

# Active Learning for Reward Estimation in Inverse Reinforcement Learning

Manuel Lopes<sup>1</sup>, Francisco Melo<sup>2</sup>, and Luis Montesano<sup>3</sup>

<sup>1</sup> Instituto de Sistemas e Robótica - Instituto Superior Técnico  
Lisboa, Portugal

`macl@isr.ist.utl.pt`

<sup>2</sup> Carnegie Mellon University  
Pittsburgh, PA, USA

`fmelo@cs.cmu.edu`

<sup>3</sup> Universidad de Zaragoza  
Zaragoza, Spain

`lmontesa@unizar.es`

# Contents

- Introduction
- Preliminaries
- Algorithm
- Experiments
- Conclusions

# Introduction

- **Inverse reinforcement learning** addresses the general problem of recovering a reward function from samples of a policy provided by an expert/demonstrator.
- We introduce **active learning** for **inverse reinforcement learning** to query the demonstrator for samples at specific states, instead of relying only on samples provided at “arbitrary” states.
- Similar accuracy
- Less queries from the expert

# Preliminaries - Bayesian Inverse Reinforcement Learning

- State-action pairs:

$$\mathcal{D} = \{(x_1, a_1), (x_2, a_2), \dots, (x_n, a_n)\}$$

- The likelihood of a pair  $(x, a)$  and  $\mathcal{D}$ :

$$L_r(x, a) = \mathbb{P}[(x, a) \mid r] = \frac{e^{\eta Q_r^*(x, a)}}{\sum_{b \in A} e^{\eta Q_r^*(x, b)}} \quad L_r(\mathcal{D}) = \prod_{(x_i, a_i) \in \mathcal{D}} L_r(x_i, a_i)$$

- Get reward function:

$$\mathbb{P}[r \mid \mathcal{D}] \propto L_r(\mathcal{D}) \mathbb{P}[r] \quad r^* = \max_r \mathbb{P}[r \mid \mathcal{D}]$$

# Preliminaries - Two Methods for Bayesian IRL – Gradient-based IRL

- The log-likelihood of the  $r$ :

$$\Lambda_r(\mathcal{D}) = \sum_{(x_i, a_i) \in \mathcal{D}} \log(L_r(x_i, a_i))$$

- Compute gradient:

$$[\nabla_r \Lambda_r(\mathcal{D})]_{xa} = \sum_{(x_i, a_i) \in \mathcal{D}} \frac{1}{L_r(x_i, a_i)} \frac{\partial L_r(x_i, a_i)}{\partial r_{xa}}$$

- Update  $r$ :

$$\mathbf{r}_{t+1} = \mathbf{r}_t + \alpha_t \nabla_r \Lambda_{r_t}(\mathcal{D})$$

# Algorithm - Active Learning for Reward Estimation

- Define :
  - $\mathcal{R}_{xa}(p)$  as the set of reward functions  $r$ , s.t.  $\pi_r(x, a) = p$
  - Density:  $\bar{\mu}_{xa}(p) = \mathbb{P} [\pi(x, a) = p \mid \mathcal{D}] = \mathbb{P} [r \in \mathcal{R}_{xa}(p) \mid \mathcal{D}]$
  - Differential entropy:  $-\int \bar{\mu}_{xa}(p) \log \bar{\mu}_{xa}(p) dp$
- Define a new discrete probability distribution:
  - $I = [0, 1]$        $I_k = \left(\frac{k}{K}, \frac{k+1}{K}\right]$        $I_0 = \left[0, \frac{1}{K}\right]$
  - $\mu_{xa}(k) = \mathbb{P} [\pi(x, a) \in I_k \mid \mathcal{D}] = \int_{I_k} \bar{\mu}_{xa}(p) dp, \quad k = 1, \dots, K$

