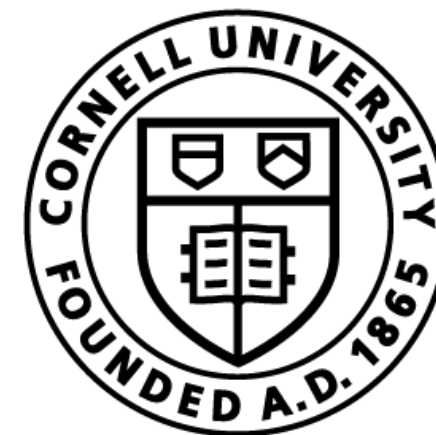


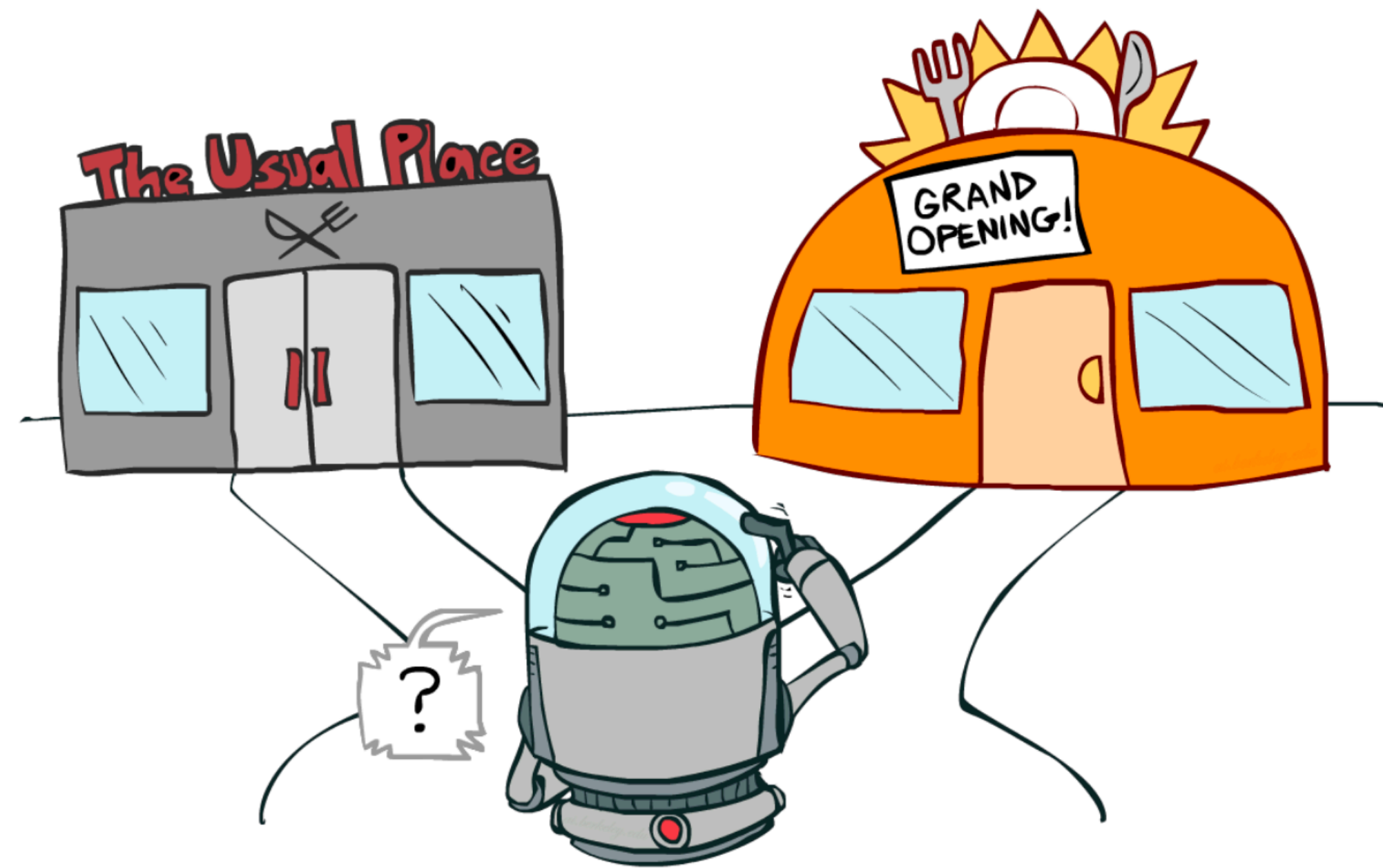
Approximate Dynamic Programming

Sanjiban Choudhury

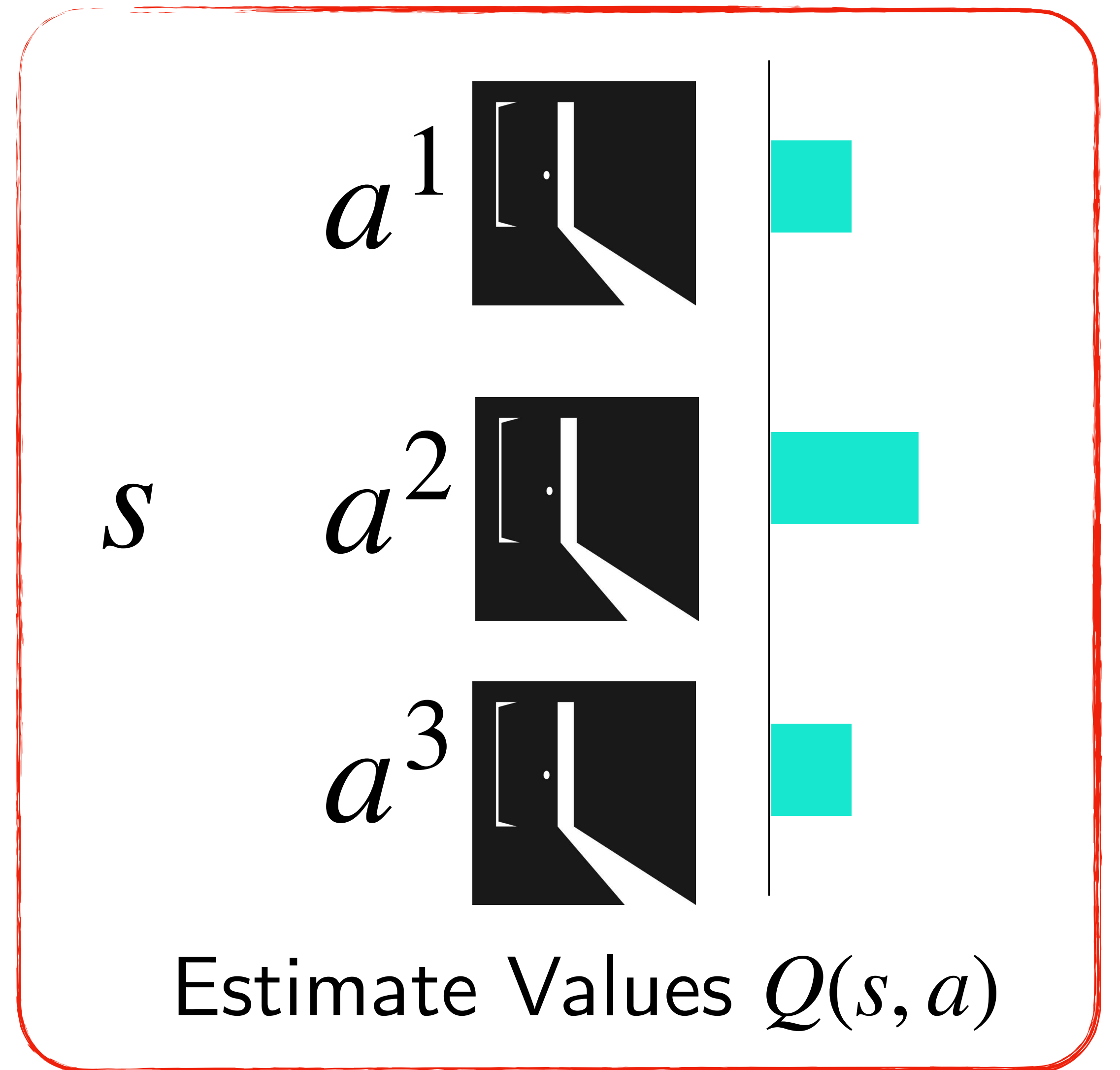


Cornell Bowers CIS
Computer Science

Recap: Two Ingredients of RL

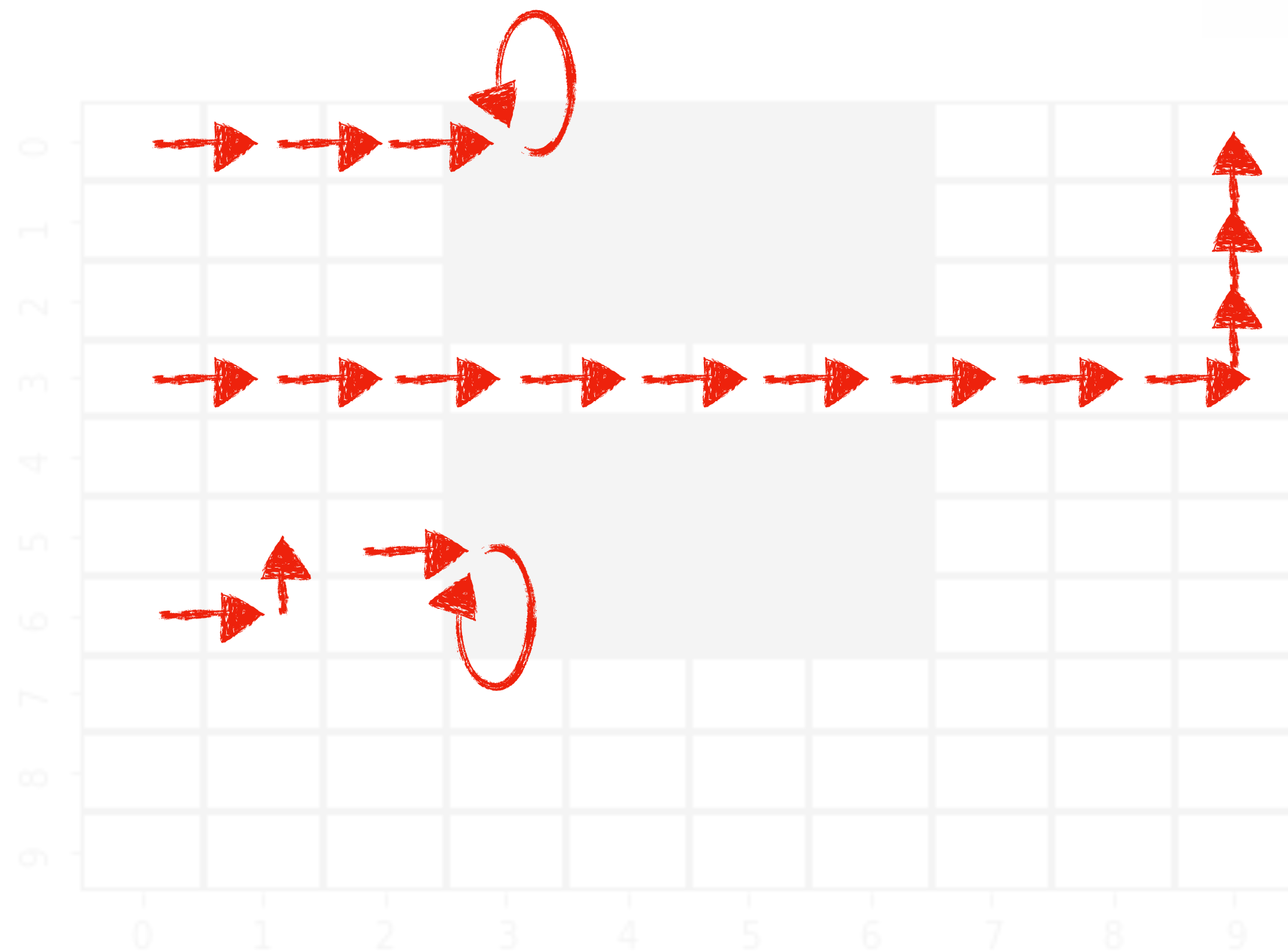


Exploration Exploitation

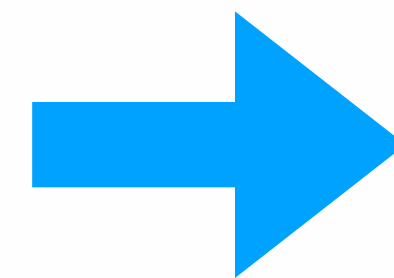


Estimate Values $Q(s, a)$

Estimate the value of policy from sample rollouts

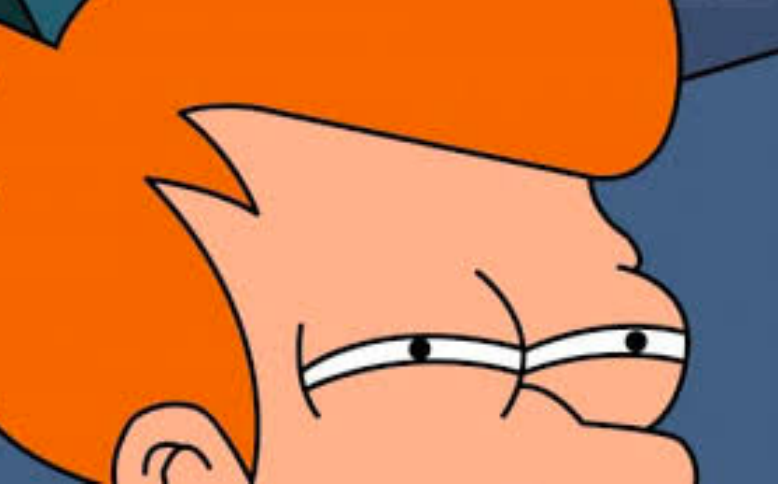


Roll outs



0	74	75	76	77	77	77	77	2	1	0
1	74	75	76	77	77	77	77	3	2	1
2	74	75	76	77	77	77	77	3.9	3	2
3	55	56	56	57	50	40	26	4.9	3.9	3
4	74	75	76	77	77	77	77	5.9	4.9	3.9
5	74	75	76	77	77	77	77	6.8	5.9	4.9
6	74	75	76	77	77	77	77	7.7	6.8	5.9
7	15	14	13	12	11	10	9.6	8.6	7.7	6.8
8	16	15	14	13	12	11	10	9.6	8.6	7.7
9	17	16	15	14	13	12	11	10	9.6	8.6
	0	1	2	3	4	5	6	7	8	9

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

Last lecture, we looked at
tabular values ...

Is it trivial to upgrade to a
Neural Network?



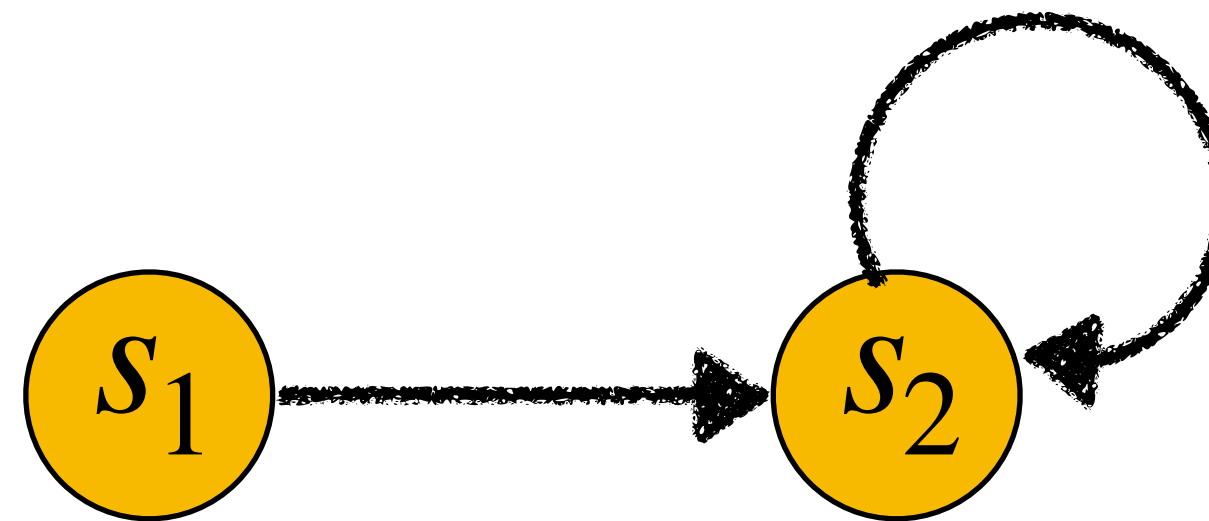
Activity!



