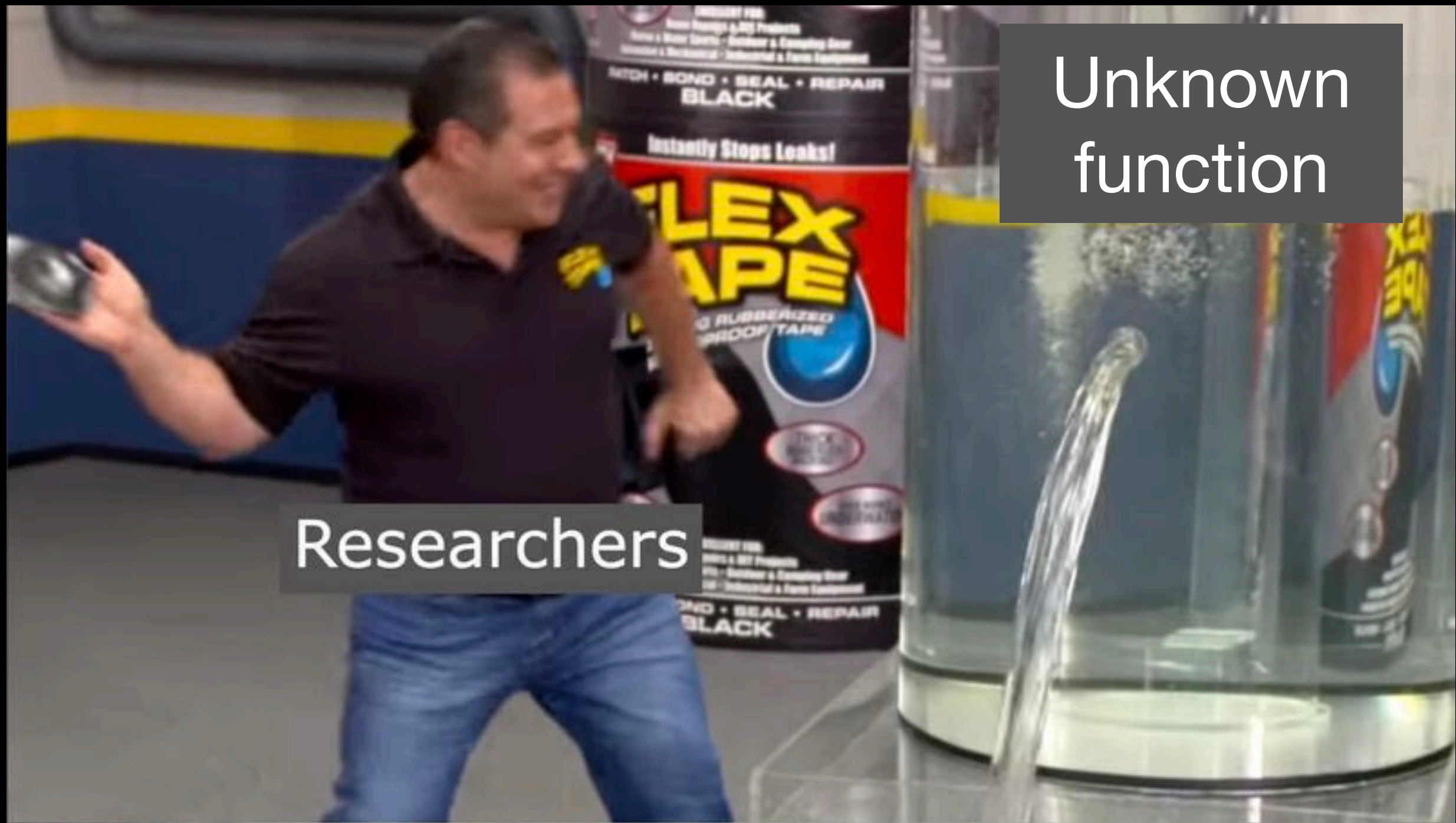
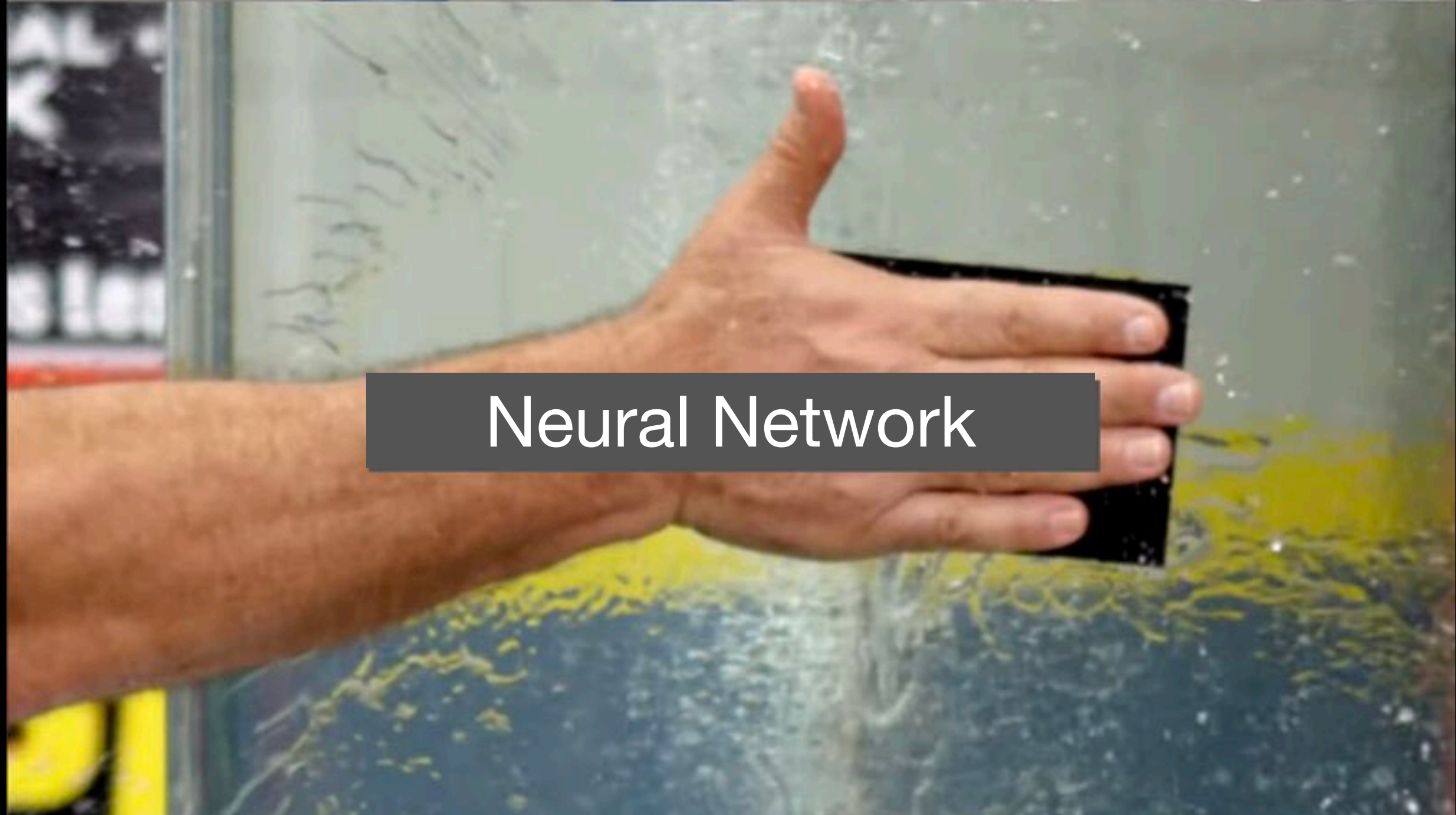Tabular setting is cute.

