

Proximal Meta-Policy Optimization: ProMP

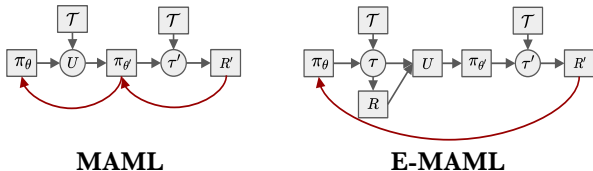
Jonas Rothfuss*, Dennis Lee*, Ignasi Clavera*,
Tamim Asfour, and Pieter Abbeel



Goal

1. Analyze **credit assignment** in meta-reinforcement learning
2. Develop a **new objective** that trains for the pre-update sampling distribution

Credit Assignment Sampling Distribution



Low Variance Curvature Estimator (LVC)

$$J^{\text{LVC}}(\tau) = \sum_{t=0}^{H-1} \frac{\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}{\perp(\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t))} \left(\sum_{t'=t}^{H-1} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) \quad \tau \sim P_{\mathcal{T}}(\tau)$$

- Meta-gradient with **low variance**
- **Unbiased** closed to local optima

Proximal Meta-Policy Optimization: ProMP

ProMP Objective:

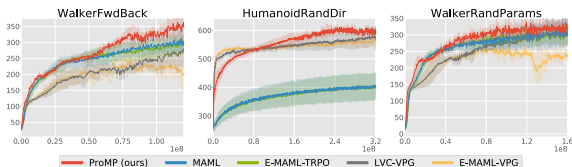
$$J_{\mathcal{T}}^{\text{ProMP}}(\theta) = J_{\mathcal{T}}^{\text{CLIP}}(\theta') - \eta \bar{\mathcal{D}}_{KL}(\pi_{\theta_o}, \pi_{\theta}) \quad \text{s.t.} \quad \theta' = \theta + \alpha \nabla_{\theta} J_{\mathcal{T}}^{\text{LR}}(\theta)$$

Incorporates the benefits of:

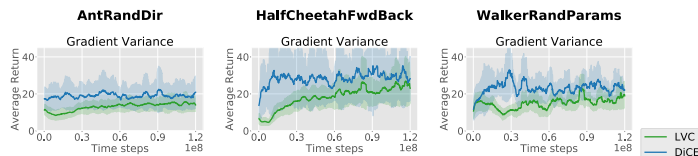
- Proximal Policy Optimization
- LVC Estimator

Experiments

Performance Comparison



Variance Comparison



Exploration – Exploitation