# Algorithm

$$\bar{H}(x) = \frac{1}{|\mathcal{A}|} \sum_a H(\mu_{xa}) = -\frac{1}{|\mathcal{A}|} \sum_{a,k} \mu_{xa}(k) \log \mu_{xa}(k)$$

$$x^* = \arg \max_{x \in \mathcal{X}} \bar{H}(x)$$

---

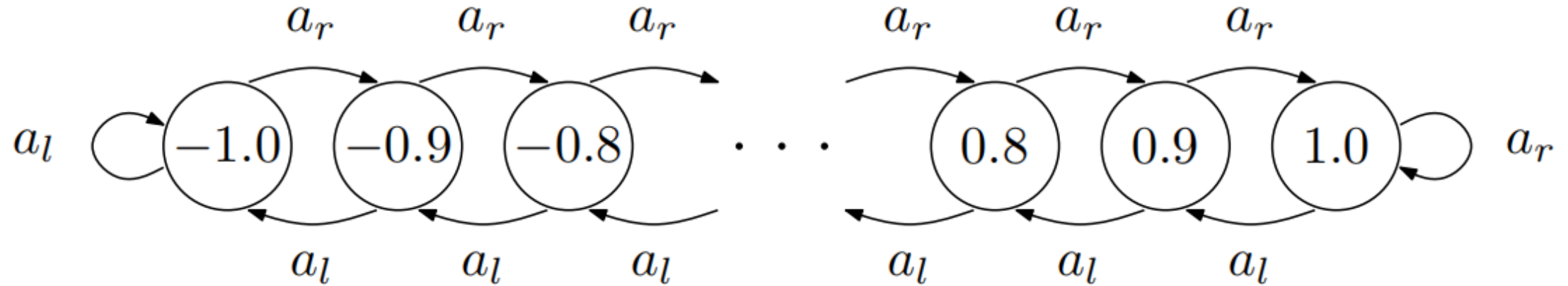
## Algorithm 1. General active IRL algorithm

---

**Require:** Initial demo  $\mathcal{D}$

- 1: Estimate  $\mathbb{P}[r \mid \mathcal{D}]$  using general MC algorithm
  - 2: **for all**  $x \in \mathcal{X}$  **do**
  - 3:   Compute  $\bar{H}(x)$
  - 4: **end for**
  - 5: Query action for  $x^* = \arg \max_x \bar{H}(x)$
  - 6: Add new sample to  $\mathcal{D}$
  - 7: Return to □
-

