

# Deep Inverse Reinforcement Learning

Matthew Alger

October 21, 2015

# Inverse reinforcement learning

**Given:** Observations of an agent's behaviour in an environment

# Inverse reinforcement learning

**Given:** Observations of an agent's behaviour in an environment

**Find:** The agent's motivations

# Outline

- 1 Preliminaries
- 2 Linear programming formulation
- 3 Maximum entropy formulation
- 4 Deep maximum entropy

## Markov decision process

$$\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}_{ss'}^a, R, \gamma)$$

- $\mathcal{S}$ , finite set of states
- $\mathcal{A}$ , finite set of actions
- $\{\mathcal{P}_{ss'}^a\}$ , set of transition probabilities
- $R : \mathcal{S} \rightarrow \mathbb{R}$ , reward function
- $\gamma \in [0, 1)$ , discount factor

# Inverse reinforcement learning

**Given:** Observations of an agent's behaviour in an MDP  $\mathcal{M}$

**Find:** The agent's motivations

## Policy

$$\pi : \mathcal{S} \rightarrow \mathcal{A}$$

## Policy

$$\pi : \mathcal{S} \rightarrow \mathcal{A}$$

## Trajectory/path

$$\zeta = [(s_1, a_1), (s_2, a_2), \dots, (s_n, a_n)]$$



## Policy

$$\pi : \mathcal{S} \rightarrow \mathcal{A}$$

## Trajectory/path

$$\zeta = [(s_1, a_1), (s_2, a_2), \dots, (s_n, a_n)]$$

## Trajectory reward

$$R(\zeta) = \sum_{s \in \zeta} R(s)$$

## Value function

$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R(s_{t+k+1}) \mid s_t = s \right]$$

## Value function

$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R(s_{t+k+1}) \mid s_t = s \right]$$

## Bellman equation for $V^\pi$

$$V^\pi(s) = \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^{\pi(s)} (R(s') + \gamma V^\pi(s'))$$

## Value function

$$V^\pi(s) = \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R(s_{t+k+1}) \mid s_t = s \right]$$

## Bellman equation for $V^\pi$

$$V^\pi(s) = \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^{\pi(s)} (R(s') + \gamma V^\pi(s'))$$

## Action-value function

$$Q^\pi(s, a) = \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a (R(s') + \gamma V^\pi(s'))$$

## Optimal policy

$$\pi^* \text{ s. t. } V^{\pi^*}(s) \geq V^{\pi'}(s)$$

# Inverse reinforcement learning

Given: Ideally  $\pi$  and  $\mathcal{P}_{ss'}^a$ , or minimally  $\{\zeta\}$

Find:  $R(s)$  such that  $\pi$  is an optimal policy or  $\{\zeta\}$  is generated by an optimal policy

# Outline

- 1 Preliminaries
- 2 Linear programming formulation**
- 3 Maximum entropy formulation
- 4 Deep maximum entropy

$$(\mathbf{V}^\pi)_i = V^\pi(s_i)$$

$$(\mathbf{R})_i = R(s_i)$$

$$(\mathbf{P}^a)_{ij} = \mathcal{P}_{s_i s_j}^a$$

$$\mathbf{a}^\pi \equiv \pi(s)$$



$$(\mathbf{V}^\pi)_i = V^\pi(s_i)$$

$$(\mathbf{R})_i = R(s_i)$$

$$(\mathbf{P}^a)_{ij} = \mathcal{P}_{s_i s_j}^a$$

$$a^\pi \equiv \pi(s)$$

Can now write the value as a matrix equation:

$$\mathbf{V}^\pi = (\mathbf{R} + \gamma \mathbf{P}^{a^\pi} \mathbf{V}^\pi)$$

## LP — Reward constraint

$$\mathbf{V}^\pi = (\mathbf{R} + \gamma \mathbf{P}^{a^\pi} \mathbf{V}^\pi)$$

## LP — Reward constraint

$$\begin{aligned}\mathbf{V}^\pi &= (\mathbf{R} + \gamma \mathbf{P}^{a^\pi} \mathbf{V}^\pi) \\ &= (\mathbf{I} - \gamma \mathbf{P}^{a^\pi})^{-1} \mathbf{R}\end{aligned}$$

An explicit solution for  $\mathbf{V}^\pi$ !

