

Meta-Inverse Reinforcement Learning with Probabilistic Context Variables

Lantao Yu*, Tianhe Yu*, Chelsea Finn, Stefano Ermon

Department of Computer Science, Stanford University



Highlights

- We aim at addressing two key limitations of existing inverse reinforcement learning (IRL) methods:
 - Learning reward functions from scratch and requiring large numbers of demonstrations to correctly infer the reward for each task.
 - Assuming demos are for one isolated task, while in practice it is more natural and scalable to obtain heterogeneous demos.
- We propose a new meta-inverse reinforcement learning framework based on latent probabilistic context variables termed PEMIRL.
- PEMIRL is capable of learning rewards from unstructured, multi-task demonstration data, and critically, use this experience to infer robust rewards for new, structurally-similar tasks from a single demonstration.
- We demonstrate the effectiveness of our approach compared to state-of-the-art imitation and inverse reinforcement learning methods on multiple continuous control tasks.

Backgrounds

Markov Decision Process (MDP): time horizon T ; state space \mathcal{S} ; action space \mathcal{A} ; transition dynamics $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$; reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$; initial state distribution $\eta : \mathcal{S} \rightarrow [0, 1]$; trajectory τ , a sequence of state action pairs for one episode.

Inverse RL Basic Principle: find a reward function r_ω that explains the expert behaviors. (*ill-defined problem*)

Maximum Entropy Inverse RL (MaxEnt IRL) (Ziebart et al., 2008) provides a general probabilistic framework that solves the reward ambiguity problem:

$$p_\omega(\tau) \propto \left[\eta(s_1) \prod_{t=1}^T P(s_{t+1}|s_t, a_t) \right] \exp \left(\sum_{t=1}^T r_\omega(s_t, a_t) \right), \max_{\omega} \mathbb{E}_{\pi_E} [\log p_\omega(\tau)] = \mathbb{E}_{\tau \sim \pi_E} \left[\sum_{t=1}^T r_\omega(s_t, a_t) \right] - \log Z_\omega$$

where Z_ω is the *intractable* partition function, i.e., an integral over all possible trajectories.

Adversarial Inverse RL (AIRL) (Fu et al., 2017) provides an efficient sampling-based approximation to MaxEnt IRL, with a special parameterization for discriminator that allows us to extract reward functions at optimality:

$$D_{\omega, \phi}(s, a, s') = \frac{\exp(f_{\omega, \phi}(s, a, s'))}{\exp(f_{\omega, \phi}(s, a, s')) + \pi(a|s)}, f_{\omega, \phi}(s, a, s') = r_\omega(s, a) + \gamma h_\phi(s') - h_\phi(s)$$

Under certain conditions, $r_\omega(s, a)$ is guaranteed to recover the ground-truth reward up to a constant.

Context-based Meta-Learning & Inverse Reinforcement Learning

Generalizing the notion of MDP with a probabilistic context variables $m \in \mathcal{M}$, where \mathcal{M} is the (discrete or continuous) value space of m . We use $p(m)$ to denote the prior distribution over m .

