

Robust Inverse Reinforcement Learning under Transition Dynamics Mismatch

L Viano, YT Huang, P Kamalaruban, A Weller, V Cevher

lions@epfl

EPFL

The
Alan Turing
Institute



UNIVERSITY OF
CAMBRIDGE



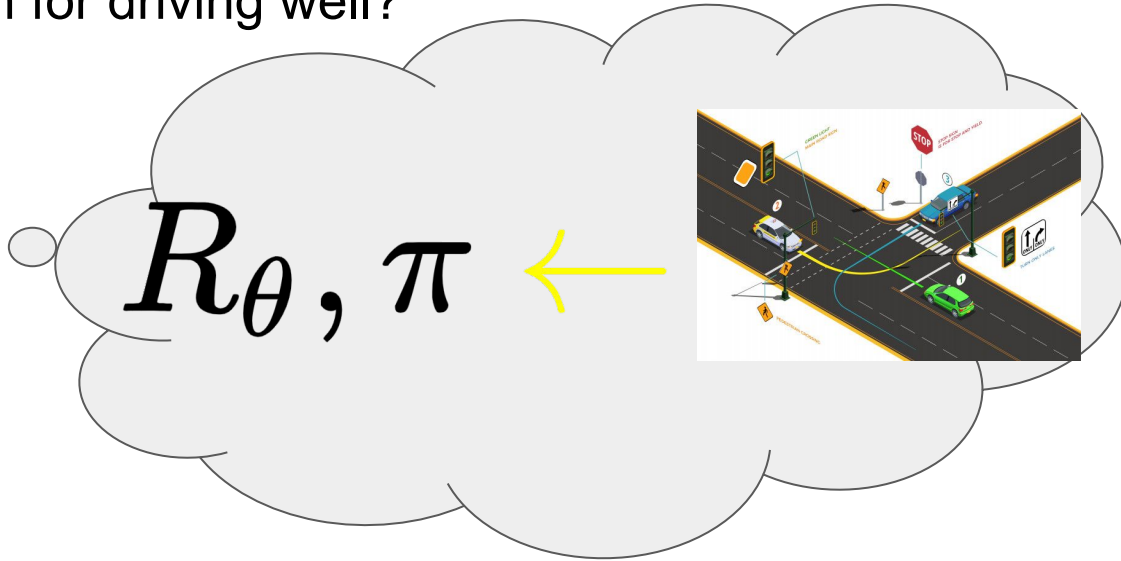
European
Research
Council



HASLERSTIFTUNG

Why Inverse Reinforcement Learning (IRL)?

Which is the reward function for driving well?



Maximum Causal Entropy IRL (MCE-IRL)



Are features meant to capture purposeful characteristics of the observed expert behaviour

Linear Reward Function

$$R_{\theta}(s) = \langle \theta, \phi(s) \rangle \quad \phi(s) \in \mathbb{R}^d$$

If the transition dynamics are known, then $P[S_t = s \mid \pi, M]$ is known and we can compute:

$$\rho_M^{\pi} := (1 - \gamma) \sum_{s \in \mathcal{S}} \sum_{t=0}^{\infty} \gamma^t P[S_t = s \mid \pi, M] \phi(s)$$

known as [occupancy measure](#)

$$V_{M_{\theta}}^{\pi} = \frac{1}{1-\gamma} \langle \theta, \rho_M^{\pi} \rangle$$

$$M_{\theta} = \{\mathcal{S}, \mathcal{A}, T, \gamma, P_0, R_{\theta}\}$$

$$M = M_{\theta} \setminus R_{\theta}$$

Assumptions

- Linear Reward Function : $R_{\theta}(s) = \langle \theta, \phi(s) \rangle$ $\phi(s) \in \mathbb{R}^d$
- We introduce the **occupancy measure** as: $\rho_M^{\pi} := (1 - \gamma) \sum_{s \in \mathcal{S}} \sum_{t=0}^{\infty} \gamma^t P[S_t = s \mid \pi, M] \phi(s)$
- The Value function is also linear: $V_{M_{\theta}}^{\pi} = \frac{1}{1-\gamma} \langle \theta, \rho_M^{\pi} \rangle$

For any θ , policies with same occupancy measure achieves the same value function.