A tiny MDP

Reward for being at any state is 0.0

Discount factor $\gamma = 0.9$



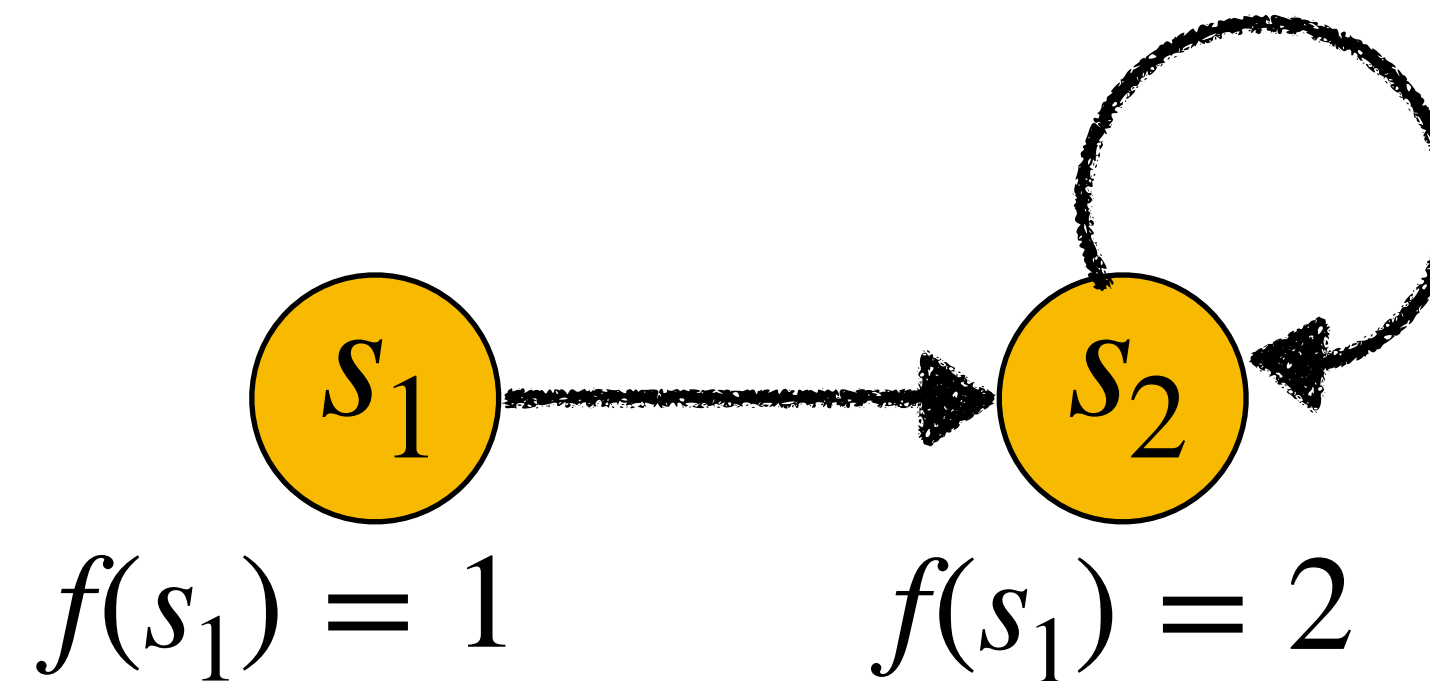
What happens when you run value iteration?

(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 $\gamma = 0.9$



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

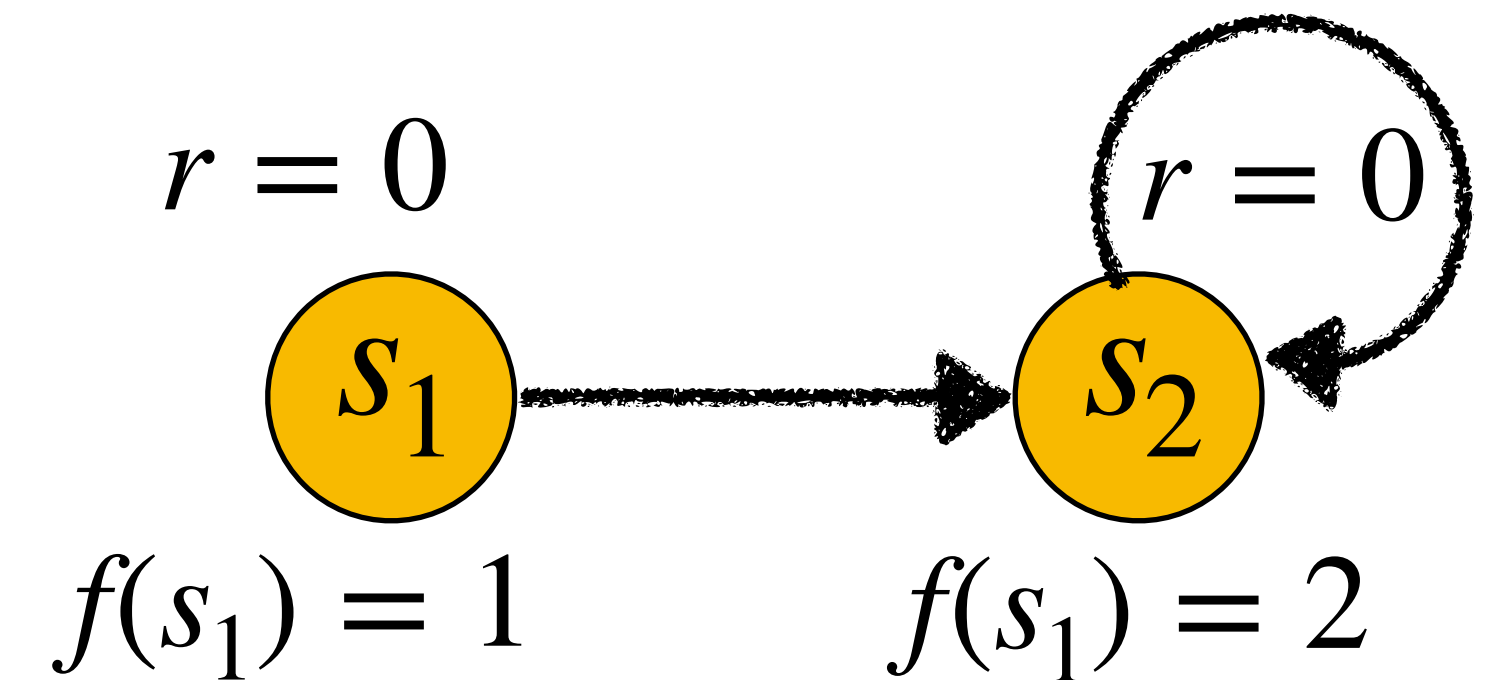
$$V(s) = wf(s) = w * \begin{cases} 1 & \text{if } s = s_1 \\ 2 & \text{if } s = s_2 \end{cases}$$

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

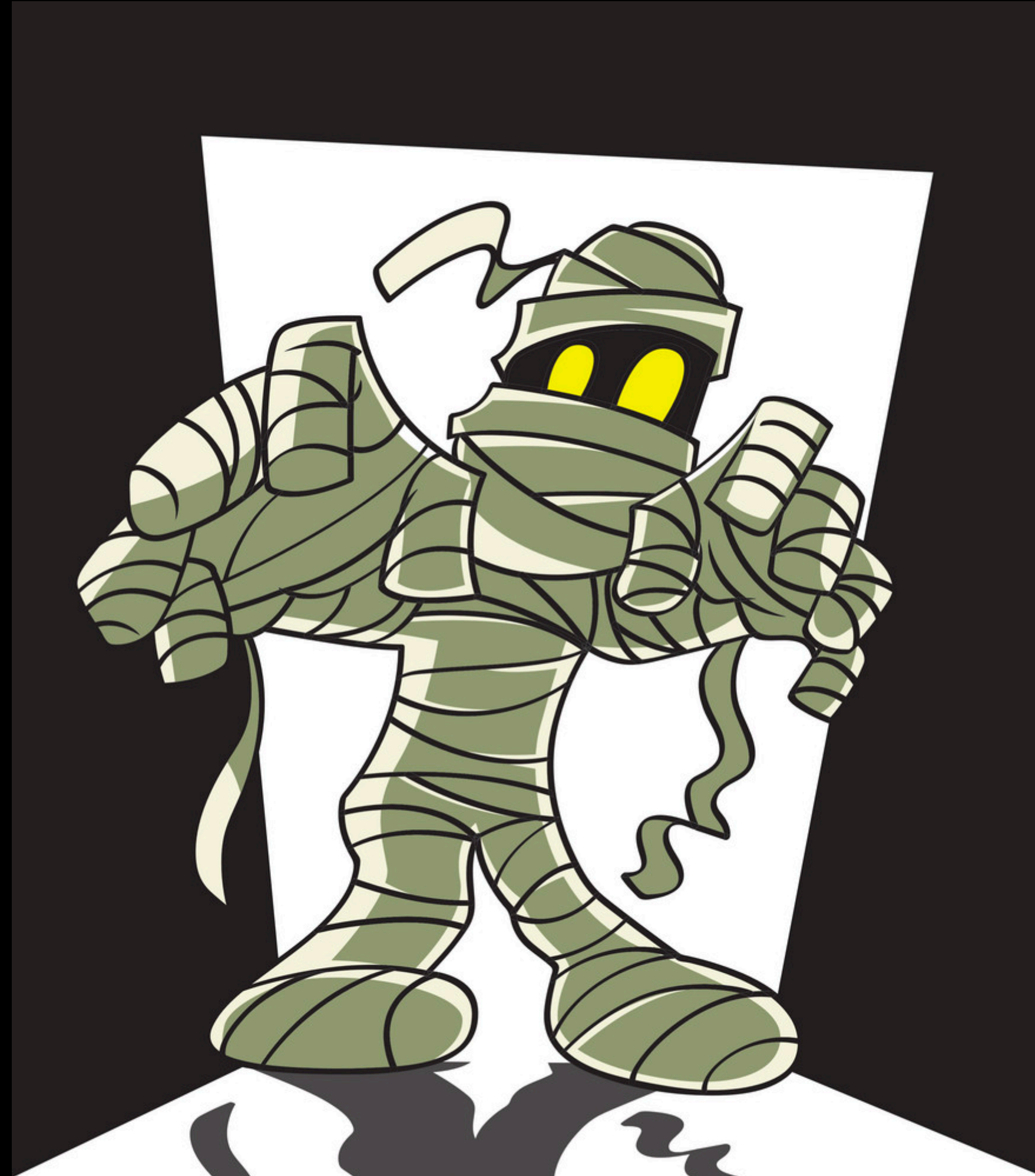


Share (45 sec): Partners exchange ideas

$$V(s) = wf(s)$$

Init with $w = 1$

CURSE OF APPROXIMATION!

