| latent dim. | return |
|---|---|
| 1 | $-10.58 \pm 1.27$ |
| 3 | $-14.13 \pm 1.21$ |
| 5 | $-15.41 \pm 1.40$ |

| method | return |
|---|---|
| PEMIRL w/o MI | $-39.24 \pm 3.48$ |
| PEMIRL | $-14.13 \pm 1.21$ |

Table 1: PEMIRL is robust to various latent dimensions.

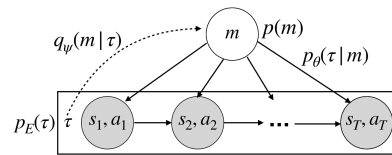Table 2: The MI term is important for training PEMIRL.



Figure 1: Graphical model underlying PEMIRL

1 We thank all the reviewers for the constructive feedback. We will incorporate the valuable suggestions in the revised
2 version. We have conducted more experiments and addressed all of the comments below:

3 **To Reviewer #2:**

4 **Q1: The importance of mutual information (MI) term?** We conducted an ablation study on the MI term with the
5 Point-Maze-Shift environment. The reward function learned without MI failed to induce a good policy in the reward
6 adaptation setting. Results in Table 2 (on the top) demonstrates the importance of our MI term.
7 Theoretically, without MI regularization, the resulting method indeed resembles a VAE. As analyzed in [36], the ELBO
8 of VAE can be interpreted as enforcing consistency between $p(m)p_\theta(\tau|m)$ and $p_E(\tau)q_\psi(m|\tau)$ by minimizing the KL
9 divergence between these joint distributions. Without maximizing the MI between $m$ and $\tau$, a simple degenerate case is
10 $p_\theta(\tau|m) = p_E(\tau)$ and $q_\psi(m|\tau) = p(m)$, which satisfies the consistency constraints, yet completely fails to capture the
11 dependencies between $m$ and $\tau$.

12 **Q2: What if latent dimension is mis-specified?** We conducted additional experiments with the Point-Maze-Shift
13 environment (where the ground-truth latent dimension is 3). See the results in Table 1 (on the top). We can observe
14 that PEMIRL with various latent dimension specifications all outperform the best baseline (return -28.61) stably and is
15 hence robust to dimension mis-specifications.

16 **Q3: Performance on a stochastic environment?** We create a stochastic version of Point-Maze-Shift (maze size:
17 $60 \times 100$ cm) by changing its deterministic transition dynamics into a stochastic one. Specifically, $p(s_{t+1}|s_t, a_t)$ is
18 now realized as a Gaussian with standard deviation being 1 cm. The average return of PEMIRL in reward adaptation is
19 $-17.39 \pm 0.84$, which outperforms the best baseline (average return $-30.58$) by a large margin.

20 **Q4: Test generalization in more realistic environments?** We will add an experiment with a simulated Sawyer robot
21 button pressing task to the revised version, which we were unable to complete during the rebuttal period.

22 **To Reviewer #3:**

23 **Q1: Discussion on data efficiency?** We would like to clarify that in reward adaptation, we use the inferred reward
24 function to train a policy from scratch rather than finetuning the learned policy. Although efficiency is not the focus of
25 this work, we are happy to provide more discussions on this aspect in the revised version. The sample complexity of
26 PEMIRL at meta-testing phase is comparable to RL training with the oracle ground-truth reward, *e.g.* (PEMIRL vs RL
27 with oracle reward): Point-Maze-Shift: 5.4M vs 4M simulation steps; Disabled-Ant: 15M vs 18M simulation steps.

28 **Q2: Can the tasks also change the dynamics during training?** In principle, our algorithm can also handle changes
29 in dynamics during meta-training. We leave this as an interesting avenue for future work.

30 **Q3: The meaning of unstructured demonstrations?** As described in line 58-59, "unstructured" means the demon-
31 strations are not grouped according to the task or labeled by task-specific variables. To elaborate, as discussed in line
32 196-199, previous Meta-IRL methods [12, 32] make simplifying assumptions that each provided expert demonstration
33 contains its corresponding task information (hence "structured"), while PEMIRL has to learn to infer the underlying
34 task corresponding to each demonstration. We will rephrase corresponding parts to clarify it in the revised version.

35 **Q4: Minor comments (1)** We will revise the captions to make them more informative. Policy generalization examines
36 if the policy learned by Meta-IL is able to generalize to new tasks with new dynamics. **(2 & 3)** [11, 23] focus on
37 standard IRL and meta-RL respectively rather than Meta-IRL as in PEMIRL. Although [32] focuses on Meta-IRL, their
38 method derivation (*e.g.* Eq 5) requires a tabular MDP. We will rephrase corresponding parts to make this clear.

39 **To Reviewer #5**

40 **Q1: Discussion on the efficiency of the proposed method?** Although efficiency is not the focus of this work, we are
41 happy to provide more discussions on this aspect in the revised version. During meta-training, for the Point-Maze
42 environment, it takes about 32M simulation steps to converge (similar to other methods such as Meta-InfoGAIL
43 that takes 28M), which amounts to about 2 hours on one Nvidia Titan-Xp GPU; for the Ant environment, it takes
44 about 13.8M simulation steps (Meta-InfoGAIL takes 12M) and about 40 hours on the same hardware (the state-action
45 dimension is much larger than that of Point-Maze). For the sample complexity of meta-testing phase, please refer to the
46 response to Q1 for reviewer #3.

47 **Q2: Graphical model illustration?** We will add the graphical model illustration in Figure 1 to the revised version.