Maximum Causal Entropy IRL (MCE-IRL)

$$V_{M_\theta}^\pi = \frac{1}{1-\gamma} \langle \boldsymbol{\theta}, \boldsymbol{\rho}_M^\pi \rangle$$

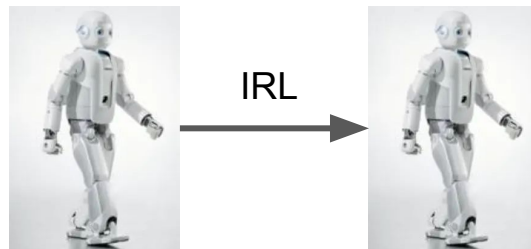
$$\operatorname{argmax}_{\pi \in \Pi} \mathbb{E} \left[\sum_{t=0}^{\infty} -\gamma^t \log \pi(a_t | s_t) \mid \pi, M^L \right]$$

subject to

$$\boldsymbol{\rho}_{M^L}^\pi = \boldsymbol{\rho}$$

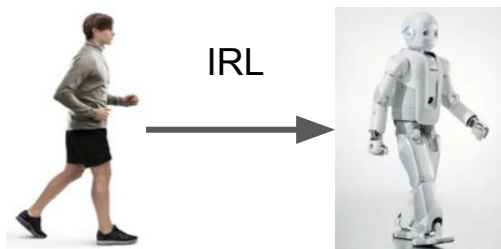
Teacher's Occupancy Measure

Learner's Occupancy Measure

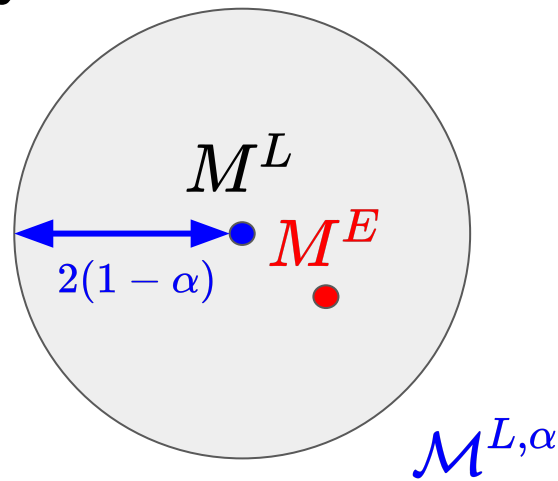


Absence of mismatch

Robust MCE-IRL under transition dynamics mismatch



Mismatch

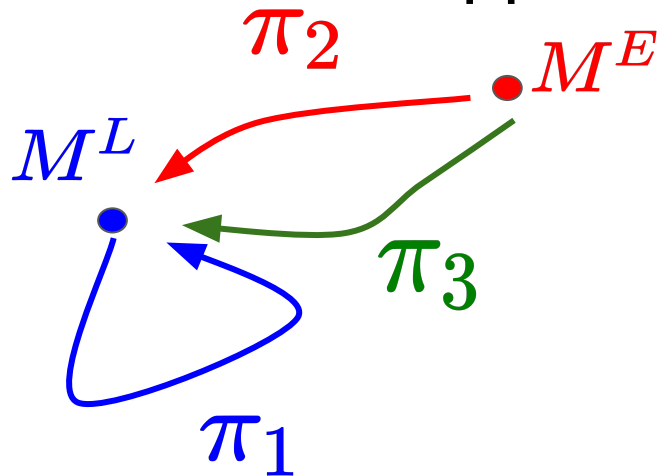


Solve the dual program under the worst case environment **under a fixed** α :

$$\operatorname{argmax}_{\pi \in \Pi} \min_{M \in \mathcal{M}^{L, \alpha}} \mathbb{E} \left[\sum_{t=0}^{\infty} -\gamma^t \log \pi(a_t | s_t) \mid \pi, M \right]$$