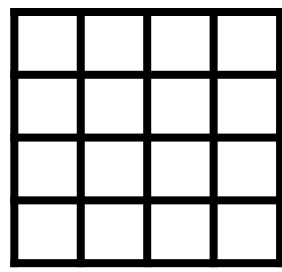


From dynamic
programming to
Fitted
dynamic
programming



Approximate (Fitted) Value Iteration

Q-iteration



$$Q(s, a) \leftarrow 0$$

while *not converged* **do**

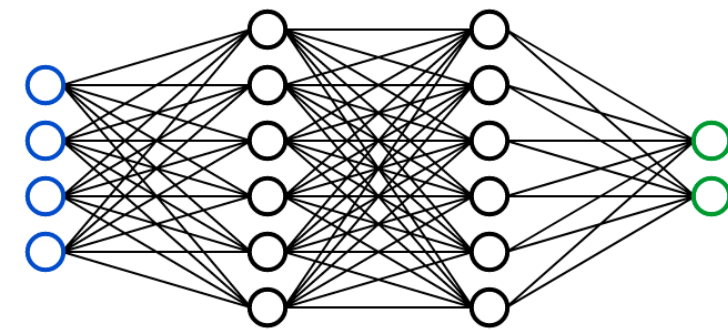
for $s \in S, a \in A$

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

$$Q \leftarrow Q^{new}$$

return Q

Fitted Q-iteration



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

$$\text{Init } Q_\theta(s, a) \leftarrow 0$$

while *not converged* **do**

$$D \leftarrow \emptyset$$

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

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

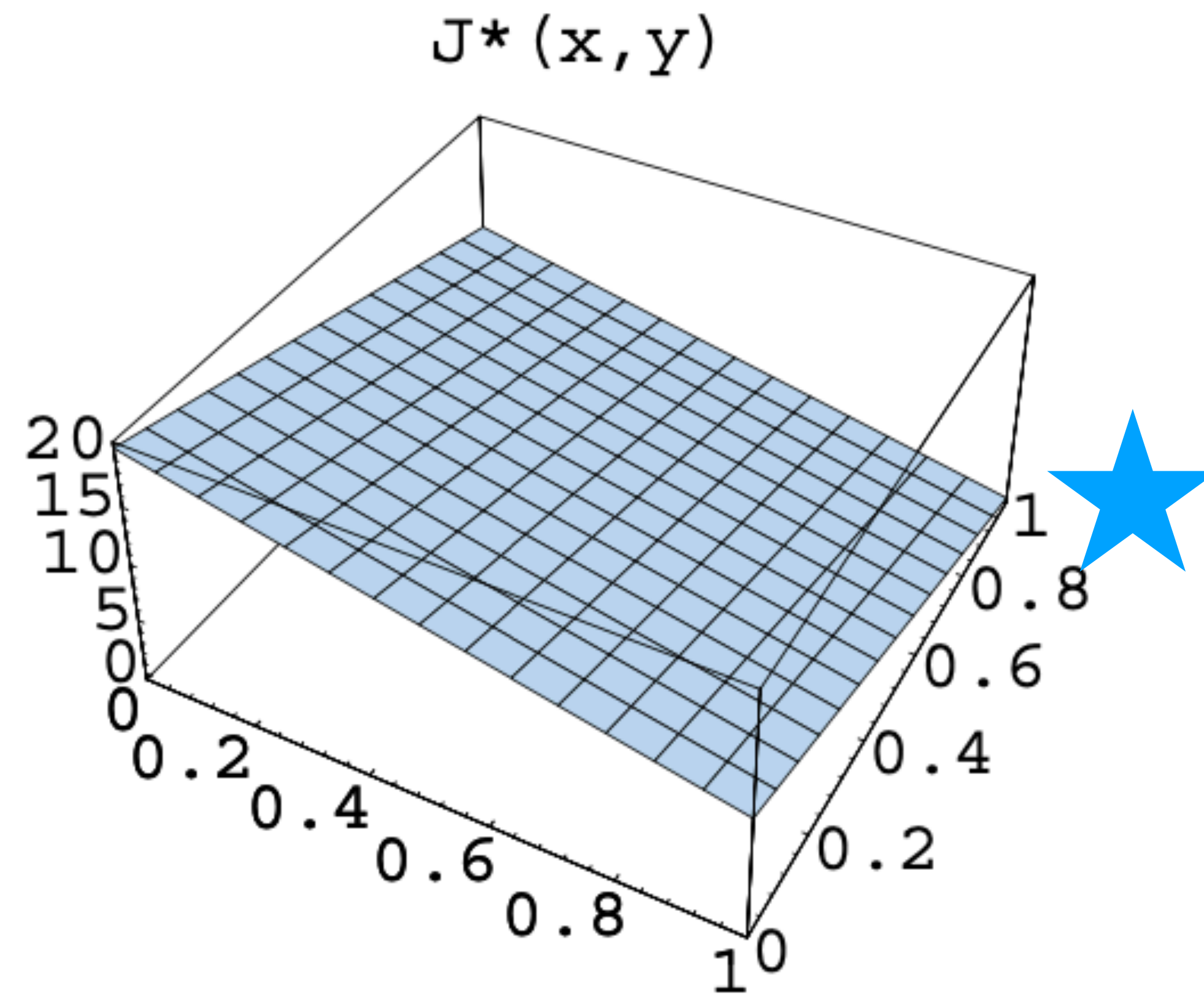
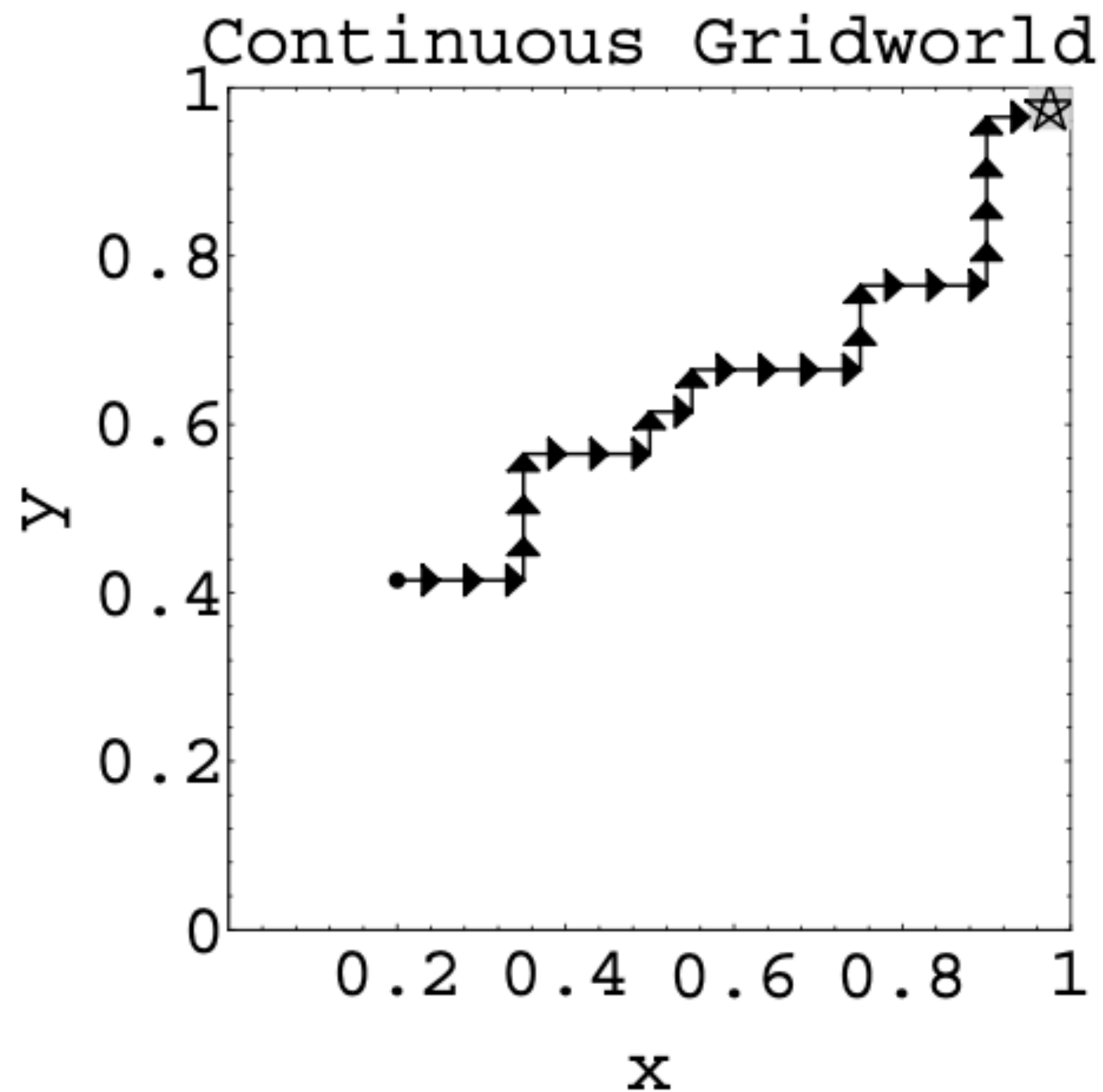
$$\text{target} \leftarrow c_i + \gamma \min_{a'} Q_\theta(s'_i, a')$$

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

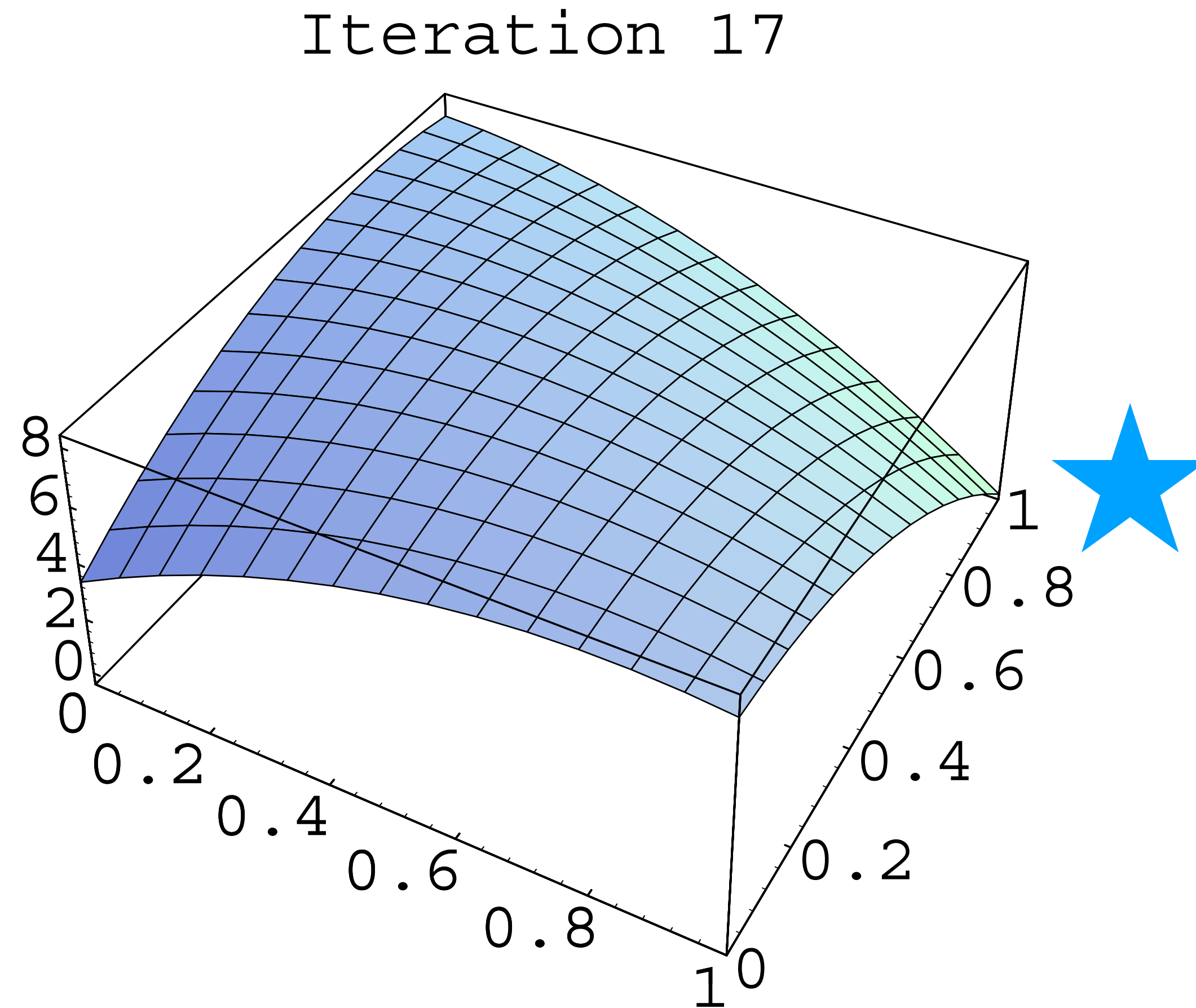
$$Q_\theta \leftarrow \text{Train}(D)$$

return Q_θ

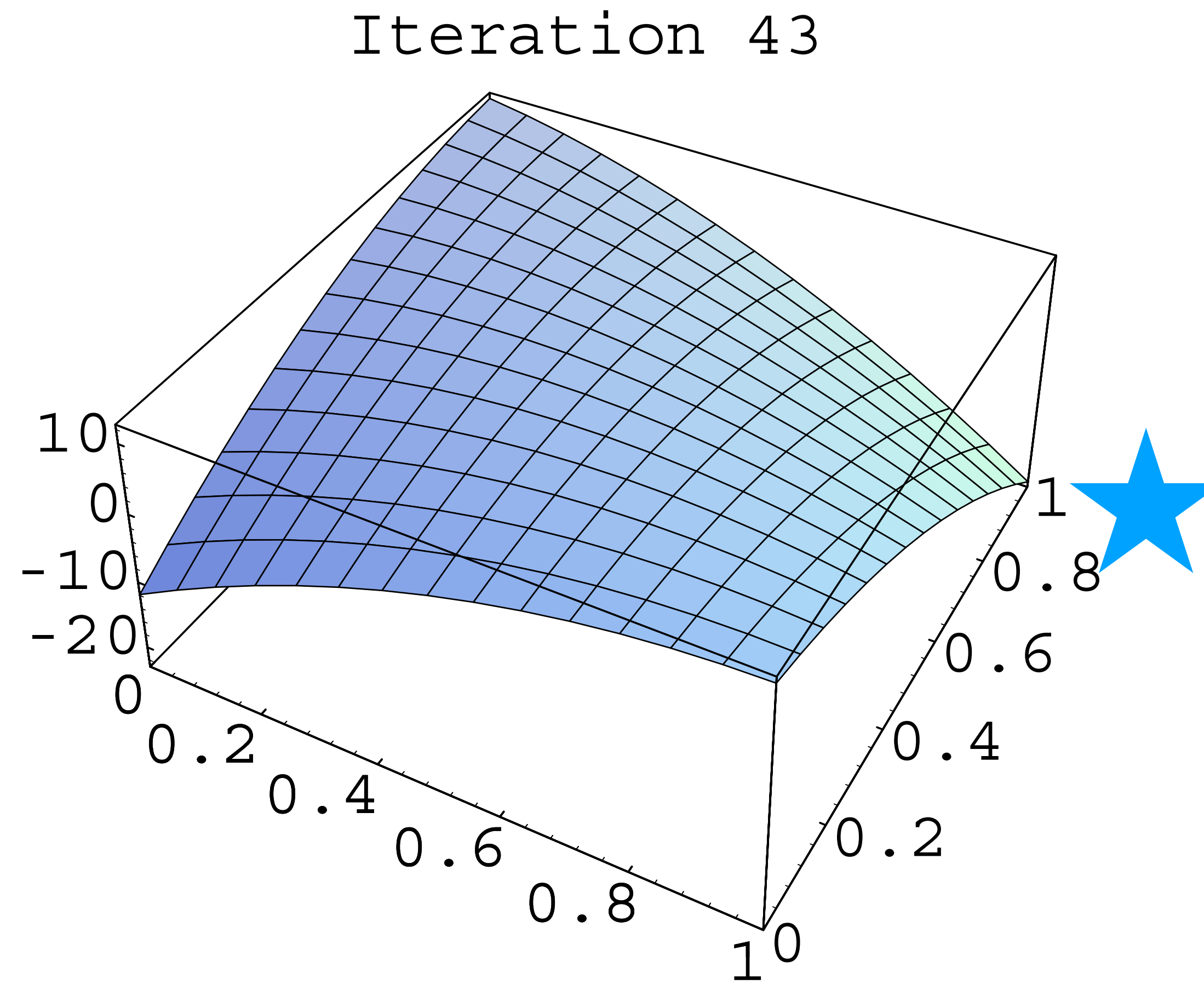
A simple example: Gridworld



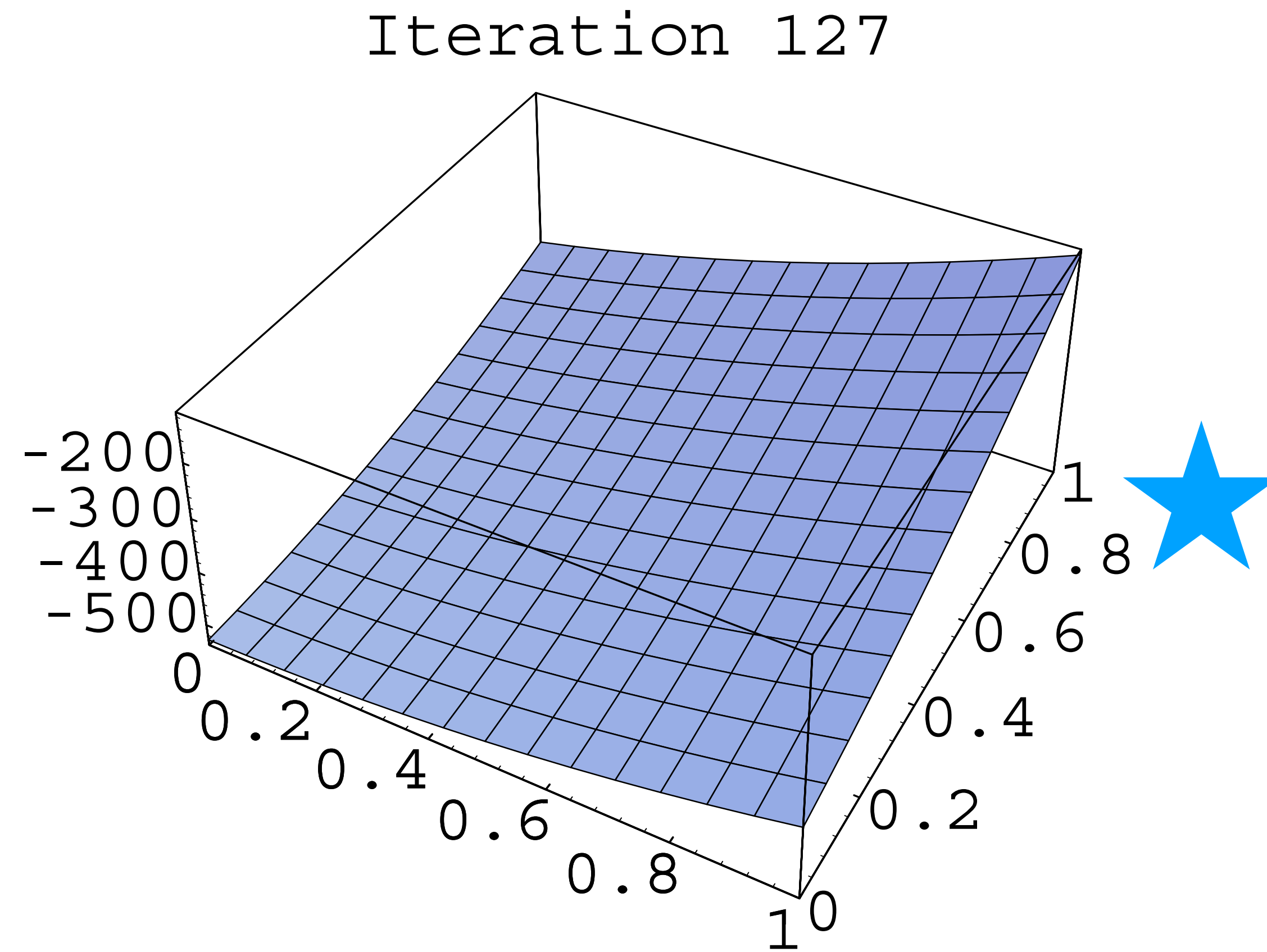
What happens when we run value iteration with a
quadratic?



