

## Introduction

Policy-gradient RL algorithms struggle in environments with delayed or episodic rewards.





### Our contributions:

- Self-imitation (SI): Exploiting useful agent behavior from the past, to improve temporal credit assignment.
- Exploration via a diverse ensemble of Self-imitation agents, using Stein Variational Gradient Descent (SVGD).

## SI with GAIL

- GAIL (Ho & Ermon, 2016) frames imitation learning as matching the state-action visitation distributions of the expert and the policy.
- **Exploiting the past:** Maintain a priority buffer  $\mathcal{B}$  of high return trajectories from the previous policy rollouts, and optimize the policy to match the state-action visitation distribution in the buffer.

 $\min_{\pi} \mathbb{D}(\rho_{\pi}, \{s_i, a_i\}_{\mathcal{B}})$  $\mathbb{D}$ : distance metric





• Similar to GAIL, using Jenson-Shanon divergence as the distance metric gives a min-max objective,

 $\min_{\pi} \max_{D} \mathbb{E}_{(s,a) \sim \rho_{\mathcal{B}}}[\log D(s,a)] + \mathbb{E}_{(s,a) \sim \rho_{\pi}}[\log(1 - D(s,a))]$ 

• Combining RL and SI:

$$\nabla_{\theta} \eta(\pi_{\theta}) = \nabla_{\theta} \mathbb{E}_{(s,a) \sim \rho_{\pi}} [r(s,a)] + \beta \nabla_{\theta} \mathbb{E}_{(s,a) \sim \rho_{\pi}} [-\log(1 - D(s,a))]$$

- Environmental rewards r(s,a) can be sparse or delayed.
- -Synthetic rewards  $-\log(1 D(s, a))$  are available for each (s, a), helping with temporal credit assignment.

# Learning Self-Imitating Diverse Policies

Tanmay Gangwani<sup>1</sup> Qiang Liu<sup>2</sup> Jian Peng<sup>1</sup>

<sup>1</sup>University of Illinois, Urbana-Champaign



<sup>2</sup>The University of Texas at Austin

## Results

**Experiment I:** Compare single-agent Self-imitation to standard policy-gradient RL in MuJoCo environments with episodic rewards.

	SI	PPO	
Walker	2996	252	(a) 3500 3000 2500 1500 1500 1000
Humanoid	3602	532	
H-Standup (× $10^4$ )	18.1	4.4	
Hopper	2618	354	
Swimmer	173	21	e 500
Invd.Pendulum	8668	344	<u> </u>

**Experiment II:** Elucidate failure of Self-imitation in harder tasks, and measure the benefit of training an ensemble of SI agents with explicit  $D_{JS}$ repulsion between policies (labeled SI-interact-JS). The SI-independent baseline trains isolated SI policies and selects the best among them.

• 2D-Maze with deceptive rewards.



center-of-mass is beyond a certain threshold distance. - SI-interact-RBF uses RBF kernel for SVGD



a. Diversification with SVGD helps in *discovery* of sparse rewards . Self-imitation then helps to efficiently *exploit* the discovered rewards

Code - github.com/tgangwani/selfImitationDiverse Paper Link - https://arxiv.org/abs/1805.10309 Correspondence - gangwan2@illinois.edu





• More in paper: Evaluation on noisy environments (each reward  $r_t$  masked to zero with some probability); comparison to CEM and ES.

• MuJoCo locomotion tasks: Forward velocity reward is only provided if the

## Takeaways