subject to $\rho = \rho_M^\pi$

Worst case upper bound for MCE-IRL



$$d(T^E, T^L) := \max_{s,a} \|T^E(\cdot|s,a) - T^L(\cdot|s,a)\|_1$$

π_2 MCE IRL

π_3 ROBUST MCE IRL

$$|V_{T^L}^{\pi_1} - V_{T^L}^{\pi_2}| \leq \mathcal{O}(d(T^E, T^L))$$

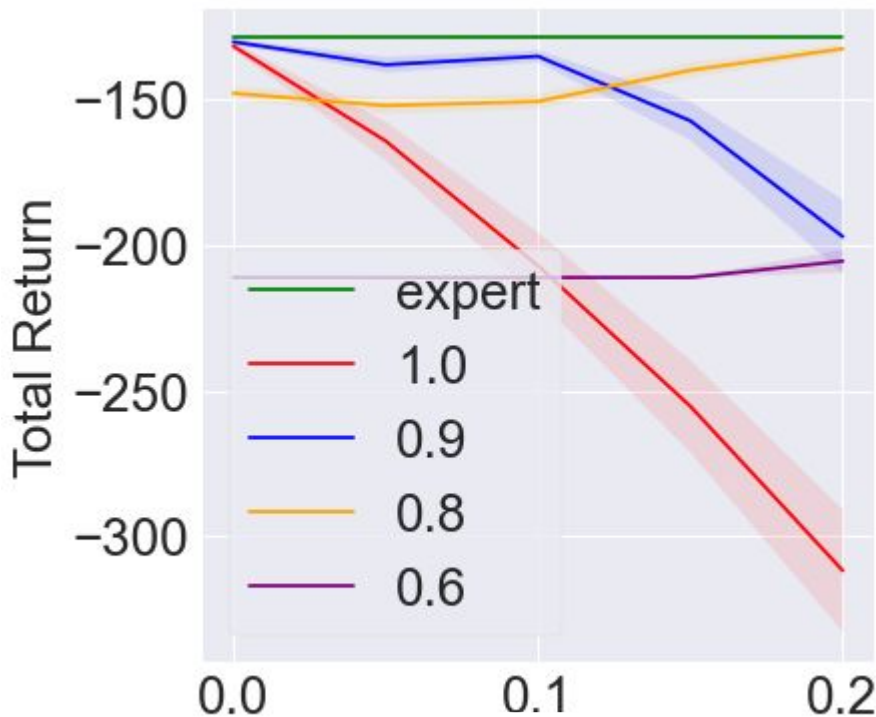
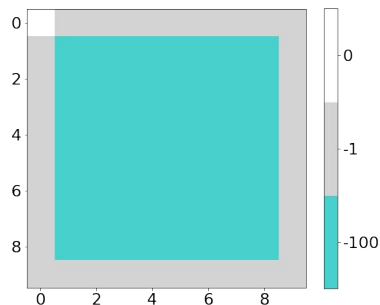
$$|V_{T^L}^{\pi_1} - V_{T^L}^{\pi_3}| \leq \mathcal{O}((1 + \alpha)d(T^E, T^L) + (1 - \alpha)d(T^E, T^*) + 2(1 - \alpha)^2)$$