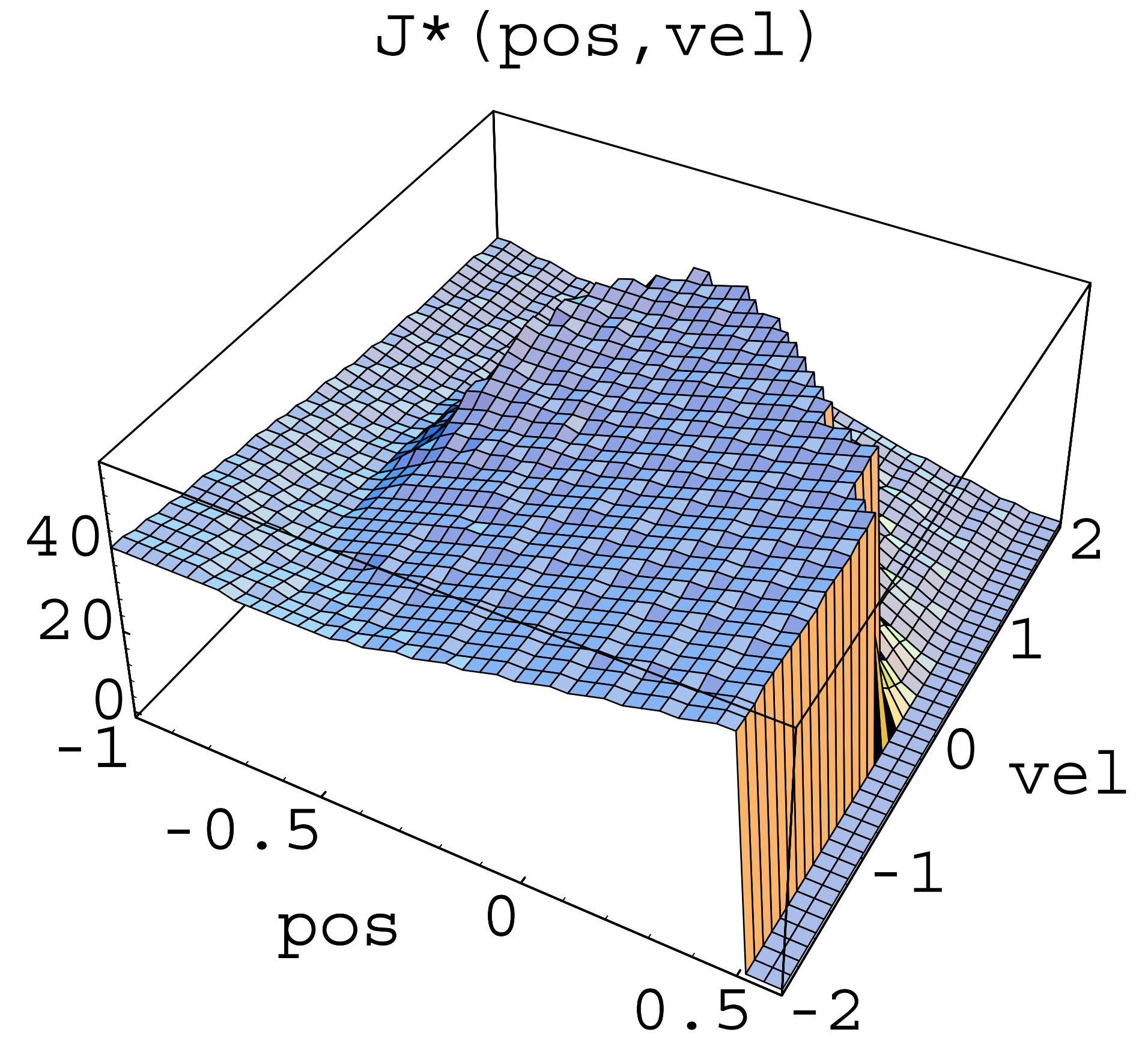
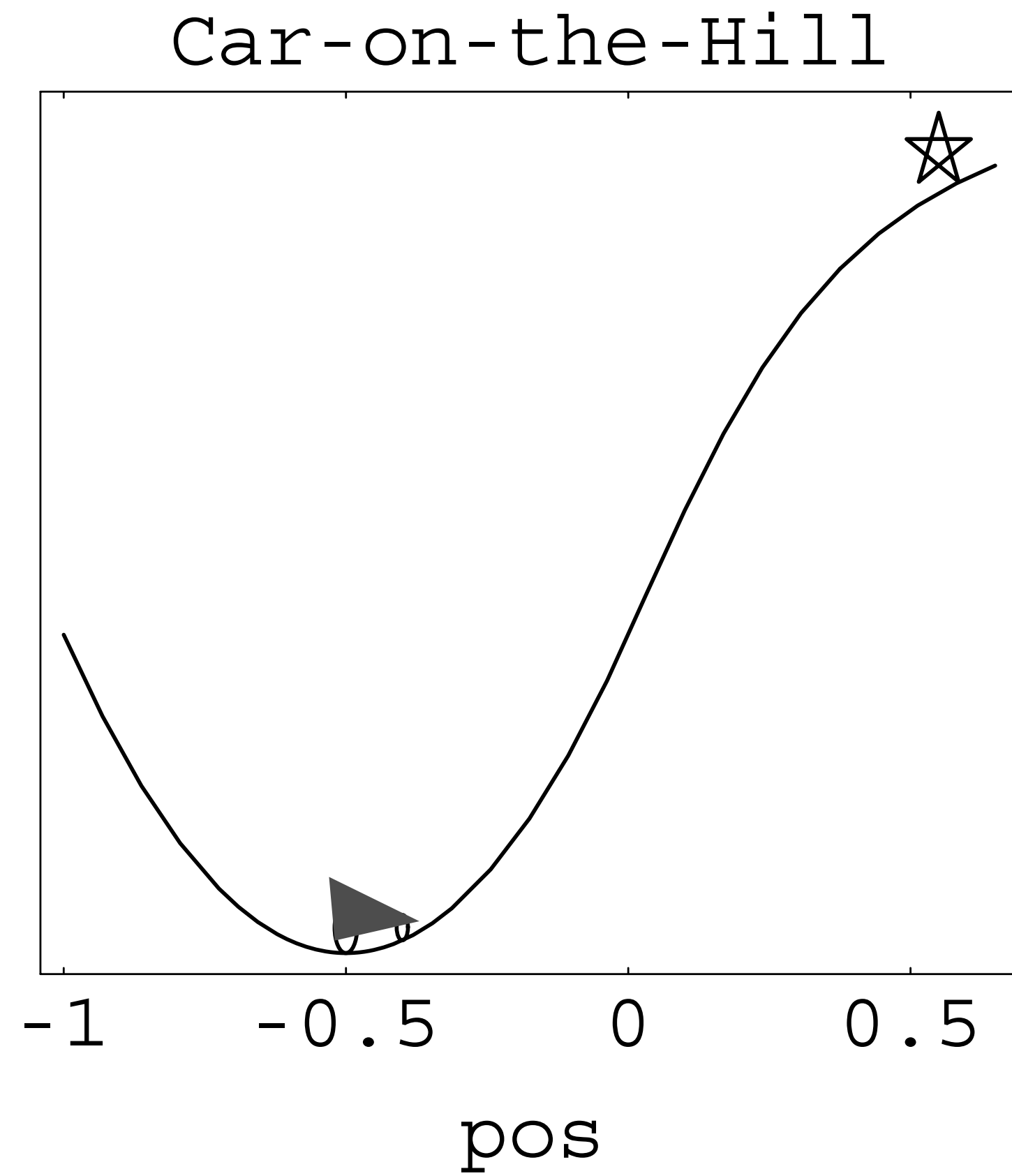
What happens when we run value iteration with a
quadratic?



What happens when we run value iteration with a
quadratic?

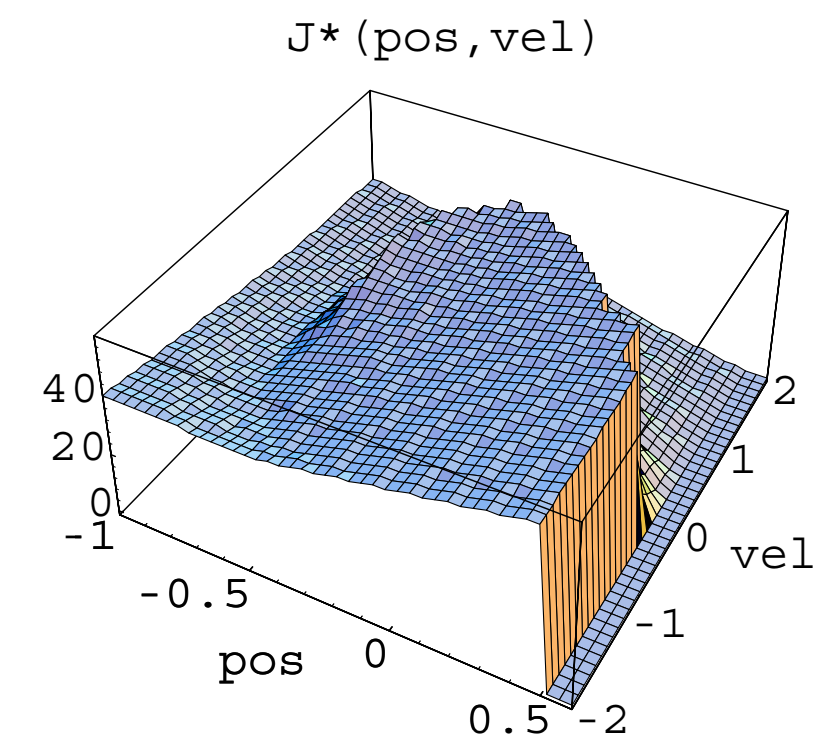
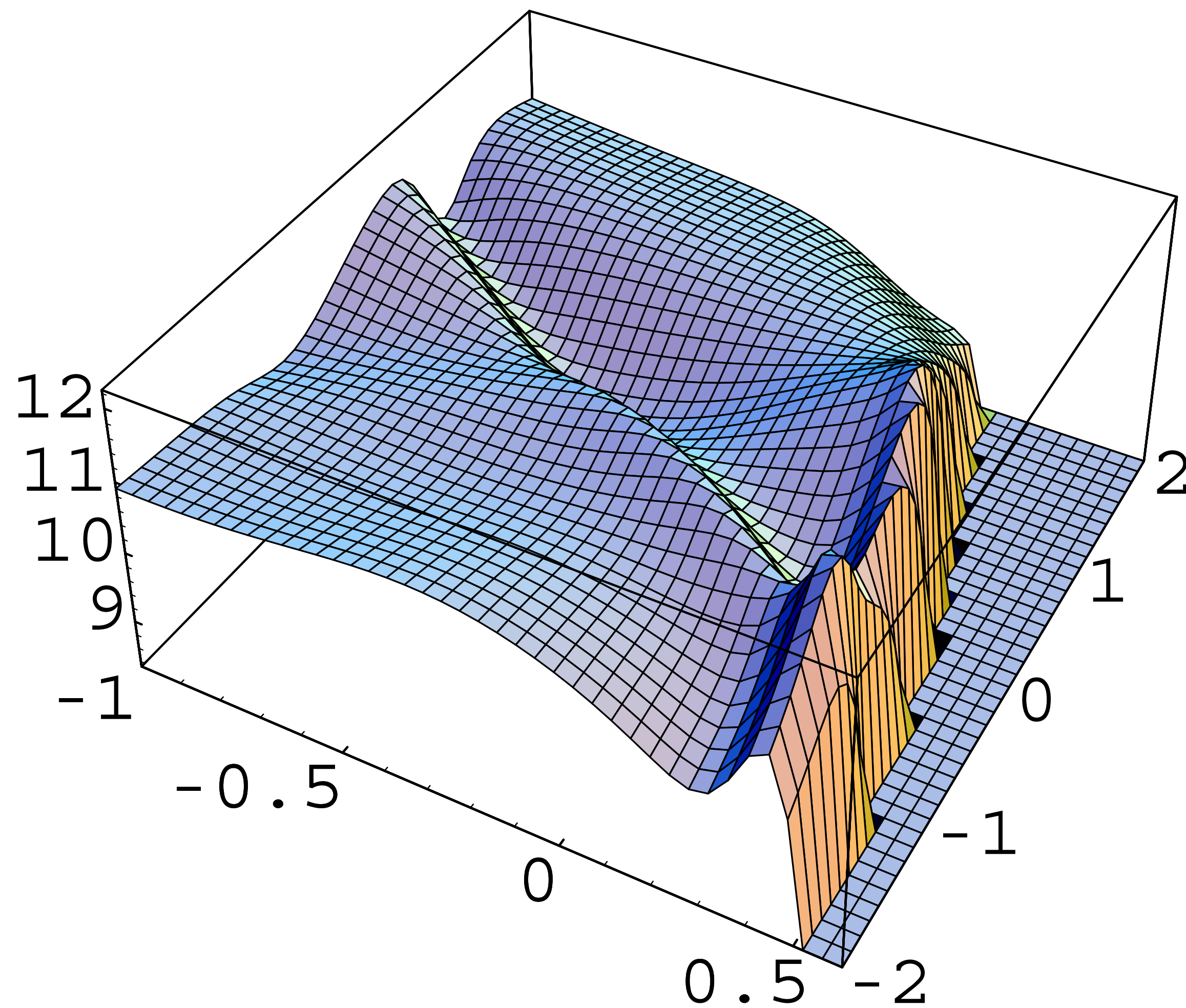


Another Example: Mountain Car!