## LP — Reward constraint

$$\begin{aligned}\mathbf{V}^\pi &= (\mathbf{R} + \gamma \mathbf{P}^{a^\pi} \mathbf{V}^\pi) \\ &= (\mathbf{I} - \gamma \mathbf{P}^{a^\pi})^{-1} \mathbf{R}\end{aligned}$$

An explicit solution for  $\mathbf{V}^\pi$ !

Assuming  $\pi$  is optimal:

$$\forall a \in \mathcal{A} \setminus a^\pi. \quad \mathbf{P}^{a^\pi} \mathbf{V}^\pi \geq \mathbf{P}^a \mathbf{V}^\pi$$

$$\begin{aligned}\mathbf{V}^\pi &= (\mathbf{R} + \gamma \mathbf{P}^{a^\pi} \mathbf{V}^\pi) \\ &= (\mathbf{I} - \gamma \mathbf{P}^{a^\pi})^{-1} \mathbf{R}\end{aligned}$$

An explicit solution for  $\mathbf{V}^\pi$ !

Assuming  $\pi$  is optimal:

$$\forall a \in \mathcal{A} \setminus a^\pi. \quad \mathbf{P}^{a^\pi} \mathbf{V}^\pi \geq \mathbf{P}^a \mathbf{V}^\pi$$

Using our  $\mathbf{V}^\pi$  solution:

$$\forall a \in \mathcal{A} \setminus a^\pi. \quad (\mathbf{P}^{a^\pi} - \mathbf{P}^a)(\mathbf{I} - \gamma \mathbf{P}^{a^\pi})^{-1} \mathbf{R} \geq \mathbf{0}$$

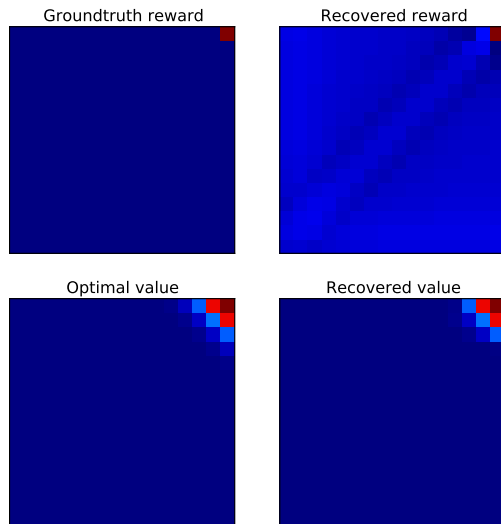
Lots of solutions, so employ a heuristic:

$$\sum_{s \in \mathcal{S}} \left( Q^\pi(s, a^\pi) - \max_{a \in \mathcal{A} \setminus a^\pi} Q^\pi(s, a) \right)$$

# LP — Linear programming problem

$$\begin{aligned} & \text{maximise} \sum_{i=1}^{|\mathcal{S}|} \min_{a \in (\mathcal{A} \setminus a^\pi)} ((\mathbf{P}^{a^\pi})_i - (\mathbf{P}^a)_i) (\mathbf{I} - \gamma \mathbf{P}^{a^\pi})^{-1} \cdot \mathbf{R} - \lambda \|\mathbf{R}\|_1 \\ & \text{s.t. } \forall i \in 1, \dots, |\mathcal{S}|. \forall a \in (\mathcal{A} \setminus a^\pi). -((\mathbf{P}^{a^\pi})_i - (\mathbf{P}^a)_i) (\mathbf{I} - \gamma \mathbf{P}^{a^\pi})^{-1} \cdot \mathbf{R} \leq 0 \\ & \quad \text{and } \forall i \in 1, \dots, |\mathcal{S}|. |R_i| \leq R_{\max} \end{aligned}$$

