# Approximate Value and Policy Iteration

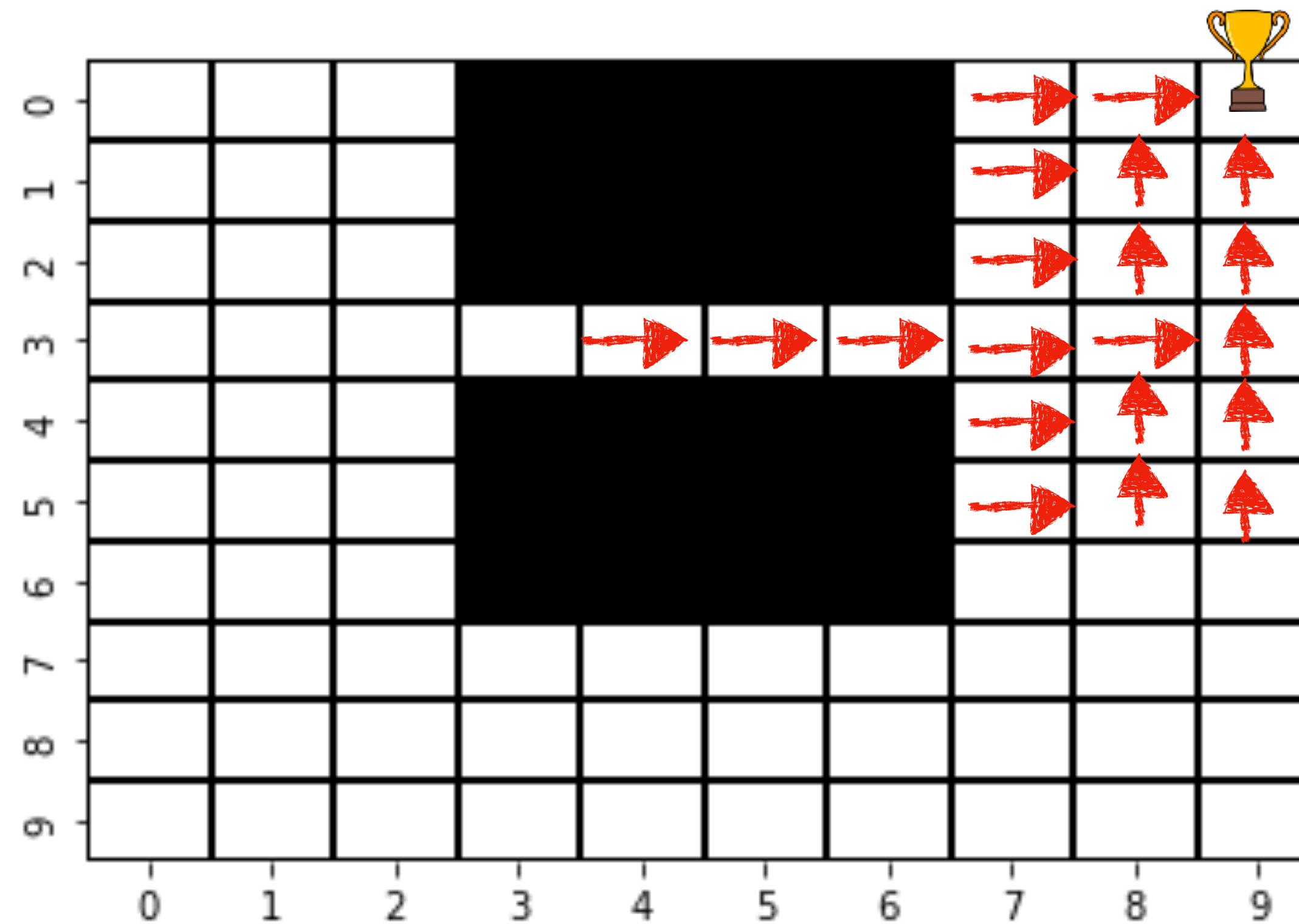Sanjiban Choudhury

# The story thus far ...

# The story thus far

We know how to define an MDP

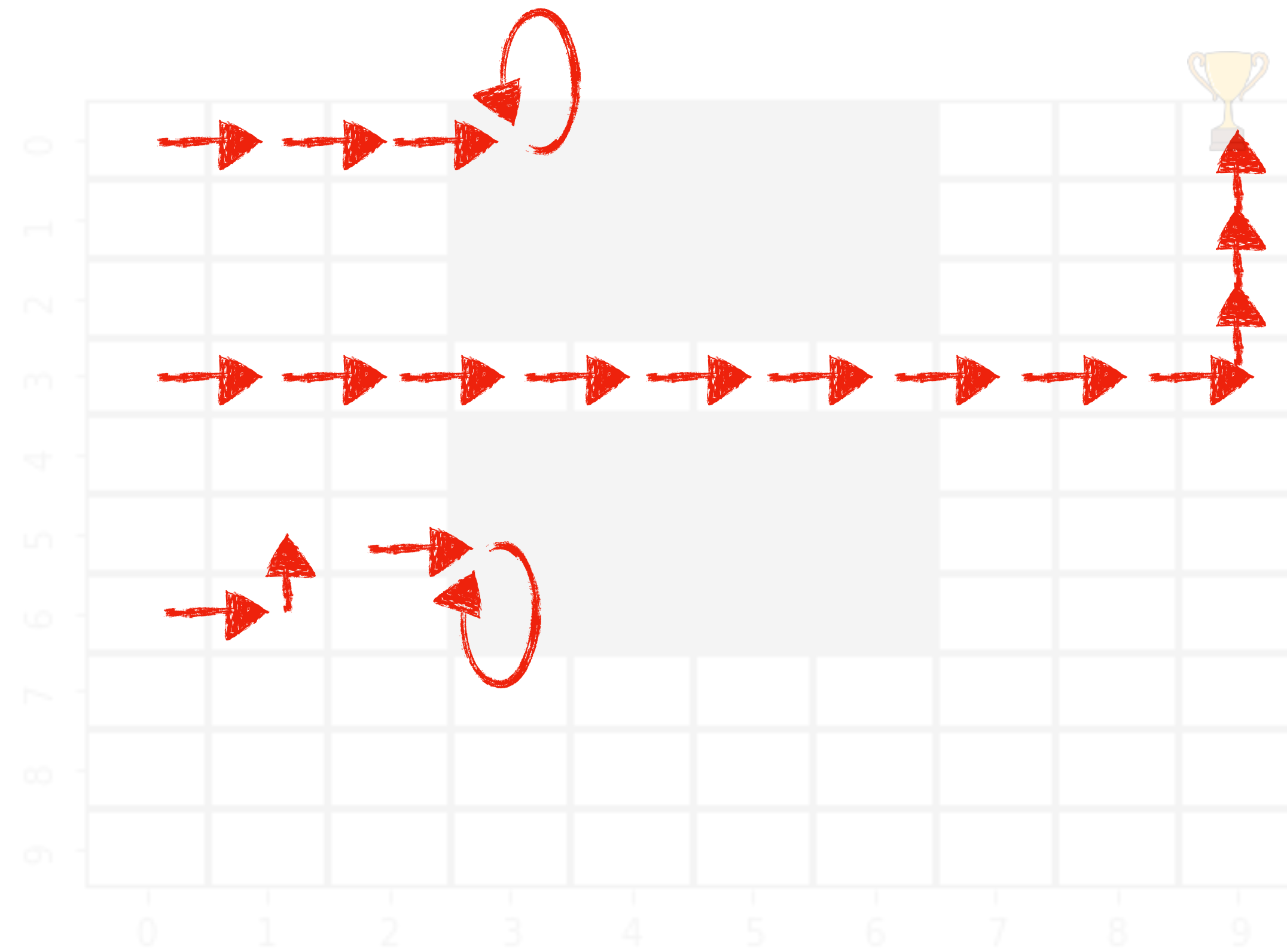If the MDP is known (i.e. I know my costs and my transition)

We know how to solve a MDP

What happens if the MDP is unknown?

# Known MDP



*If I know the transition function, I could teleport to any state, try any action and know the next state*

# *Unknown* MDP



*I don't know the transition, I can only roll-out from start state, and see where I end up*