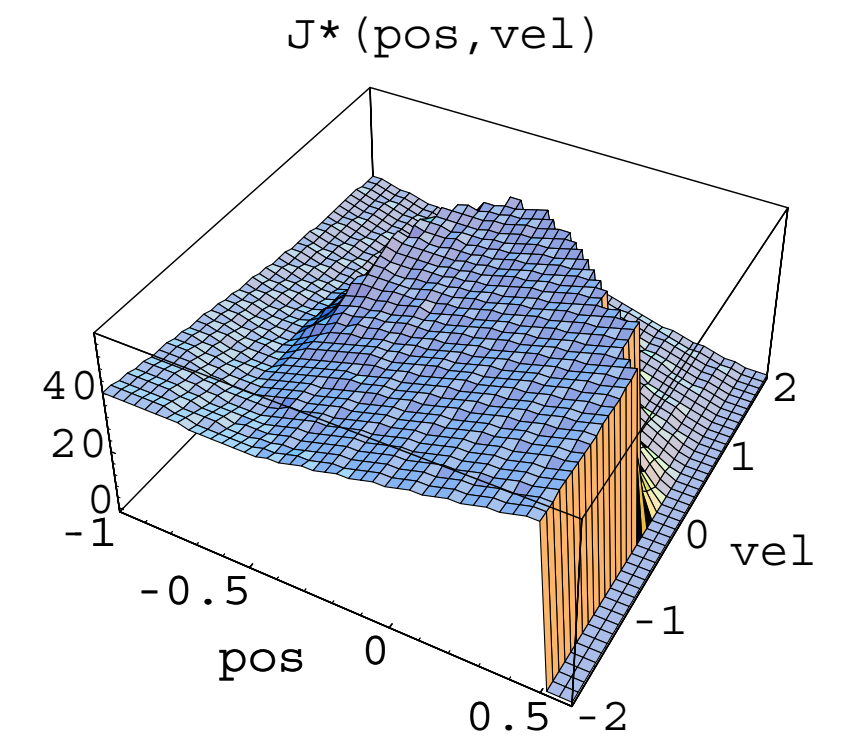
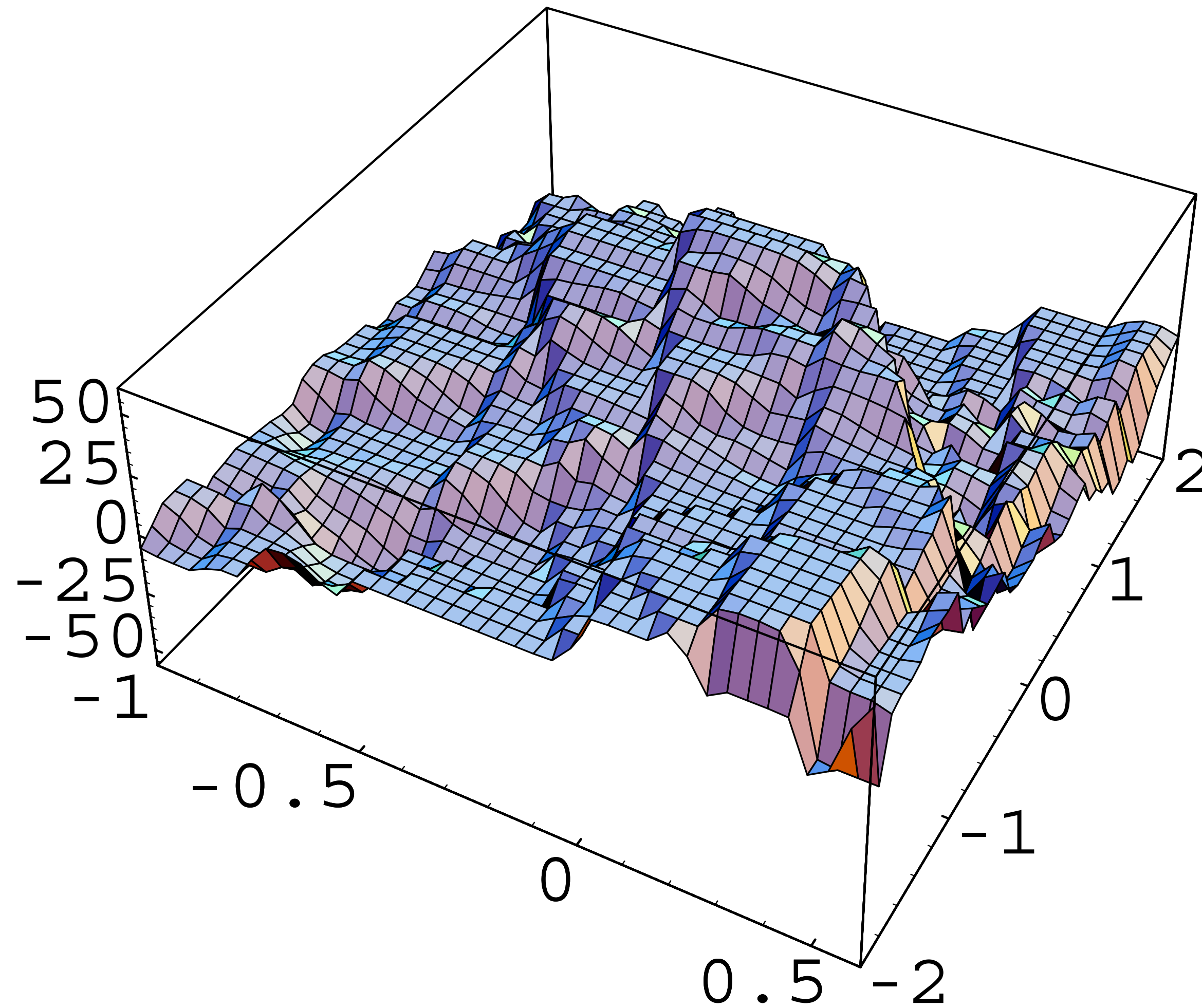
What happens when we run value iteration with a *2 Layer MLP?*

Iteration 11



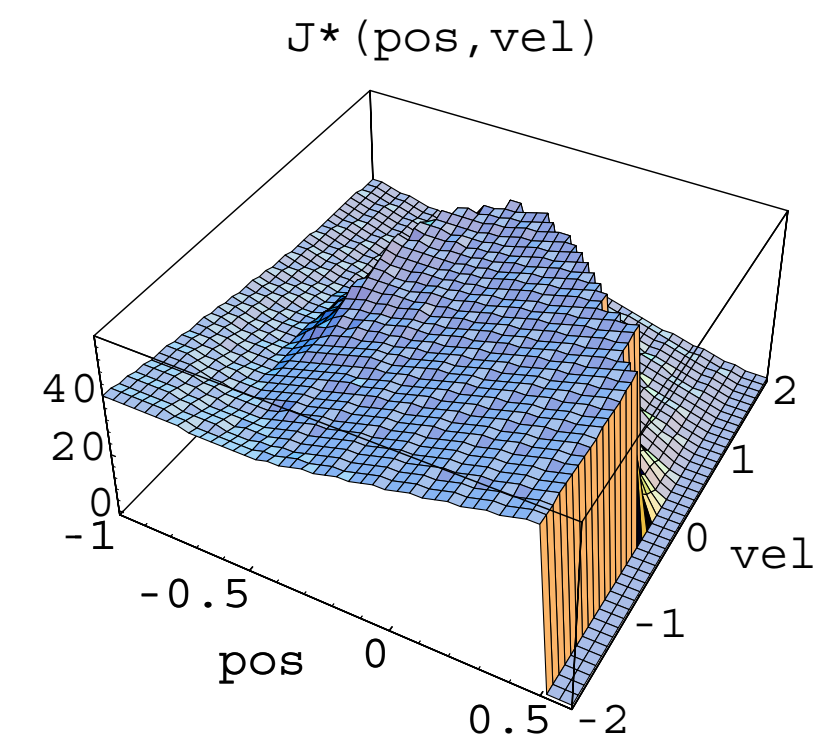
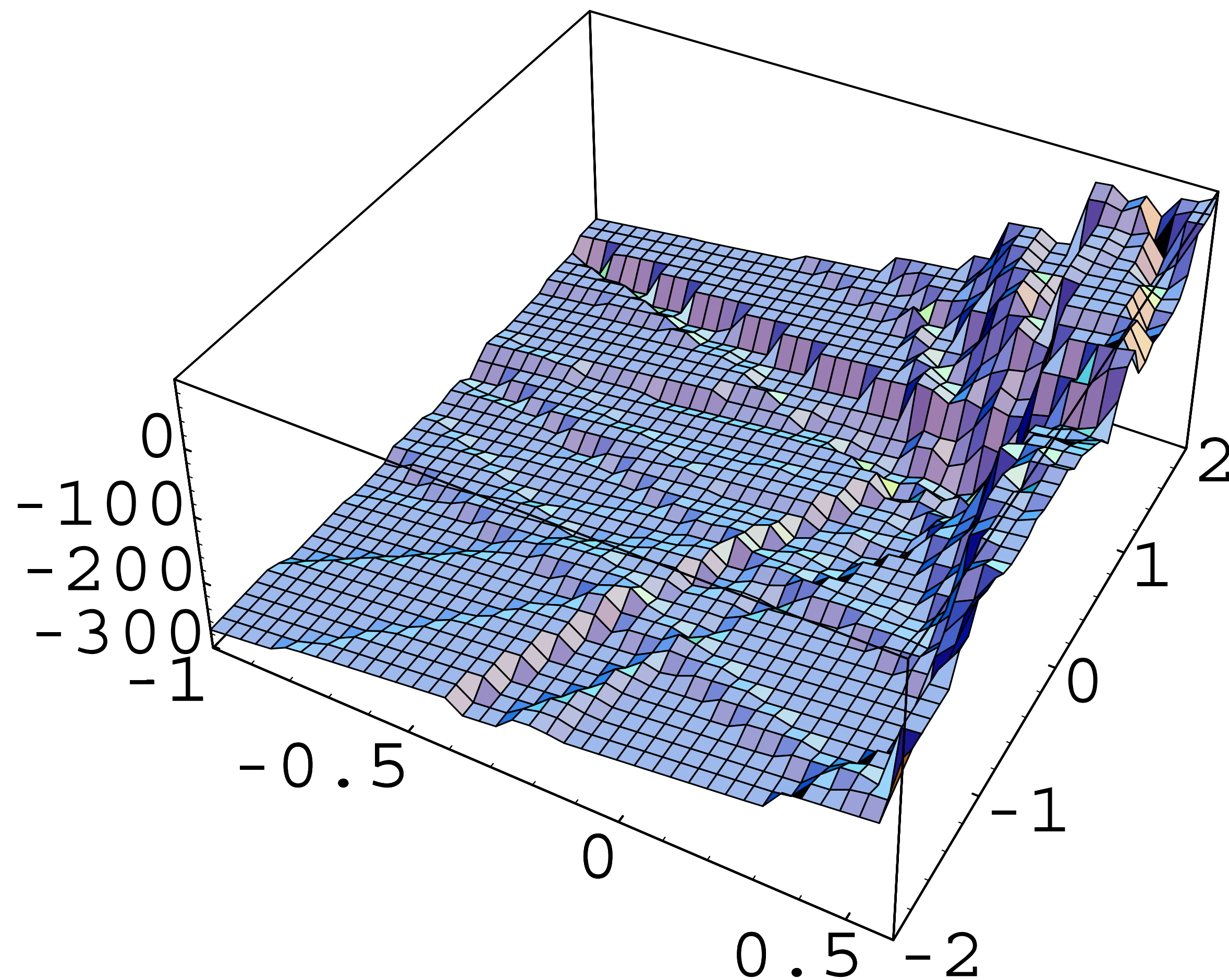
What happens when we run value iteration with a *2 Layer MLP?*

Iteration 101



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

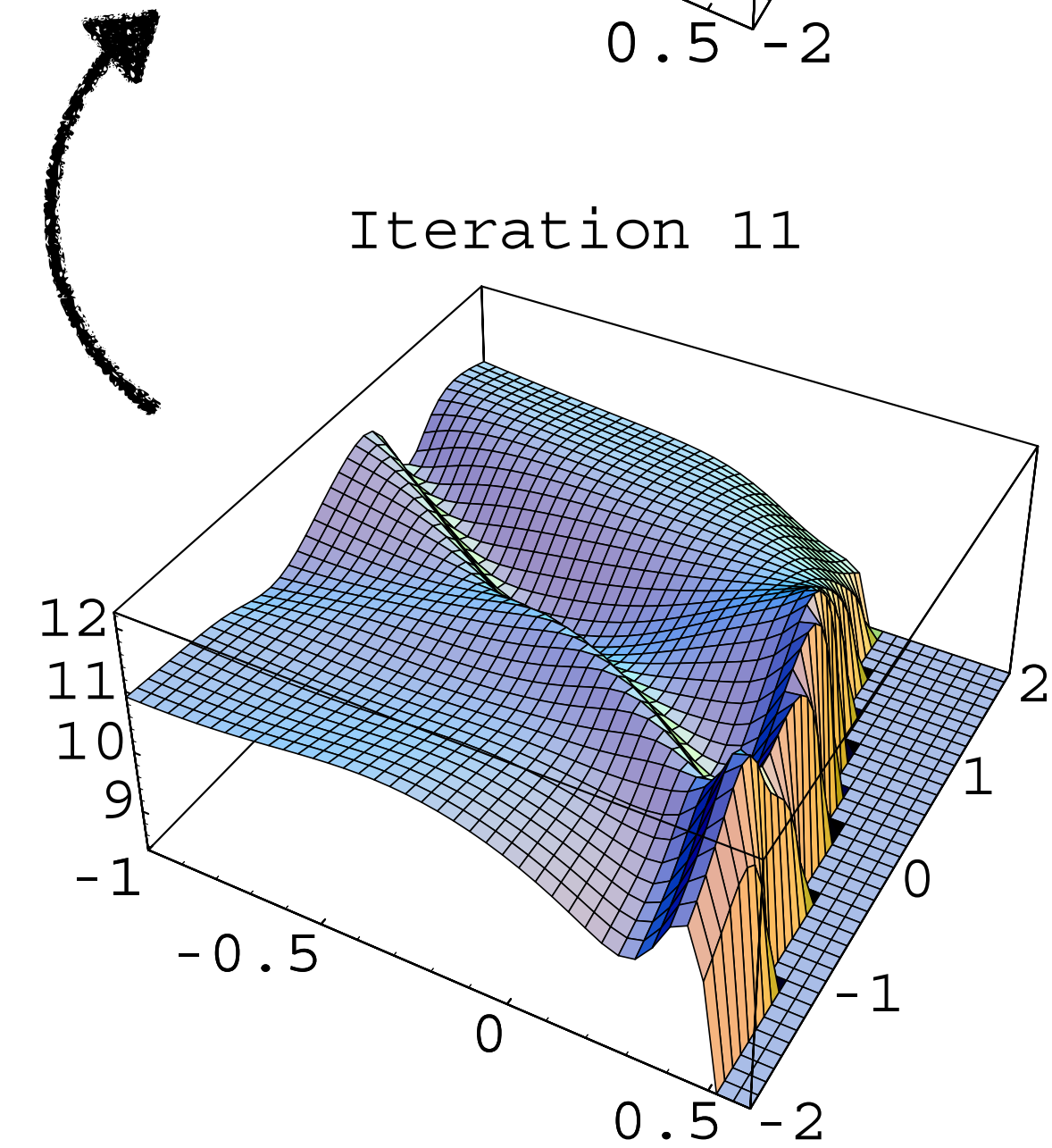
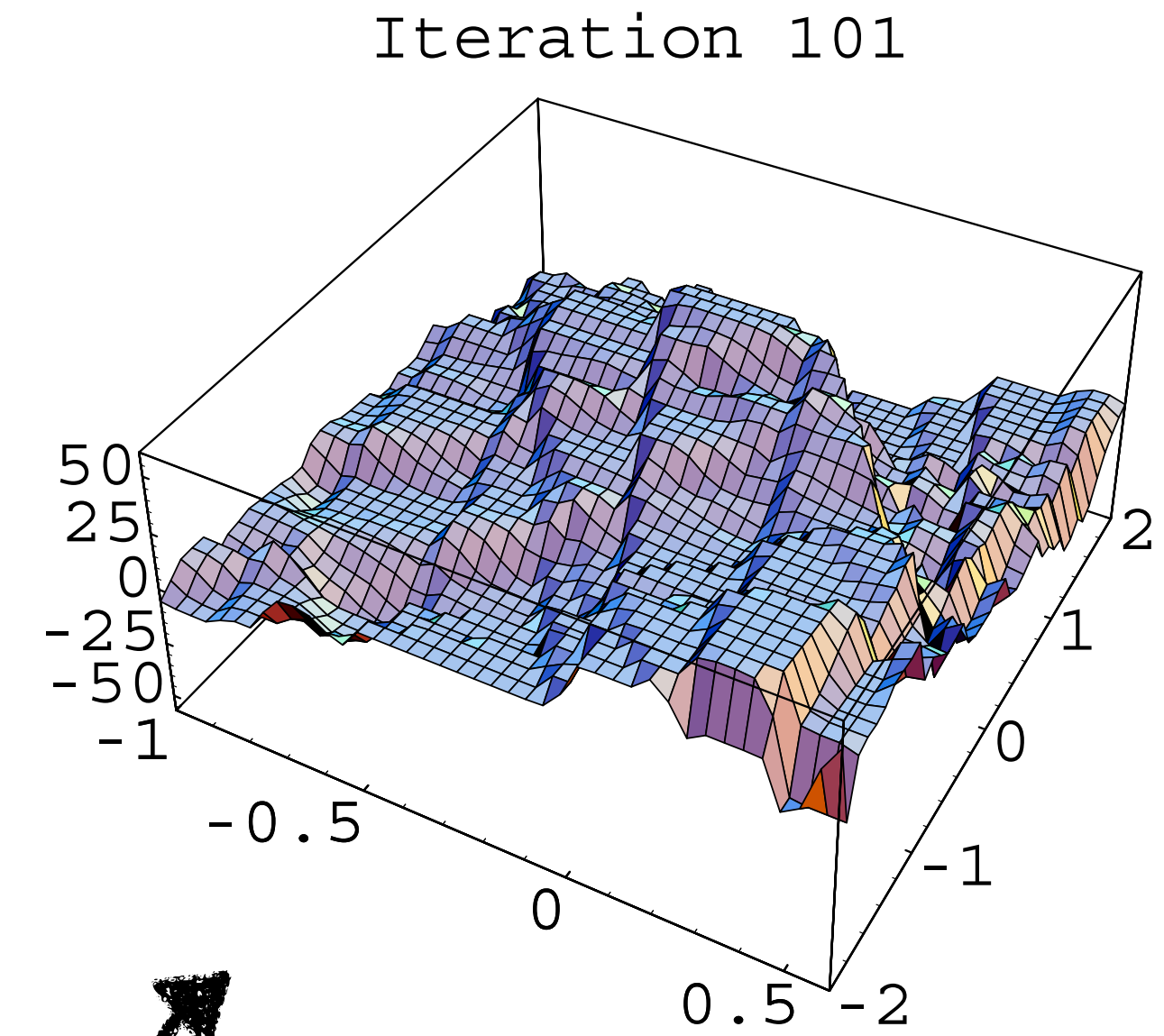
Iteration 201



The problem of Bootstrapping!



