

Bayesian Inverse Reinforcement Learning

Presented by Monica Dinculescu

McGill University
February 7, 2007

Markov Decision Processes(MDP)

Definition

An MDP is a tuple (S, A, T, γ, R)

- S is a finite set of states
- A is a finite set of actions
- $T : S \times A \times S \rightarrow [0, 1]$ is the transition probability function
- $\gamma \in [0, 1)$ is the discount factor
- $R : S \rightarrow \mathbb{R}$ is the reward function

Bellman Equations

Definition

A **policy** is a map $\pi : S \rightarrow A$

Definition

The **value** of a policy π under a reward function R

$$V^\pi(s) = R(s) + \gamma \sum_{s'} T(s, \pi(s), s') V^\pi(s')$$

Definition

The **Q-function** of a policy π under a reward function R

$$Q^\pi(s, a) = R(s) + \gamma \sum_{s'} T(s, a, s') V^\pi(s')$$

Motivation

Problem

Learn the reward function of an underlying Markov Decision Process given

- 1 behaviour of an expert
- 2 dynamics of the system

Tasks

Reward Learning

Estimate the reward function as accurately as possible

- 1 modelling opponents in competitive games
- 2 preference elicitation

Apprenticeship Learning

Use observations of expert behaviour to decide own behaviour

- 1 policy learning
- 2 better generalization of tasks

Bayesian Inference

- evidence is used to infer the probability that a hypothesis may be true
- given a hypothesis H and evidence E , we define:

Definition

Prior probability $P(H)$

Definition

Posterior probability $P(H|E) = \frac{P(E|H)P(H)}{P(E)}$

Main Idea

- the hypothesis: is the reward function that explains the agent's behaviour
- the evidence: observations about expert behaviour
- the evidence is used to infer the probability that a hypothesis may be true (i.e. the posterior distribution of the rewards, from a prior distribution)
- important how to choose the prior:
 - uniform distribution, if there is no information given
 - Laplacian or Gaussian, if most states have negligible rewards
 - Beta, if the MDP is a planning-type problem (most states have low rewards, but a few, goal states have high rewards)

Reward Learning

- learn a reward function

Definition

Error loss function : $L(R, \bar{R}) = \| R - \bar{R} \|^2$, where:

- R is the actual reward
- \bar{R} is the estimated reward

- if R is drawn from the posterior distribution, then $L(R, \bar{R})$ is minimized by setting \bar{R} to the mean of the posterior
- use a maximum a posteriori estimator (MAP) as the estimator

Apprenticeship Learning

- want to learn a policy (i.e. how to act), given expert behaviour

Definition

Class of policy loss functions : $L(R, \pi) = \| V^*(R) - V^\pi(R) \|^P$,
where:

- V^* is the vector of optimal values for each state, under the optimal policy π for R
- want to find a π that minimizes the policy loss over the posterior distribution for R
- direct minimization is hard; instead, find optimal policy π for the mean reward function

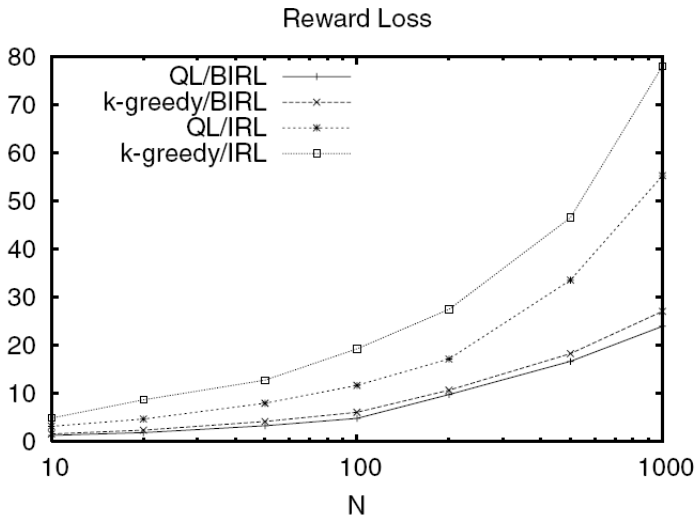
PolicyWalk Sampling Algorithm

Algorithm PolicyWalk(Distribution P , MDP M , Step Size δ)

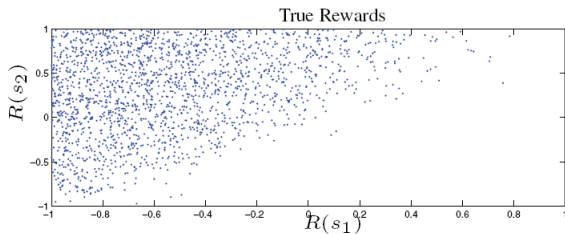
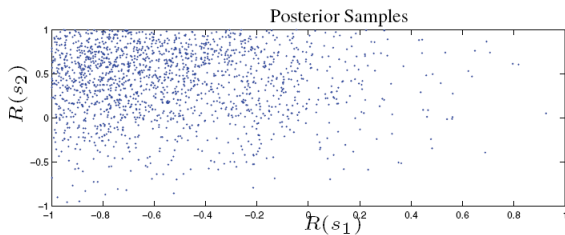
1. Pick a random reward vector $\mathbf{R} \in \mathbb{R}^{|S|}/\delta$.
2. $\pi := \text{PolicyIteration}(M, \mathbf{R})$
3. Repeat
 - (a) Pick a reward vector $\tilde{\mathbf{R}}$ uniformly at random from the neighbours of \mathbf{R} in $\mathbb{R}^{|S|}/\delta$.
 - (b) Compute $Q^\pi(s, a, \tilde{\mathbf{R}})$ for all $(s, a) \in S, A$.
 - (c) If $\exists (s, a) \in (S, A), Q^\pi(s, \pi(s), \tilde{\mathbf{R}}) < Q^\pi(s, a, \tilde{\mathbf{R}})$
 - i. $\tilde{\pi} := \text{PolicyIteration}(M, \tilde{\mathbf{R}}, \pi)$
 - ii. Set $\mathbf{R} := \tilde{\mathbf{R}}$ and $\pi := \tilde{\pi}$ with probability $\min\{1, \frac{P(\tilde{\mathbf{R}}, \tilde{\pi})}{P(\mathbf{R}, \pi)}\}$
 - Else
 - i. Set $\mathbf{R} := \tilde{\mathbf{R}}$ with probability $\min\{1, \frac{P(\tilde{\mathbf{R}}, \pi)}{P(\mathbf{R}, \pi)}\}$
4. Return \mathbf{R}

- generate samples from the prior distribution (of rewards)
- the sample mean is the estimate of the true mean of the distribution

Reward Loss



Posterior Distributions



Conclusions

- reward function can be learned by posing the problem as a Bayesian learning task
- algorithm yields a probability distribution over reward functions that is close to the true one
- policy learning is possible

References

- D. Ramachandran and E. Amir, *Bayesian Inverse Reinforcement Learning*, In IJCAI, 2006
- P. Abbeel and A. Y. Ng. *Apprenticeship learning via inverse reinforcement learning*. In ICML, 2004