



University of Stuttgart
Germany

Maximum Entropy Inverse Reinforcement Learning

Algorithms for Imitation Learning

Maximilian Luz

Summer Semester 2019

MLR/IPVS

Nomenclature

Basis

- Feature Expectation Matching

- Principle of Maximum Entropy

Maximum Entropy IRL

- Algorithm and Derivation

- Extensions

Demonstration

NOMENCLATURE

Markov Decision Process (MDP)

$S = \{s_i\}_i$ States

$A = \{a_i\}_i$ Actions

$T = p(s_{t+1} | s_t, a_t)$ Transition dynamics

$R : S \rightarrow \mathbb{R}$ Reward

Trajectories & Demonstrations

$\tau = ((s_1, a_1), (s_2, a_2), \dots, s_{|\tau|})$ Trajectory

$\mathcal{D} = \{\tau_i\}_i$ Demonstrations

Features

$$\phi : S \rightarrow \mathbb{R}^d \quad \text{with} \quad \phi(\tau) = \sum_{s_t \in \tau} \phi(s_t)$$

Policies

$\pi(a_j | s_i)$ Policy (stochastic)
 π^L Learner Policy
 π^E Expert Policy

BASIS

Feature Expectation Matching

Idea: Learner should visit same features as expert (in expectation).

Feature Expectation Matching [Abbeel and Ng 2004]

$$\mathbb{E}_{\pi^E} [\phi(\tau)] = \mathbb{E}_{\pi^L} [\phi(\tau)]$$

Note: We want to find reward $R : S \rightarrow \mathbb{R}$ defining $\pi^L(a | s)$ and thus $p(\tau)$.

$$\mathbb{E}_{\pi^L} [\phi(\tau)] = \sum_{\tau \in \mathcal{T}} p(\tau) \cdot \phi(\tau)$$

Observation: Optimality for *linear* (unknown) reward [Abbeel and Ng 2004].

\Rightarrow $R(s) = \omega^\top \phi(s), \quad \omega \in \mathbb{R}^d$: Reward parameters

Problem: Multiple (infinite) solutions \Rightarrow ill-posed (Hadamard).

- Reward shaping [Ng et al. 1999]:
 - Multiple reward functions R lead to same policy π .

Idea (Ziebart et al. 2008):

- Regularize by maximizing entropy $H(p)$.
 - But why?

Shannon's Entropy

Entropy $H(p)$

$$H(p) = - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x)$$

$x \in \mathcal{X}$: Event

$p(x)$: Probability of occurrence

$-\log_2 p(x)$: Optimal encoding length

Expected information received when observing $x \in \mathcal{X}$.

⇒ Measure of uncertainty.



No uncertainty, $H(p)$ minimal.



Uniformly random, $H(p)$ maximal.

Principle of Maximum Entropy [Jaynes 1957]

Consider: A problem with solutions p, q, \dots

(e.g. feature expectation matching)

$\Rightarrow p, q$ represent *partial* information.



\Rightarrow Maximizing entropy minimizes bias.

MAXIMUM ENTROPY IRL

Problem Formulation

$$\begin{array}{lll} \arg \max_p & H(p) & \text{(entropy)} \\ \text{subject to} & \mathbb{E}_{\pi^E} [\phi(\tau)] = \mathbb{E}_{\pi^L} [\phi(\tau)], & \text{(feature matching)} \\ & \sum_{\tau \in \mathcal{T}} p(\tau) = 1, \quad \forall \tau \in \mathcal{T} : p(\tau) > 0 & \text{(prob. distr.)} \end{array}$$

Solution: Deterministic Dynamics

Solution via Lagrange multipliers [Ziebart et al. 2008]:

$$p(\tau) \propto \exp(R(\tau))$$

where

$$R(\tau) = \boldsymbol{\omega}^\top \boldsymbol{\phi}(\tau)$$

Lagrange multipliers for feature matching

Deterministic transition dynamics:

$$p(\tau | \boldsymbol{\omega}) = \frac{1}{Z(\boldsymbol{\omega})} \exp(\boldsymbol{\omega}^\top \boldsymbol{\phi}(\tau))$$

normalization

reward

with

$$Z(\boldsymbol{\omega}) = \sum_{\tau \in \mathcal{T}} \exp(\boldsymbol{\omega}^\top \boldsymbol{\phi}(\tau))$$

partition function

Stochastic transition dynamics:

$$p(\tau | \omega) = \underbrace{\frac{1}{Z_S(\omega)} \exp(\omega^\top \phi(\tau))}_{\propto \text{deterministic}} \underbrace{\prod_{S_t, a_t, S_{t+1} \in \tau} p(S_{t+1} | S_t, a_t)}_{\text{combined transition probability}}$$

assumes limited transition randomness

via adaption of deterministic solution [Ziebart et al. 2008].

Problem: Adaption introduces bias [Osa et al. 2018; Ziebart 2010]:

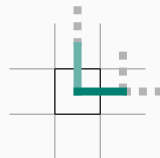
$$\tilde{R}(\tau) = \omega^\top \phi(\tau) + \sum_{S_t, a_t, S_{t+1} \in \tau} \log p(S_{t+1} | S_t, a_t)$$

Solution: Maximum Causal Entropy IRL (Ziebart 2010, not covered here).

Observation:

can be computed via R

$$\pi_{ME}(a_j | s_i, \omega) \propto \sum_{\tau \in \mathcal{T}: s_i, a_j \in \tau_{t=1}} p(\tau | \omega)$$



Idea: Split into sub-problems.