Results with Linear Reward Function

Recall: Standard MCE-IRL (Ziebart, 2010) is recovered with

$$\alpha = 1$$

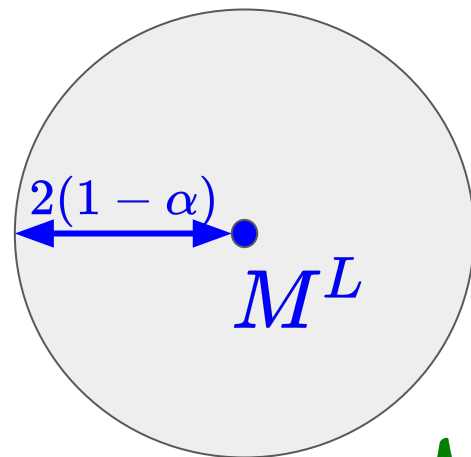
The legend reports the value for α



Small Mismatch

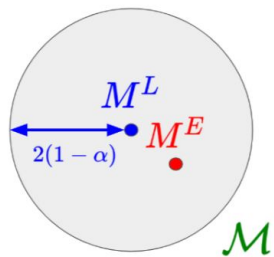
Large Mismatch

MDP uncertainty region

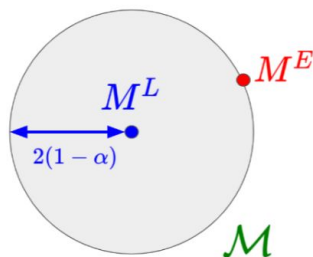


M

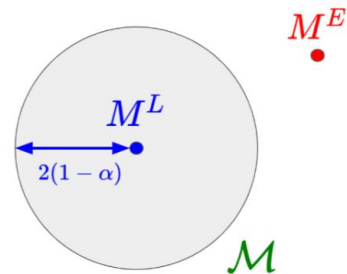
Choosing the right uncertainty set



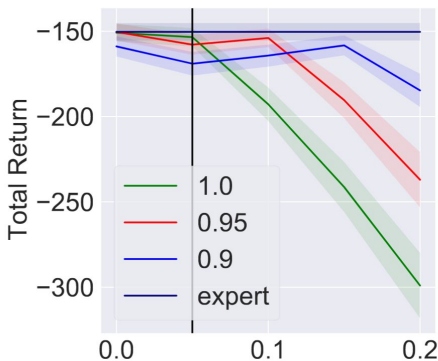
(a) Overestimating $1 - \alpha$



(b) Perfect estimation of $1 - \alpha$



(c) Underestimating $1 - \alpha$

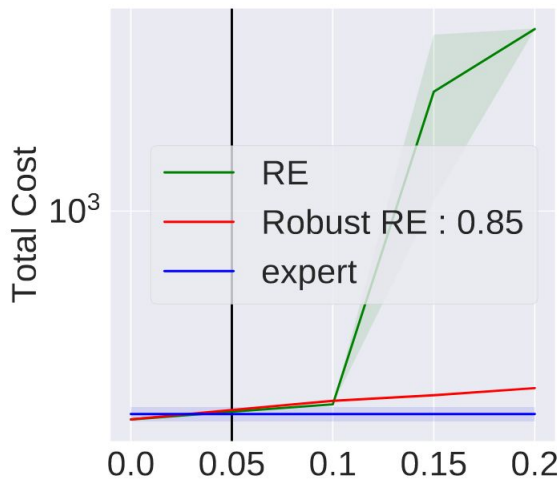
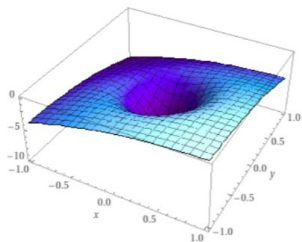


Continuous States and Actions

We propose an extension to continuous states and actions based on Relative Entropy IRL.

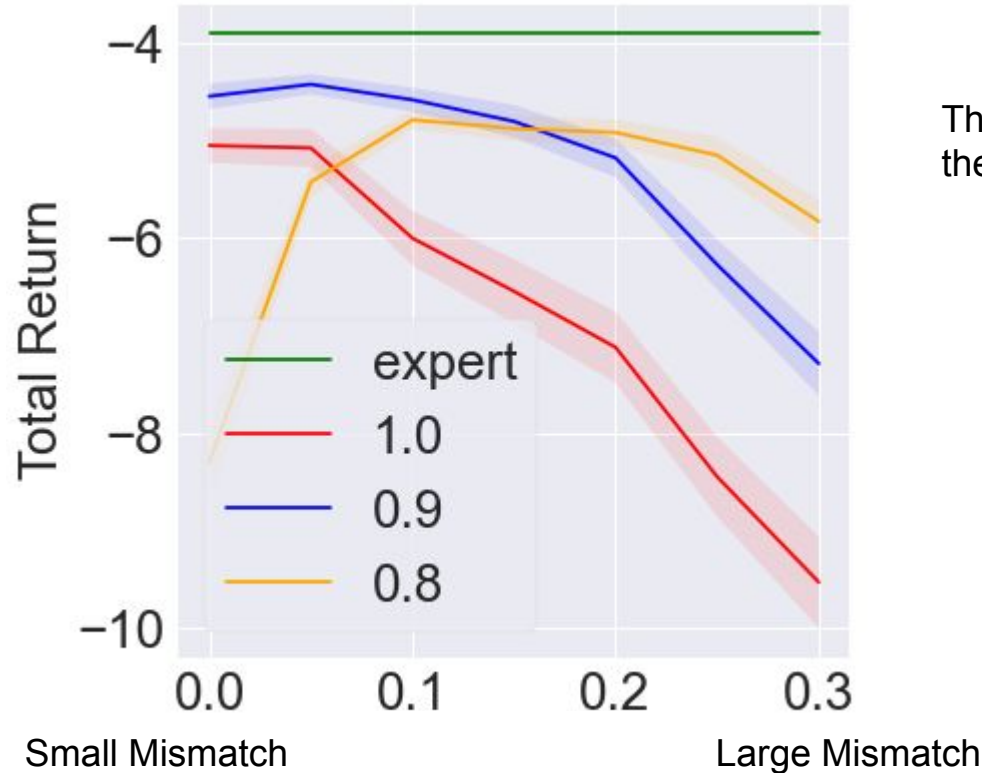
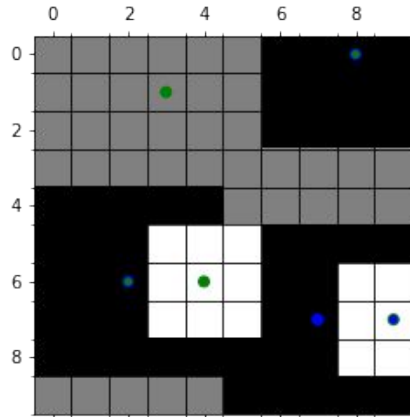
Relative Entropy-IRL
(Boularias et al.,
2011) is recovered
with

