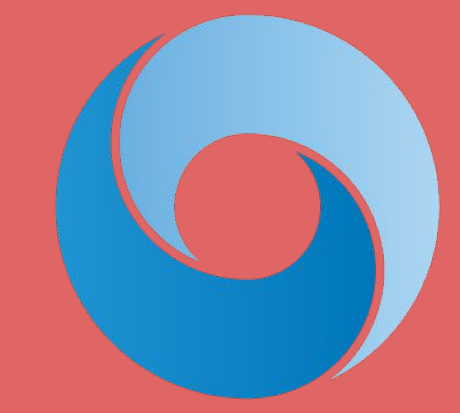


ADAPTIVE TEMPORAL-DIFFERENCE LEARNING FOR POLICY EVALUATION WITH PER-STATE UNCERTAINTY ESTIMATES



Hugo Penedones*, Carlos Riquelme*, Damien Vincent, Hartmut Maennel, Timothy Mann, Andre Barreto, Sylvain Gelly, Gergely Neu



Google DeepMind

In 20 seconds. Temporal Differences learning (TD) propagates approximation errors for policy evaluation. We adaptively choose between the TD and MC targets via MC confidence intervals.

BASIC DEFINITIONS

Value Function for Policy π :

$$v^\pi(s) = \mathbf{E} \left[\sum_{t=0}^{\infty} \gamma^t r(S_t, \pi(S_t)) \mid S_0 = s \right]$$

Mean Squared Value Error (MSVE):

$$\text{MSVE}(\hat{V}) = \mathbf{E}_{S_0 \sim \mu_0} \left[\left(v(S_0) - \hat{V}(S_0) \right)^2 \right]$$

Monte Carlo Target:

$$T_{\text{MC}}(s_t^{(i)}) := \sum_{k=0}^{n_i-t-1} \gamma^k r(s_{t+k}^{(i)})$$

TD₀ Target:

$$T_{\text{TD}(0)}(s_t^{(i)}) := r(s_t^{(i)}) + \gamma \hat{V}(s_{t+1}^{(i)})$$

CONFIDENCE INTERVALS

We fit m networks V_i on the MC target data:

$$\mathbf{D} = \{(s, T_{\text{MC}}(s)) \mid s \in S\}.$$

Given a new state s , get m predictions $v_i = V_i(s)$.

Assumption: $v_1, v_2, \dots, v_m \sim \mathbf{F}$ for some \mathbf{F} .

Then, compute its **predictive interval** at level $1-\alpha$.

In the paper we assume $\mathbf{F} = \mathbf{N}(\mu, \sigma^2)$, for unknown (μ, σ) leading to:

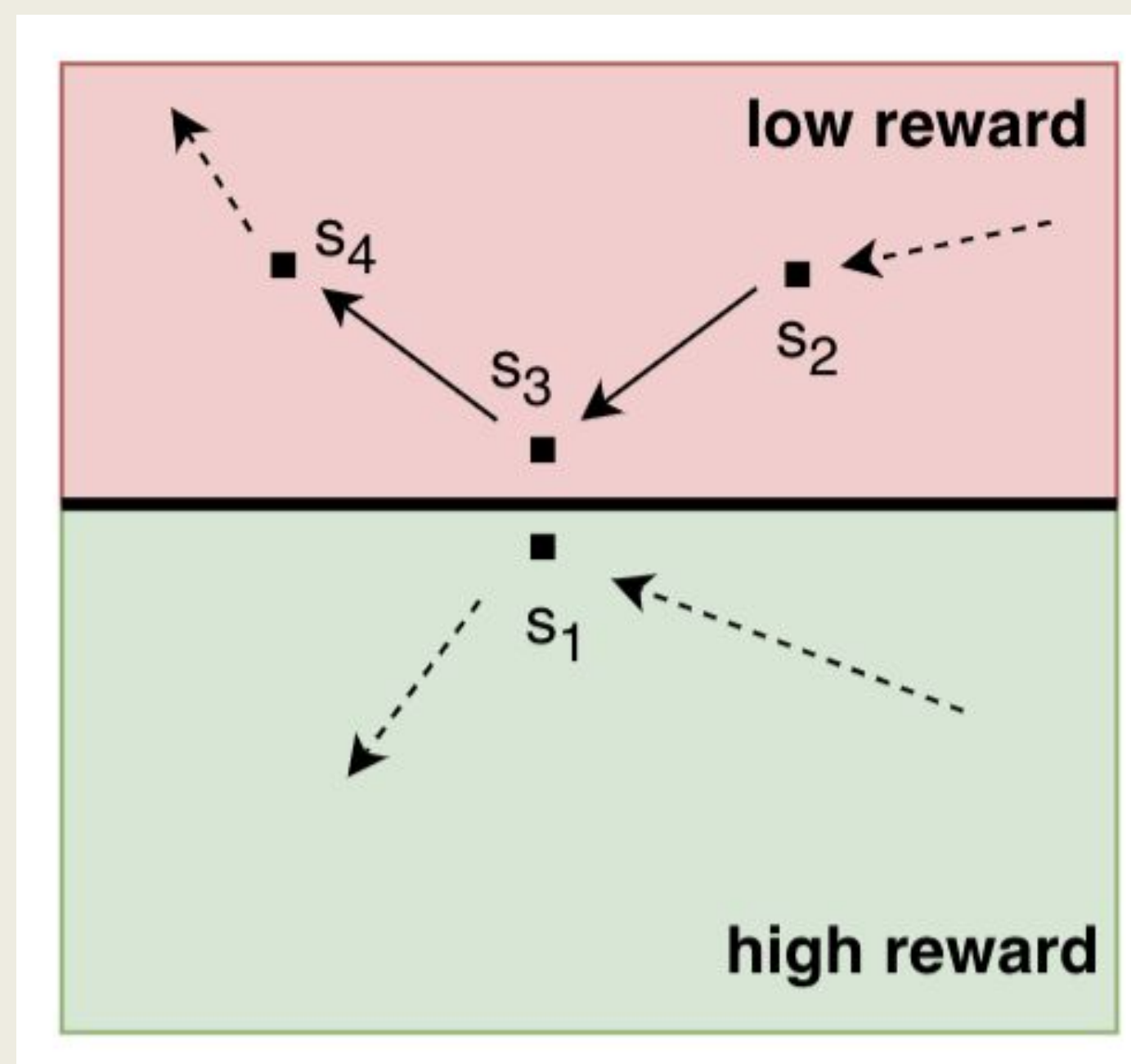
$$\bar{v} - z_\alpha \hat{\sigma}_m \sqrt{1+1/m} \leq v_{m+1} \leq \bar{v} + z_\alpha \hat{\sigma}_m \sqrt{1+1/m}$$

where $\bar{v} = \sum_i v_i / m$, and $\hat{\sigma}_m^2 = \sum_i (v_i - \bar{v})^2 / (m-1)$