But how do we estimate V(s)
in the continuous setting

Unknown function

Researchers

Neural Network
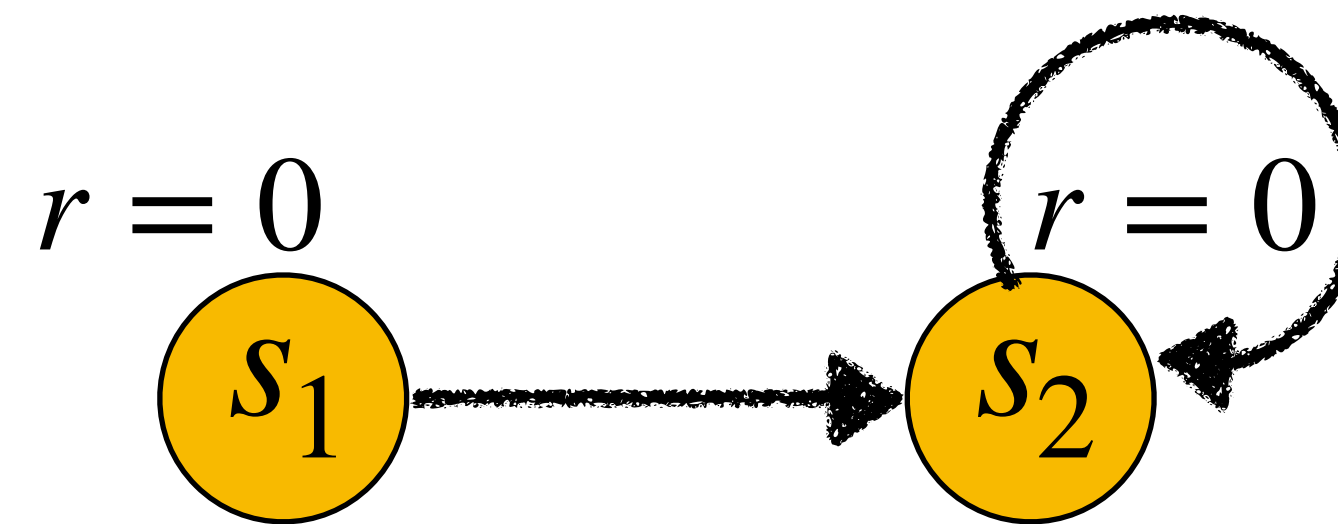
# Activity!

# A *tiny* MDP

Reward for being at
any state is 0.0

Discount factor $\gamma = 0.9$



$r = 0$          $r = 0$

$s_1$ → $s_2$

What happens when you run value iteration?

(Initialize with random values, say $V(s_1) = 1$ and 2)

# A *tiny* MDP

Reward for being at any state is 0.0

Discount factor $\gamma = 0.9$

$r = 0$      $r = 0$

$s_1 \longrightarrow s_2$

$f(s_1) = 1$      $f(s_2) = 2$

Let's say we want to use a *linear value function approximator*

$$V(s) = wf(s) = w * \begin{cases} 1 & \textit{if } s = s_1 \\ 2 & \textit{if } s = s_2 \end{cases}$$

What happens if you run value iteration? (Initialize with w=1)

# Think-Pair-Share

Think (30 sec): Initialize value iteration with w=1. What happens? What's the explanation?
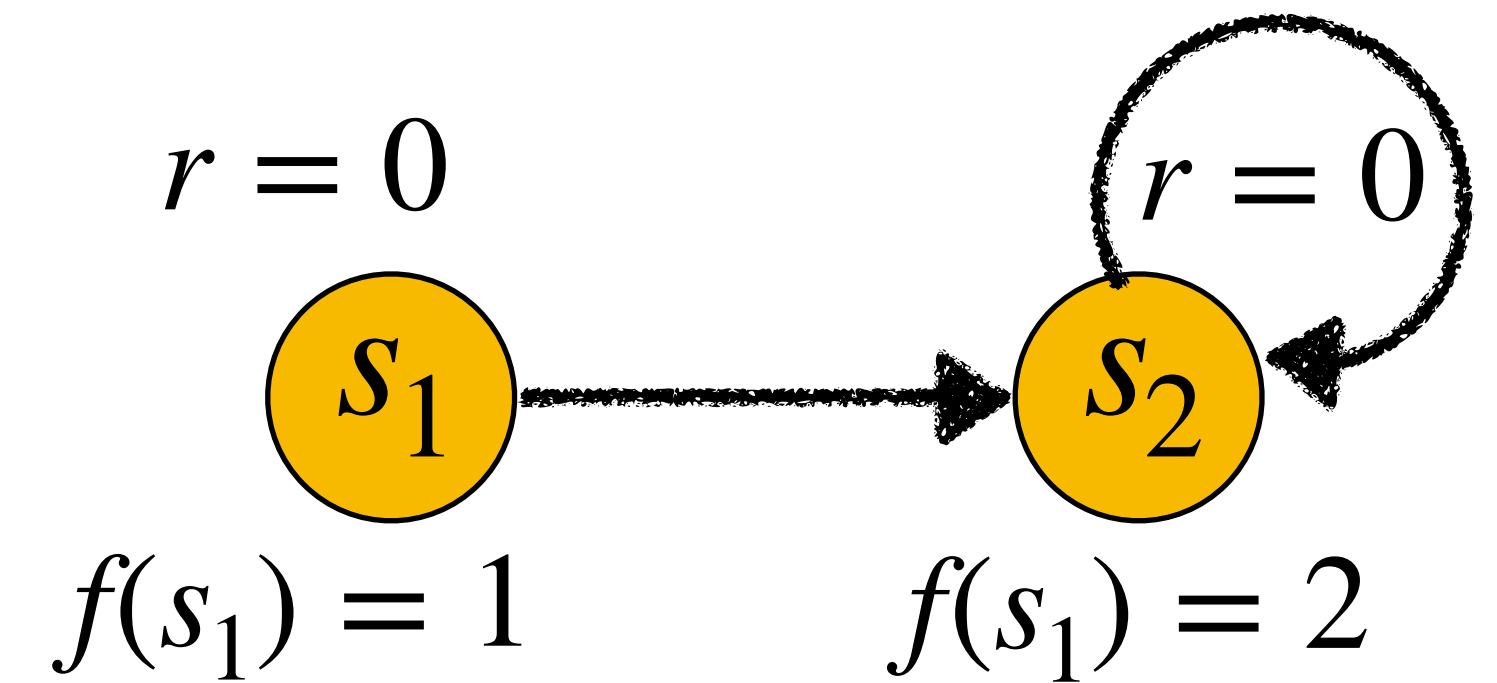
Pair: Find a partner
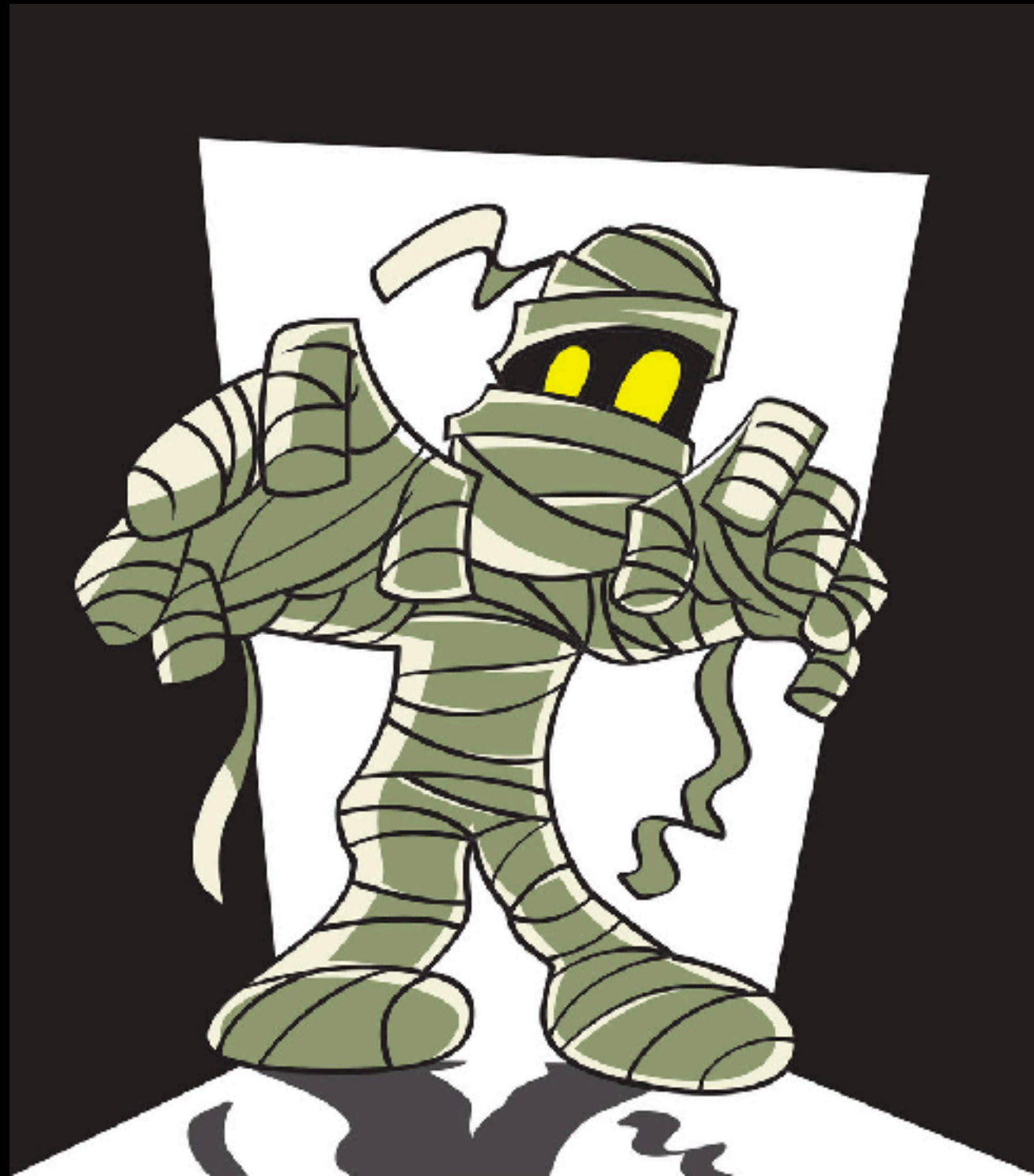
Share (45 sec): Partners exchange
ideas