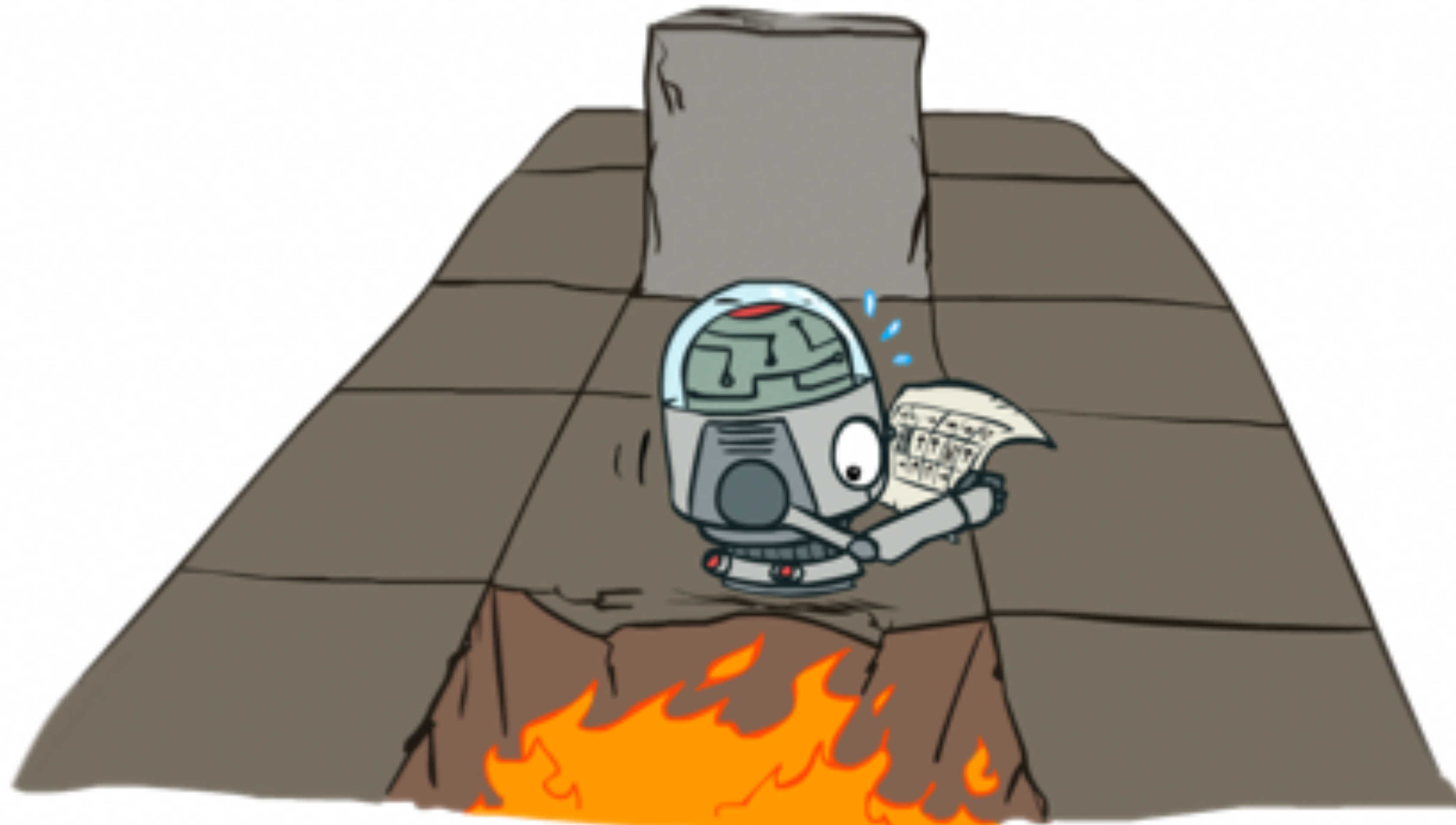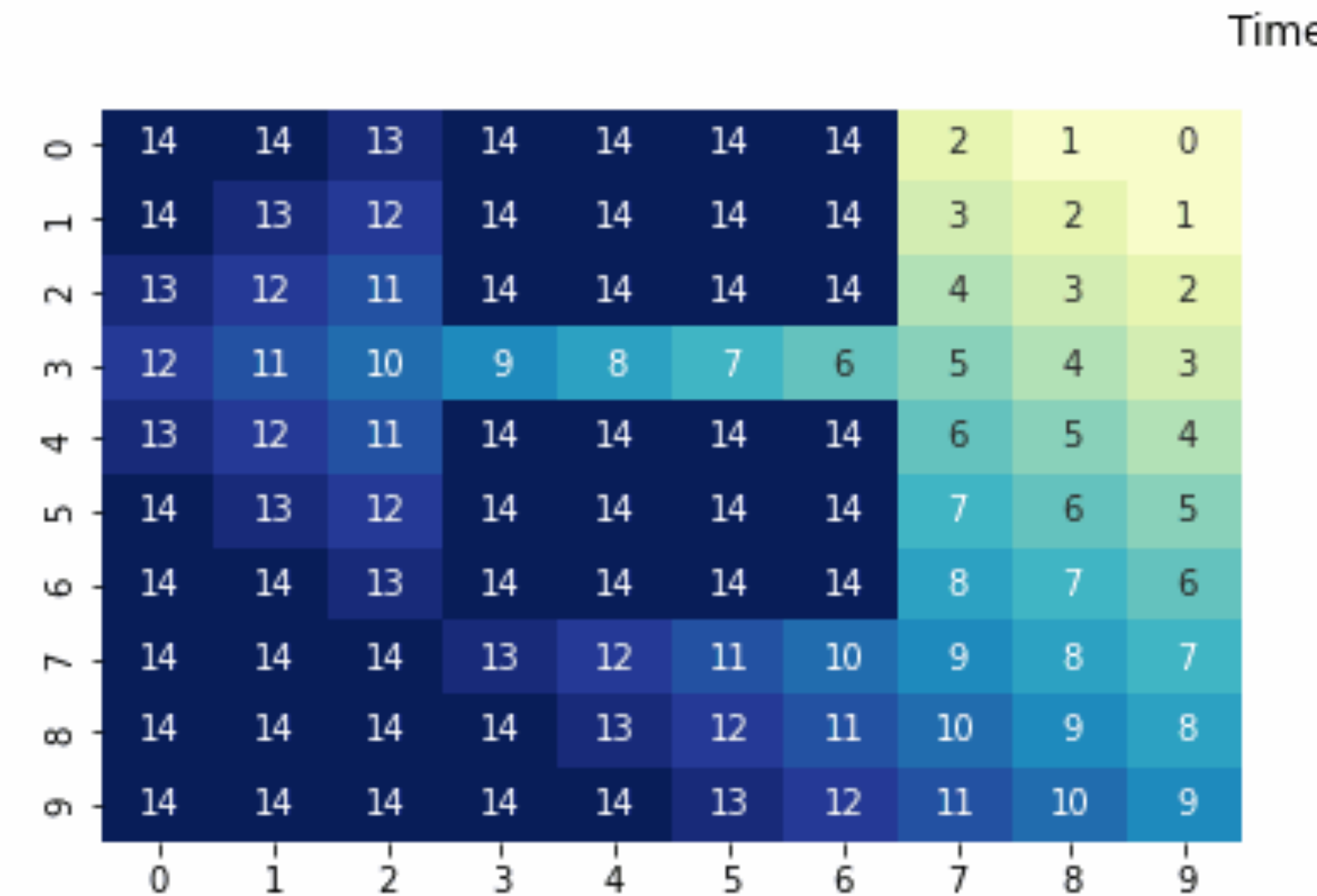# Recall: How do we solve a known MDP?

# Value Iteration

Initialize value function at last time-step

$$V^*(s, T-1) = \min_a c(s, a) \quad \forall s$$

for $t = T-2, \ldots, 0$

Compute value function at time-step t

$$V^*(s, t) = \min_a \left[ c(s, a) + \gamma \sum_{s'} \mathcal{T}(s'|s, a) V^*(s', t+1) \right] \forall s$$

# Q-Value Iteration

Initialize value function at last time-step

$$Q^*(s, a, T-1) = c(s, a) \quad \forall (s, a)$$

for $t = T-2, \ldots, 0$

Compute value function at time-step t

$$Q^*(s, a, t) = c(s, a) + \gamma \sum_{s'} \mathcal{T}(s' \,|\, s, a) \min_{a'} \; Q^*(s', a', t+1) \quad \forall s, a$$

# Q-Value Iteration (Infinite horizon)

Initialize value function at last time-step

$$Q^*(s, a) = c(s, a) \quad \forall (s, a)$$

While not converged

Update value function

$$Q^*(s, a) = c(s, a) + \gamma \sum_{s'} \mathcal{T}(s' \mid s, a) \min_{a'} Q^*(s', a') \quad \forall (s, a)$$

# Two Problems

Initialize value function at last time-step

$$Q^*(s, a) = c(s, a) \quad \forall (s, a)$$

*1) What happens when states are continuous?*

While not converged

*2) What happens when I don't know the MDP?*

Update value function

$$Q^*(s, a) = \boxed{c(s, a)} + \gamma \sum_{s'} \boxed{\mathcal{T}(s'|s, a)} \min_{a'} Q^*(s', a') \quad \boxed{\forall (s, a)}$$