# Experiments - Finding the Maximum of a Quadratic Function

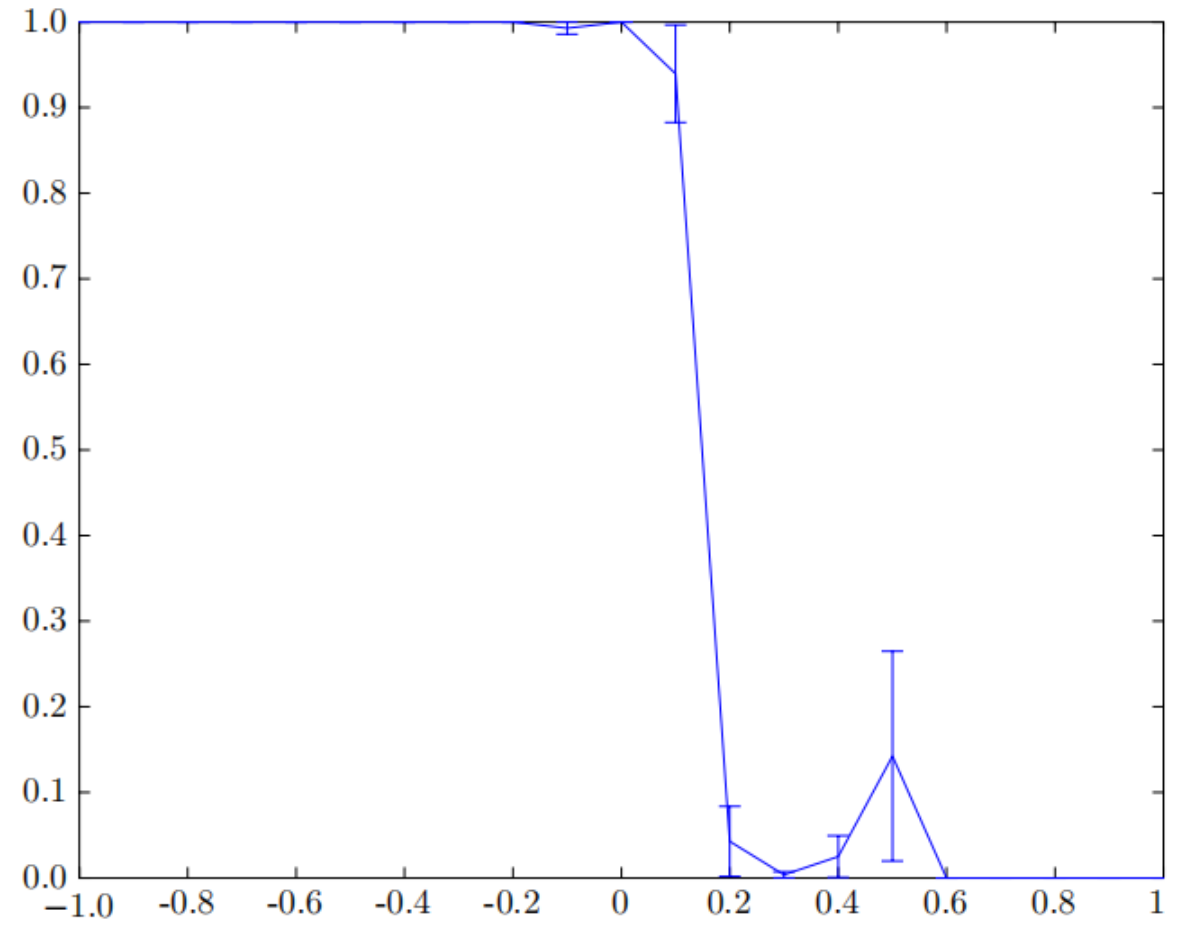
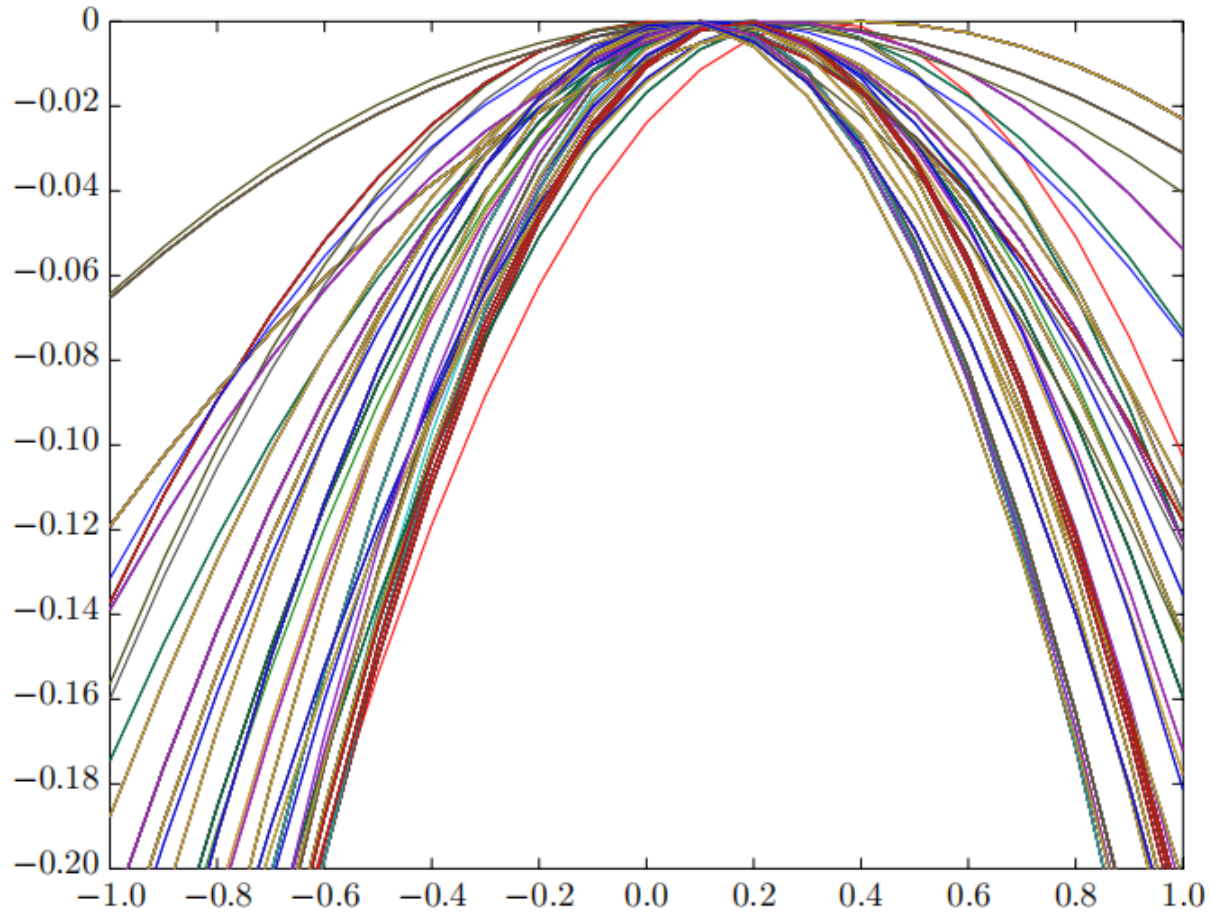