$r = 0$        $r = 0$

$s_1 \longrightarrow s_2$

$f(s_1) = 1$      $f(s_1) = 2$

$V(s) = wf(s)$

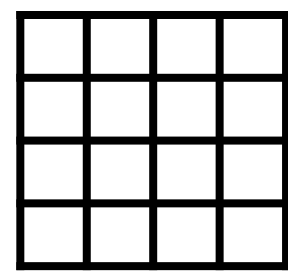Init with $w = 1$

From dynamic programming to *Fitted* dynamic programming

# Approximate (Fitted) Value Iteration

## Q-iteration



$Q(s, a) \leftarrow 0$
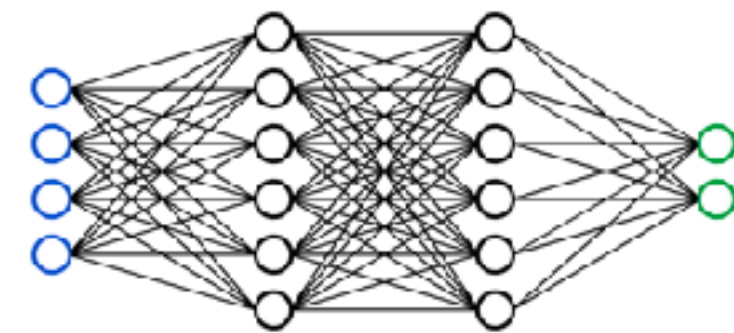
**while** *not converged* **do**