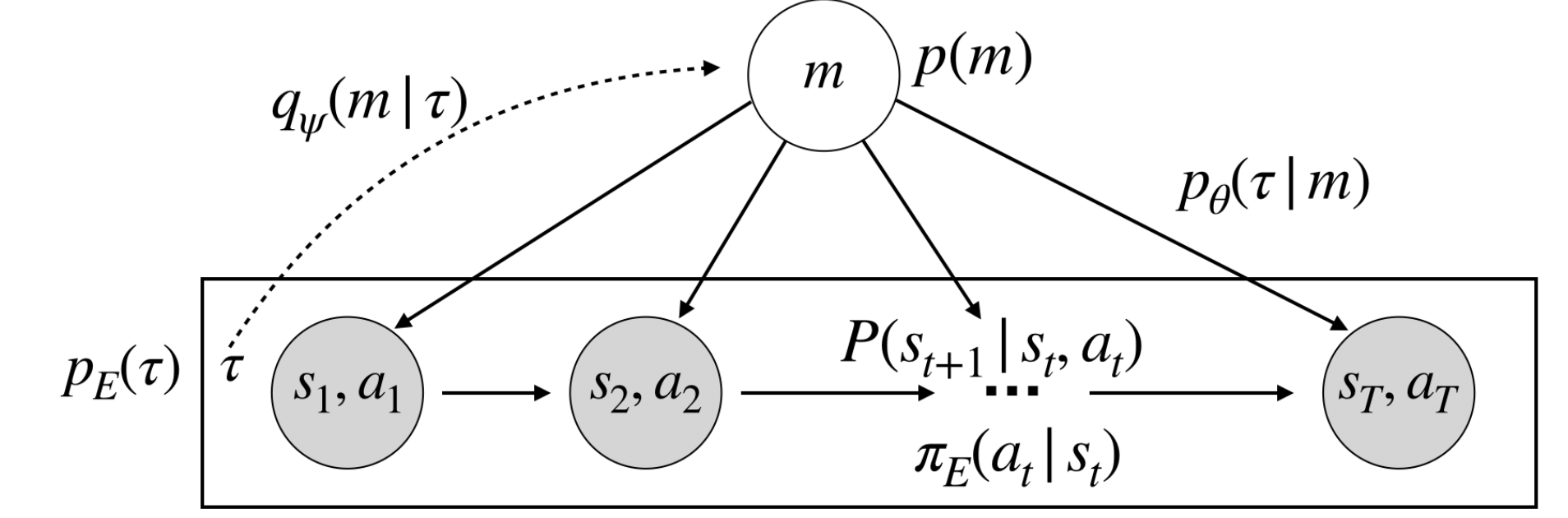
- Context-dependent reward function $r : \mathcal{S} \times \mathcal{A} \times \mathcal{M} \rightarrow \mathbb{R}$; Context-dependent policy $\pi : \mathcal{S} \times \mathcal{M} \rightarrow \mathcal{P}(\mathcal{A})$.
- Expert policy: $\pi_E = \arg \max_{\pi} \mathbb{E}_{m \sim p(m), (s_{1:T}, a_{1:T}) \sim \pi(\cdot|m)} \left[\sum_{t=1}^T r(s_t, a_t, m) - \log \pi(a_t|s_t, m) \right]$
- Marginal trajectory distribution of expert: $p_{\pi_E}(\tau) = \int_{\mathcal{M}} p(m) \prod_{t=1}^T \pi_E(a_t|s_t, m) P(s_{t+1}|s_t, a_t) dm$

Meta-Inverse Reinforcement Learning (Meta-IRL): Given a set of **unstructured** demonstrations *i.i.d.* sampled from the expert policy-induced marginal distribution $p_{\pi_E}(\tau)$, the goal is to meta-learn an inference model $q(m|\tau)$ and a reward function $f(s, a, m)$, such that given some new demonstration τ_E generated by sampling $m' \sim p(m)$, $\tau_E \sim p_{\pi_E}(\tau|m')$, with \hat{m} being inferred as $\hat{m} \sim q(m|\tau_E)$, the learned reward function $f(s, a, \hat{m})$ and the ground-truth reward $r(s, a, m')$ will induce the same set of optimal policies.

Meta-IRL with Probabilistic Context Variables

Under the framework of MaxEnt IRL, we first parametrize two components:

- Context variable inference model $q_\psi(m|\tau)$.
- Context-dependent reward function $f_\theta(s, a, m)$.