# LP — Gridworld





# LP — Large state space approximation

$$R(s) = \alpha \cdot \phi(s)$$

# LP — Large state space approximation

$$R(s) = \alpha \cdot \phi(s)$$

$$V^\pi(s) = \sum_i \alpha_i V_i^\pi(s)$$

# LP — Large state space approximation

$$R(s) = \alpha \cdot \phi(s)$$

$$V^\pi(s) = \sum_i \alpha_i V_i^\pi(s)$$

$$\forall a \in \mathcal{A} \setminus a^\pi. \quad P_s^{a^\pi} \cdot V^\pi \geq P_s^a \cdot V^\pi$$

# LP — Large state space approximation

$$R(s) = \alpha \cdot \phi(s)$$

$$V^\pi(s) = \sum_i \alpha_i V_i^\pi(s)$$

$$\forall a \in \mathcal{A} \setminus a^\pi. \quad \mathbf{P}_s^{a^\pi} \cdot \mathbf{V}^\pi \geq \mathbf{P}_s^a \cdot \mathbf{V}^\pi$$

$$\text{maximise } \sum_{s \in \mathcal{S}_0} \min_{a \in (\mathcal{A} \setminus a^*)} p(\mathbf{P}_s^{a^\pi} \cdot \mathbf{V}^\pi - \mathbf{P}_s^a \cdot \mathbf{V}^\pi)$$

$$\text{s.t. } \forall i \in 1, \dots, D. \quad |\alpha_i| \leq 1$$

$$p(x) = \begin{cases} x, & x \geq 0 \\ 2x, & x < 0 \end{cases}$$

# LP — Summary

Given:  $\pi$  and  $\mathcal{P}_{ss'}^a$

Maximise:  $\sum_{s \in \mathcal{S}_0} \min_{a \in (\mathcal{A} \setminus a^*)} p (P_s^{a^\pi} \cdot V^\pi - P_s^a \cdot V^\pi)$

Such that:  $\forall i \in 1, \dots, D. |\alpha_i| \leq 1$

Result:  $R(s) = \alpha \cdot \phi(s)$

# Outline

- 1 Preliminaries
- 2 Linear programming formulation
- 3 Maximum entropy formulation**
- 4 Deep maximum entropy

## Feature counts

$$\phi_{\zeta} = \sum_{s \in \zeta} \phi(s)$$

## Feature expectations

$$\tilde{\phi} = \frac{1}{N} \sum_{i=1}^N \phi_{\zeta}$$

Good reward functions have optimal policies that generate matching feature expectations.

$$P(\zeta) = \frac{1}{Z} \exp(R(\zeta)) \quad (\text{deterministic})$$



$$P(\zeta) = \frac{1}{Z} \exp(R(\zeta)) \quad (\text{deterministic})$$

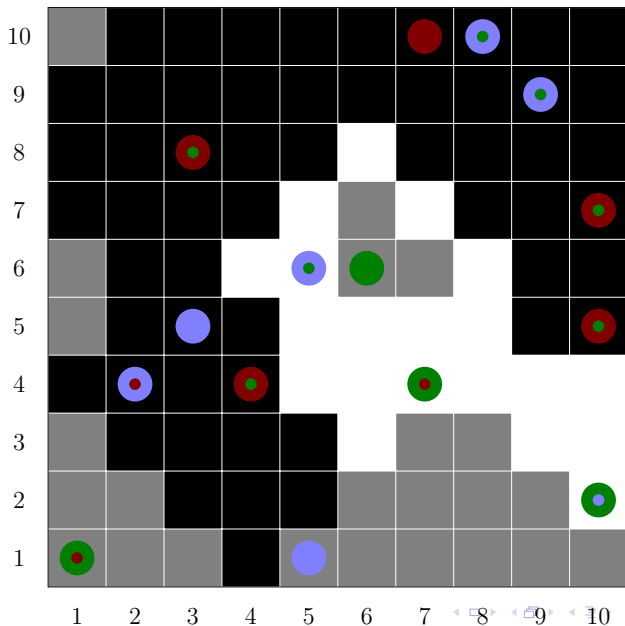
$$P(\zeta \mid \alpha) \approx \frac{\exp(\alpha \cdot \phi_\zeta)}{Z(\alpha)} \prod_{s_{t+1}, a_t, s_t \in \zeta} \mathcal{P}_{s_t s_{t+1}}^{a_t} \quad (\text{non-deterministic})$$