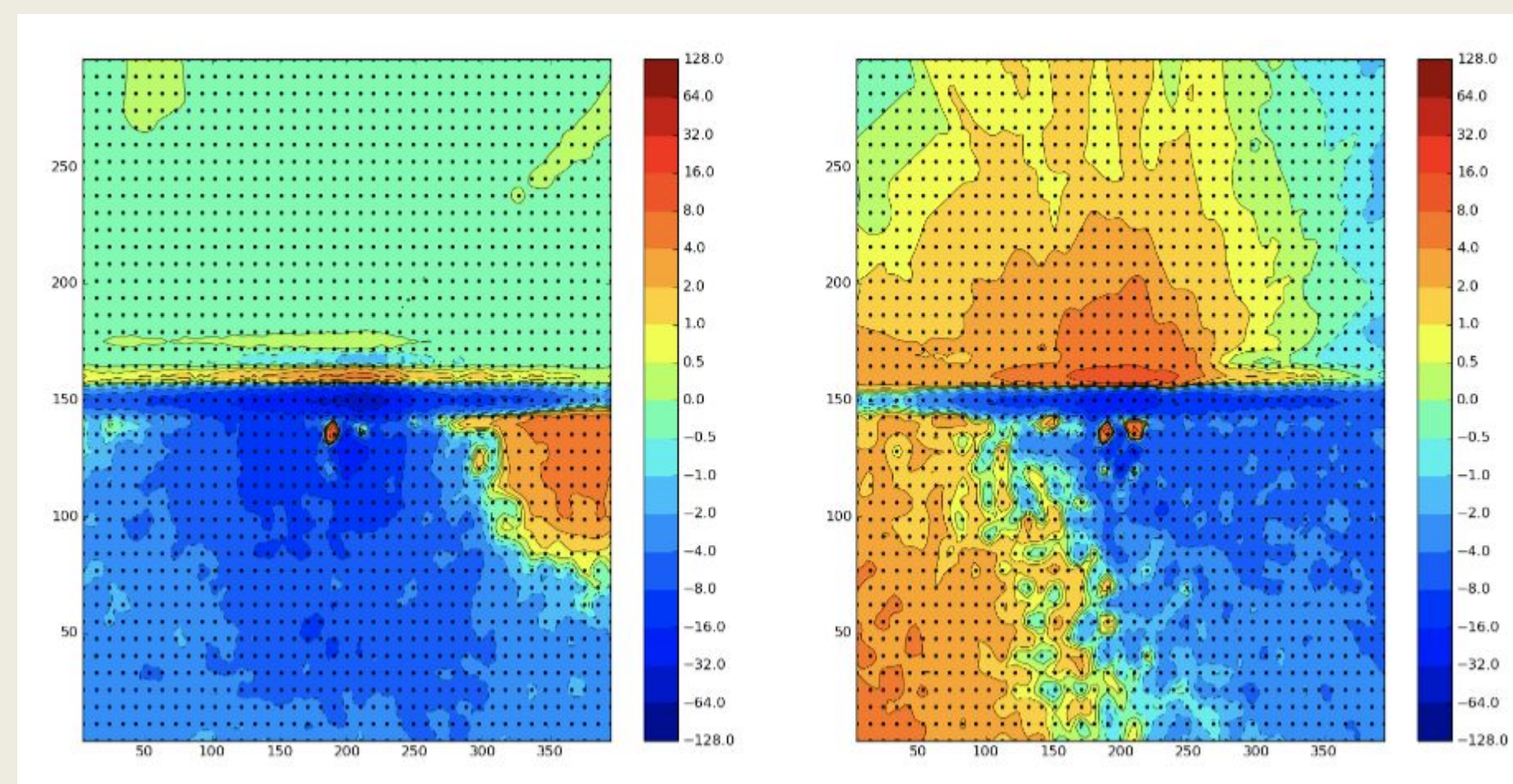
Other ways to estimate uncertainty can be used.

THE PROBLEM

Value Leakage with Function Approximation



Function approximation tradeoffs (s_1 and s_3) will be propagated through TD updates to s_2 .



MC

TD

MC versus TD on a simple environment with 2 rooms completely separated by a wall with a single target with positive reward right under the wall.

For each state s , the heatmaps show:

$$\hat{V}(s) - V(s)$$

The true value on the upper half of the plane is zero. MC overestimates the values of a narrow region right above the wall, due to function approximation limitations. With TD, these unavoidable approximation errors also occur, but things get worse when **bootstrap updates propagate errors** to much larger regions.

CONFIDENCE IN PRACTICE

THE ADAPTIVE TD ALGORITHM

Input: Confidence level $\alpha \in (0, 1)$. Trajectories τ_1, \dots, τ_n generated by policy π .