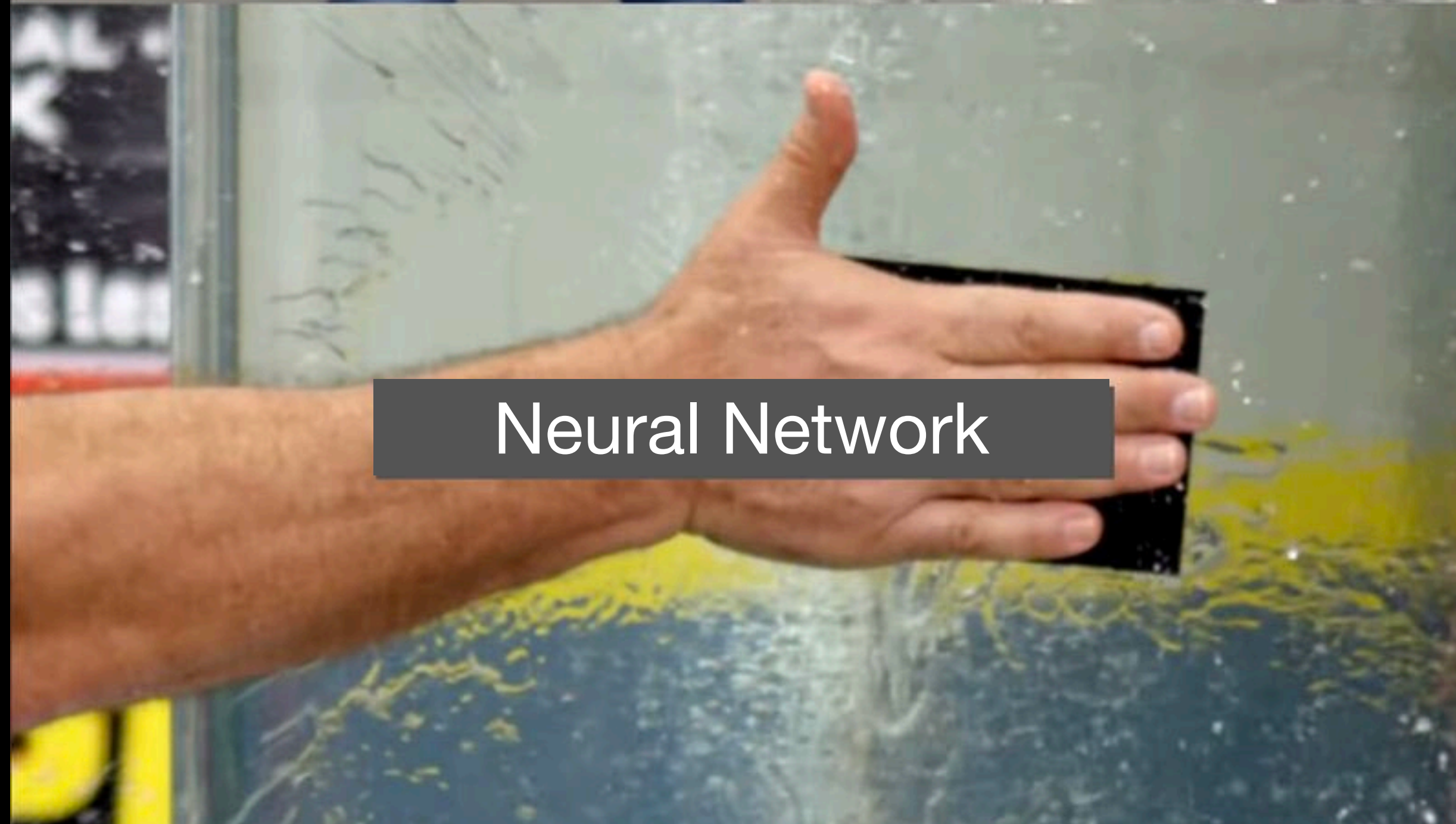*Are these known?*

*Can I do this?*

# Simple Idea

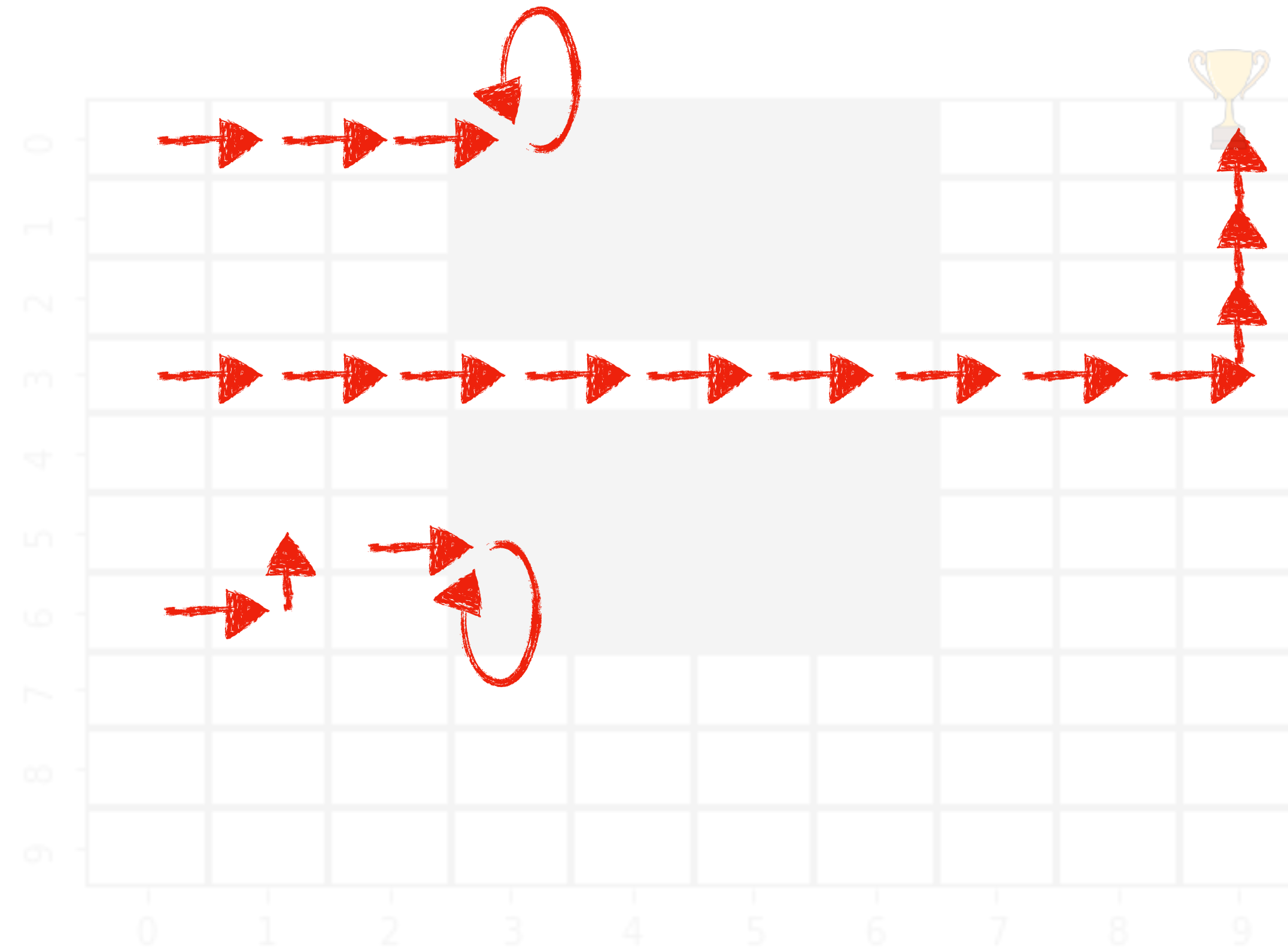Can I collect roll-out data from the real world and just fit a Q function?

Unknown function

Researchers

Neural Network

# Step 1: First collect roll-out data
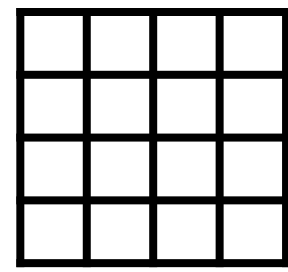


*Data is a tuple of state, action, cost, next state*

$$\mathcal{D} = \{(s_i, a_i, c_i, s_{i+1})\}_{i=1}^n$$

# Step 2: Fitted Q-Iteration

Regular Q-iteration

$Q(s, a) \leftarrow c(s, a)$

**while** *not converged* **do**
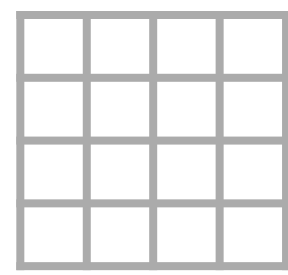
    **for** $s \in S, a \in A$

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

    $Q \leftarrow Q^{new}$

**return** $Q$

# Step 2: Fitted Q-Iteration

## Regular Q-iteration

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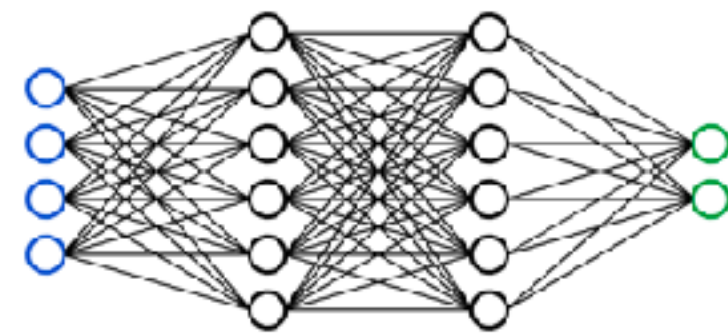
**while** *not converged* **do**

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

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

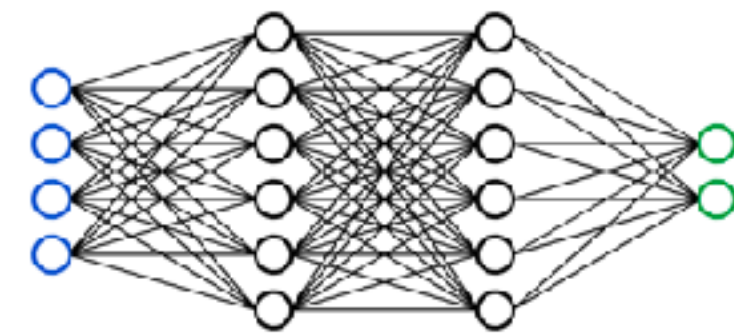     target $\leftarrow c_i + \gamma \min_{a} Q_\theta(s_i', a')$

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

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

  **return** $Q_\theta$

# Step 2: Fitted Q-Iteration

*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** $i \in 1, \ldots, n$    *Use old copy of Q*

      input $\leftarrow \{s_i, a_i\}$    *to set target*

      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** $Q_\theta$

*Training is a regression problem*

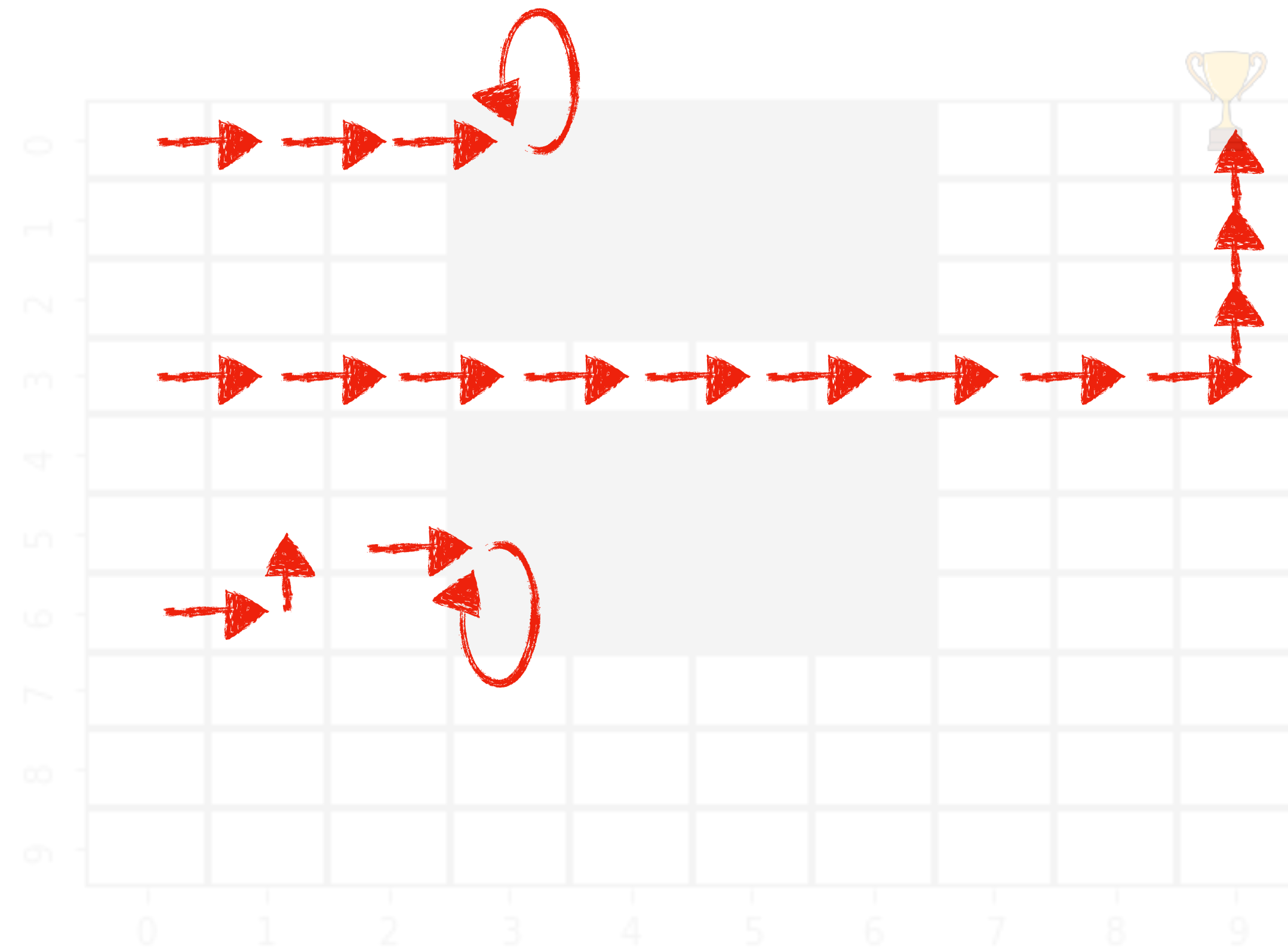$$\ell(\theta) = \sum_{i=1}^n (Q_\theta(s_i, a_i) - target)^2$$

16

# Temporal Difference Error (TD Error)

Penalize violation of Bellman Equation

$$\ell(\theta) = \left( c(s,a) + \gamma \min_{a'} Q_{\theta_{old}}(s', a') - Q_{\theta}(s_t, a_t) \right)^2$$

$$\theta = \theta_{\text{old}} - \alpha \nabla_{\theta} l(\theta)$$

# What policy do I use to collect data?



Do I explore randomly? Do I use my learnt Q function?
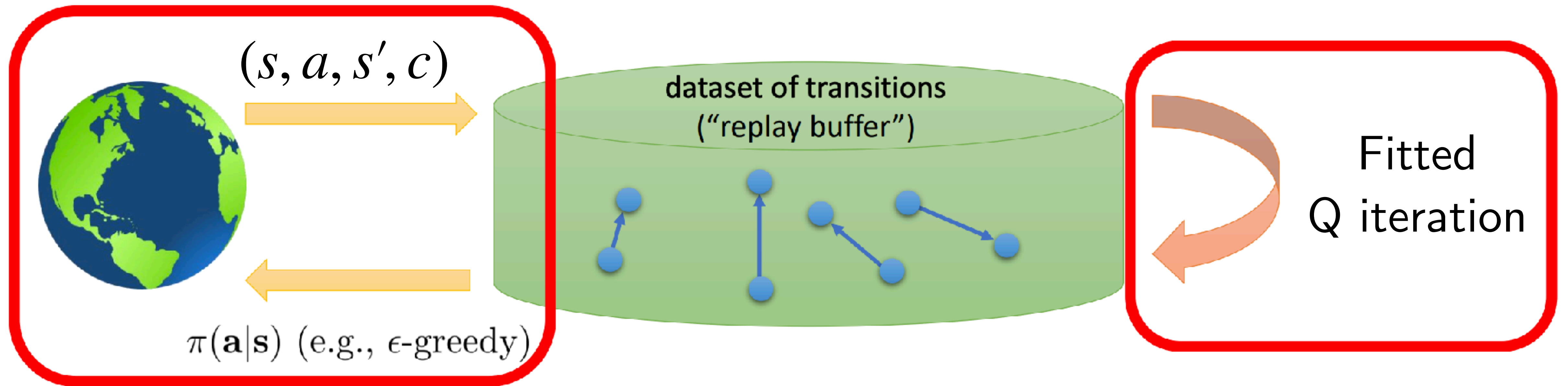
# What policy do I use to collect data?

When poll is active respond at **PollEv.com/sc2582**

Send **sc2582** to **22333**

Do I explore randomly? Do I use my learnt Q function?