$$\alpha = 1$$



Non Linear Reward Function

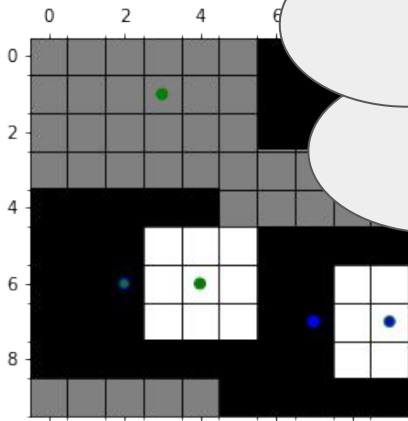
Notice: Standard Deep MCE-IRL (Wulfmeier et al. 2015) is recovered with $\alpha = 1$



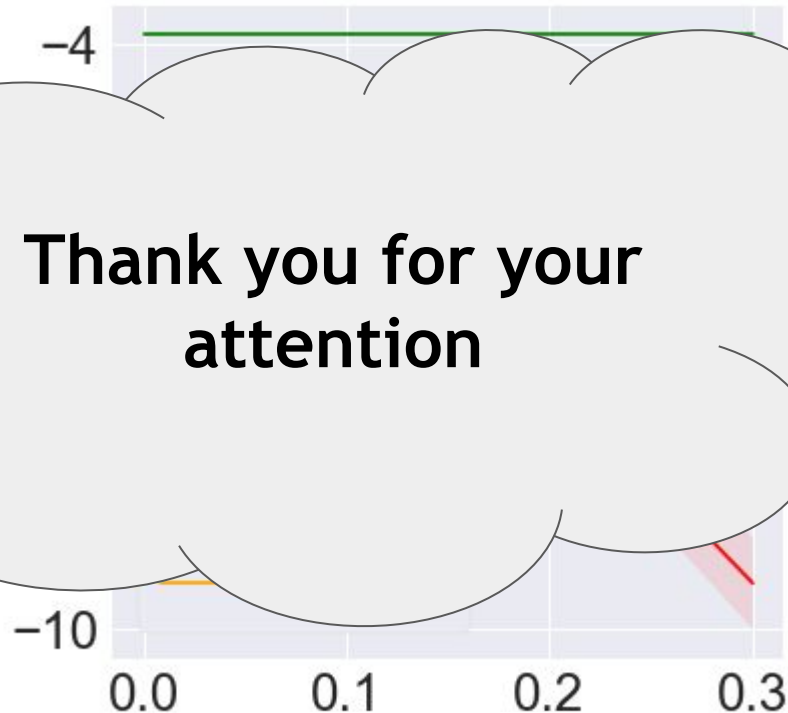
The legend reports the value for α

Non Linear Reward Function

Notice: Standard Deep MCE-IRL (Wulfmeier et al. 2015) is recovered with $\alpha = 1$



The legend reports value for α



Small Mismatch

Large Mismatch