

## Abstract

We study MDPs evolving over time (Definition 1) and consider planning in this setting. We make two hypotheses:

1. Continuous, bounded, evolution (Definition 2);
2. Snapshot model (Definition 3) known at each decision epoch.

Our contribution can be presented in three points:

1. Proposal of a planning method robust to the environment's evolution;
2. Introduction of a zero-shot model-based algorithm, Risk-Averse Tree-Search (RATS), computing the best worst-case action;
3. Illustration of the benefits of the approach in experiments.

## Non-Stationary Markov Decision Processes

### Definition (1)

An **NSMDP** is an MDP whose transition and reward functions depend on the decision epoch. It is defined by a 5-tuple  $\{\mathcal{S}, \mathcal{T}, \mathcal{A}, (p_t)_{t \in \mathcal{T}}, (r_t)_{t \in \mathcal{T}}\}$  where  $\mathcal{S}$  is a state space;  $\mathcal{T} \equiv \{1, 2, \dots, N\}$  is the set of decision epochs,  $N \leq +\infty$ ;  $\mathcal{A}$  is an action space;  $p_t(s' | s, a)$  is the probability of reaching state  $s'$  with action  $a$  at decision epoch  $t$  in state  $s$ ;  $r_t(s, a, s')$  is the reward associated to the transition from  $s$  to  $s'$  with action  $a$  at decision epoch  $t$ .

### Definition (2)

An  $(L_p, L_r)$ -**LC-NSMDP** is an NSMDP whose transition and reward functions are respectively  $L_p$ -LC and  $L_r$ -LC w.r.t. time, i.e.,

$$\forall (t, \hat{t}, s, s', a) \in \mathcal{T}^2 \times \mathcal{S}^2 \times \mathcal{A}, \begin{cases} W_1(p_t(\cdot | s, a), p_{\hat{t}}(\cdot | s, a)) & \leq L_p |t - \hat{t}| \\ |r_t(s, a, s') - r_{\hat{t}}(s, a, s')| & \leq L_r |t - \hat{t}|. \end{cases}$$

## Risk-Averse Tree-Search algorithm

### Definition (3)

The **snapshot of an NSMDP**  $\{\mathcal{S}, \mathcal{T}, \mathcal{A}, (p_t)_{t \in \mathcal{T}}, (r_t)_{t \in \mathcal{T}}\}$  at decision epoch  $t_0$ , denoted by  $\text{MDP}_{t_0}$ , is the stationary MDP defined by the 4-tuple  $\{\mathcal{S}, \mathcal{A}, p_{t_0}, r_{t_0}\}$  where  $p_{t_0}(s' | s, a)$  and  $r_{t_0}(s, a, s')$  are the transition and reward functions of the NSMDP at  $t_0$ .

### Proposition (1)

**Set of admissible snapshot models.** Consider an  $(L_p, L_r)$ -**LC-NSMDP**,  $s, t, a \in \mathcal{S} \times \mathcal{T} \times \mathcal{A}$ . The transition and expected reward functions  $(p_t, R_t)$  of the snapshot  $\text{MDP}_t$  respect

$$(p_t, R_t) \in \Delta_t := \mathcal{B}_{W_1}(p_{t-1}(\cdot | s, a), L_p) \times \mathcal{B}_{|\cdot|}(R_{t-1}(s, a), L_p + L_r)$$

where  $\mathcal{B}_d(c, r)$  is the ball of centre  $c$ , defined with metric  $d$  and radius  $r$ .