 **for** $s \in S, a \in A$

  $Q^{new}(s, a) = c(s, a) + \gamma \mathbb{E}_{s'} \min_{a'} Q(s', a')$

 $Q \leftarrow Q^{new}$

**return** $Q$

## *Fitted* Q-iteration



**Given** $\{s_i, a_i, c_i, s_i'\}_{i=1}^{N}$

**Init** $Q_\theta(s, a) \leftarrow 0$

**while** *not converged* **do**

 $D \leftarrow \varnothing$
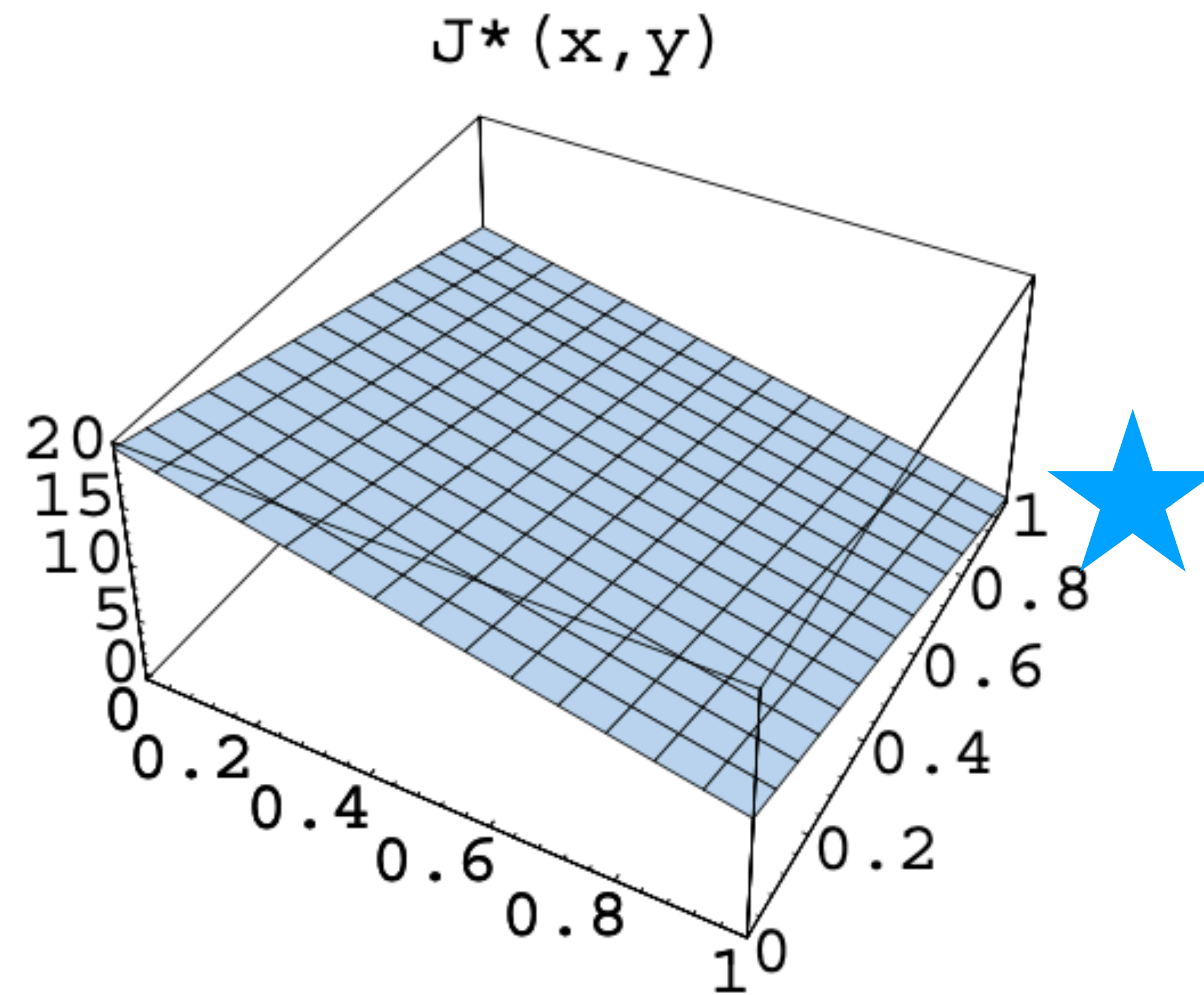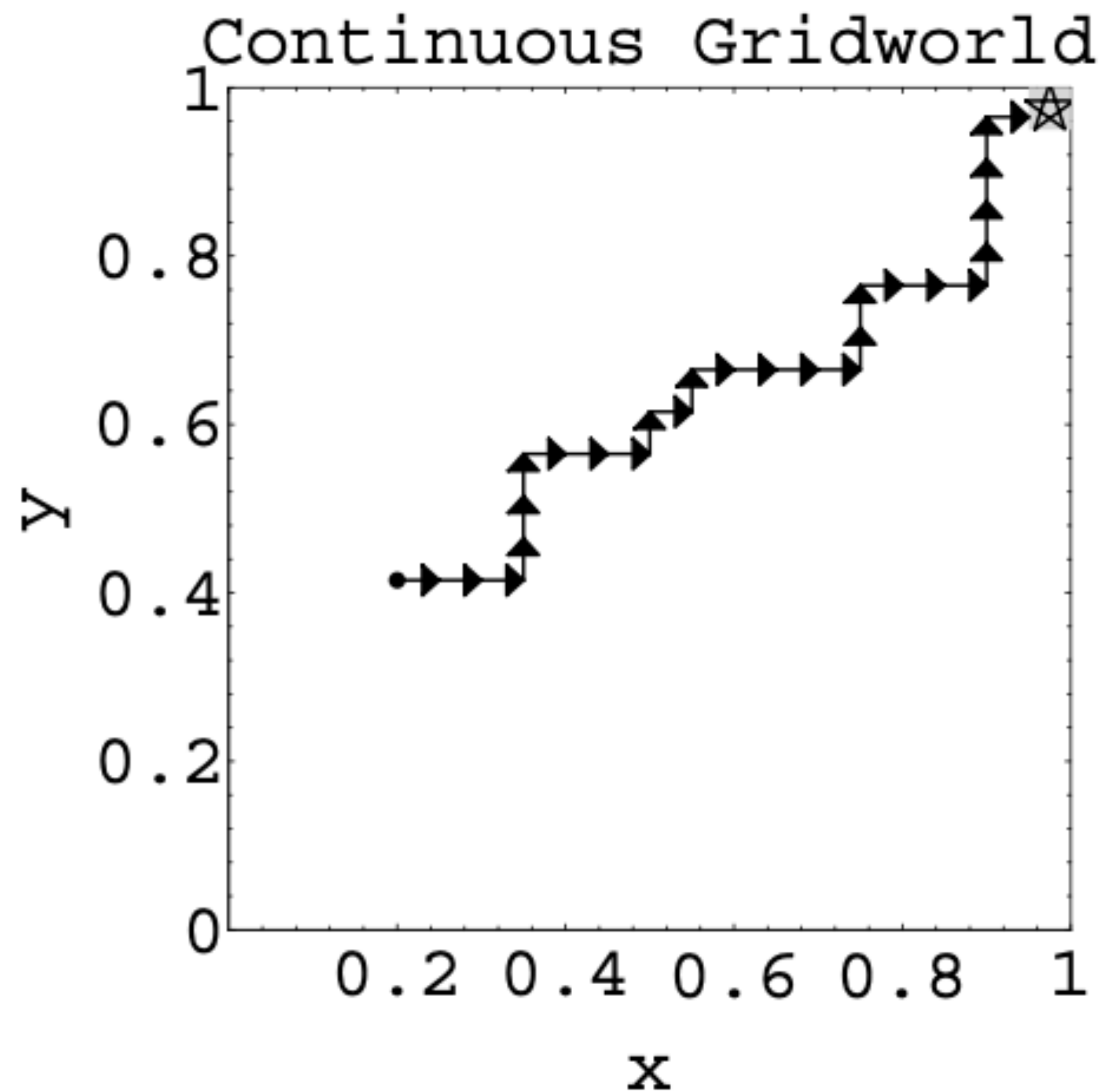
 **for** $i \in 1, \ldots, n$

  input $\leftarrow \{s_i, a_i\}$

  target $\leftarrow c_i + \gamma \min_{a} Q_\theta(s_i', a')$

  $D \leftarrow D \cup \{\text{input, output}\}$

 $Q_\theta \leftarrow \textbf{Train}(D)$

**return** $Q_\theta$

# A simple example: Gridworld



Optimal path

True value function

Boyan,Justin A and Moore, Andrew W, Generalization in Reinforcement Learning: Safely Approximating the Value Function. NeurIPS 1994.

