



MIT



RSS Paper

Active Preference-Based Learning of Reward Functions [1]

Nikita Jaipuria

Aerospace Controls Laboratory
Department of Mechanical Engineering
Massachusetts Institute of Technology

August 11, 2017

► Objective:

- model a **human's preference** for how a dynamical system should act

- learn

$$R_H(\xi) = R_H(x^0, \mathbf{u}_R, \mathbf{u}_H) = \sum_{t=0}^N r_H(x^t, u_R^t, u_H^t) = \sum_{t=0}^N \mathbf{w}^T \phi(x^t, u_R^t, u_H^t) = \mathbf{w}^T \Phi(\xi)$$

► Problem Domain:

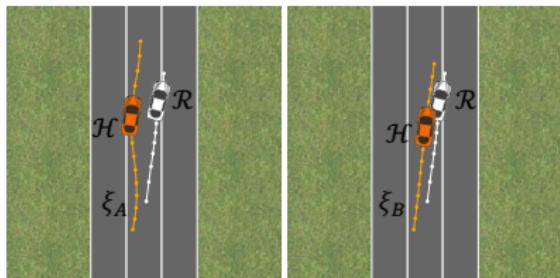
- difficult to provide demonstrations of **desired** system trajectory (IRL)
- assign numerical reward to an action/trajectory

► Main Idea: active preference-based learning

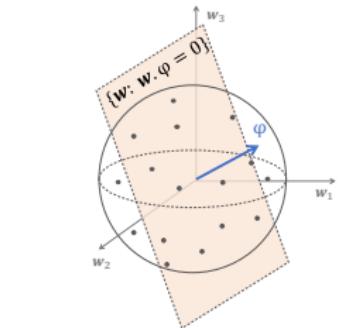
- system decides on what preference queries to make (**active**)
- build on label ranking; learn from preferences/comparisons (**preference-based**)

► Challenges/Contribution

- complexity and continuous nature of **queries**
- **active synthesis** of queries satisfying system dynamics: $x^{t+1} = f_{HR}(x^t, u_R^t, u_H^t)$ using continuous optimization
- **maximize volume removed** from continuous hypothesis space of reward functions by each query

$\xi_A \text{ or } \xi_B \rightarrow I_t$


(a) Preference query.



(b) Query response effect.

Algorithm 1 Preference-Based Learning of Reward Functions

- 1: **Input:** Features ϕ , horizon N , dynamics f , $iter$
- 2: **Output:** Distribution of \mathbf{w} : $p(\mathbf{w})$
- 3: Initialize $p(\mathbf{w}) \sim \text{Uniform}(B)$, for a unit ball B
- 4: **While** $t < iter$:
- 5: $W \leftarrow M$ samples from AdaptiveMetropolis($p(\mathbf{w})$)
- 6: $(x^0, \mathbf{u}_R, \mathbf{u}_H^A, \mathbf{u}_H^B) \leftarrow \text{SynthExps}(W, f)$
- 7: $I_t \leftarrow \text{QueryHuman}(x^0, \mathbf{u}_R, \mathbf{u}_H^A, \mathbf{u}_H^B)$
- 8: $\varphi = \Phi(x^0, \mathbf{u}_R, \mathbf{u}_H^A) - \Phi(x^0, \mathbf{u}_R, \mathbf{u}_H^B)$
- 9: $f_\varphi(\mathbf{w}) = \min(1, I_t \exp(\mathbf{w}^\top \varphi))$
- 10: $p(\mathbf{w}) \leftarrow p(\mathbf{w}) \cdot f_\varphi(\mathbf{w})$
- 11: $t \leftarrow t + 1$
- 12: **End for**

Results

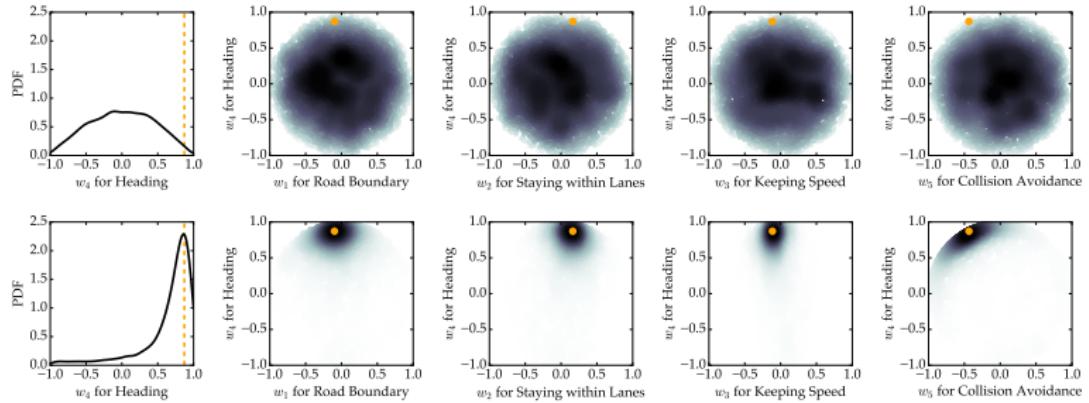
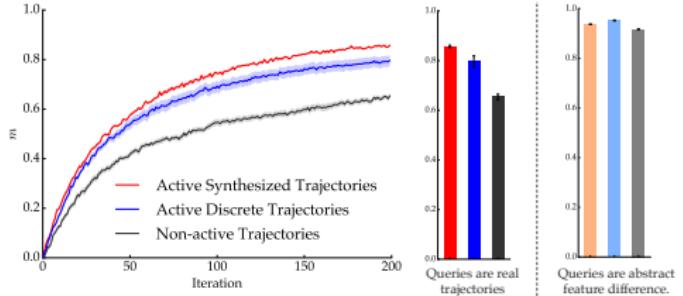


Fig. 2: Distribution of w_4 , the weight for the heading feature, relative to the other features. The top plots shows the starting distribution, and the bottom plot shows the distribution at convergence. The orange dot and dotted line show the ground truth for the weights.



► What is novel/interesting?

- preference-based approach allows users to compare two trajectories in various scenarios without having to demonstrate a full trajectory (of how they would like to drive, not how they actually drive)
- active enables choosing informative test cases otherwise difficult to encounter in driving scenarios; addresses the limitation of training data sparsity/informativeness
- In-depth investigation of the effect of **synthesis** of real trajectories for dynamical systems compared to relying on a discrete set

► Limitations

- Markov assumption is inherently built into system dynamics
- Feature selection is not active/based on expert knowledge
- Reward function is constrained to be linear in the feature space
- The entire formulation is for a human H with only one other robot R
- Formulation of w_{true} was unclear
- Experiments:
 - potential functions' based reward function
 - f_{HR} assumed as a simple point mass-dynamics model

Questions?

References I



MIT



- [1] Anca D Dragan Dorsa Sadigh, Shankar Sastry, and Sanjit A Seshia. Active preference-based learning of reward functions.