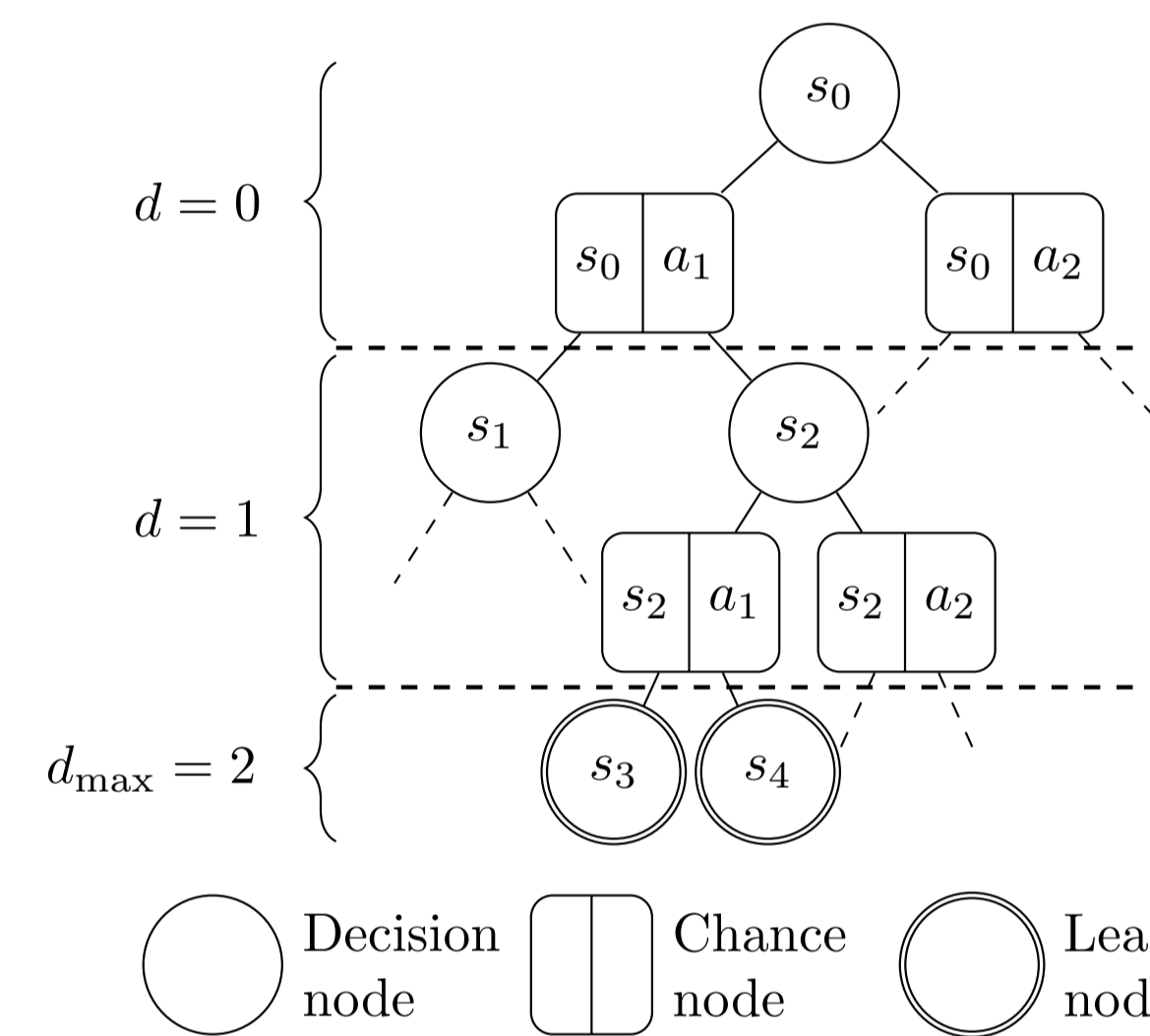


Figure: Tree structure: alternation between **decision nodes** labelled by a unique state and **chance nodes** labelled by a state-action pair. Maximum depth  $d_{\max} = 2$ , action space  $\mathcal{A} = \{a_1, a_2\}$ . The tree is entirely developed until  $d_{\max}$  which makes the per-time-step complexity of the RATS algorithm  $\mathcal{O}(|\mathcal{S}|^{3.5} |\mathcal{A}|^2)^{d_{\max}}$ .

### Algorithm 1: RATS algorithm

**RATS** ( $s_0, t_0, \text{maxDepth}$ )

$\nu_0 = \text{rootNode}(s_0, t_0)$

Minimax( $\nu_0$ )

$\nu^* = \arg \max_{\nu'} \text{in } \nu.\text{children } \nu'.\text{value}$

**return**  $\nu^*.\text{action}$

**Minimax** ( $\nu, \text{maxDepth}$ )

**if**  $\nu$  is *DecisionNode* **then**

**if**  $\nu.\text{state}$  is *terminal* **or**  $\nu.\text{depth} = \text{maxDepth}$  **then**

**return**  $\nu.\text{value} = \text{heuristicValue}(\nu.\text{state})$

**else**

**return**  $\nu.\text{value} = \max_{\nu' \in \nu.\text{children}} \text{Minimax}(\nu', \text{maxDepth})$

**else**

**return**  $\nu.\text{value} = \min_{(p, R) \in \Delta_{t_0}^i} R(\nu) + \gamma \sum_{\nu' \in \nu.\text{children}} p(\nu' | \nu) \text{Minimax}(\nu', \text{maxDepth})$