# Q-learning: Learning off-policy



$(s, a, s', c)$

dataset of transitions
("replay buffer")

Fitted
Q iteration

$\pi(\mathbf{a}|\mathbf{s})$ (e.g., $\epsilon$-greedy)

# QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation
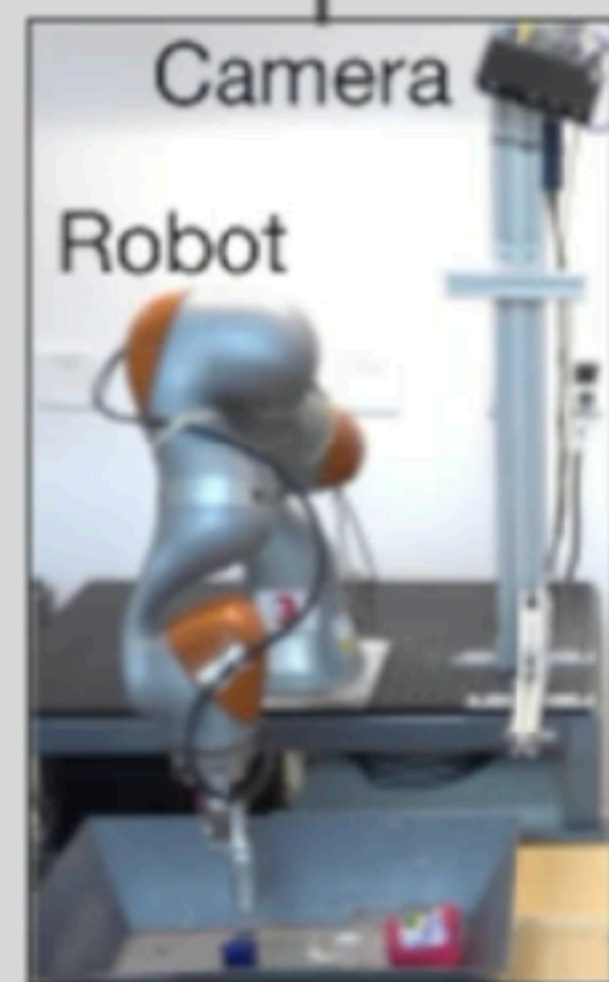
**Training time**

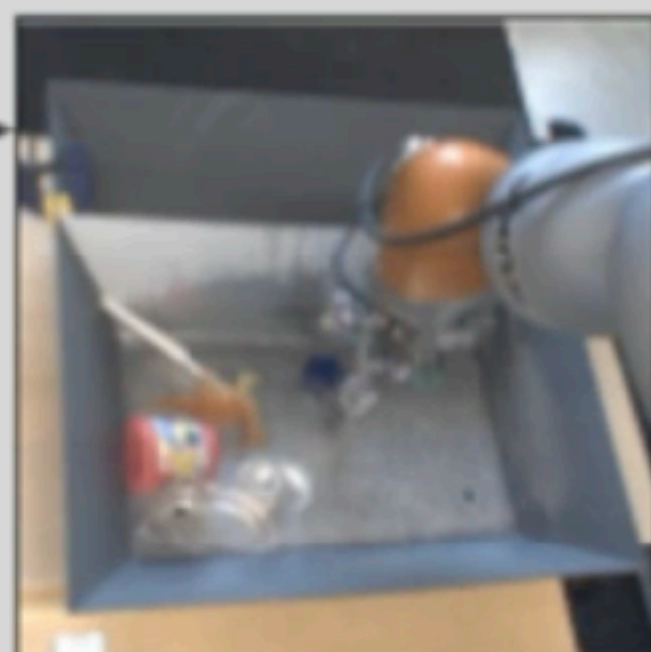Reward: Grasp success determined by subtracting pre and post-drop images

Distributed RL

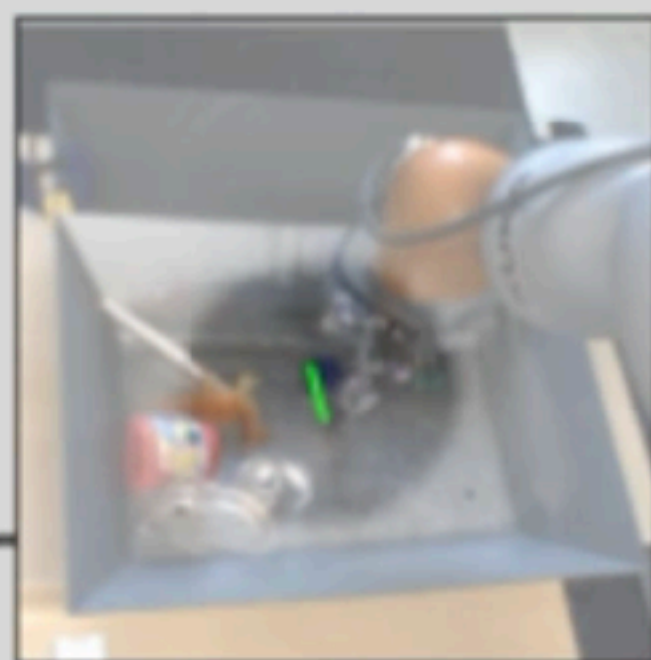State, Action, Reward

Learned weights

**Inference time**

Camera

Robot

State: 472x472 Image and gripper aperture

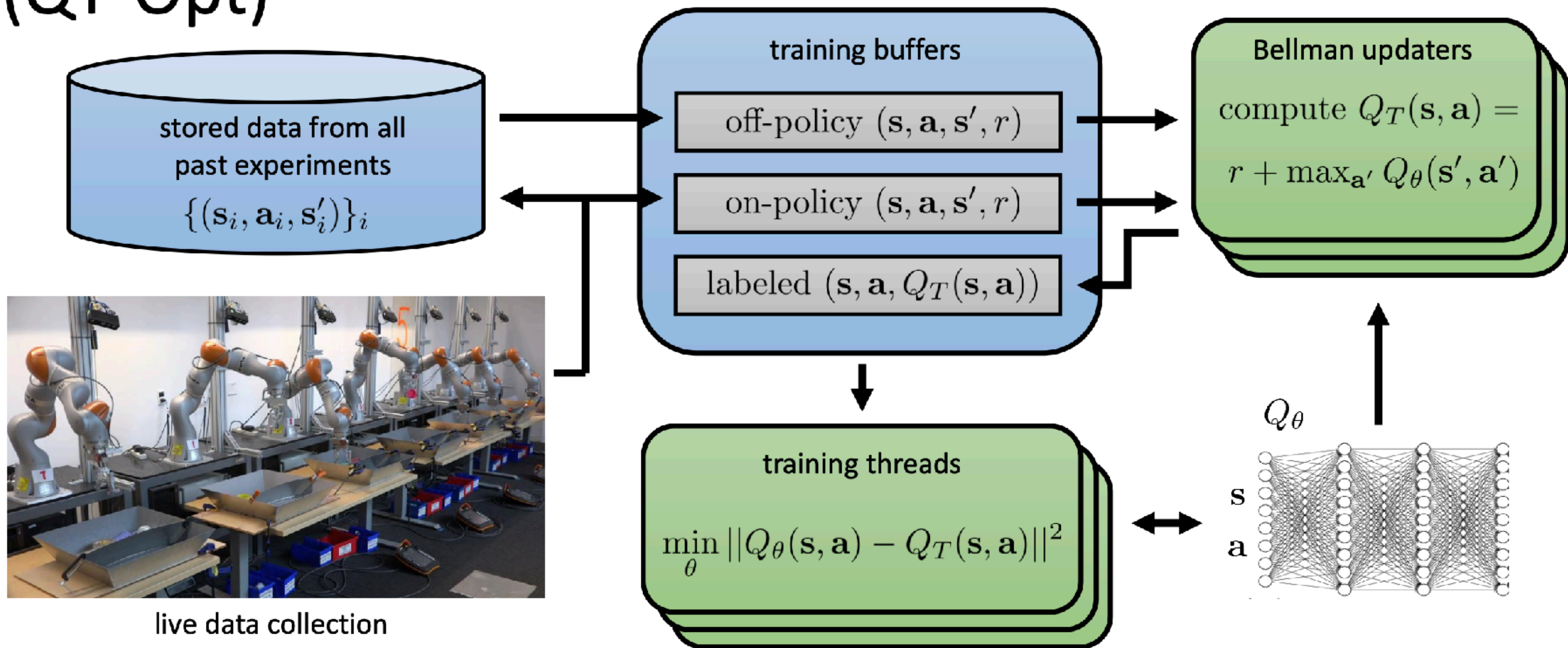Action: Gripper displacement and aperture

Action proposals

Critic Function
$Q(\text{State}, \text{Action})$

Q-Values

Cross-Entropy Method
$\arg\max_{\text{Action}} Q(\text{State}, \text{Action})$

# Large-scale Q-learning with continuous actions (QT-Opt)



training buffers

off-policy $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$

on-policy $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)$

labeled $(\mathbf{s}, \mathbf{a}, Q_T(\mathbf{s}, \mathbf{a}))$

Bellman updaters

$$\text{compute } Q_T(\mathbf{s}, \mathbf{a}) = r + \max_{\mathbf{a}'} Q_\theta(\mathbf{s}', \mathbf{a}')$$

stored data from all past experiments

$$\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}_i')\}_i$$

live data collection

training threads

$$\min_\theta ||Q_\theta(\mathbf{s}, \mathbf{a}) - Q_T(\mathbf{s}, \mathbf{a})||^2$$

$Q_\theta$

$\mathbf{s}$

$\mathbf{a}$

Kalashnikov, Irpan, Pastor, Ibarz, Herzong, Jang, Quillen, Holly, Kalakrishnan, Vanhoucke, Levine. **QT-Opt: Scalable Deep Reinforcement Learning of Vision-Based Robotic Manipulation Skills**
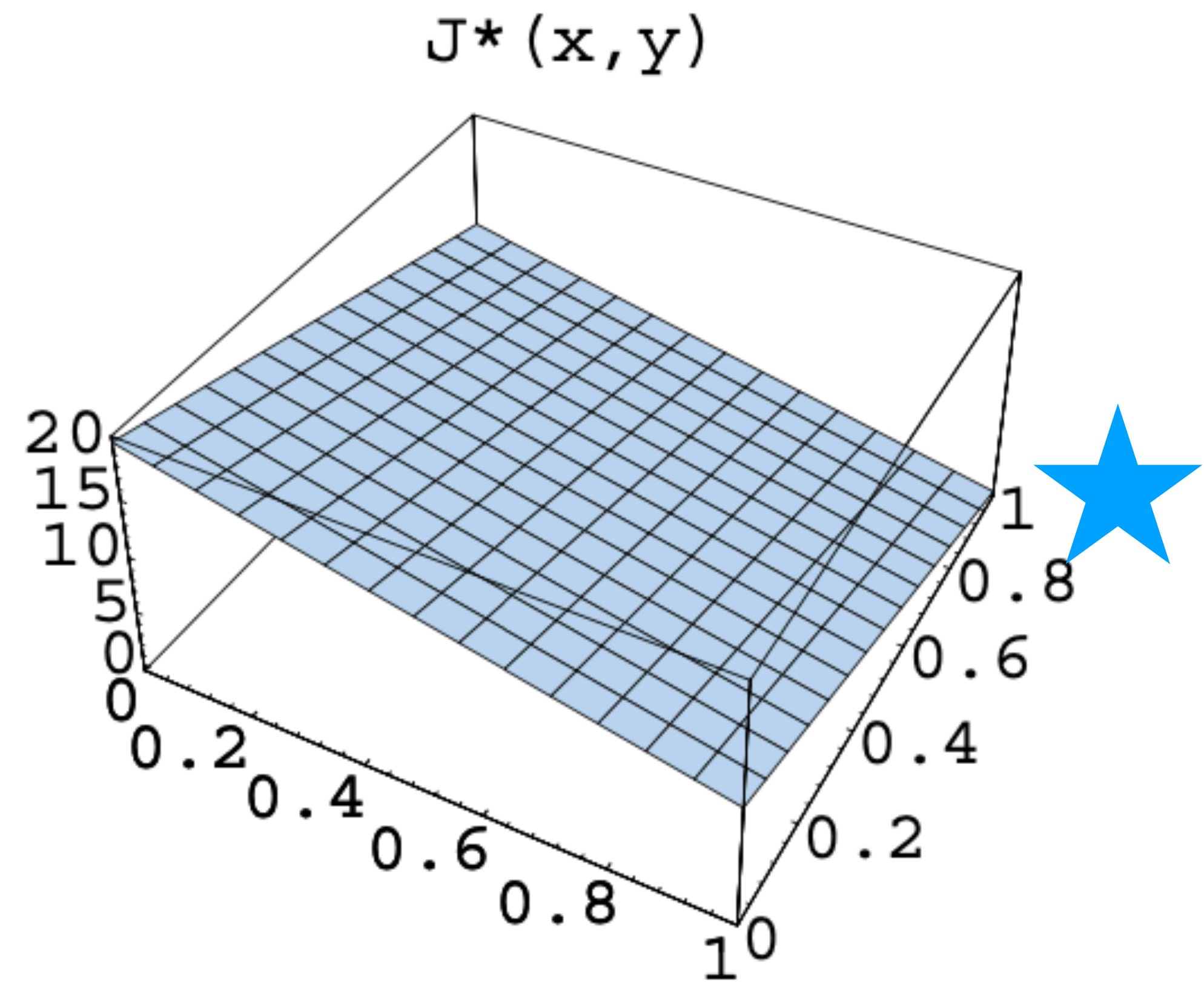
So does approximate
value iteration work?