We would like to maximize the mutual information between two random variables m and τ under joint distribution $p_\theta(m, \tau) = p(m)p_\theta(\tau|m)$:

$$I_{p_\theta}(m; \tau) = \mathbb{E}_{m \sim p(m), \tau \sim p_\theta(\tau|m)} [\log p_\theta(m|\tau) - \log p(m)]$$

subject to two consistency constraints:

- Desideratum 1. Matching conditional distributions: $\mathbb{E}_{p(m)} [D_{\text{KL}}(p_{\pi_E}(\tau|m) || p_\theta(\tau|m))] = 0$
- Desideratum 2. Matching posterior distributions: $\mathbb{E}_{p_\theta(\tau)} [D_{\text{KL}}(p_\theta(m|\tau) || q_\psi(m|\tau))] = 0$

With Lagrangian duality and Lagrangian multipliers taking specific values, we have the relaxed problem:

$$\begin{aligned} \min_{\theta, \psi} \mathbb{E}_{p(m)} [D_{\text{KL}}(p_{\pi_E}(\tau|m) || p_\theta(\tau|m))] + \mathbb{E}_{p_\theta(m, \tau)} \left[\log \frac{p(m)}{p_\theta(m|\tau)} + \log \frac{p_\theta(m|\tau)}{q_\psi(m|\tau)} \right] \\ \equiv \max_{\theta, \psi} -\mathbb{E}_{p(m)} [D_{\text{KL}}(p_{\pi_E}(\tau|m) || p_\theta(\tau|m))] + \mathbb{E}_{m \sim p(m), \tau \sim p_\theta(\tau|m)} [\log q_\psi(m|\tau)] \end{aligned}$$

We can leverage adversarial reward learning (AIRL) and sampling-based gradient estimation to achieve tractability for optimizing the first and second term in above objective respectively (formally derived in **Section 3** in the paper).

Experiments

Empirical evaluations in various continuous control tasks demonstrate the effectiveness of our framework:

- It can learn a policy with competitive few-shot generalization abilities compared to one-shot imitation learning methods using only unstructured demonstrations.
- It can efficiently infer robust reward functions of new continuous control tasks where one-shot imitation learning fails to generalize, enabling an agent to continue to improve with more trials.



Figure: **Experimental domains** (left to right): Point-Maze, Ant, Sweeper, and Sawyer Pusher.

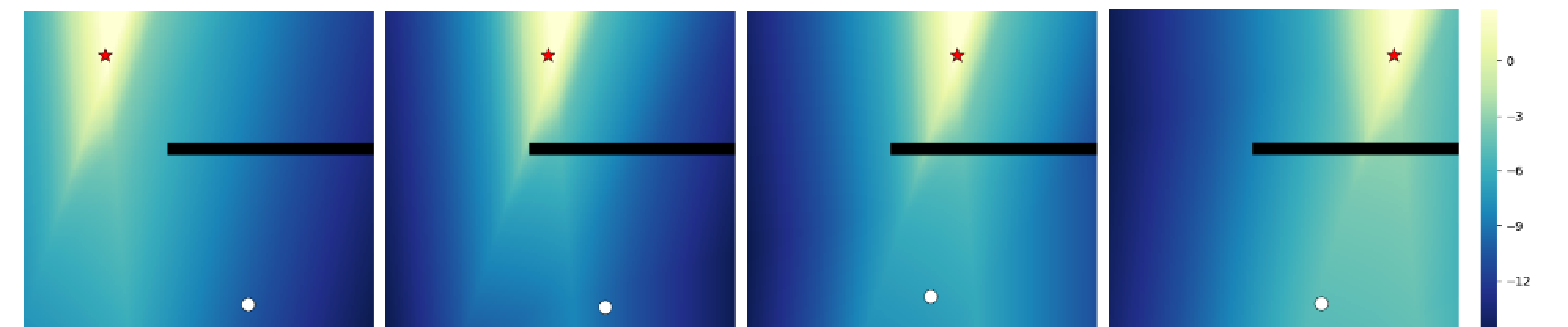


Figure: Visualizations of meta-learned reward functions for Point-Maze Navigation.