## Structured Apprenticeship Learning

Abdeslam Boularias<sup>1</sup>, Oliver Kroemer<sup>2</sup> and Jan Peters<sup>1,2</sup>



Max Planck Institute for Intelligent Systems, Germany
Technische Universität Darmstadt, Germany

## Aim: Learning to grasp unknown objects



(a) Barrett robot hand

(b) Unknown object

- A grasp quality (or reward) depends on the configuration of the hand, and the shape of the object in the contact region.
- Defining shapes, e.g. handles, using local geometrical features is not a trivial task.
- Features acquired through sensors are often subject to noise.

#### Apprenticeship Learning via Inverse Reinforcement Learning



Apprenticeship Learning via Inverse Reinforcement Learning

## Given:

- A Markov Decision Process without a reward function
- A demonstration  $\{(s_i, a_i)\}$
- 1 Learn a reward function R such that the demonstrations are generated by an optimal policy
- 2 Use the learned reward to find an optimal policy

#### Notations

A Markov Decision Process (MDP) is defined by:

- $\mathcal{S}$ : states set
- $\mathcal{A}$ : actions set
- T: transition function defined as T(s, a, s') = P(s'|s, a)
- R: reward function where R(s, a) is the reward given for executing action a in state s
- $\mu_0$ : initial state distribution
- $\gamma \in [0,1[:$  discount factor

A policy  $\pi$  is a function that selects an action for each state.

#### Reward features

• The reward is defined as a function of state-action features

$$R(s,a) \stackrel{def}{=} \sum_{k=1}^{n} \theta_k \phi_k(s,a).$$

• The expected value of feature  $\phi_k$  given policy  $\pi$  is defined as

$$\phi_k^{\pi} \stackrel{def}{=} \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t \phi_k(s_t, a_t) | \mu_0, \pi, T\right]$$

 The expected value of policy π is a linear function of the expected feature values

$$V(\pi) = \sum_{k=1}^{n} \theta_k \phi_k^{\pi}.$$

Structured Apprenticeship Learning- Key Insights

- Keward features are often obtained from empirical measurements and are subject to noise.
- Features of complex reward functions, such as a grasp quality, cannot be easily defined.

Structured Apprenticeship Learning- Key Insights

- Keward features are often obtained from empirical measurements and are subject to noise.
- Features of complex reward functions, such as a grasp quality, cannot be easily defined.
- In practice, optimal policies are often structured: states that are close to each other tend to have similar optimal actions.

Structured Apprenticeship Learning- Key Insights

- Keward features are often obtained from empirical measurements and are subject to noise.
- Features of complex reward functions, such as a grasp quality, cannot be easily defined.
- In practice, optimal policies are often structured: states that are close to each other tend to have similar optimal actions.
- Create a graph that connects similar states together, by using the Euclidean distance, for example, as a similarity measure. Notation: *E* is the edge set and ψ<sub>k</sub> are the edge features.
- Constrain the learned policy to often select similar actions in neighboring states.

### Problem Statement

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

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

subject to the following constraints

$$\sum_{\pi \in \mathcal{A}^{|\mathcal{S}|}} P(\pi) = 1,$$
$$\forall \phi_k : \sum_{\pi \in \mathcal{A}^{|\mathcal{S}|}} P(\pi) \sum_{s \in \mathcal{S}} \mu^{\pi}(s) \phi_k(s, \pi(s)) = \hat{\phi}_k,$$
$$\forall \psi_k : \sum_{(s_i, s_j) \in \mathcal{E}} \psi_k(s_i, s_j) \sum_{\pi, \pi(s_i) = \pi(s_j)} P(\pi) = \hat{\psi}_k.$$

 $\hat{\phi}$  and  $\hat{\psi}$  are the empirical values of the reward and edge features respectively, calculated from the demonstration.  $\mu^{\pi}(s) = \mu_0(s) + \gamma \sum_{s'} T(s, \pi(s), s') \mu^{\pi}(s').$ 

#### Solution



- $\theta$  and  $\lambda$  are learned by maximizing the likelihood of the demonstration.
- This distribution corresponds to a Markov Random Field

$$P(\pi) = \frac{1}{Z} \prod_{s \in \mathcal{S}} f_{\pi}(s) \prod_{(s_i, s_j) \in \mathcal{S}^2} g_{\pi}(s_i, s_j).$$

• A policy  $\pi^* = \arg \max_{\pi} P(\pi)$  is found by reducing the planning problem to a sequence of inference in MRFs

#### Grasping as a Markov Decision Process



(c) Barrett hand

(d) Object seen through a depth camera



## Results



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

### Results



Percentage of successful grasps

Figure : Percentage of grasps labeled as successful (out of 7 objects).

- Efficient inference algorithms
- Other applications: autonomous mobile manipulation

# Thank you