## Experiments

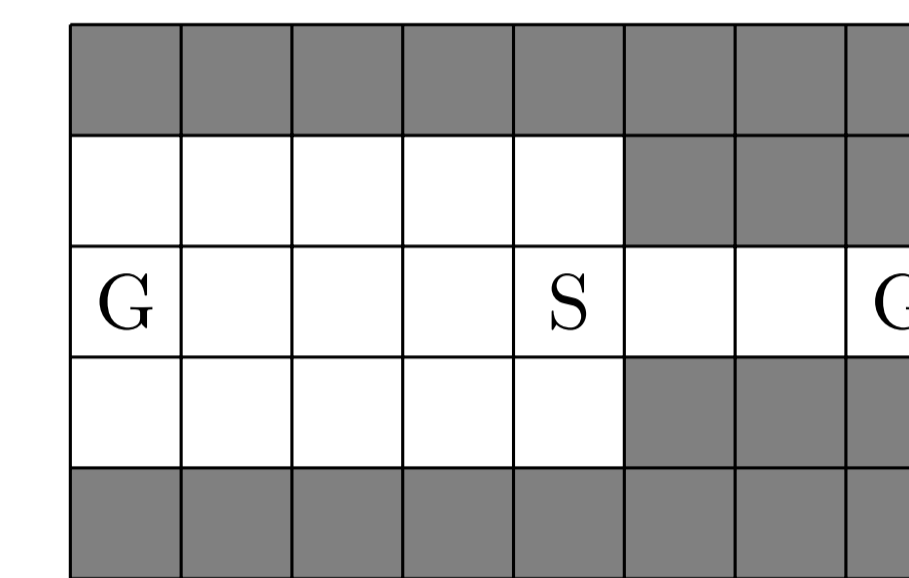


Figure: Non-Stationary bridge MDP.

The value of  $\epsilon \in [0, 1]$  defines different possible evolutions:

- ▶  $\epsilon = 0$  left cells are slippery;
- ▶  $\epsilon = 1$  right cells are slippery;
- ▶  $\epsilon \in (0, 1)$  linear balance between extreme cases.