# 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

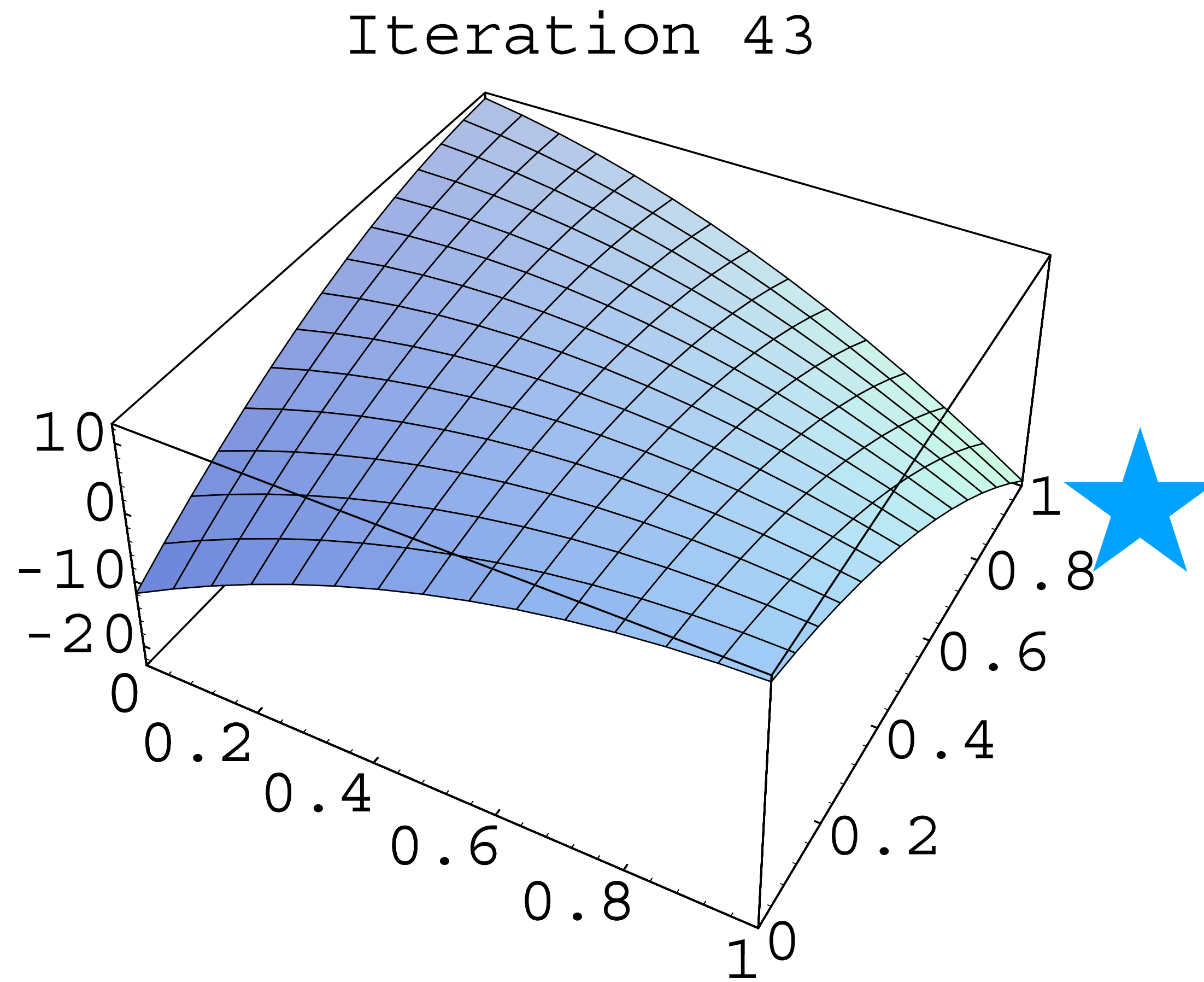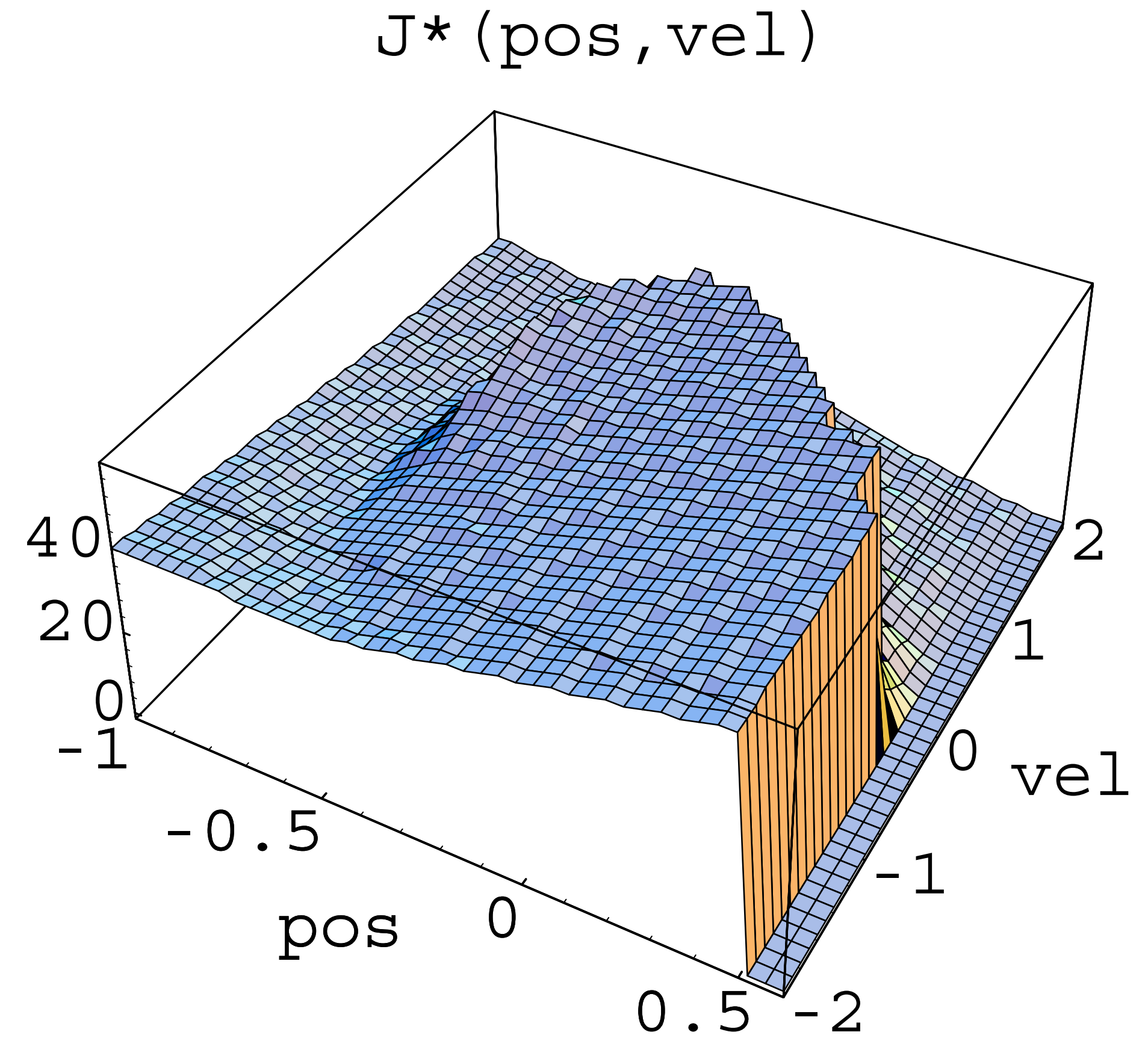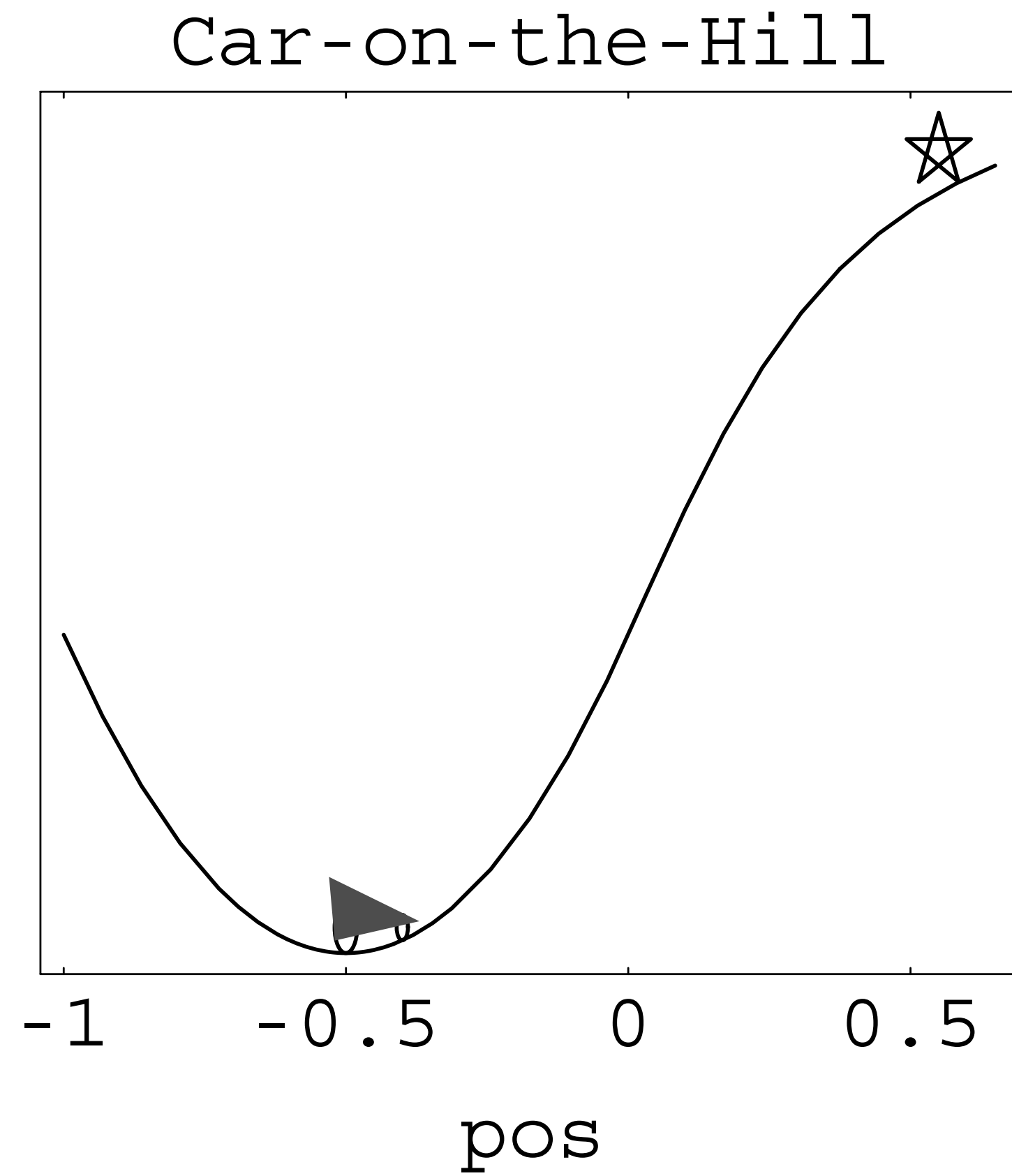# What happens when we run value iteration with a *quadratic?*