$$\boldsymbol{\alpha}^* = \operatorname{argmax}_{\boldsymbol{\alpha}} \sum_{i=1}^N \log P(\zeta_i \mid \boldsymbol{\alpha})$$

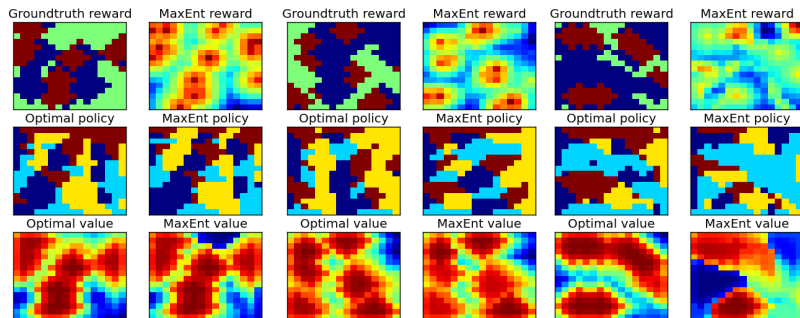
$$\alpha^* = \operatorname{argmax}_{\alpha} \sum_{i=1}^N \log P(\zeta_i | \alpha)$$

$$\frac{\partial}{\partial \alpha} \sum_{i=1}^N \log P(\zeta_i | \alpha) = \tilde{\phi}_{\text{obs}} - \sum_{i=1}^N P(\zeta_i | \alpha) \phi_{\zeta_i} = \tilde{\phi}_{\text{obs}} - \sum_{s \in \mathcal{S}} D(s) \phi(s)$$

# Objectworld



# MaxEnt — Objectworld



# Outline

- 1 Preliminaries
- 2 Linear programming formulation
- 3 Maximum entropy formulation
- 4 Deep maximum entropy**

$$R(s) = \alpha \cdot \phi(s)$$

## DeepMaxEnt — Removing linearity

$$R(s) = \alpha \cdot \phi(s)$$

$$R(s) = \alpha \cdot \varphi(s)$$

$$\varphi(s) = \sigma(\mathbf{W} \cdot \phi(s))$$



$$R(s) = \alpha \cdot \phi(s)$$

$$R(s) = \alpha \cdot \varphi(s)$$

$$\varphi(s) = \sigma(\mathbf{W} \cdot \phi(s))$$

$$R(s) = \alpha \cdot \varphi_n(s)$$

$$\varphi_n(s) = \sigma(\mathbf{W}_n \cdot \varphi_{n-1}(s))$$

⋮

$$\varphi_1(s) = \sigma(\mathbf{W}_1 \cdot \phi(s))$$

$$\partial_{\alpha} \sum_{i=1}^N \log P(\zeta_i | \alpha) = \tilde{\varphi}_{\text{obs},n} - \sum_{s \in \mathcal{S}} D(s) \varphi_n(s)$$

$$\partial_{\alpha} \sum_{i=1}^N \log P(\zeta_i | \alpha) = \tilde{\varphi}_{\text{obs},n} - \sum_{s \in \mathcal{S}} D(s) \varphi_n(s)$$

$$\partial_{\alpha} \sum_{i=1}^N \log P(\zeta_i | \alpha) = \sum_{s \in \mathcal{S}} (D_{\text{obs}}(s) - D(s)) \varphi_n(s)$$

$$\partial_{\alpha} \sum_{i=1}^N \log P(\zeta_i | \alpha) = \tilde{\varphi}_{\text{obs},n} - \sum_{s \in \mathcal{S}} D(s) \varphi_n(s)$$