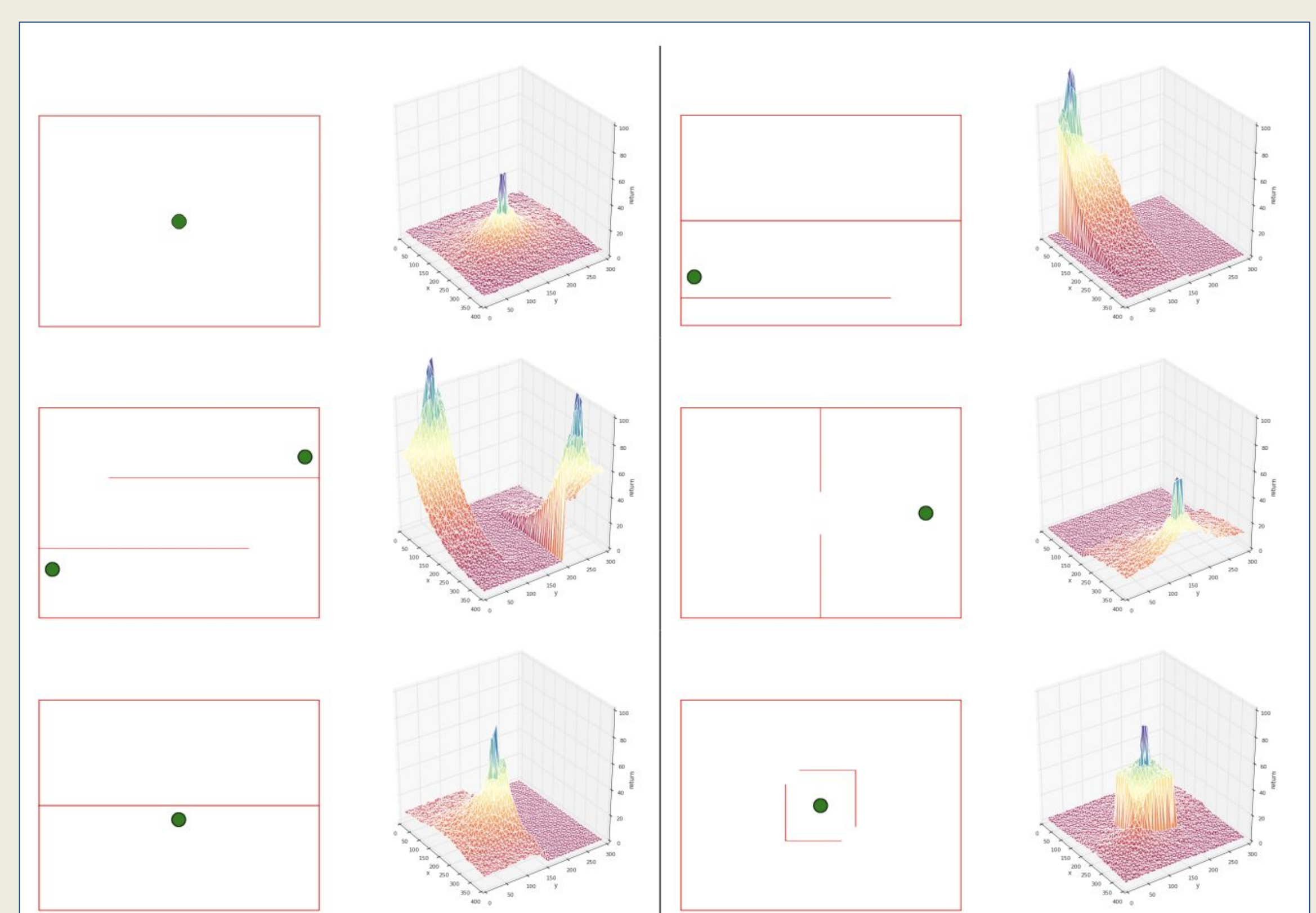
Let S be the set of visited states in τ_1, \dots, τ_n . Initialize $\hat{V}(s) = 0$, for all s . Compute Monte-Carlo returns dataset as in (3): $D_{\text{MC}} = \{(s, T_{\text{MC}}(s)) \mid s \in S\}$. Fit confidence function to D_{MC} : $\text{CI}_{\text{MC}}^\alpha(s) := (L_{\text{MC}}^\alpha(s), U_{\text{MC}}^\alpha(s))$.

```
repeat
  for  $i = 1$  to  $n$  do
    for  $t = 1$  to  $|\tau_i| - 1$  do
       $s_t^{(i)}$  is the  $t$ -th state of  $\tau_i$ .
       $T_{\text{TD}(0)} = r(s_t^{(i)}) + \gamma \hat{V}(s_{t+1}^{(i)})$ 
      if  $T_{\text{TD}(0)} \in (L_{\text{MC}}^\alpha(s_t^{(i)}), U_{\text{MC}}^\alpha(s_t^{(i)}))$  then
         $T_{i,t} \leftarrow T_{\text{TD}(0)}$ 
      else
         $T_{i,t} \leftarrow (L_{\text{MC}}(s_t^{(i)}) + U_{\text{MC}}(s_t^{(i)})) / 2$ 
      end if
      Use target  $T_{i,t}$  to fit  $\hat{V}(s_t^{(i)})$ .
    end for
  end for
until epochs exceeded
```