`max()`



The problem of Bootstrapping!

- Since these methods rely on approximating the value function inductively, errors in approximation are propagated, and, even worse, amplified as the algorithm encourages actions that lead to states with sub-optimal values.
- The key reason behind this is the minimization operation performed when generating the target value used for the action value function. The minimization operation pushes the desired policy to visit states where the value function approximation is less than the true value of that state – that is to say, states that look more attractive than they should. This typically happens in areas of state spaces which are very few training samples and could, in fact, be quite bad places to arrive.

What about policy iteration?



Policy Iteration

Policy Evaluation

0	-	0	0	0	0	0	0	0	0	0	0
1	-	0	0	0	0	0	0	0	0	0	0
2	-	0	0	0	0	0	0	0	0	0	0
3	-	0	0	0	0	0	0	0	0	0	0
4	-	0	0	0	0	0	0	0	0	0	0
5	-	0	0	0	0	0	0	0	0	0	0
6	-	0	0	0	0	0	0	0	0	0	0
7	-	0	0	0	0	0	0	0	0	0	0
8	-	0	0	0	0	0	0	0	0	0	0
9	-	0	0	0	0	0	0	0	0	0	0
		0	1	2	3	4	5	6	7	8	9

Iter: 0

Policy Improvement

0	-	→	→	→	→	→	→	→	→	→	↑
1	-	→	→	→	→	→	→	→	→	→	↑
2	-	→	→	→	→	→	→	→	→	→	↑
3	-	→	→	→	→	→	→	→	→	→	↑
4	-	→	→	→	→	→	→	→	→	→	↑
5	-	→	→	→	→	→	→	→	→	→	↑
6	-	→	→	→	→	→	→	→	→	→	↑
7	-	→	→	→	→	→	→	→	→	→	↑
8	-	→	→	→	→	→	→	→	→	→	↑
9	-	→	→	→	→	→	→	→	→	→	↑
		0	1	2	3	4	5	6	7	8	9

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

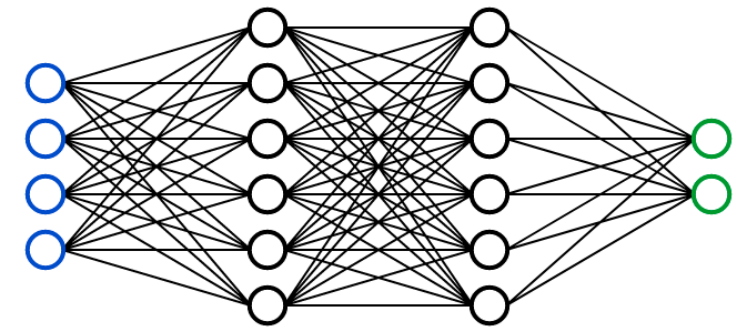
$$\pi^+(s) = \arg \min_a Q^\pi(s, a)$$

Approximate (Fitted) Policy Iteration

Fitted policy evaluation

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**

$D \leftarrow \emptyset$

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

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

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

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

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

return Q_θ

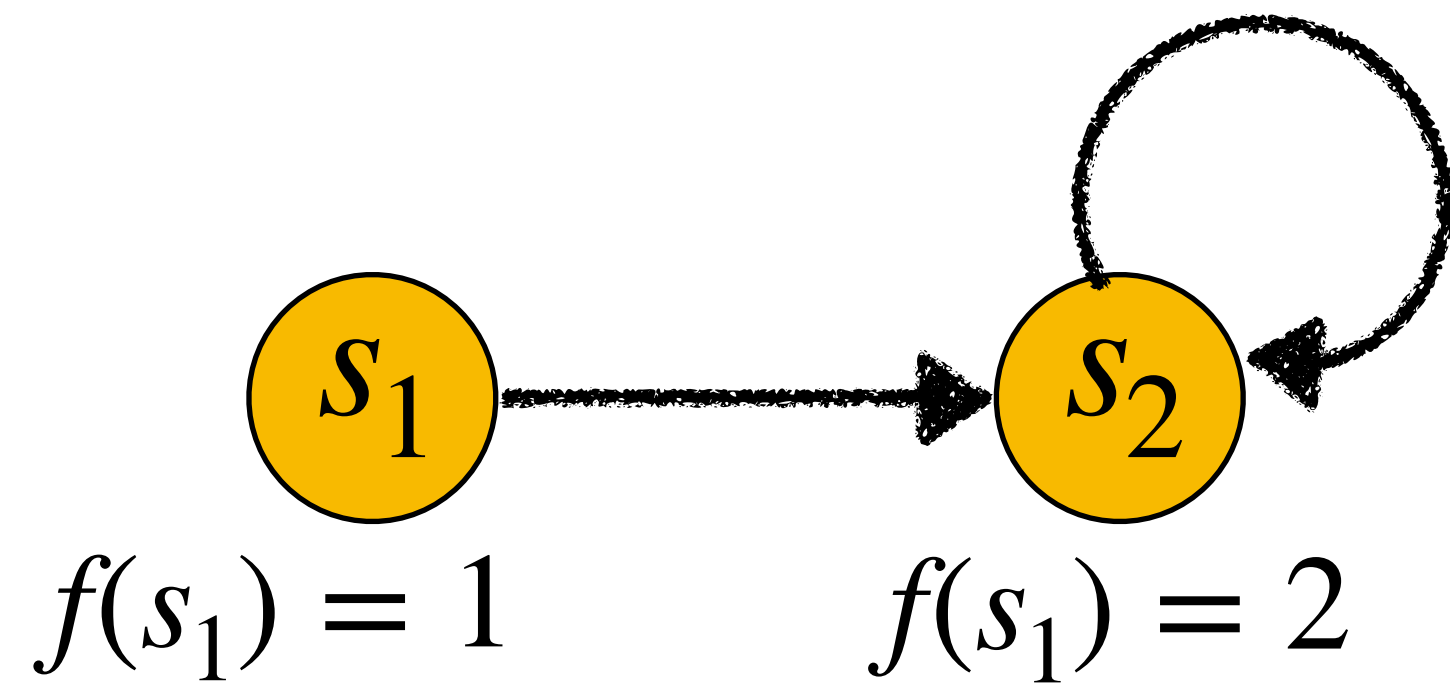
This remains
the same!

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

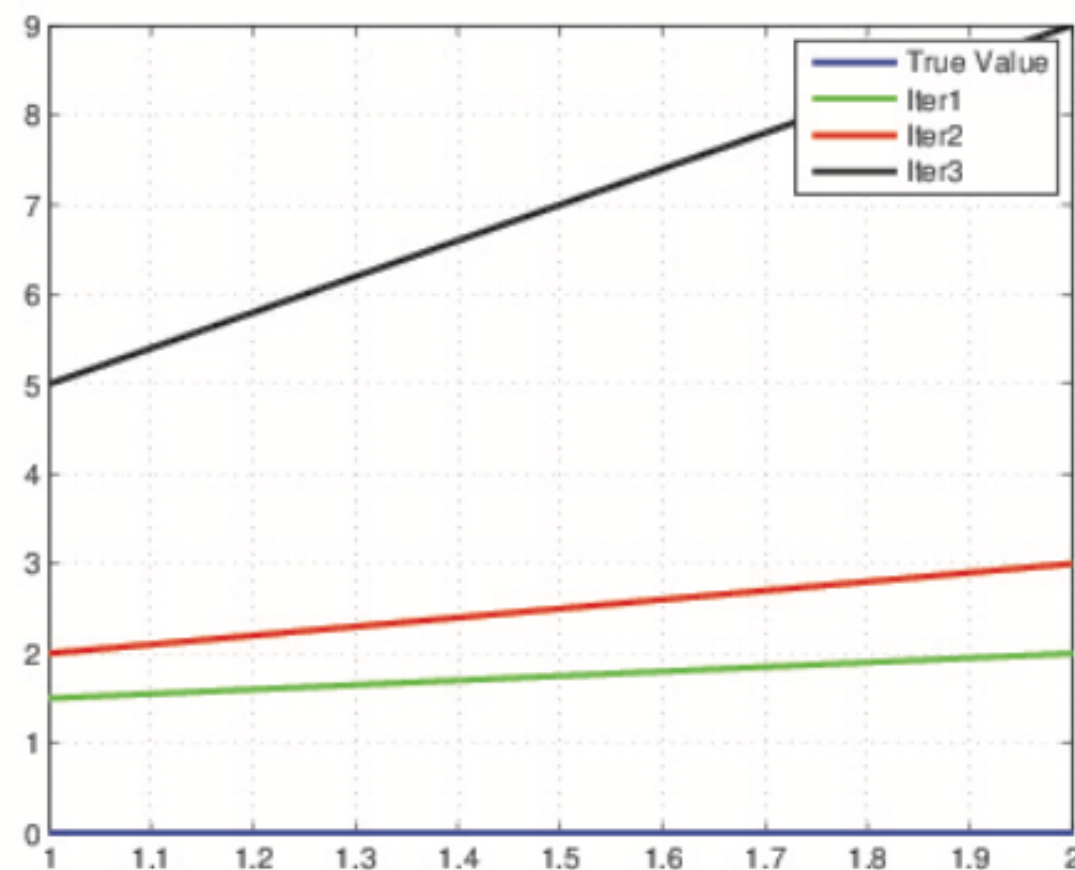
Surely approximate value
evaluation is more stable than
approximate value iteration?
(There is no `min()`!)



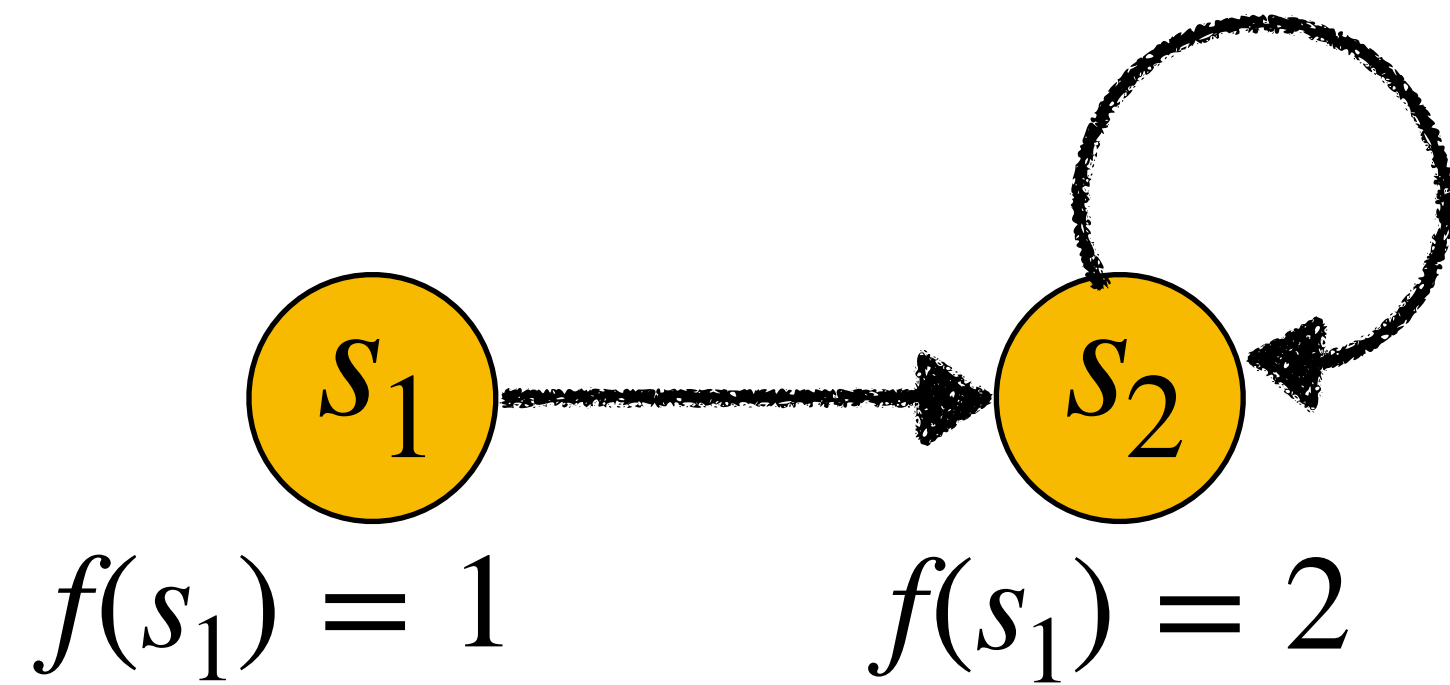
Well ... not quite



w blows up!