# What happens when we run value iteration with a *quadratic?*



Iteration 17

# What happens when we run value iteration with a *quadratic?*



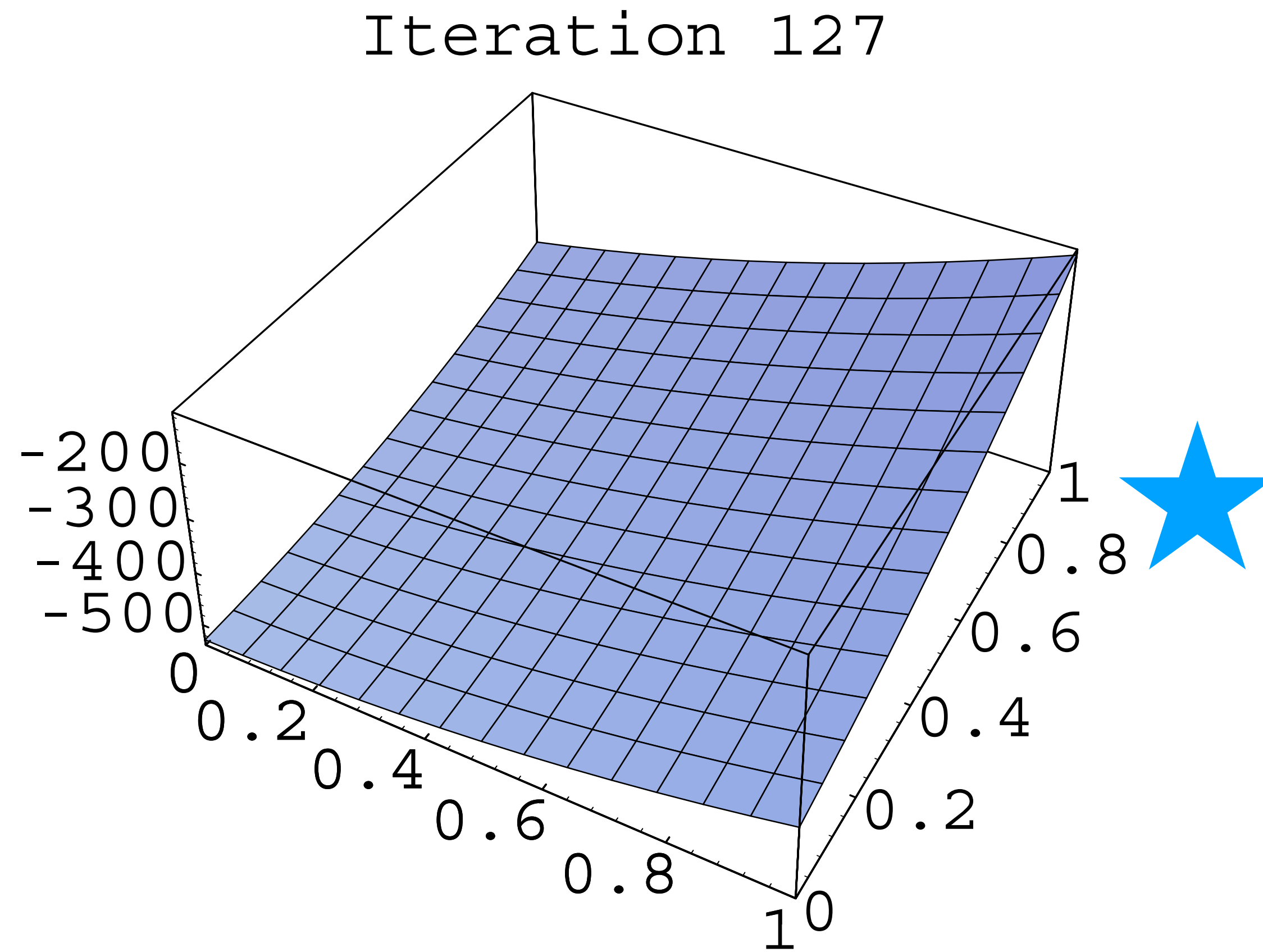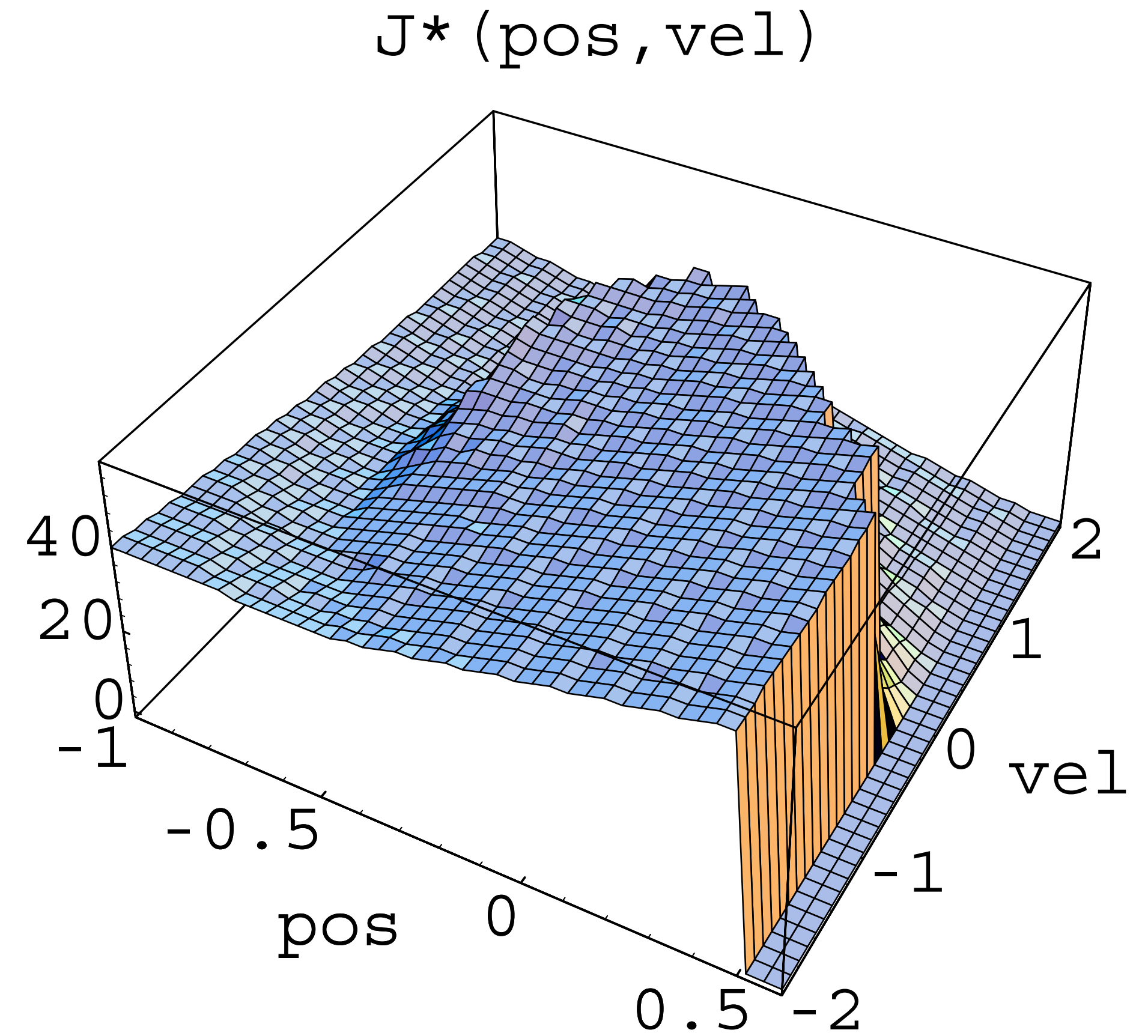Iteration 43

# What happens when we run value iteration with a *quadratic?*

# Another Example: Mountain Car!



Car-on-the-Hill

J*(pos,vel)