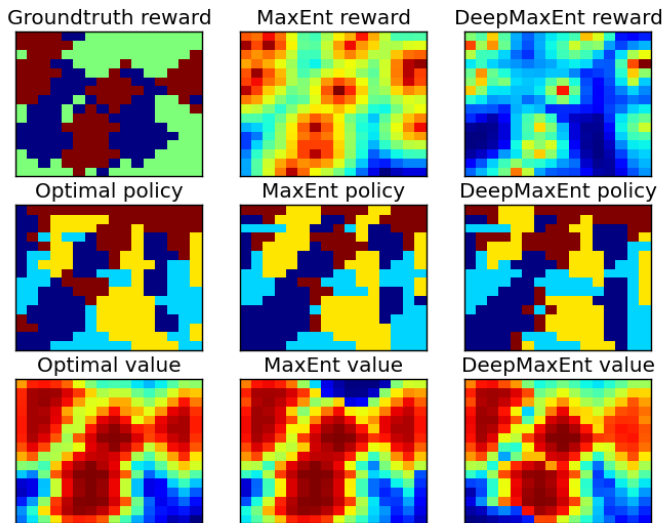
$$\partial_{\alpha} \sum_{i=1}^N \log P(\zeta_i | \alpha) = \sum_{s \in \mathcal{S}} (D_{\text{obs}}(s) - D(s)) \varphi_n(s)$$

$$R(s) = \alpha \cdot \varphi_n(s)$$

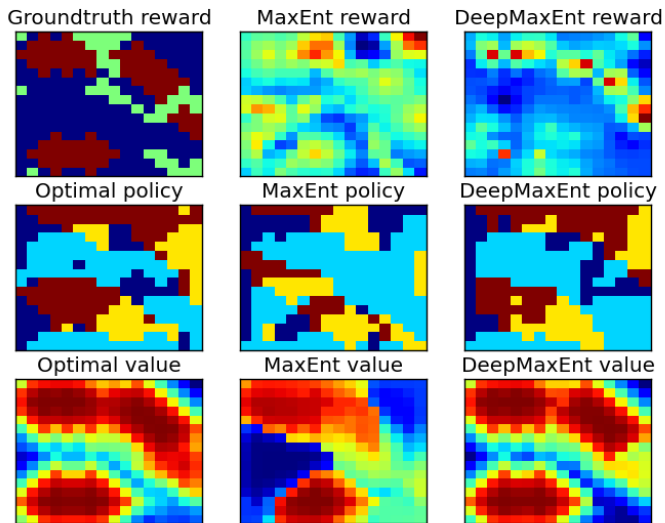
$$\partial_{\alpha} L = \sum_{s \in \mathcal{S}} (D_{\text{obs}}(s) - D(s)) \varphi_n(s)$$

$$\begin{aligned} \partial_{\mathbf{w}_i} L &= \sum_{s \in \mathcal{S}} \frac{\partial L}{\partial R(s)} \frac{\partial R(s)}{\partial \mathbf{w}_i} \\ &= \sum_{s \in \mathcal{S}} (D_{\text{obs}}(s) - D(s)) \partial_{\mathbf{w}_i} R(s) \end{aligned}$$