$$r(x) = \theta_1 (x - \theta_2)^2$$

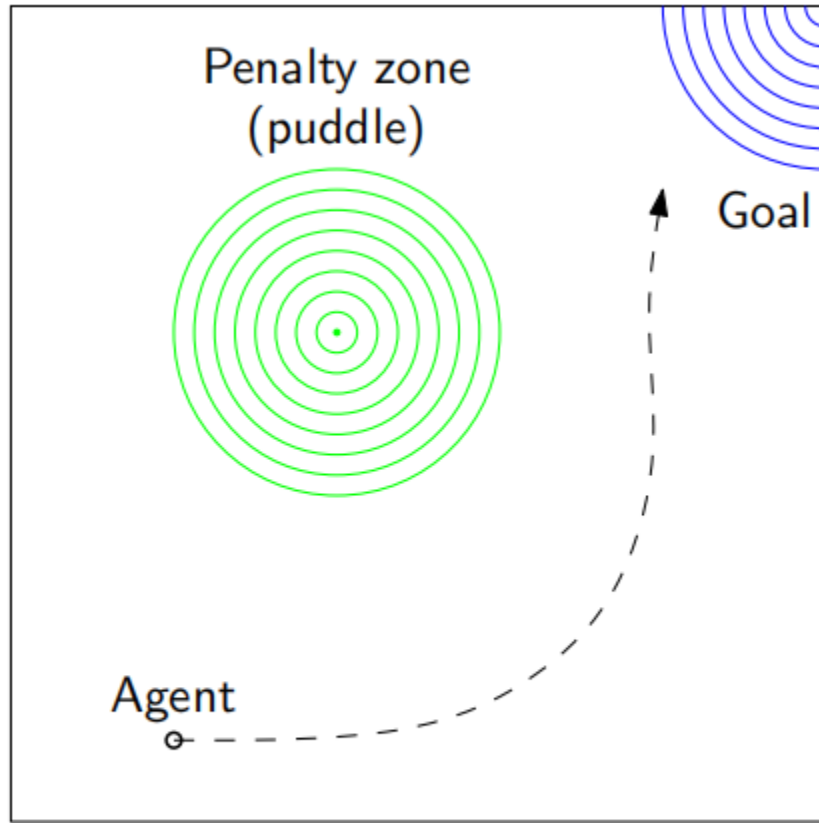
$$r(x) = -(x - 0.15)^2$$

$$\mathcal{D} = \{(-1.0, a_r), (-0.9, a_r), (-0.8, a_r), (0.8, a_l), (0.9, a_l), (1.0, a_l)\}$$