# What happens when we run value iteration with a *2 Layer MLP?*



Iteration 11

# What happens when we run value iteration with a *2 Layer MLP?*



Iteration 101

# What happens when we run value iteration with a
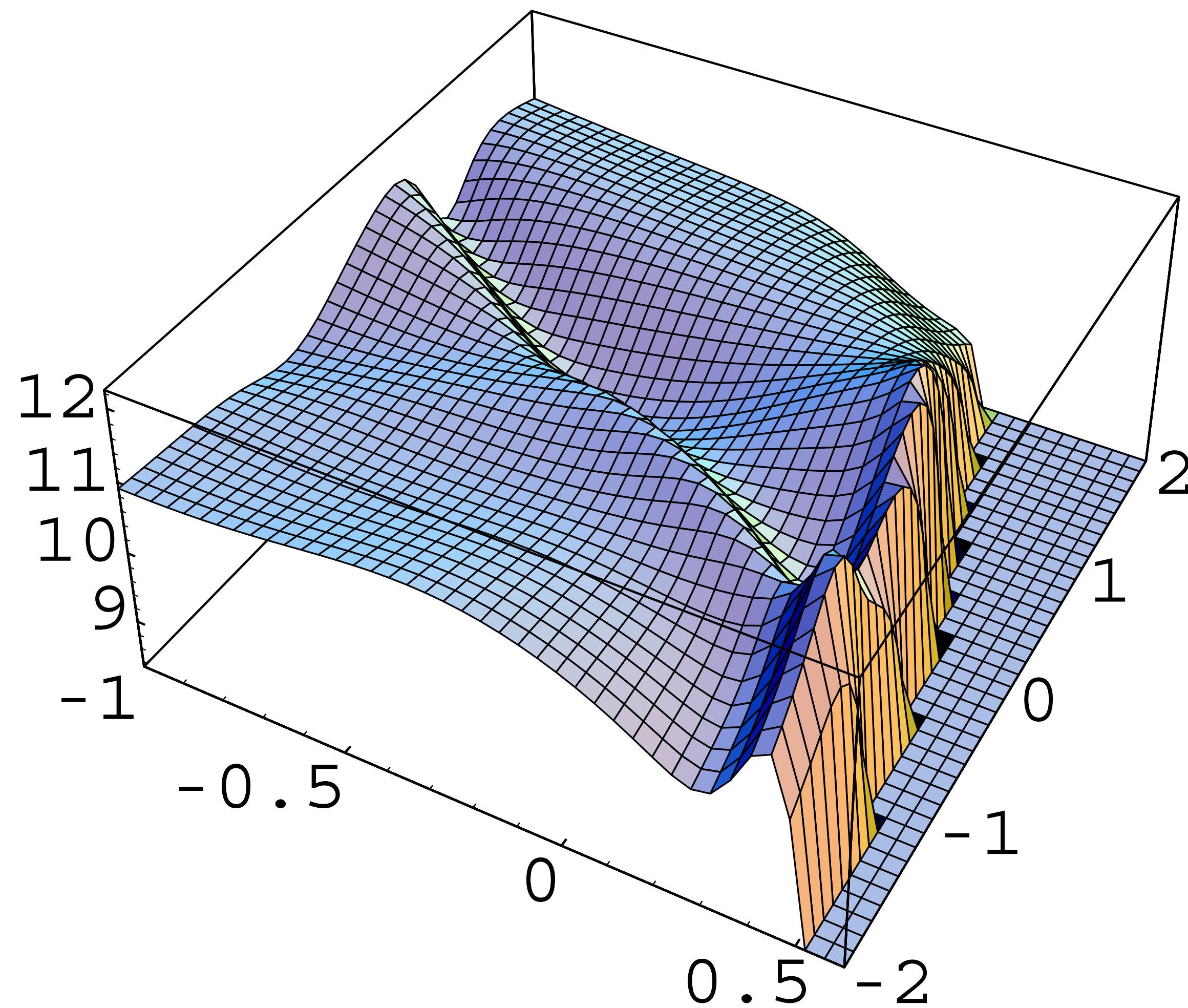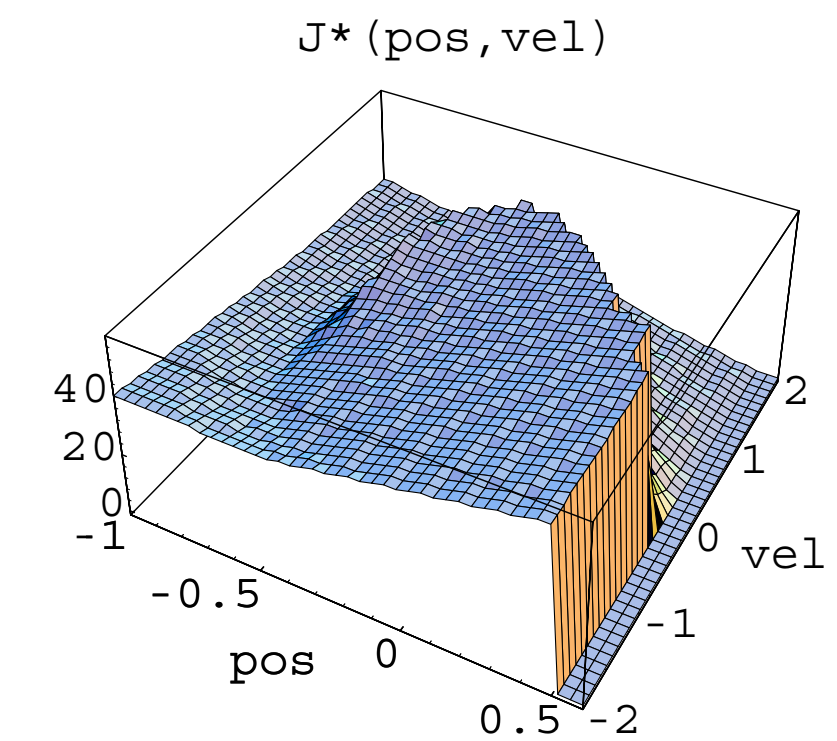## *2 Layer MLP?*

# The problem of Bootstrapping!



max()

Iteration 101

pos

Iteration 11

# The problem of Bootstrapping!

Errors in approximation are amplified

Key reason is the minimization

Minimization operation visit states where approximate values is less than the true value of that state – that is to say, states that look more attractive than they should.

Typically states where you have very few samples

max()

What about policy
iteration?

# Policy Iteration

## Policy Evaluation

## Policy Improvement



Iter: 0

$$Q^{\pi}(s, a) = c(s, a)) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s,a)} Q^{\pi}(s', \pi(s'))]$$

$$\pi^{+}(s) = \arg\min_{a} Q^{\pi}(s, a)$$

# Approximate (Fitted) Policy Iteration

### *Fitted* policy evaluation

**Given** $\{s_i, a_i, c_i, s_i'\}_{i=1}^N$



**Init** $Q_\theta(s, a) \leftarrow 0$

**while** *not converged* **do**

    $D \leftarrow \varnothing$

    **for** $i \in 1, \ldots, n$

        input $\leftarrow \{s_i, a_i\}$

        target $\leftarrow c_i + \gamma Q_\theta(s_i', \pi(s_i'))$

        $D \leftarrow D \cup \{\text{input, output}\}$

    $Q_\theta \leftarrow \textbf{Train}(D)$

**return** $Q_\theta$

### Policy Improvement

## This remains the same!

$$\pi^+(s) = \arg\min_a Q^\pi(s, a)$$

Surely approximate value
evaluation is more stable than
approximate value iteration?
(There is no min()!)

# Well ... not quite



$s_1$ → $s_2$ ↻

$f(s_1) = 1$        $f(s_1) = 2$

$w$ blows up!



37

# Well ... not quite



$f(s_1) = 1$    $f(s_1) = 2$

$w$ blows up!



But we can fix this by *on-policy weighting*

Weight each datapoint by how often the policy visits it.

# But what about policy improvement?

Policy Improvement

**Given** $\{s_i, a_i, c_i, s_i'\}_{i=1}^{N}$

**Init** $Q_\theta(s, a) \leftarrow 0$

**while** *not converged* **do**

$D \leftarrow \emptyset$

But this has
the min() step!

This is fine..