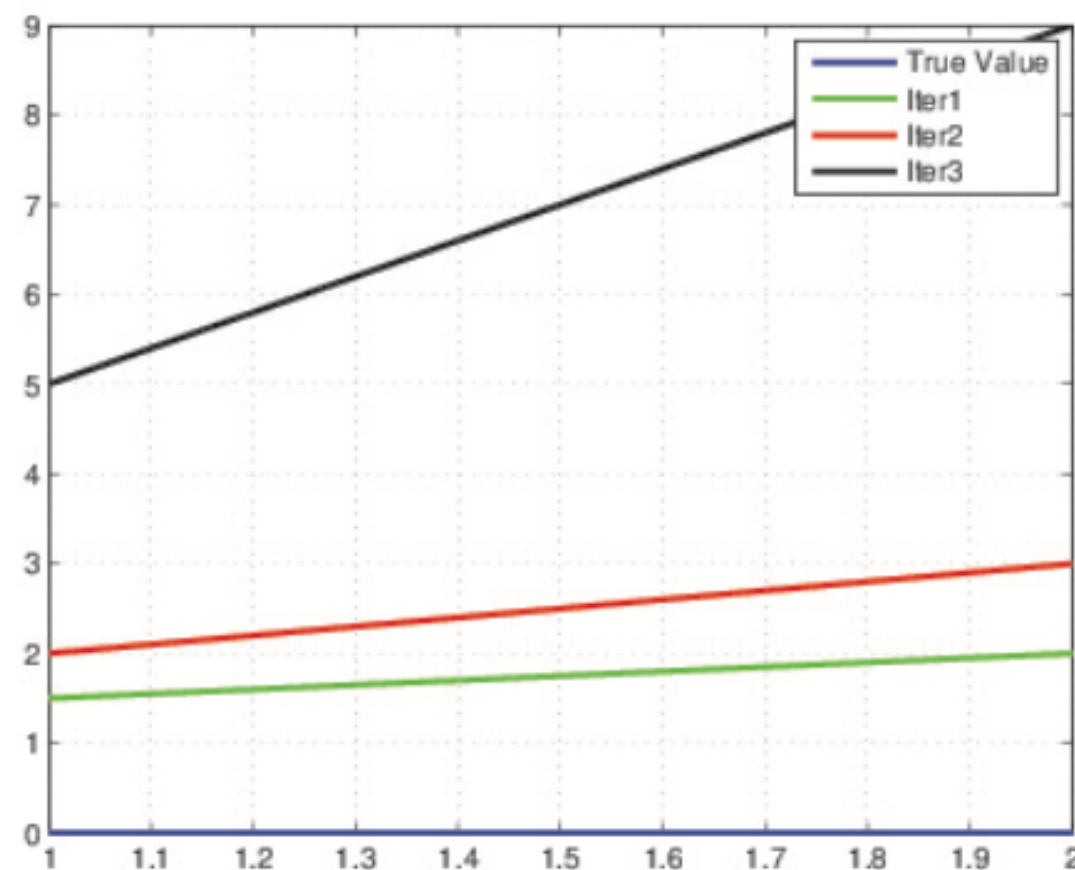


Well ... not quite



But we can fix this by
on-policy weighting

w blows
up!



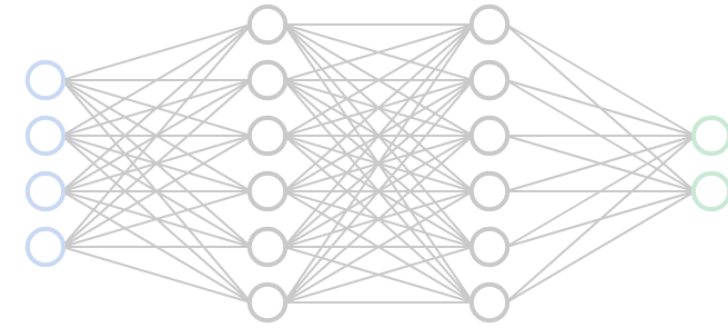
Weight each datapoint
by how often the
policy visits it.

But what about policy improvement?

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..
 $D \leftarrow \emptyset$

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

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

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

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

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

return Q_θ

But this has
the $\min()$ step!