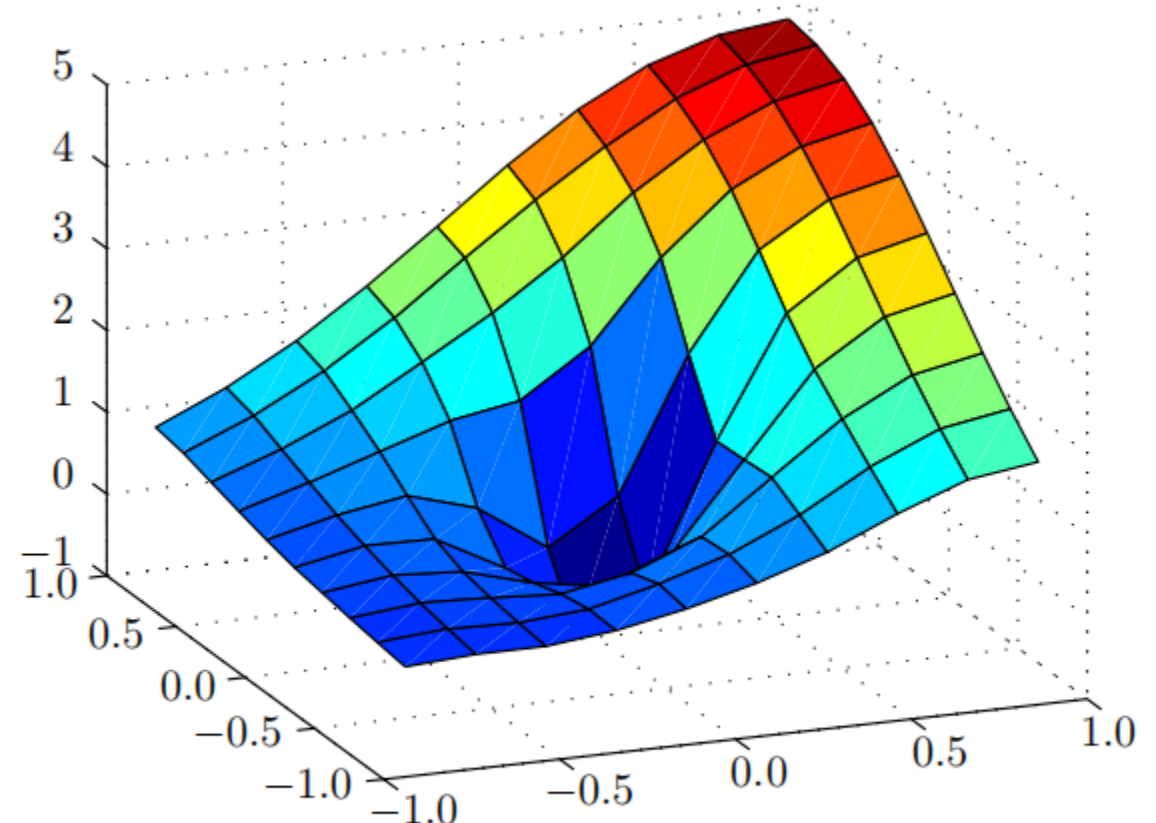
# Experiments - Finding the Maximum of a Quadratic Function