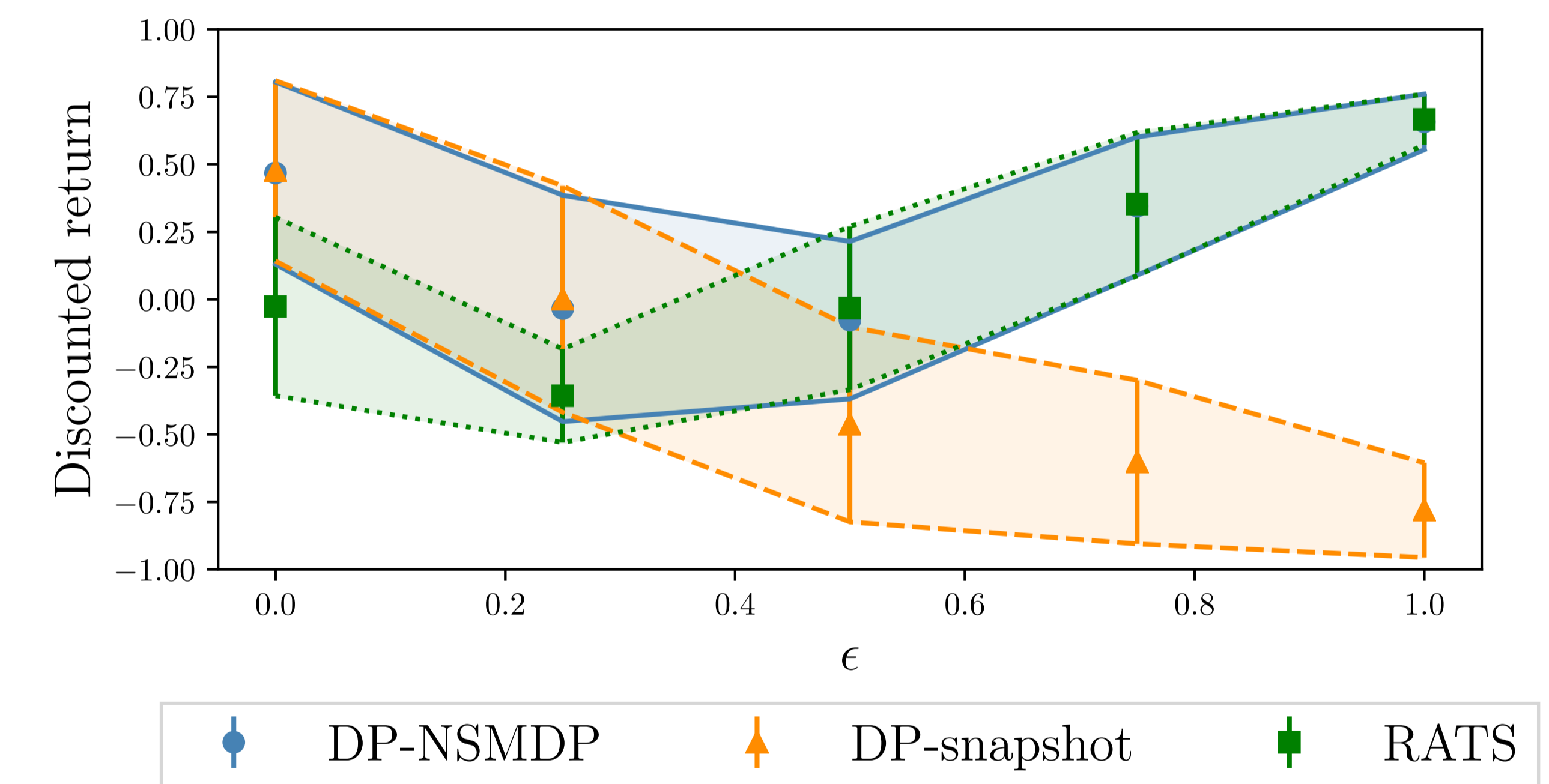


Figure: Discounted return vs  $\epsilon$ , 50% of standard deviation.

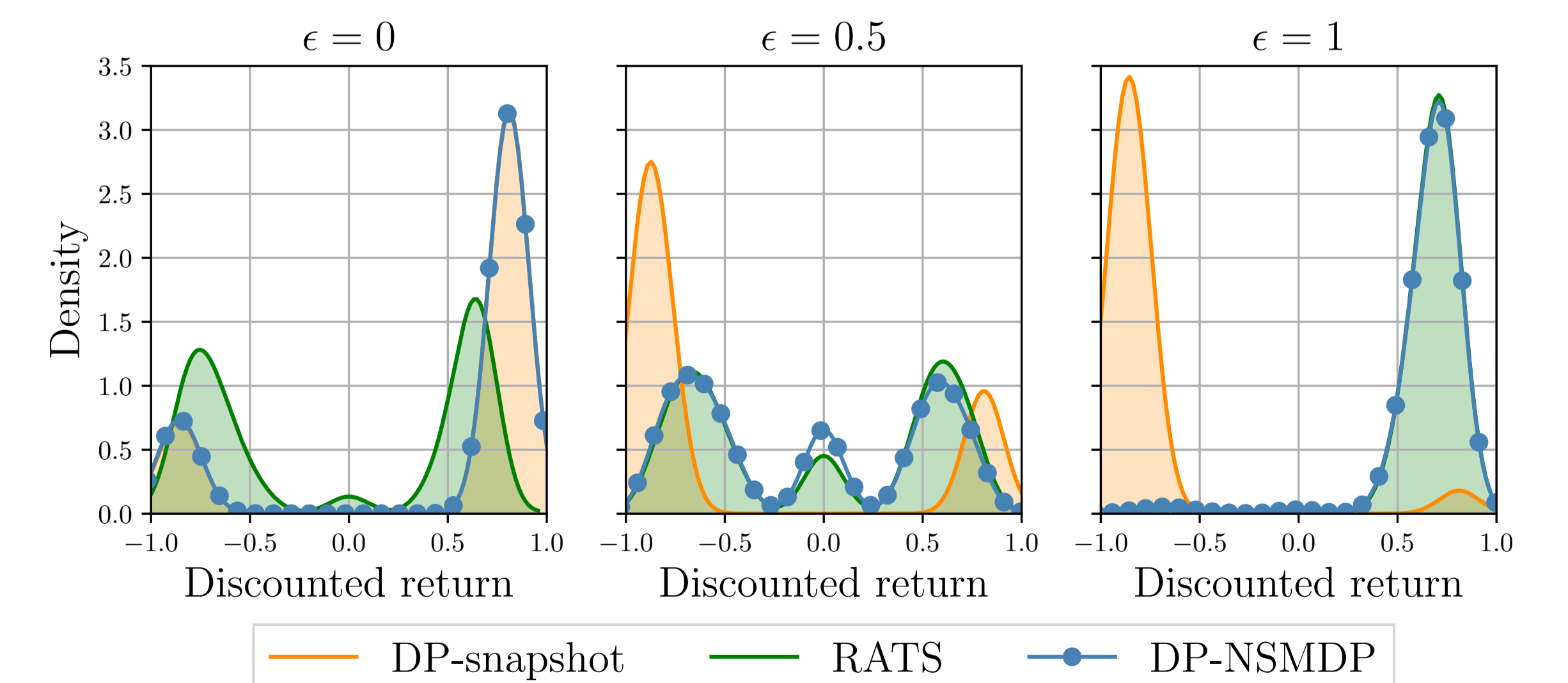


Figure: Discounted return distributions  $\epsilon \in \{0, 0.5, 1\}$ .