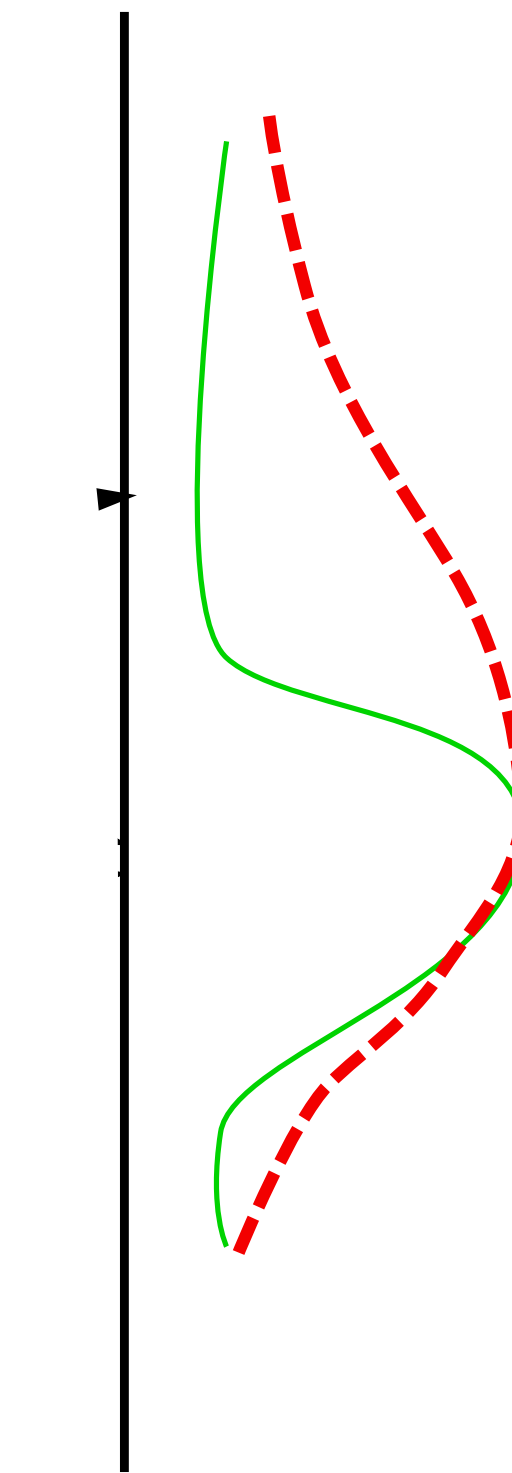
$$\pi^+(s) = \arg \min_a Q^\pi(s, a)$$

The problem of distribution shift

Upper half of state
is BAD

Lower half of state
is GOOD

T-1



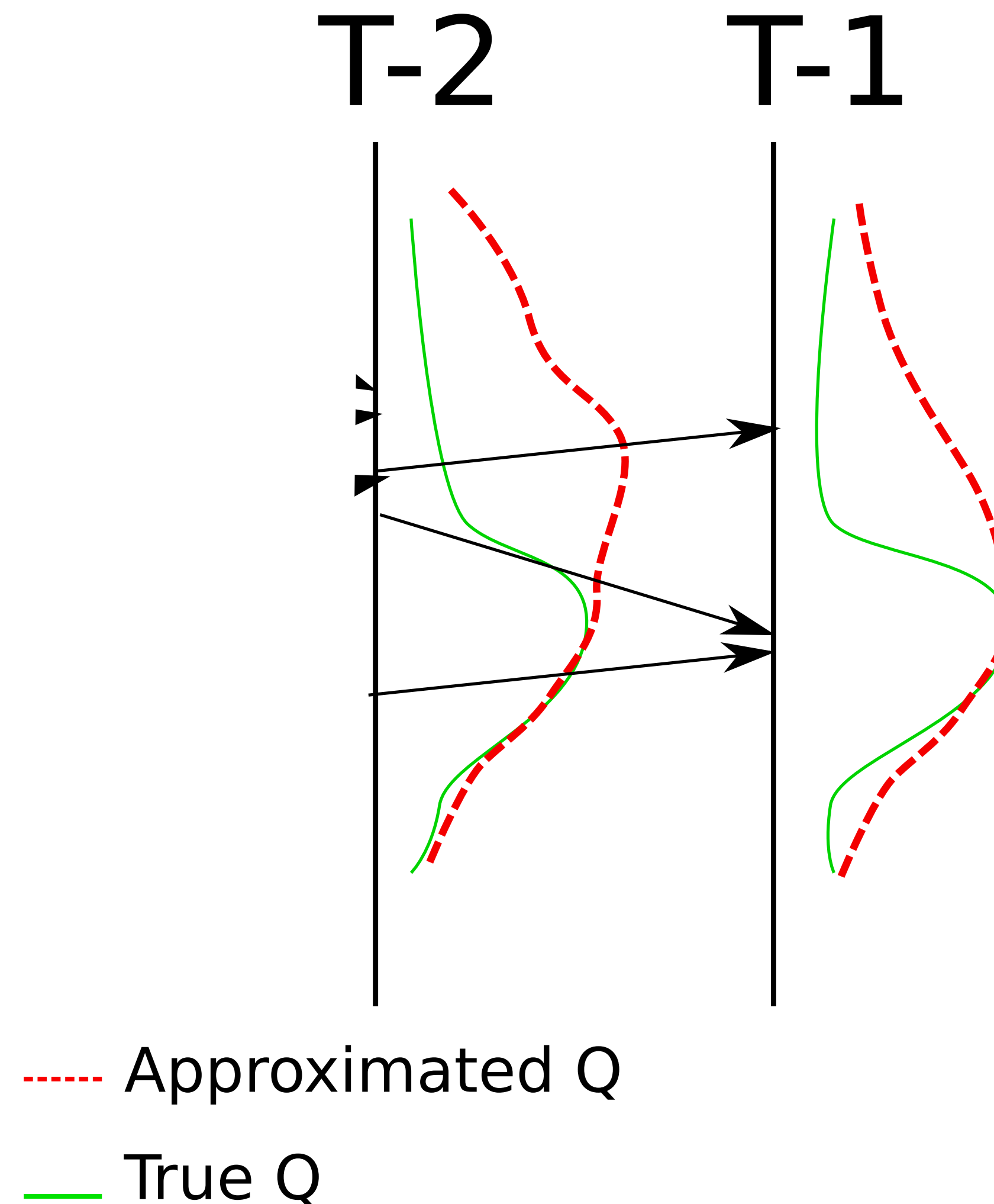
----- Approximated Q

— True Q

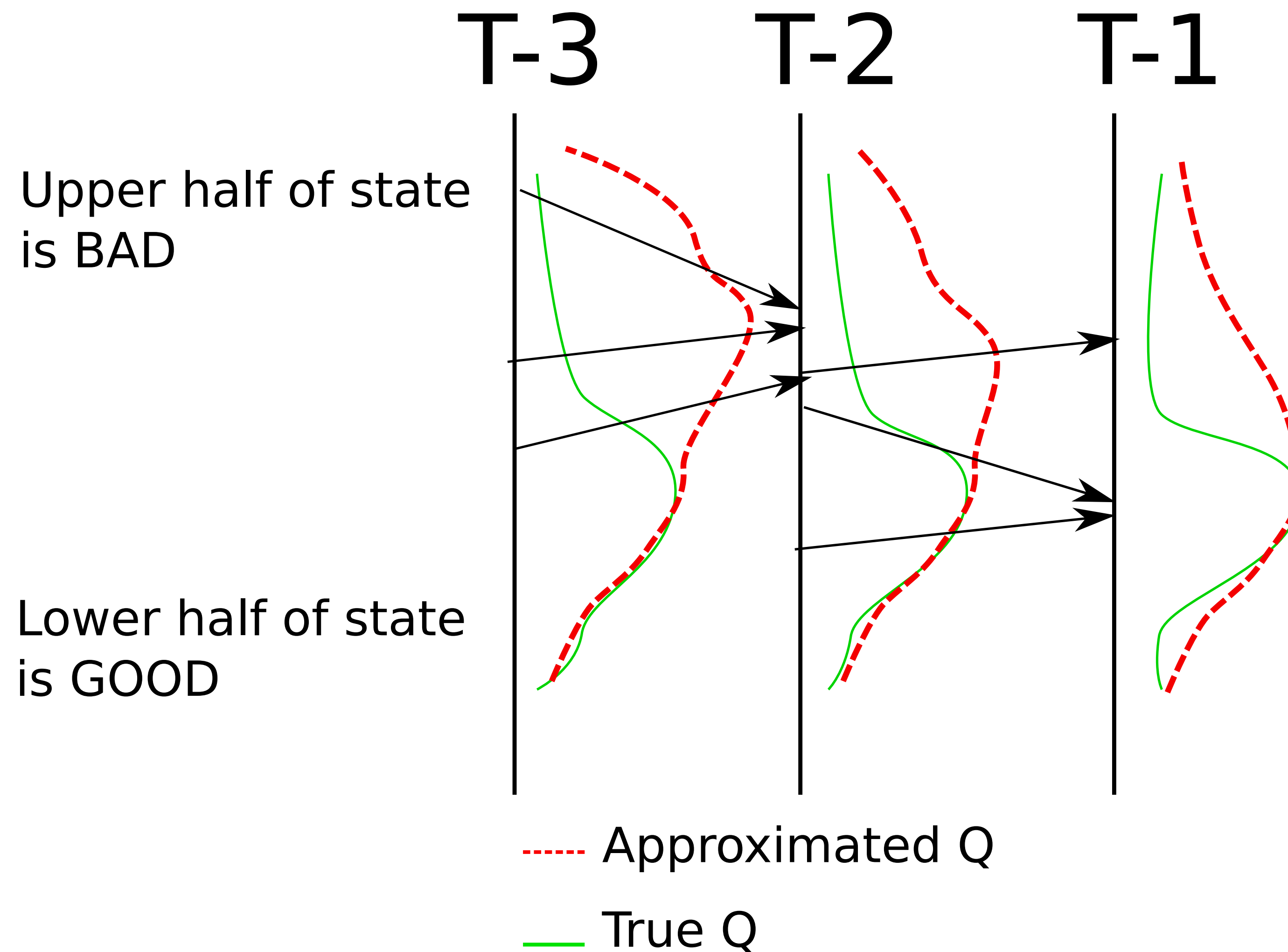
The problem of distribution shift

Upper half of state
is BAD

Lower half of state
is GOOD



The problem of distribution shift



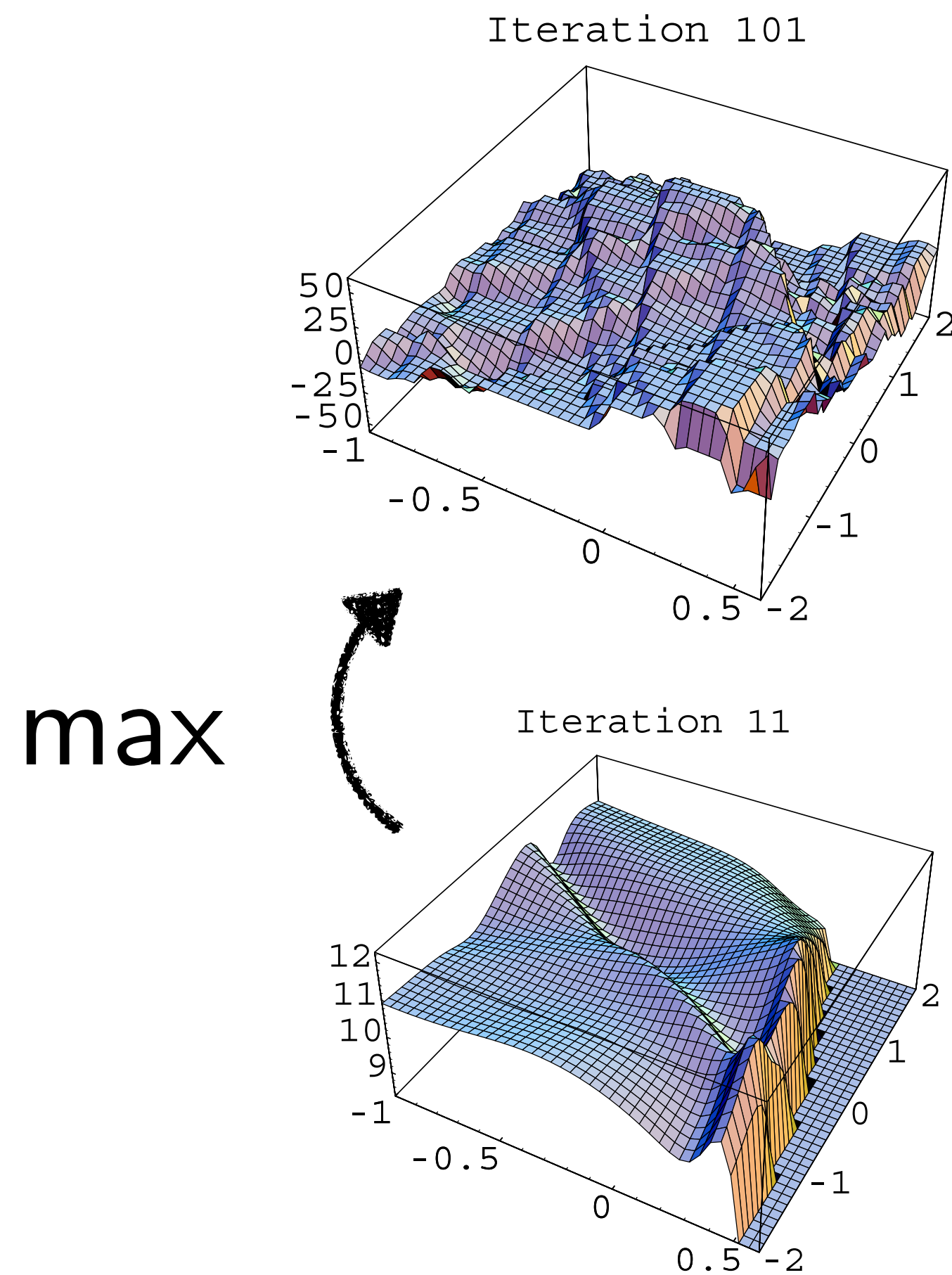
Bootstrapping

Distribution Shift

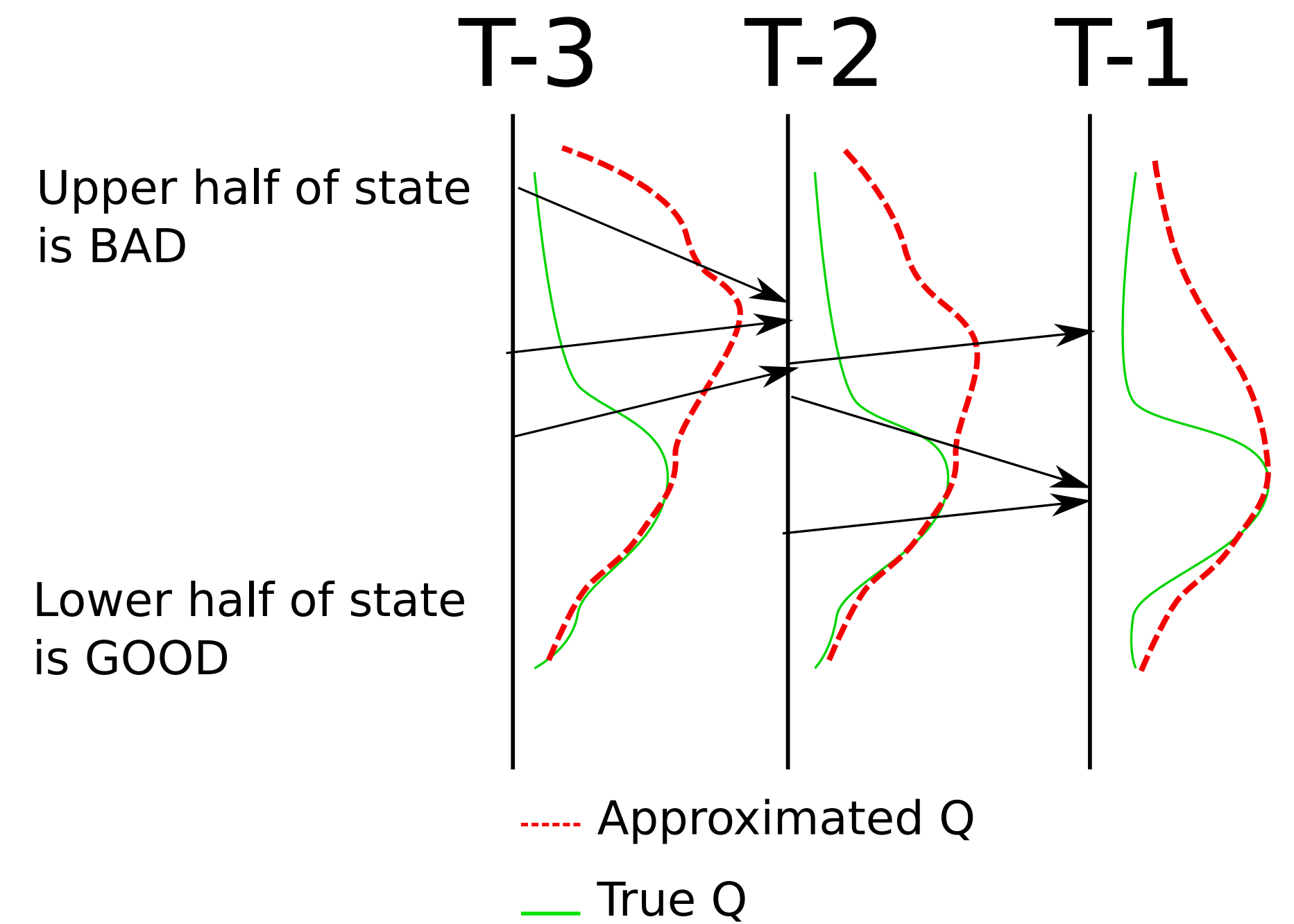


Two sides of the same coin

Bootstrapping



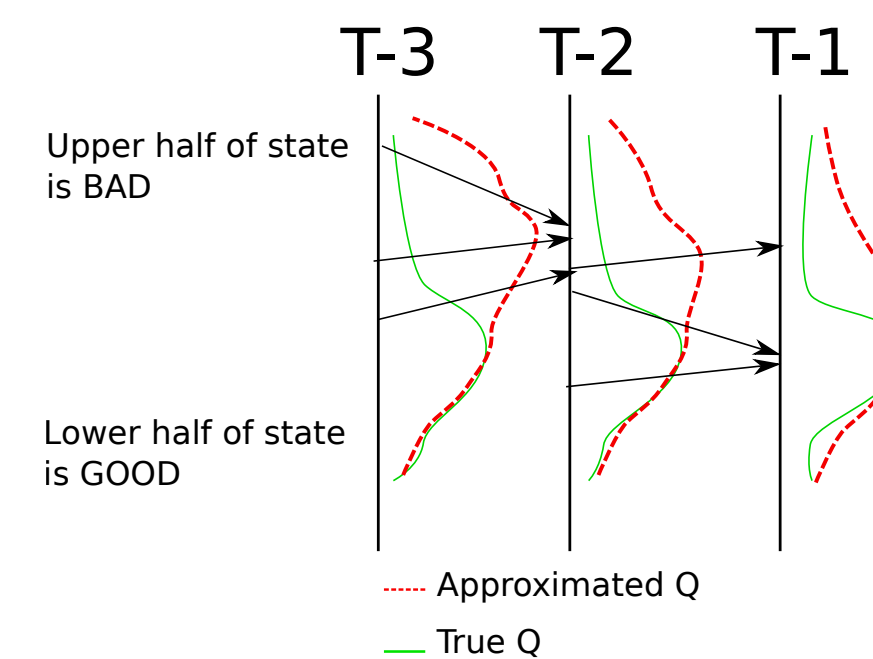
Distribution shift



Ideas for fixing
this?



Remedies



Bootstrapping



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

Execute policy
and **trust actual returns**

Distribution shift

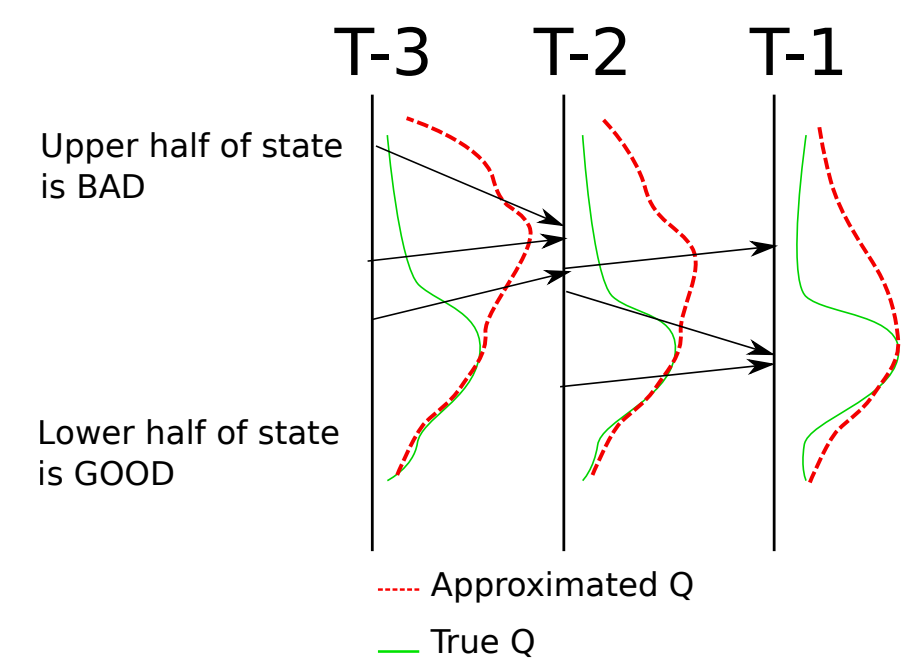


Minimize the
distribution shift

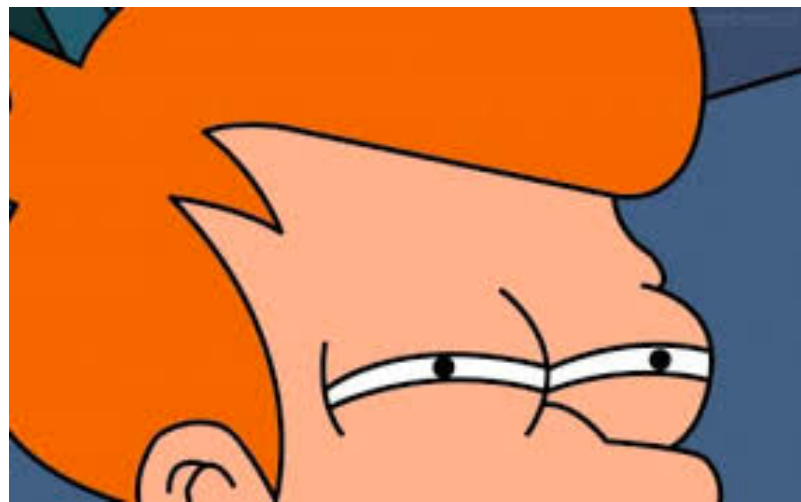
Be conservative,
change policy slowly



Remedies



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



Fixing bootstrapping

Policy Search via

Dynamic Programming

Design considerations with PSDP

Where do we get baseline distributions from?

From a human expert or from hand-crafted policies!

What do we do if T is very very large?

Just use a stationary policy

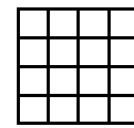
Performance Guarantees of PSDP

$$J(\pi_{PSDP}) - J(\pi^*) \leq \sum_{t=1}^T \left\| \frac{d_{\pi^*}^t}{d_{\mu}^t} \right\|_{\infty} \epsilon_t$$

tl;dr

Approximate (Fitted) Value Iteration

Q-iteration

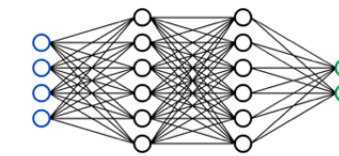


```

 $Q(s, a) \leftarrow 0$ 
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 $\{s_i, a_i, c_i, s'_i\}_{i=1}^N$

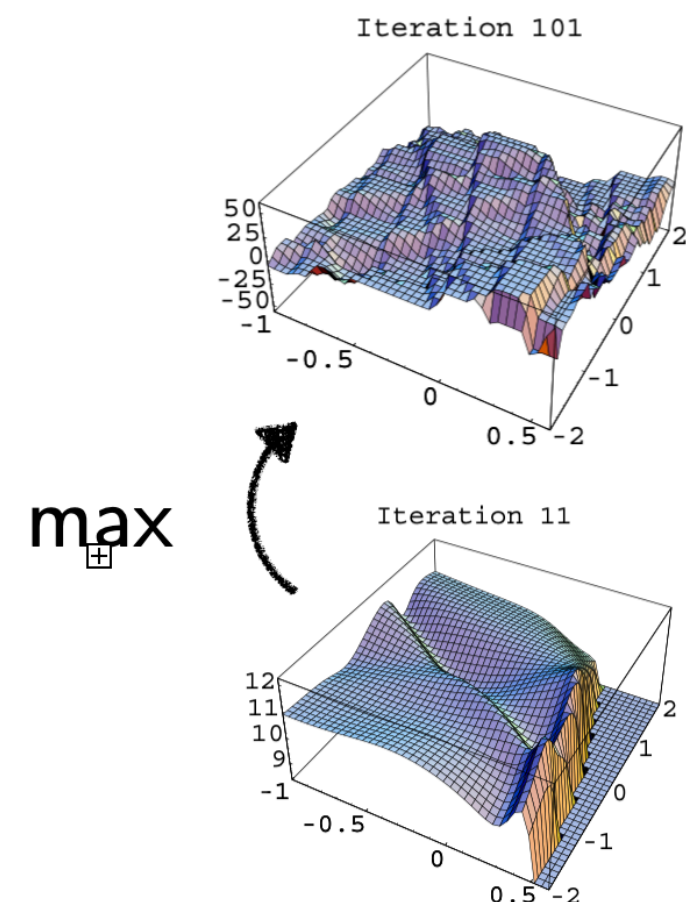
```

Init  $Q_\theta(s, a) \leftarrow 0$ 
while not converged do
   $D \leftarrow \emptyset$ 
  for  $i \in 1, \dots, n$ 
    input  $\leftarrow \{s_i, a_i\}$ 
    target  $\leftarrow c_i + \gamma \min_{a'} Q_\theta(s'_i, a')$ 
     $D \leftarrow D \cup \{\text{input}, \text{output}\}$ 
   $Q_\theta \leftarrow \text{Train}(D)$ 
return  $Q_\theta$ 

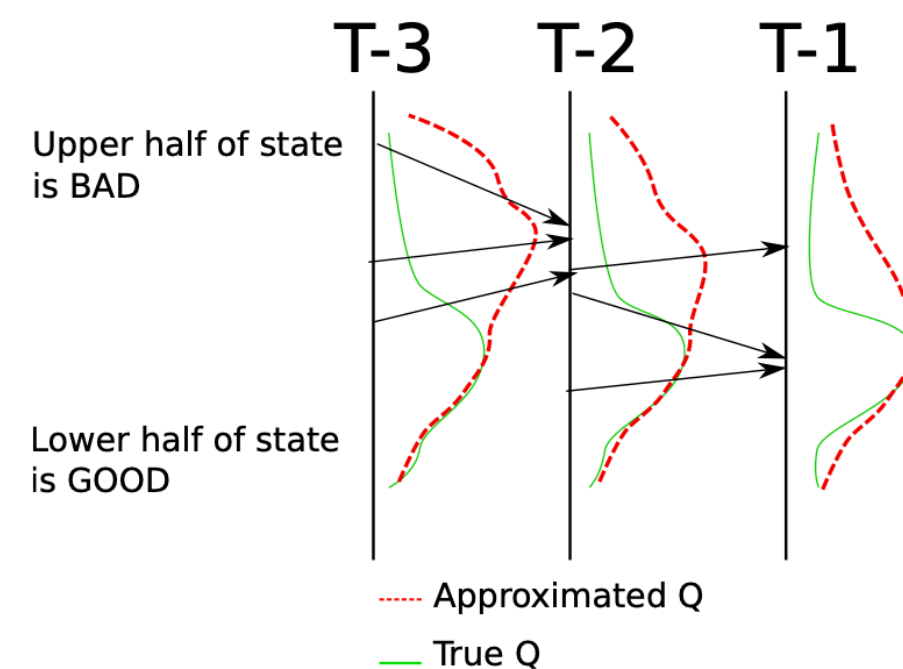
```

1

Bootstrapping



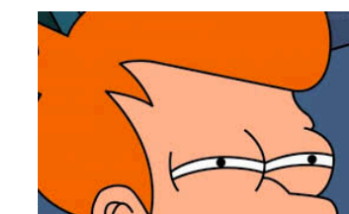
Distribution shift



Remedies



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