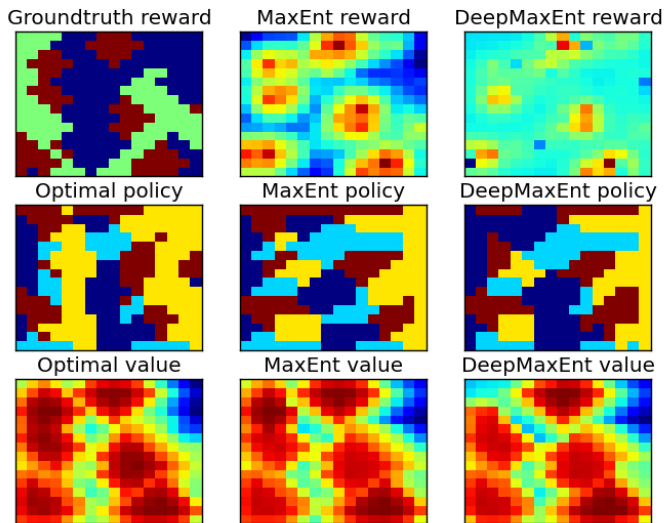
# DeepMaxEnt vs MaxEnt



# DeepMaxEnt vs MaxEnt

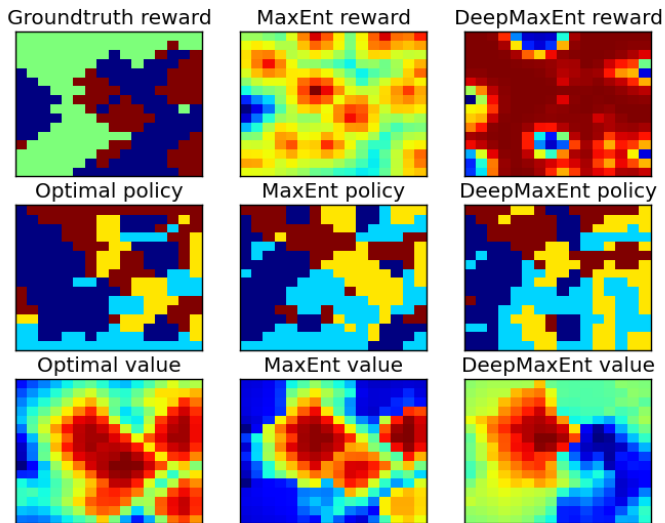


# DeepMaxEnt vs MaxEnt

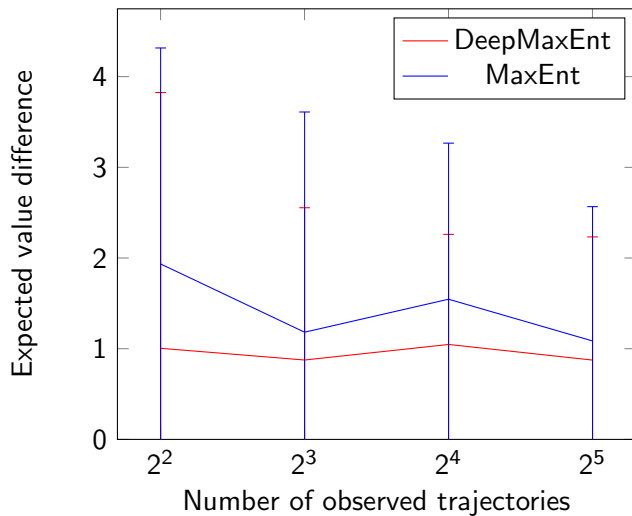




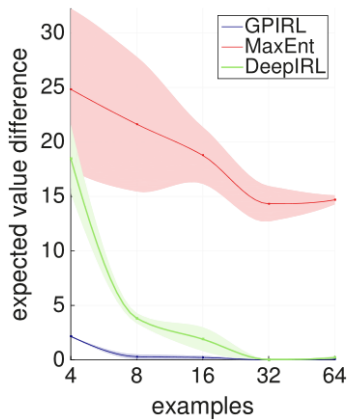
# DeepMaxEnt vs MaxEnt





# Results





# Results — Wulfmeier



 A. Y. Ng and S. Russell.  
Algorithms for inverse reinforcement learning.  
In *ICML*, pages 663–670, 2000.

 R. S. Sutton and A. G. Barto.  
Reinforcement learning: An introduction, 1998.

 M. Wulfmeier, P. Ondruska, and I. Posner.  
Deep inverse reinforcement learning.  
*arXiv preprint arXiv:1507.04888*, 2015.

 B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey.  
Maximum entropy inverse reinforcement learning.  
In *AAAI*, pages 1433–1438, 2008.