Iteration 17

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



Iteration 43

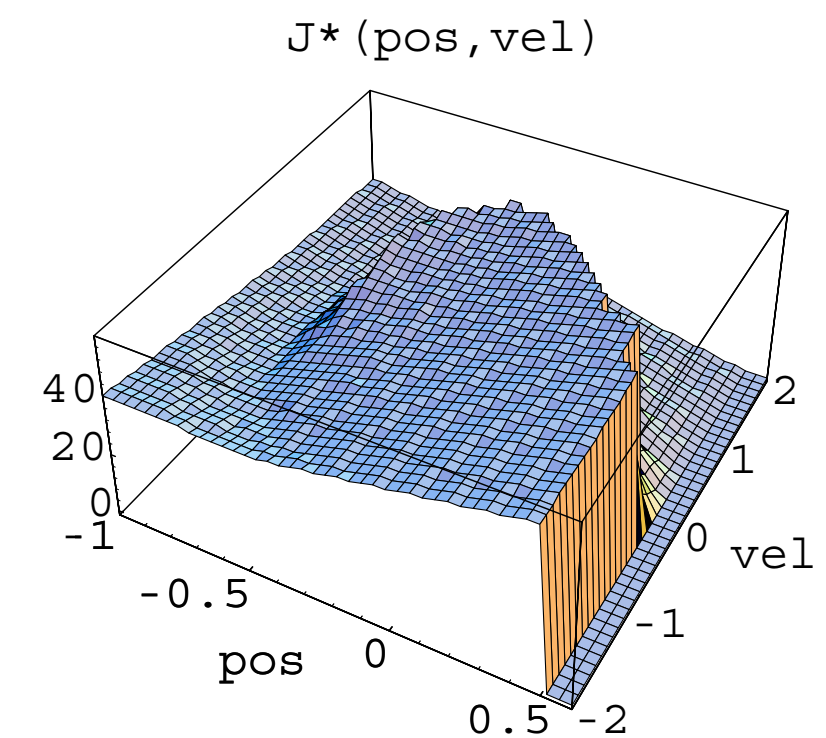# What happens when we run value iteration with a *quadratic?*



Iteration 127

# Another Example: Mountain Car!



Car-on-the-Hill

pos

J*(pos,vel)

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

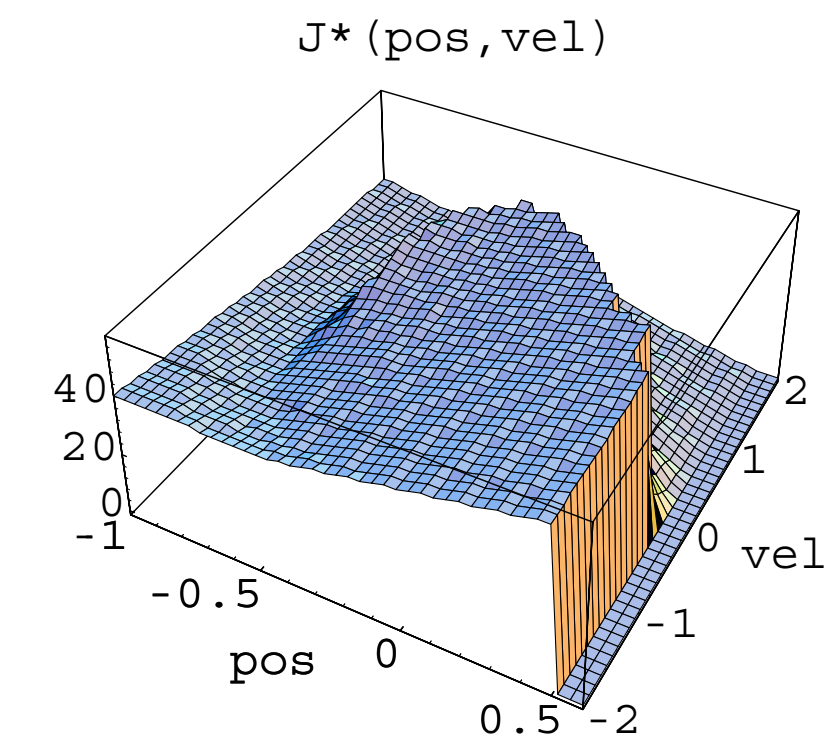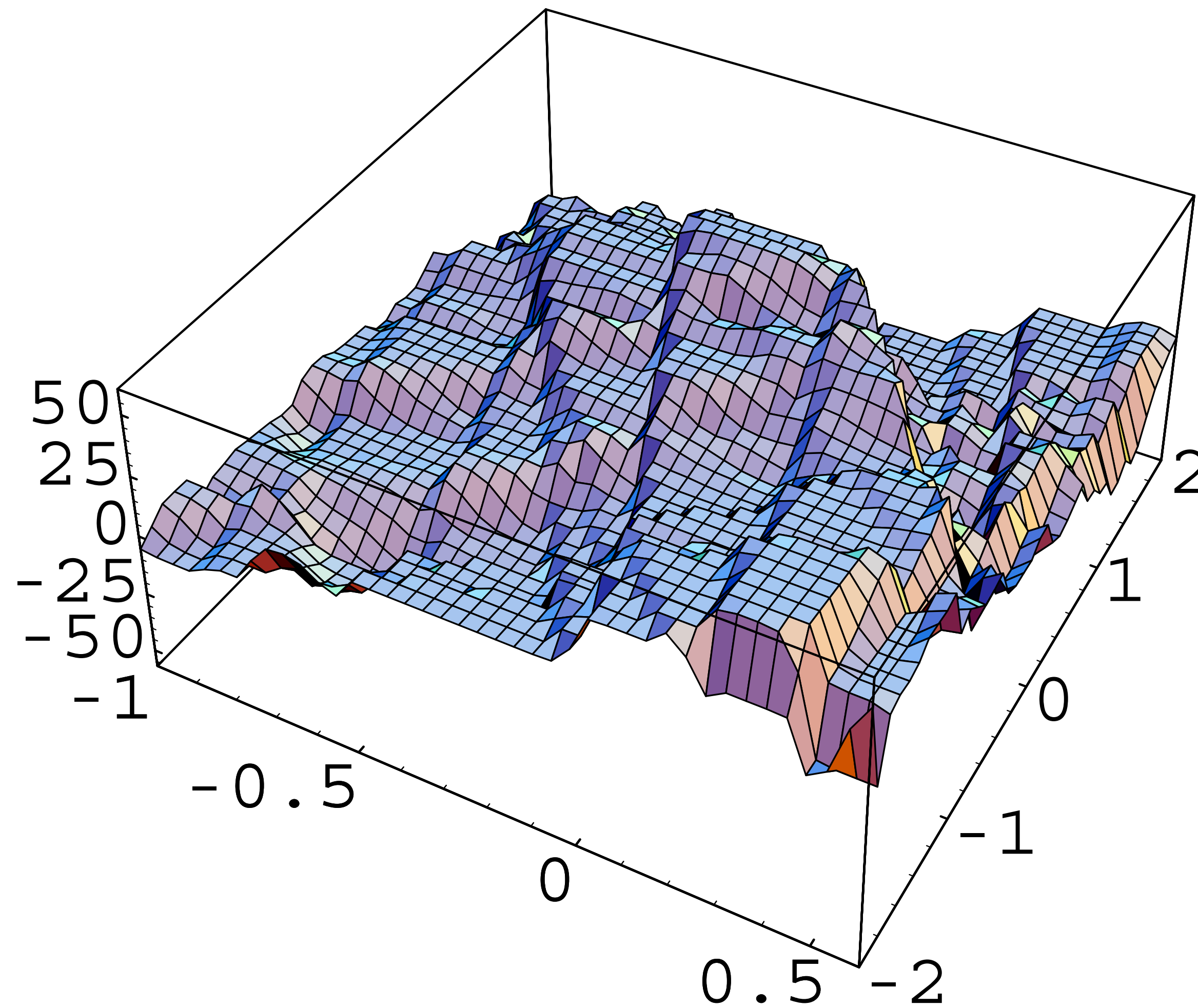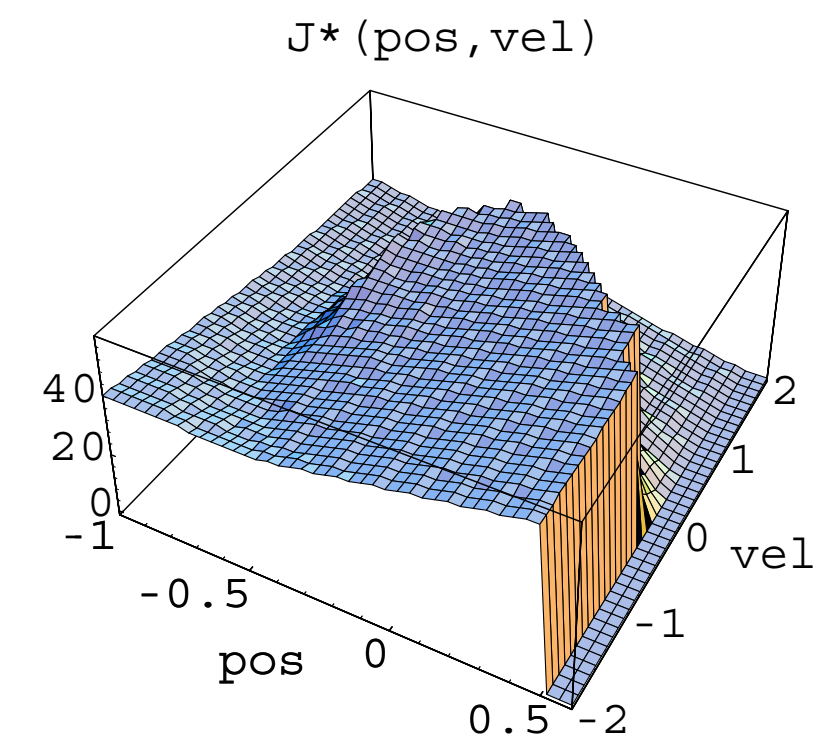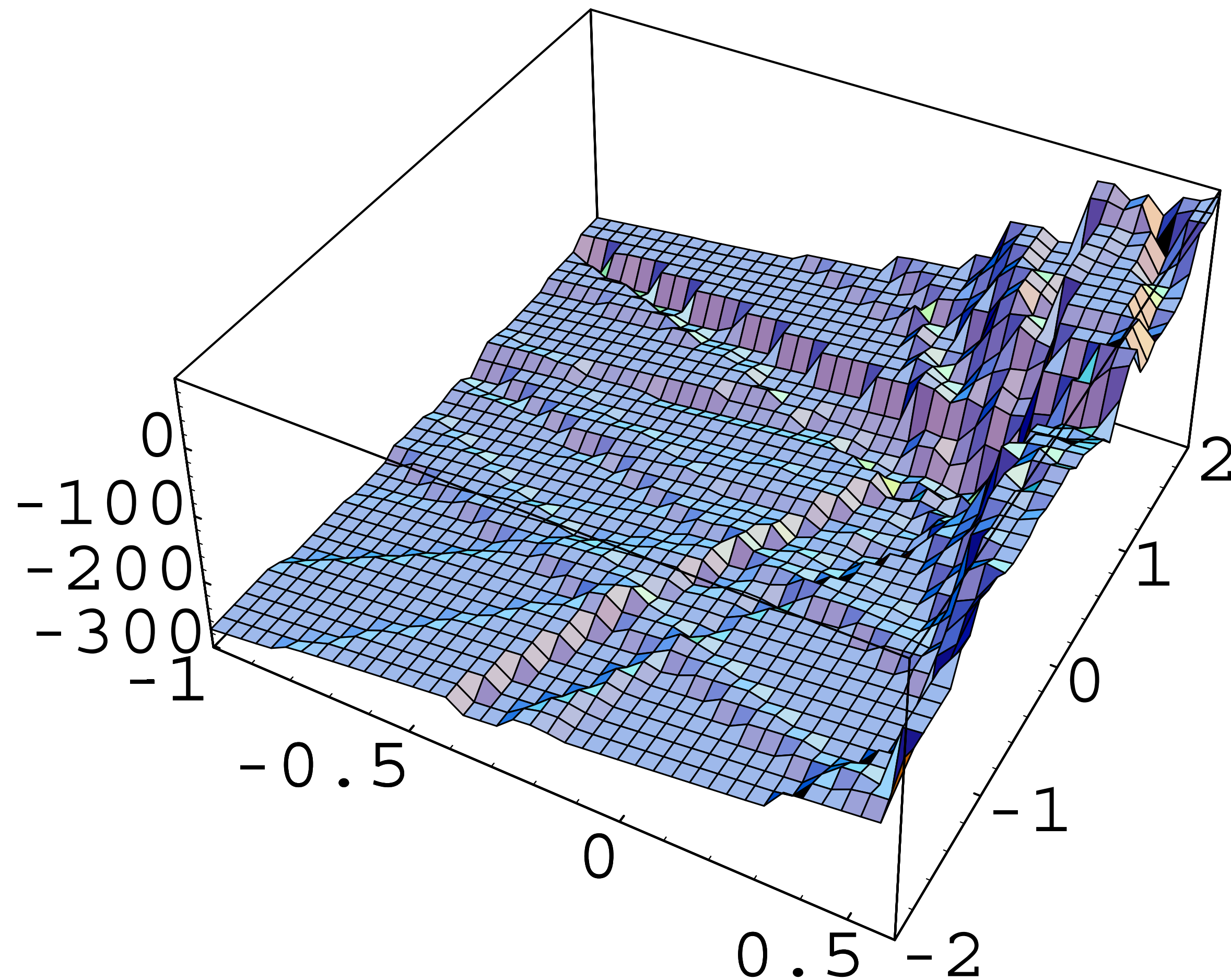# What happens when we run value iteration with a *2 Layer MLP?*