1. *Backward Pass:* Compute policy $\pi_{ME}(a | s, \omega)$.
2. *Forward Pass:* Compute state visitation frequency from $\pi_{ME}(a | s, \omega)$.

State Visitation Frequency: Backward Pass

Observation:

$$\pi_{\text{ME}}(a_j | s_i, \boldsymbol{\omega}) \propto \sum_{\tau \in \mathcal{T}: s_i, a_j \in \tau_{t=1}} p(\tau | \boldsymbol{\omega})$$

Idea:

recursively expand observation

$$\pi_{\text{ME}}(a_j | s_i, \boldsymbol{\omega}) = \frac{Z_{s_i, a_j}}{Z_{s_i}} \leftarrow \text{normalization}$$
$$Z_{s_i, a_j} = \sum_{s_k \in \mathcal{S}} p(s_k | s_i, a_j) \cdot \exp(\boldsymbol{\omega}^\top \boldsymbol{\phi}(s_i)) \cdot Z_{s_k}, \quad Z_{s_i} = \sum_{a_j \in \mathcal{A}} Z_{s_i, a_j}$$

Algorithm:

1. Initialize $Z_{s_k} = 1$ for all terminal states $s_k \in S_{\text{terminal}}$.
2. Compute Z_{s_i, a_j} and Z_{s_i} by recursively backing-up from terminal states.
3. Compute $\pi_{\text{ME}}(a_j | s_i, \boldsymbol{\omega})$.

Parallels to value-iteration.

State Visitation Frequency: Forward Pass

Idea: Propagate starting-state probabilities $p_0(s)$ forward via policy $\pi_{ME}(a | s, \omega)$.

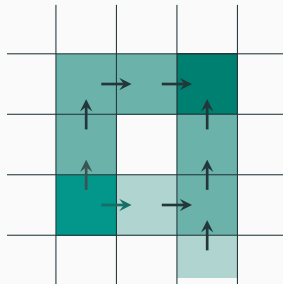
Algorithm:

1. Initialize $D_{s_i,0} = p_0(s) = p(\tau \in \mathcal{T} : s \in \tau_{t=1})$.
2. Recursively compute

$$D_{s_k,t+1} = \sum_{s_i \in \mathcal{S}} \sum_{a_j \in \mathcal{A}} D_{s_i,t} \cdot \pi_{ME}(a_j | s_i) \cdot p(s_k | a_j, s_i)$$

3. Sum up over t , i.e.

$$D_{s_i} = \sum_{t=0,\dots} D_{s_i,t}$$



Algorithm: Iterate until convergence:

1. Compute policy $\pi_{\text{ME}}(a | s, \boldsymbol{\omega})$ (*forward pass*).
2. Compute state visitation frequency D_{s_i} (*backward pass*).
3. Compute gradient $\nabla \mathcal{L}(\boldsymbol{\omega})$ of likelihood.
4. Gradient-based optimization step, e.g.: $\boldsymbol{\omega} \leftarrow \boldsymbol{\omega} + \eta \nabla \mathcal{L}(\boldsymbol{\omega})$.

Assumptions:

- Known transition dynamics $T = p(s_{t+1} | s_t, a_t)$.
- Limited transition randomness.
- Linear reward $R(s) = \boldsymbol{\omega}^\top \boldsymbol{\phi}(s)$.

Other Drawbacks:

- Need to “solve” MDP once per iteration.
- Reward bias for stochastic transition dynamics.

- Maximum Causal Entropy IRL [Ziebart 2010]
- Maximum Entropy Deep IRL [Wulfmeier et al. 2015]
- Maximum Entropy IRL in Continuous State Spaces with Path Integrals [Aghasadeghi and Bretl 2011]

DEMONSTRATION

`github.com/qzed/irl-maxent`

References

- Abbeel, Pieter and Andrew Y. Ng (2004). “Apprenticeship Learning via Inverse Reinforcement Learning”. In: *Proc. 21st Intl. Conference on Machine Learning (ICML '04)*.
- Aghasadeghi, N. and T. Bretl (Sept. 2011). “Maximum entropy inverse reinforcement learning in continuous state spaces with path integrals”. In: *Intl. Conference on Intelligent Robots and Systems (IORS 2011)*, pp. 1561–1566.
- Bishop, Christopher M. (Aug. 17, 2006). *Pattern Recognition and Machine Learning*. Springer-Verlag New York Inc.
- Jaynes, E. T. (May 1957). “Information Theory and Statistical Mechanics”. In: *Physical Review* 106.4, pp. 620–630.
- Ng, Andrew Y., Daishi Harada, and Stuart J. Russell (1999). “Policy Invariance Under Reward Transformations: Theory and Application to Reward Shaping”. In: *Proc. 16th Intl. Conference on Machine Learning (ICML '99)*, pp. 278–287.
- Osa, Takayuki et al. (2018). “An Algorithmic Perspective on Imitation Learning”. In: *Foundations and Trends in Robotics* 7.1-2, pp. 1–179.
- Wulfmeier, Markus, Peter Ondruska, and Ingmar Posner (2015). “Deep Inverse Reinforcement Learning”. In: *Computing Research Repository*. arXiv: 1507.04888.
- Ziebart, Brian D. (2010). “Modeling Purposeful Adaptive Behavior with the Principle of Maximum Causal Entropy”. PhD thesis. Carnegie Mellon University.
- Ziebart, Brian D. et al. (2008). “Maximum Entropy Inverse Reinforcement Learning”. In: *Proc. 23rd AAAI Conference on Artificial Intelligence (AAAI '08)*, pp. 1433–1438.