

MAX-PLANCK-GESELLSCHAFT



A Markov Decision Process is a tuple $(\mathcal{S}, \mathcal{A}, T, R)$, where

- \mathcal{S} is a set of states, \mathcal{A} is a set of actions,
- T are transition probabilities with T(s, a, s') = P(s'|s, a) for $s, s' \in \mathcal{S}, a \in A$,
- R is a reward function where R(s, a) is the reward given for action a in state s.

Policies and Value Functions 2.2

- A policy π is a function that maps every state into an action.
- The value of policy π is given by $V(\pi) = \mathbb{E}\left[\sum_{t=0}^{H} \gamma^t R(s_t, a_t) | \mu_0, \pi, T\right]$, where μ_0 is an initial state distribution and H is a horizon.

Apprenticeship Learning $\mathbf{2.3}$

- Designing a reward function for matching a complex behavior can be a challenging problem. It is often easier to provide examples of the desired behavior [1].
- Apprenticeship Learning via Inverse Reinforcement Learning (IRL) consists in learning a reward function that explains an observed behavior.
- The reward is usually assumed to be a linear function of state-action features ϕ_i ,

$$R(s,a) = \sum_{i} \theta_i \phi_i(s,a).$$

• The learned reward function, parameterized by θ , is used to find a policy that generalizes the observed behavior.

Structured Apprenticeship Learning Oliver Krömer² Jan Peters^{1,2}

Abdeslam Boularias¹

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany

²Darmstadt University of Technology, Darmstadt, Germany

Structured Apprenticeship Learning

3.1 Key Insight

- In practice, it is often difficult to find appropriate reward features. These features are usually obtained from empirical data and are subject to measurement errors.
- Most real-world problems have a certain structure wherein states that are close to each other tend to have the same optimal action.
- We create a graph that connects similar states together. This graph can be constructed by using the Euclidean distance, for instance, as a similarity measure.



Markov Decision Process with a Similarity Graph

• We denote the edges of this graph by \mathcal{E} and the edge features by ψ .

3.2 Problem Statement

Structured Apprenticeship Learning is formulated as the problem of maximizing the entropy of a distribution on policies P,

$$\max_{P,\mu^{\pi}} \bigg(-\sum_{\pi \in \mathcal{A}^{|\mathcal{S}|}} P(\pi) \log P(\pi) \bigg|_{\mathcal{S}} \bigg|_{\mathcal{S}}$$

subject to the following constraints

$$\begin{array}{l} & \sum_{\pi \in \mathcal{A}^{|\mathcal{S}|}} P(\pi) \\ \forall \phi_k : \sum_{\pi \in \mathcal{A}^{|\mathcal{S}|}} P(\pi) \sum_{s \in \mathcal{S}} \mu^{\pi}(s) \phi_k(s, \pi(s)) \\ \forall \psi_k : \sum_{(s_i, s_j) \in \mathcal{E}} \psi_k(s_i, s_j) \sum_{\pi, \pi(s_i) = \pi(s_j)} P(\pi) \end{array}$$

• ϕ and ψ are the empirical averages of the reward and graph features respectively, calculated from the demonstration, and $\mu^{\pi}(s) = \mu_0(s) + \gamma \sum_{s'} T(s, \pi(s), s') \mu^{\pi}(s')$.

3.3 Solution

$$P(\pi) \propto \exp\left(\sum_{s} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{usual value function}} \int_{\text{s.t.}} \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \mu^{\pi}(s) \sum_{k} \theta_{k} \phi_{k}(s, \pi(s)) + \int_{\text{s.t.}} \frac{1}{2} \left(\sum_{k} \theta_{k}(s, \pi(s)) + \int$$

- Parameters θ and λ are learned by maximizing the likelihood of the demonstration.
- An optimal policy $\pi^* = \arg \max_{\pi \in \mathcal{A}^{|\mathcal{S}|}} P(\pi)$ is found by dynamic programing.

3.4	Re	lation to other method	S
$ \mathcal{E} = 0$			
γ :	= 0	Logistic regression	Associative
γ :	$\neq 0$	Maximum Entropy IRL [2]	Structure

• • •





 $|\mathcal{E}| \neq 0$ Markov Networks (AMN) [4] ed Apprenticeship Learning



Blue indicates a low value of a feature associated with negative reward, and red indicates a higher value of that feature. In the testing domain, a noise is added to the negative-weighted feature of randomly chosen states. The figure shows the average true rewards of MaxEnt and SAL as a function of the percentage of noisy states.

4.2 Grasping New Objects

From a high-level point of view, grasping an object can be seen as an MDP with three time-steps: reaching the object, preshaping the hand, and grasping.



Learned Q-values at t = 0. Each point on an object corresponds to a reaching action. Blue indicates low values and red indicates high values. The black arrow indicates the approach direction in the optimal policy according to the learned reward function.

Notice how Structured Apprenticeship Learning improves 100 over the other methods by generally giving high values to $_{80}$ handle points only. The confusion in the other methods 60 comes from noise features and self-occlusions of the 3D $_{40}$ point clouds.

References





[1] Pieter Abbeel and Andrew Ng. Apprenticeship Learning via Inverse Reinforcement Learning. ICML 2004. [2] Ziebart, B., Maas, A., Bagnell, A., and Dey, A. Maximum Entropy Inverse Reinforcement Learning. AAAI 2008.

[3] Nathan Ratliff, J. Andrew Bagnell and Martin Zinkevich. Maximum Margin Planning. ICML 2006.

[4] Taskar, B.: Learning Structured Prediction Models: A Large Margin Approach. PhD thesis, Stanford University, 2004.