Iteration 101

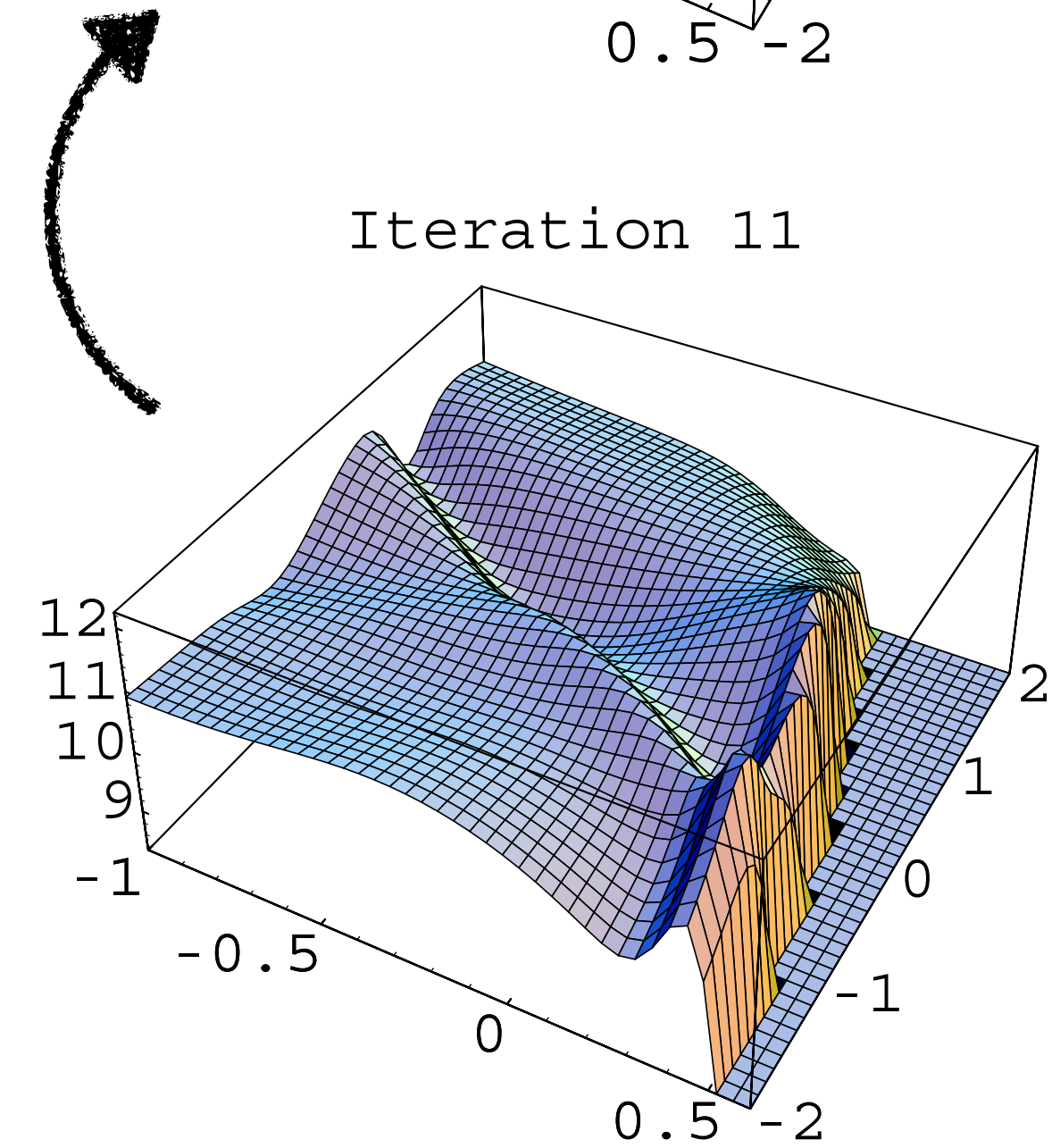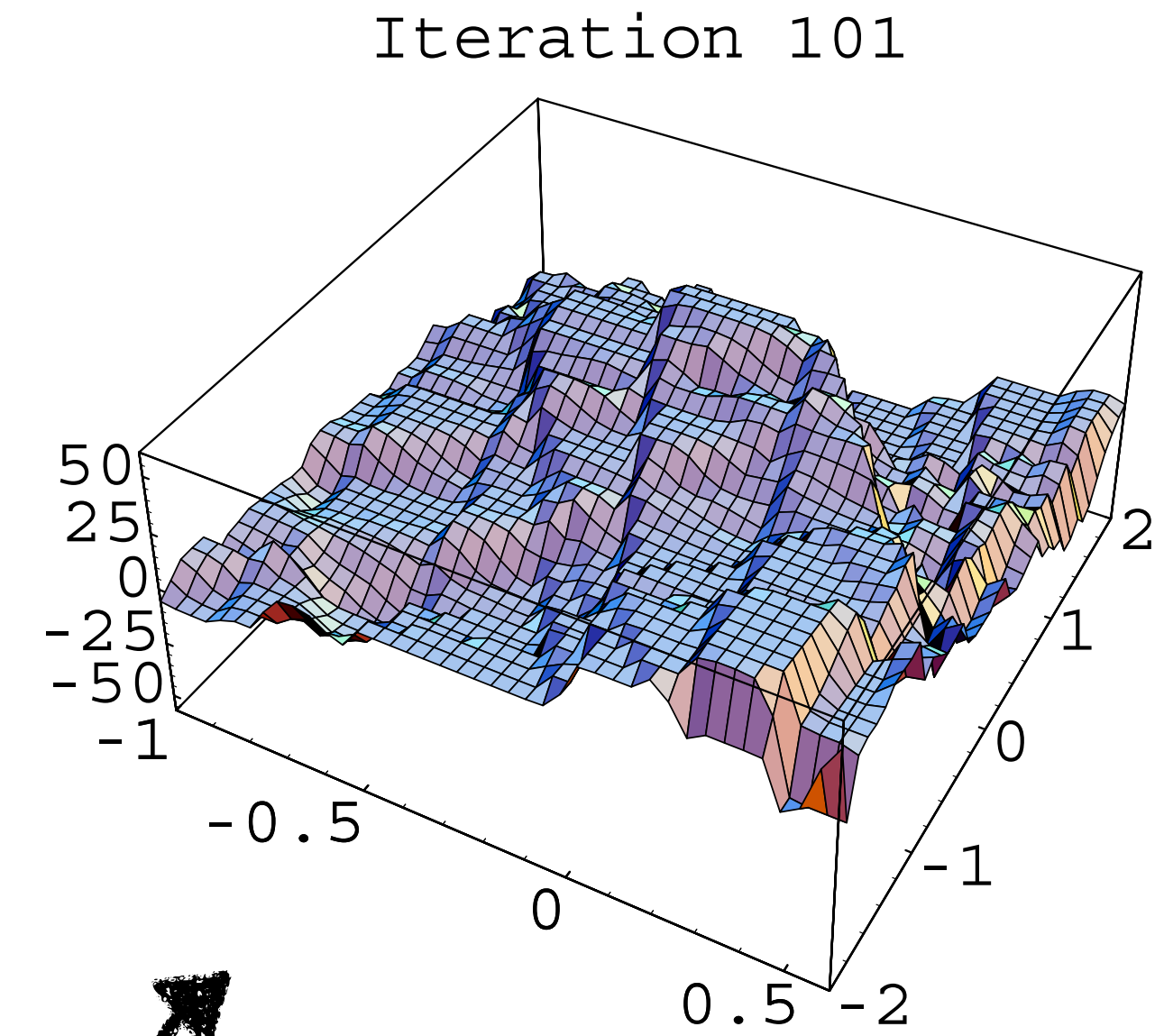# What happens when we run value iteration with a
## *2 Layer MLP?*