**for** $i \in 1, \dots, n$

    input $\leftarrow \{s_i, a_i\}$

    target $\leftarrow c_i + \gamma Q_\theta(s_i', \pi(s_i'))$

    $D \leftarrow D \cup \{\text{input, output}\}$
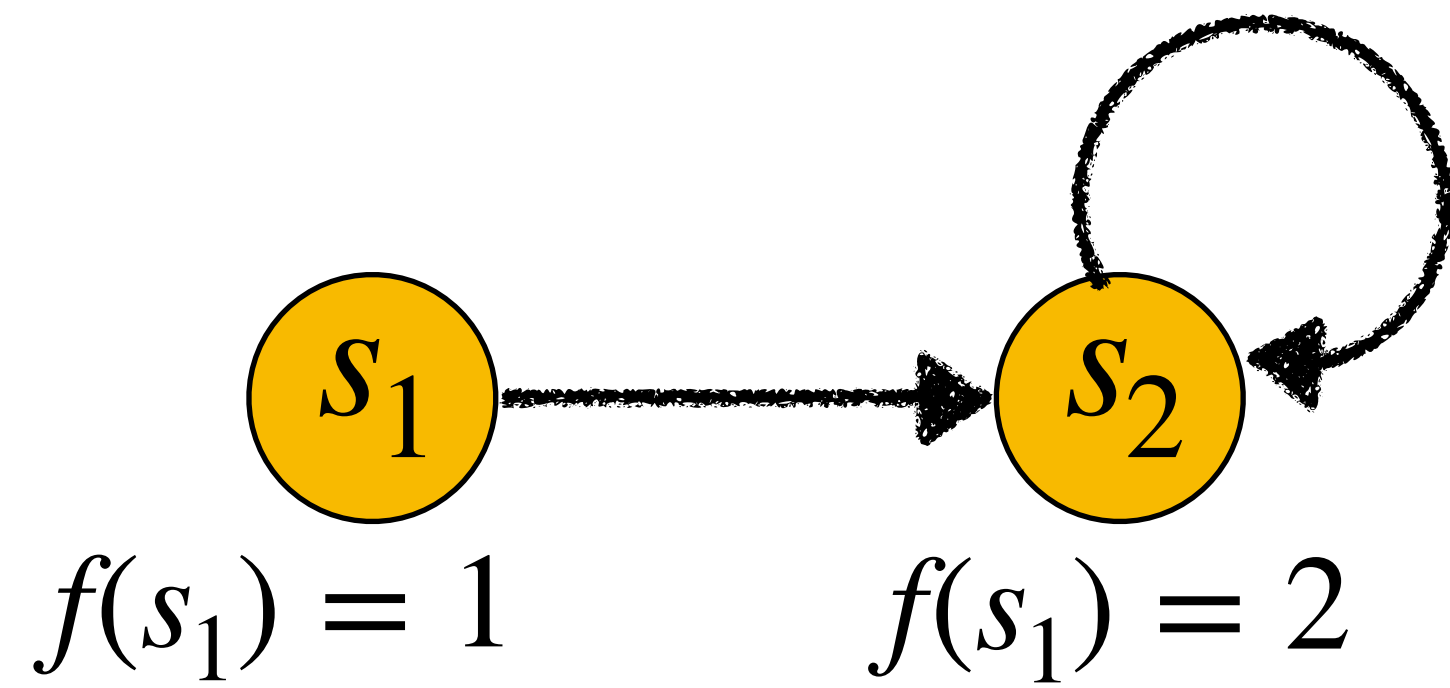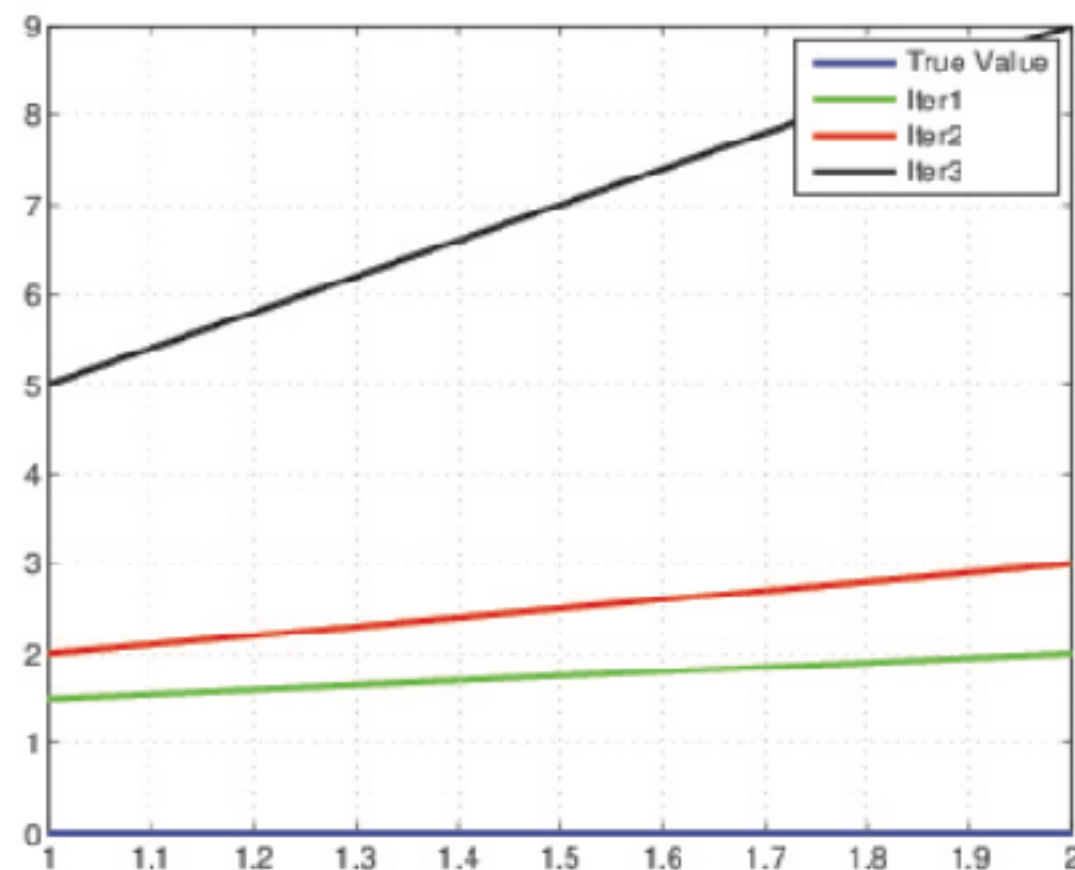
$Q_\theta \leftarrow \textbf{Train}(D)$