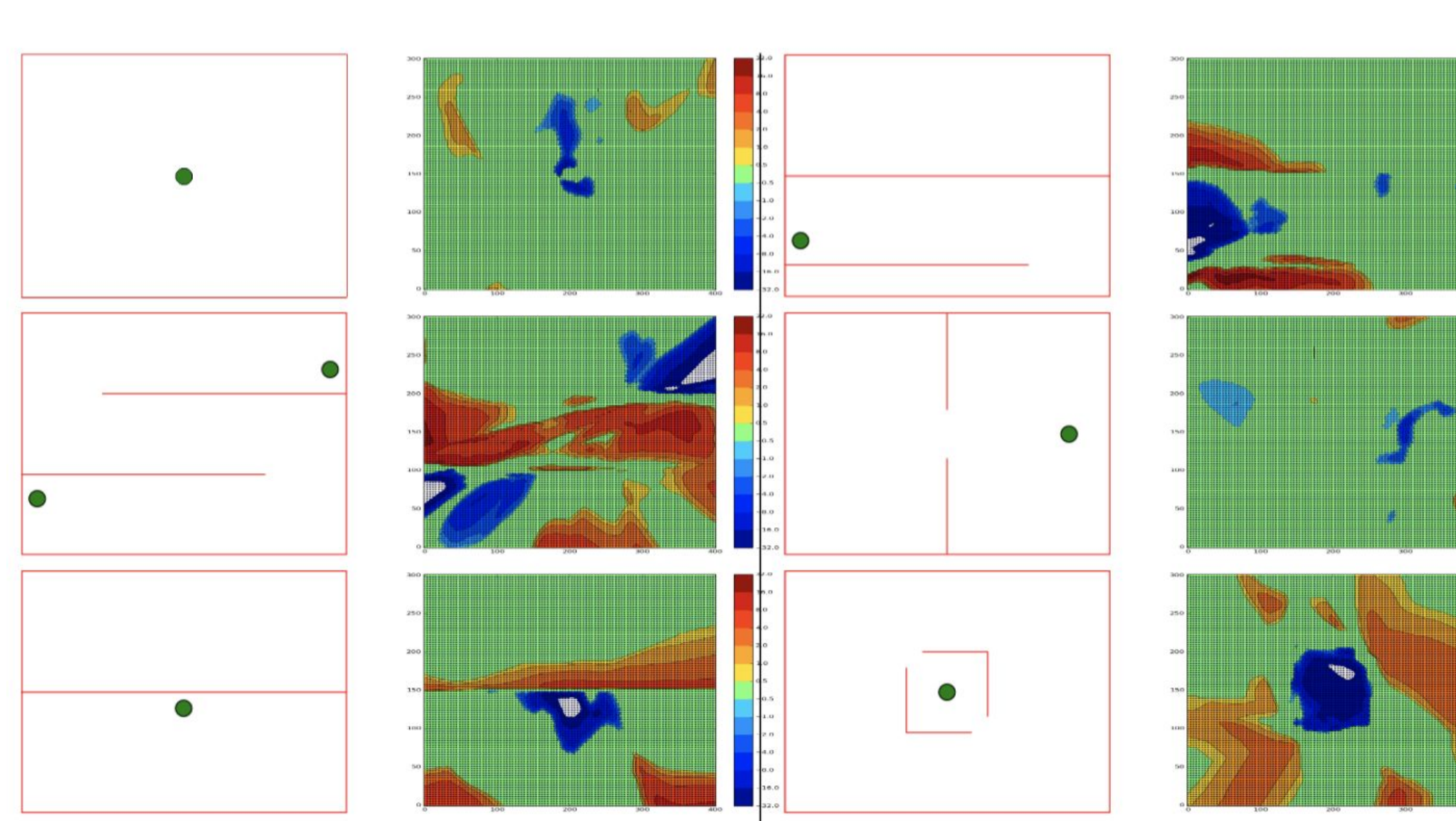
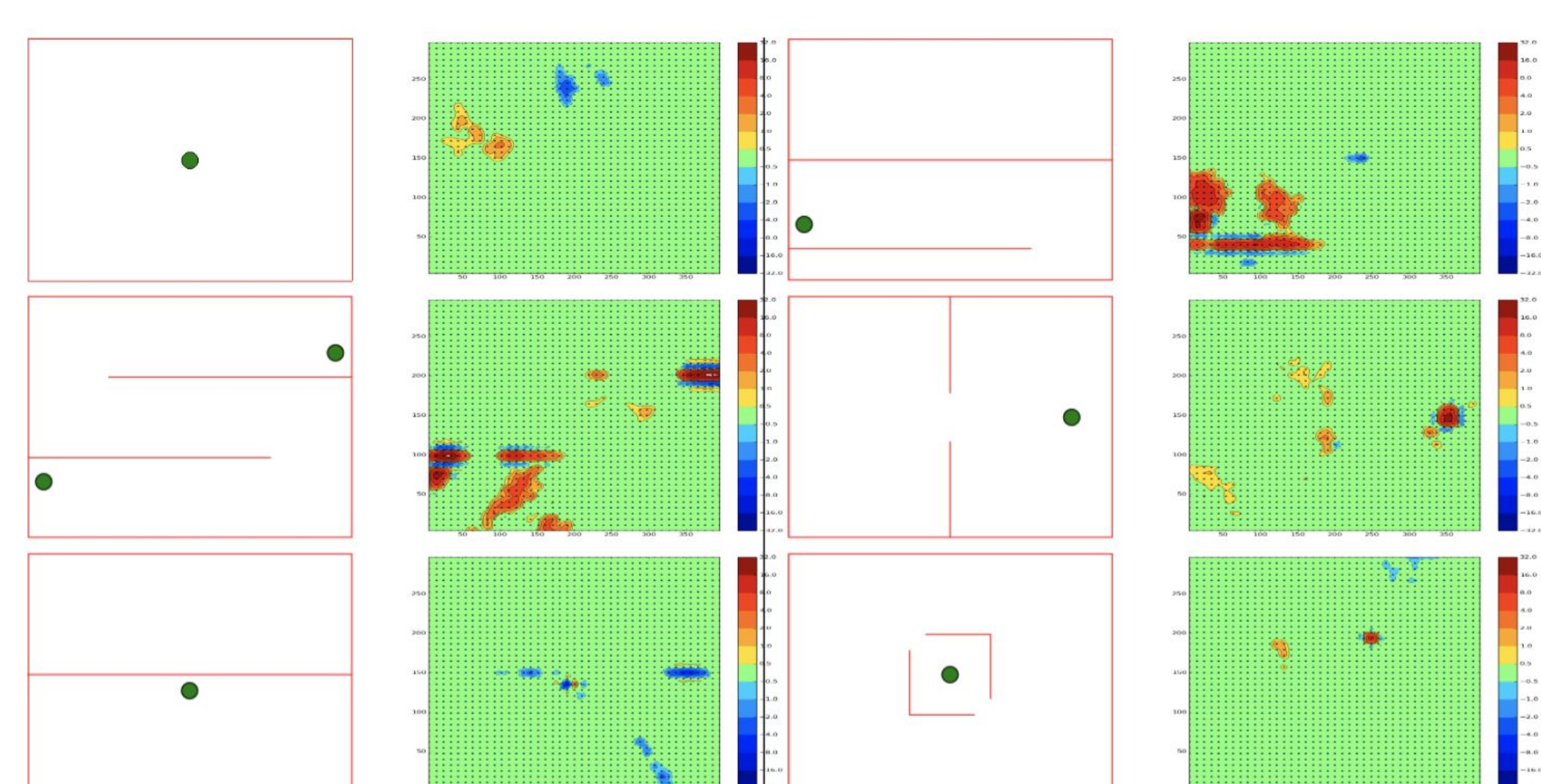
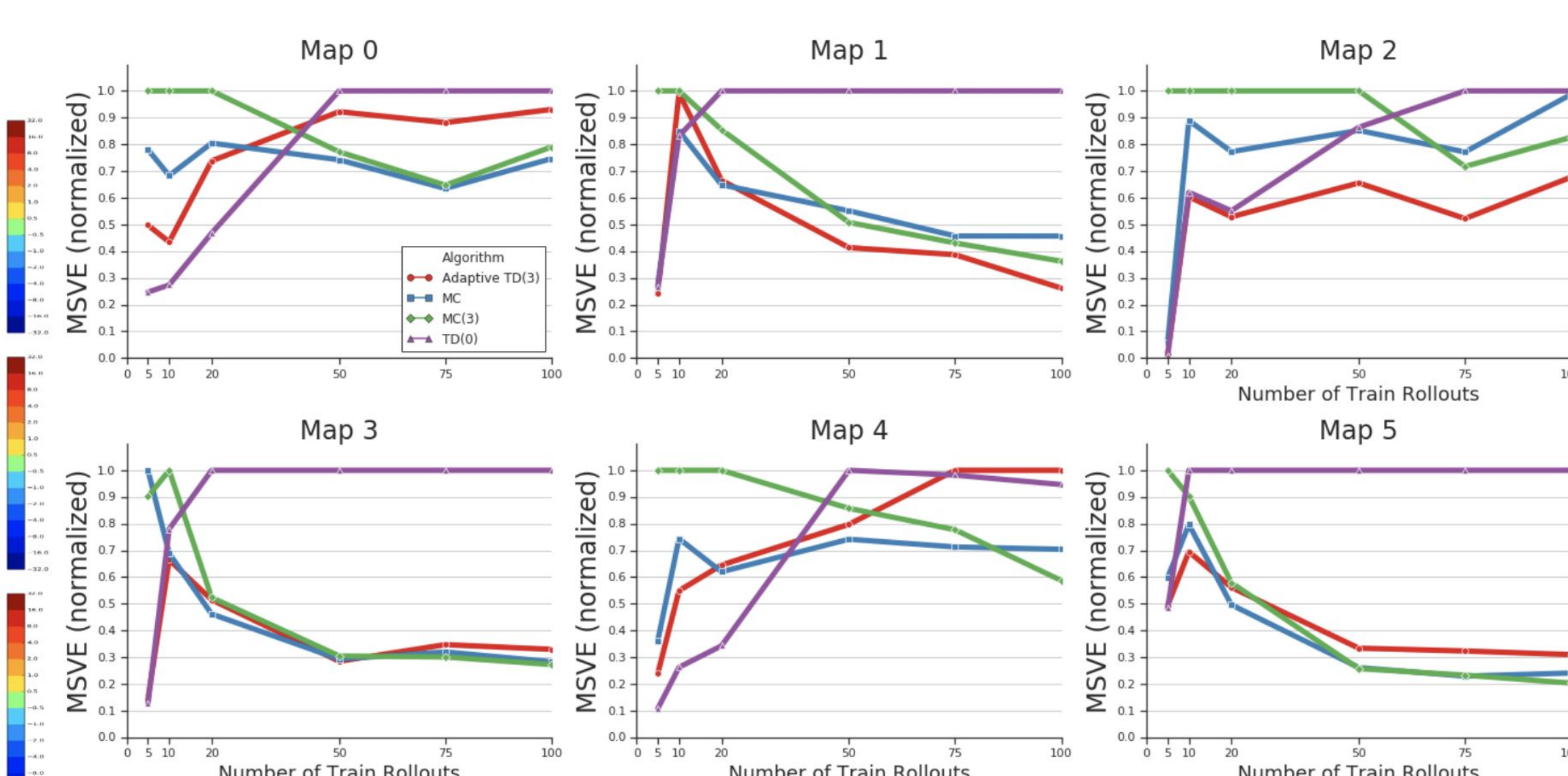
Algorithm 1: Adaptive TD

1. Train an ensemble of value networks with MC target.
2. Train a single value network with an adaptive target:
 - for each training state s :
 - a. By default, use TD target.
 - b. If TD target doesn't fall in MC interval ($\text{low}_{\text{MC}}, \text{high}_{\text{MC}}$):
 - i. Use $(\text{low}_{\text{MC}} + \text{high}_{\text{MC}}) / 2$ as target.

LAB-2D EXPERIMENTS



MSVE RESULTS



TD and MC intervals mismatch. Use MC there.

MC intervals mismatch with ground true value function