# Experiments - Puddle World

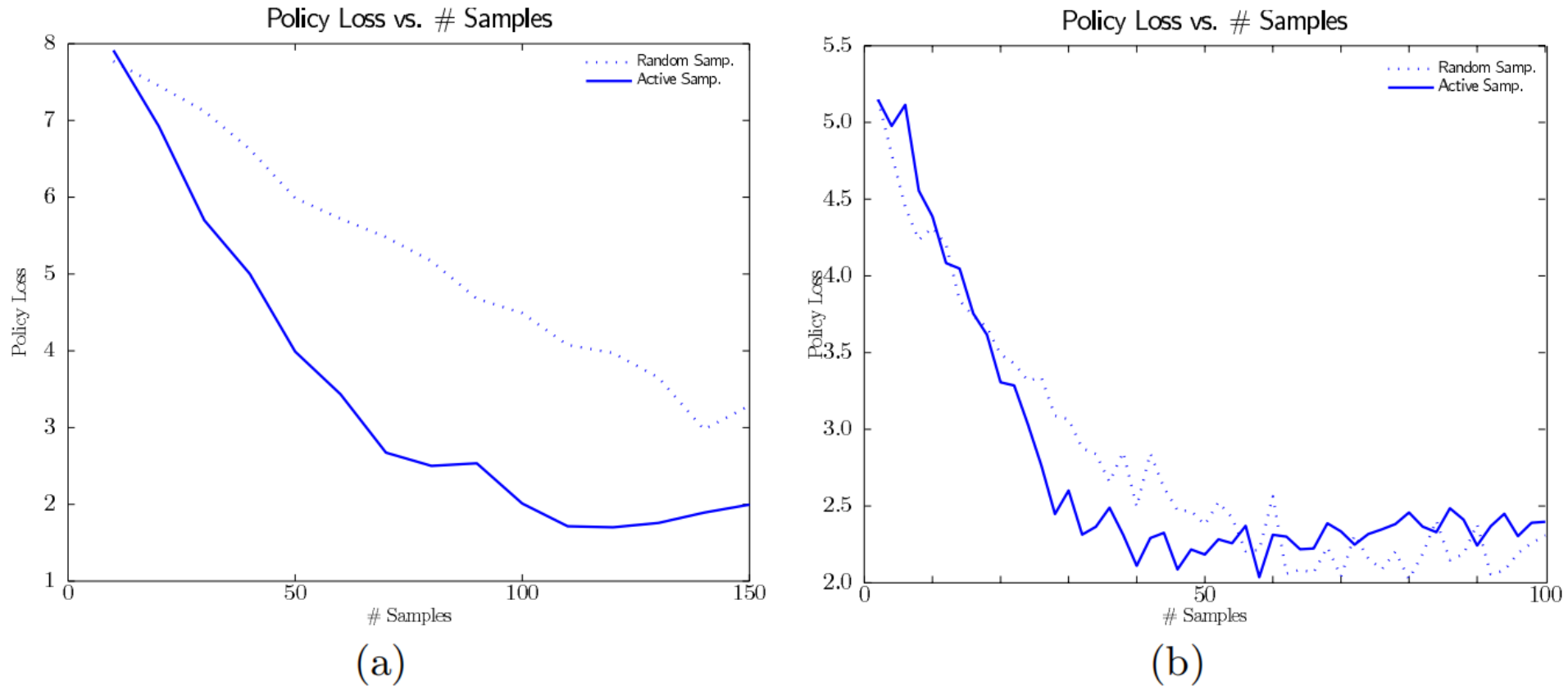


(a)



(b)

# Experiments - Active sampling vs. Random sampling – Grid-World



**Fig. 6.** Performance of Algorithm 2 comparing active sampling vs. random sampling as a function of the demonstration size. (a) Results with parameterized rewards in a  $15 \times 15$  grid-world. (b) Results with general (non-parameterized) rewards in a  $10 \times 10$  grid-world.

## Conclusions

- We introduced the first **Active Learning** algorithm explicitly designed to reward estimation in **Inverse Reinforcement Learning**.
- Successful application in **Bayesian Inverse reinforcement learning**.