**return** $Q_\theta$

$$\pi^+(s) = \arg\min_a Q^\pi(s, a)$$

# The problem of distribution shift

Upper half of state
is BAD

Lower half of state
is GOOD

- - - - Approximated Q

——— True Q

# The problem of distribution shift

Upper half of state is BAD

Lower half of state is GOOD

T-2    T-1

---- Approximated Q

___ True Q

# The problem of distribution shift

T-3        T-2        T-1

Upper half of state
is BAD

Lower half of state
is GOOD

---- Approximated Q

—— True Q

Boostrapping

Distribution Shift

# Two sides of the same coin

## Bootstrapping

## Distribution shift

**max**



| T-3 | T-2 | T-1 |
|---|---|---|

Upper half of state
is BAD

Lower half of state
is GOOD

---- Approximated Q

—— True Q

44

# Ideas for fixing this?

# Remedies



Upper half of state is BAD

Lower half of state is GOOD

----- Approximated Q
— True Q

## Bootstrapping



When doing min(),
don't trust value estimate

Execute policy
and trust actual returns

## Distribution shift



Minimize the
distribution shift

Be conservative,
change policy slowly

# Remedies



Upper half of state is BAD

Lower half of state is GOOD

----- Approximated Q

—— True Q

## Bootstrapping



When doing min(),
don't trust value estimate

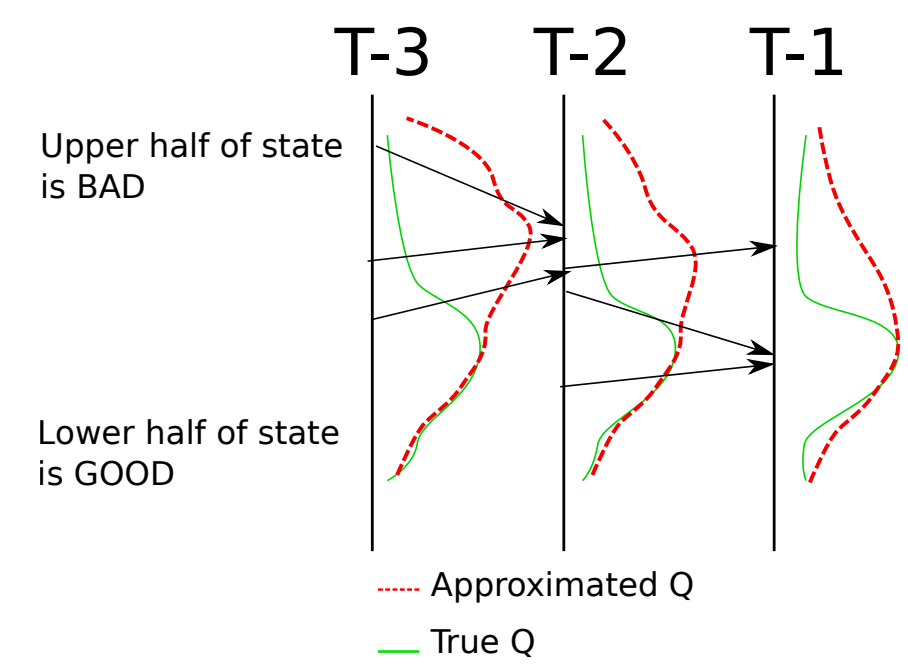Execute policy
and trust actual returns
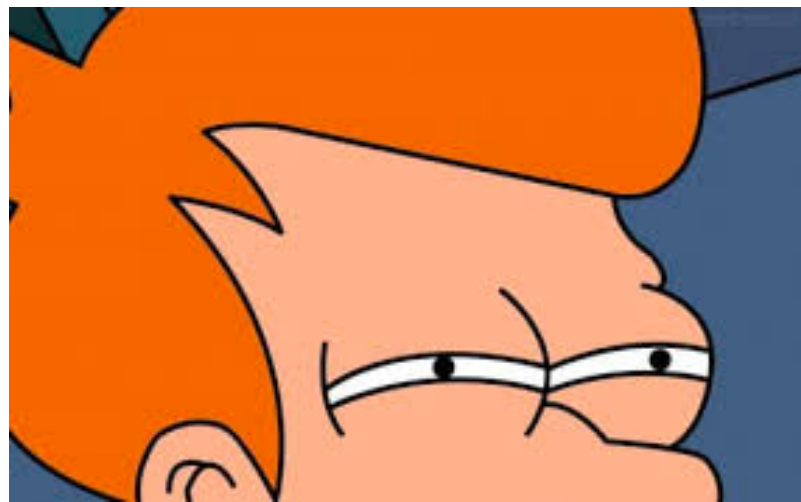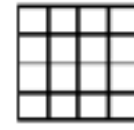
## Distribution shift



Minimize the
distribution shift

Be conservative,
change policy slowly

# tl;dr

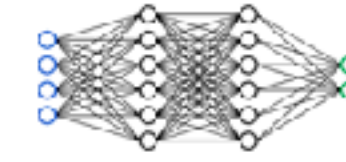## Approximate (Fitted) Value Iteration

### Q-iteration

$Q(s, a) \leftarrow 0$

**while** *not converged* **do**

    **for** $s \in S, a \in A$

        $Q^{new}(s, a) = c(s, a) + \gamma \mathbb{E}_{s'} \min_{a'} Q(s', a')$

    $Q \leftarrow Q^{new}$

**return** $Q$

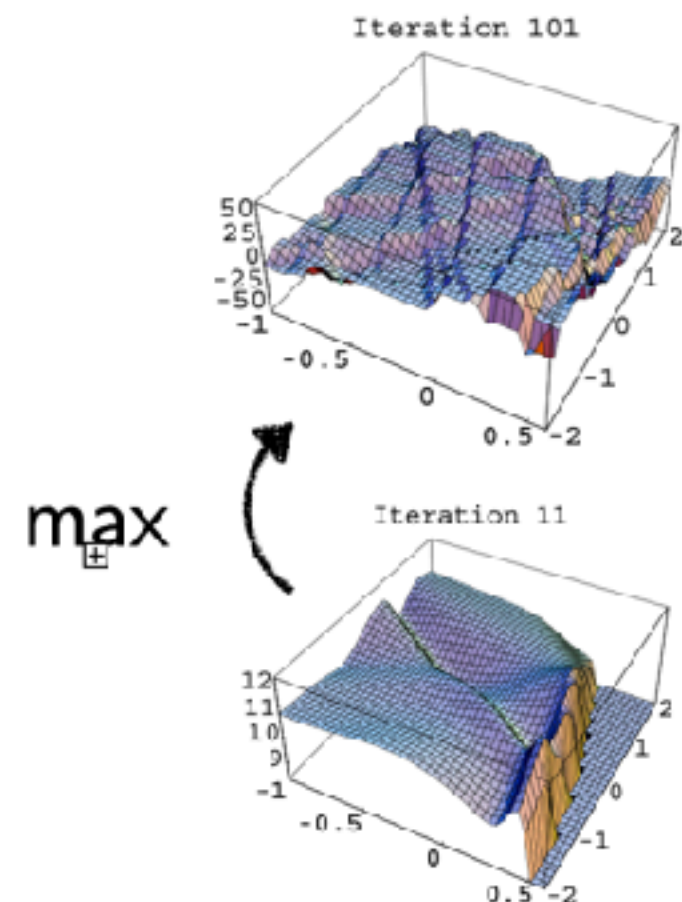### Fitted Q-iteration

**Given** $\{s_i, a_i, c_i, s_i'\}_{i=1}^N$

**Init** $Q_\theta(s, a) \leftarrow 0$

**while** *not converged* **do**

    $D \leftarrow \varnothing$

    **for** $i \in 1, \ldots, n$

        input $\leftarrow \{s_i, a_i\}$ ,

        target $\leftarrow c_i + \gamma \min_a Q_\theta(s_i', a')$

        $D \leftarrow D \cup \{input, output\}$

    $Q_\theta \leftarrow \textbf{Train}(D)$

**return** $Q_\theta$
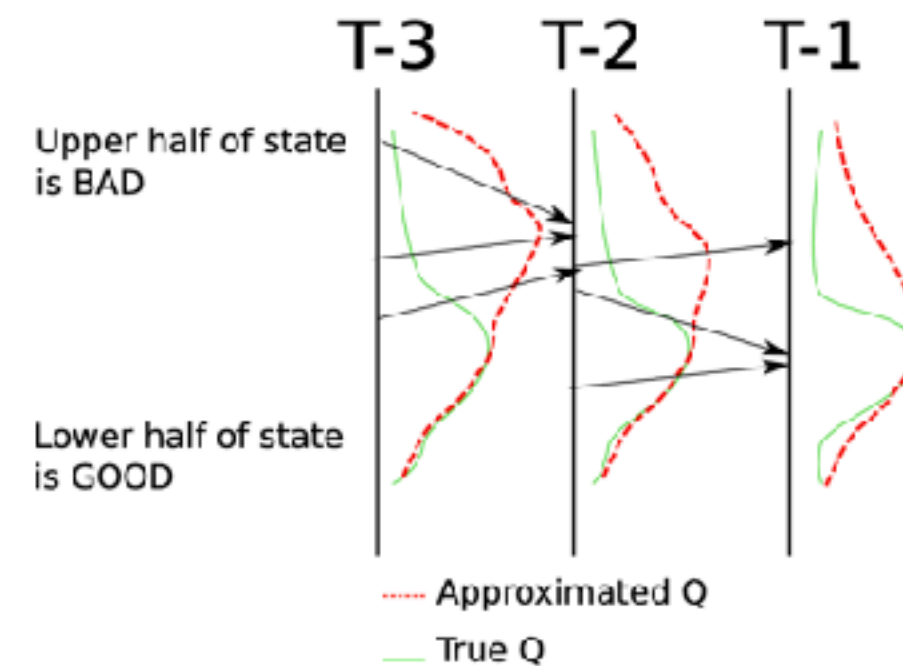
1

## Bootstrapping

Iteration 101

max

Iteration 11

## Distribution shift

        T-3        T-2        T-1

Upper half of state
is BAD

Lower half of state
is GOOD

----- Approximated Q

——— True Q

36

## Remedies

                                      T-3  T-2  T-1

Upper half of state
is BAD

Lower half of state
is GOOD

----- Approximated Q
——— True Q

### Bootstrapping

When doing min(),
don't trust value estimate

Execute policy
and trust actual returns

### Distribution shift

Minimize the
distribution shift

Be conservative,
change policy slowly