# The problem of Bootstrapping!

max()

pos

Iteration 101

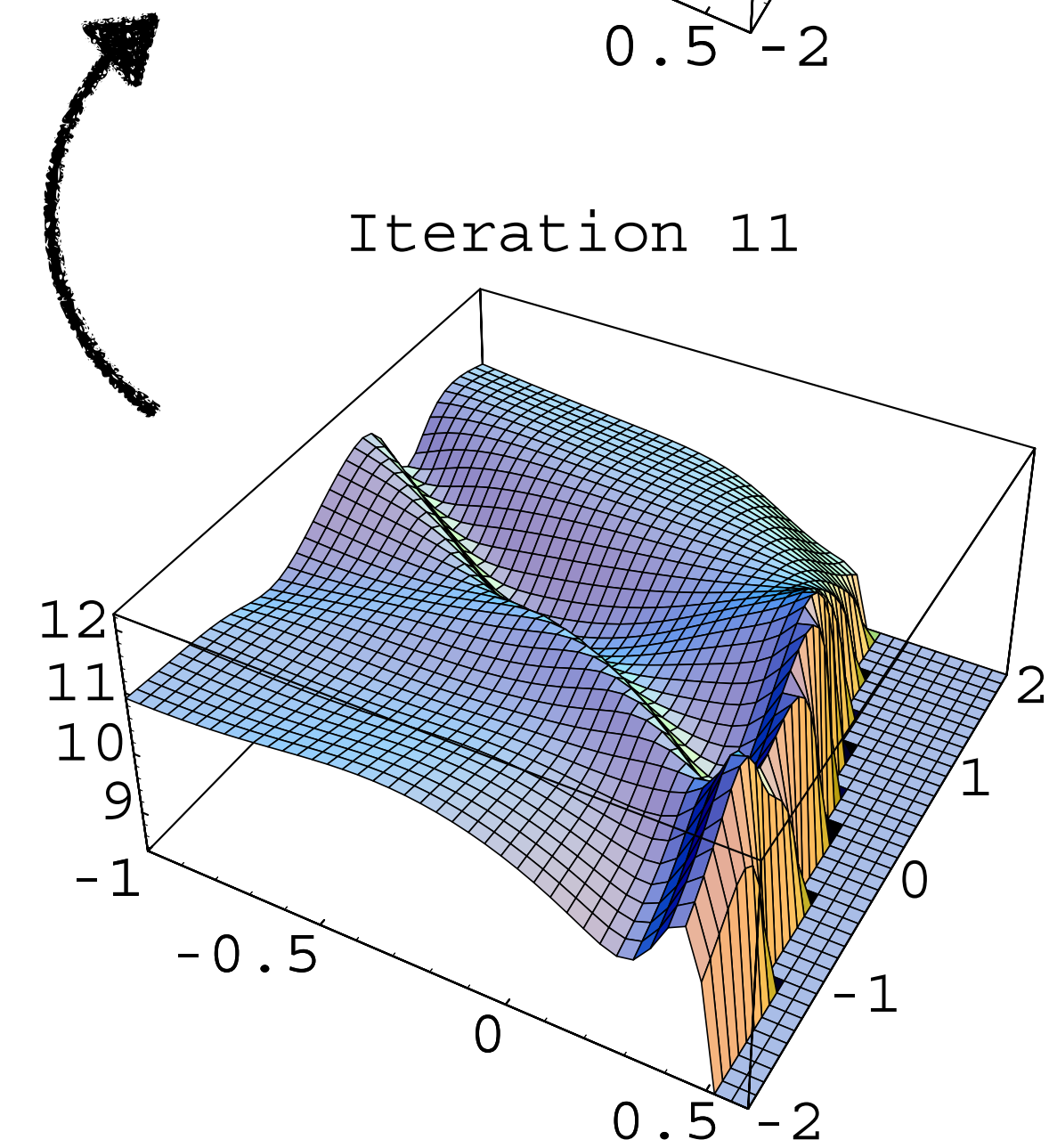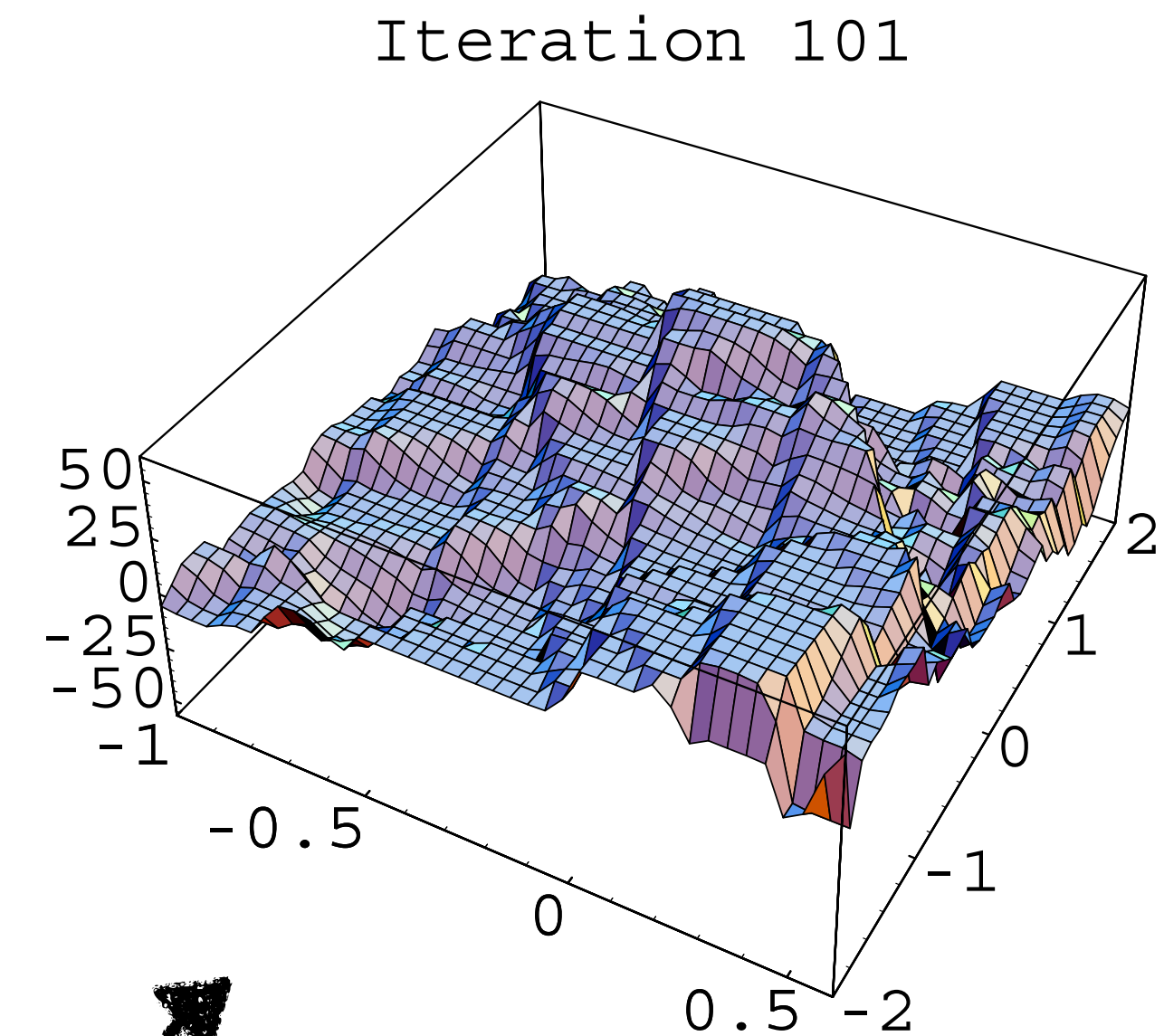Iteration 11

# 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

-1    -0.5    0    0.5

pos

50
25
0
-25
-50
-1
-0.5
0
0.5

2
1
0
-1
-2

-

0.5

min()

Iteration 11

12
11
10
9
-1
-0.5
0
0.5

2
1
0
-1
-2

-

What about policy
iteration?

# Policy Iteration

## Policy Evaluation

## Policy Improvement

# Policy Iteration

Init with some policy $\pi$

Repeat forever

    Evaluate policy

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

    Improve policy

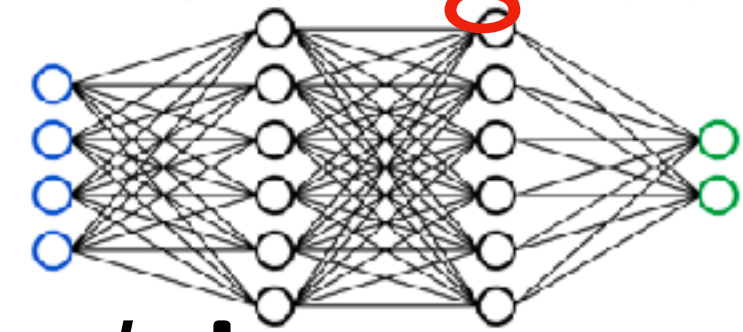$$\pi^{+}(s) = \arg\min_{a} Q^{\pi}(s, a)$$

# Fitted Policy Iteration

## *Fitted* policy evaluation

<span style="color:red">Collect data using current policy $\pi$</span>

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

        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** $Q_\theta$

## Policy Improvement

# This remains the same!

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

# Fitted Policy Iteration

*Fitted* policy evaluation                     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**

# This is fine..

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

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

<span style="color:red">No min()</span>

target $\leftarrow c_i + \gamma Q_\theta(s'_i, \pi(s'_i))$ <span style="color:red">step</span>

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

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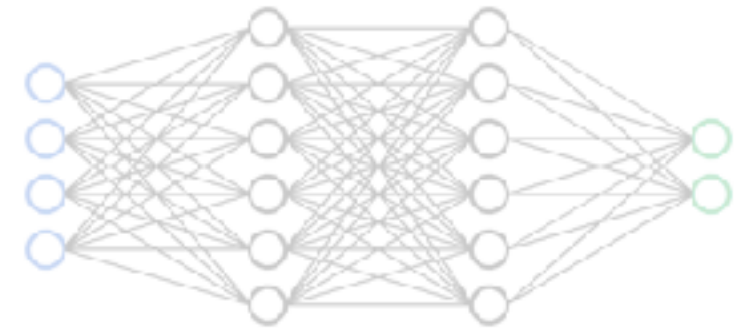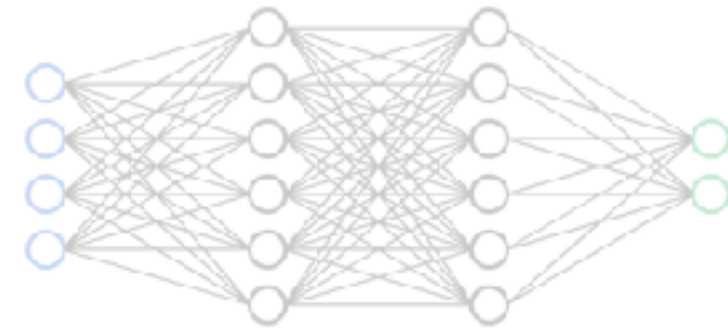
**return** $Q_\theta$

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

# Fitted Policy Iteration

## Fitted policy evaluation

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**

$D \leftarrow \emptyset$

This is fine..

But this has
the min() step!

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

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

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

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

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

**return** $Q_\theta$

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