

---

# Supplementary Material: MAP Inference for Bayesian Inverse Reinforcement Learning

---

**Jaedeug Choi and Kee-Eung Kim**

Department of Computer Science

Korea Advanced Institute of Science and Technology

Daejeon 305-701, Korea

jdchoi@ai.kaist.ac.kr, kekim@cs.kaist.ac.kr

**Corollary 1** *Given an MDP  $\langle S, A, T, \gamma, \alpha \rangle$ , policy  $\pi$  is optimal if and only if reward function  $\mathbf{R}$  satisfies*

$$\left[ \mathbf{I} - (\mathbf{I}^A - \gamma \mathbf{T})(\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1} \mathbf{E}^\pi \right] \mathbf{R} \leq \mathbf{0}, \quad (1)$$

where  $\mathbf{E}^\pi$  is an  $|S| \times |S||A|$  matrix with the  $(s, (s', a'))$  element being 1 if  $s = s'$  and  $\pi(s') = a'$ , and  $\mathbf{I}^A$  is an  $|S||A| \times |S|$  matrix constructed by stacking the  $|S| \times |S|$  identity matrix  $|A|$  times.

**Proof**

Policy  $\pi$  is optimal

$$\begin{aligned} &\Leftrightarrow \mathbf{Q}_a^\pi(\mathbf{R}) \leq \mathbf{V}^\pi(\mathbf{R}) \\ &\Leftrightarrow \mathbf{R}^a + \gamma \mathbf{T}^a \mathbf{V}^\pi(\mathbf{R}) \leq \mathbf{R}^\pi + \gamma \mathbf{T}^\pi \mathbf{V}^\pi(\mathbf{R}) \\ &\Leftrightarrow \mathbf{R}^a + \gamma \mathbf{T}^a (\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi \leq \mathbf{R}^\pi + \gamma \mathbf{T}^\pi (\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi \\ &\Leftrightarrow \mathbf{R}^a - (\mathbf{I} - \gamma \mathbf{T}^a)(\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi \leq \mathbf{R}^\pi - (\mathbf{I} - \gamma \mathbf{T}^\pi)(\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi \\ &\Leftrightarrow \mathbf{R}^a - (\mathbf{I} - \gamma \mathbf{T}^a)(\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1} \mathbf{E}^\pi \mathbf{R} \leq \mathbf{0} \end{aligned} \quad (2)$$

The third equivalence holds by  $\mathbf{V}^\pi(\mathbf{R}) = (\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1} \mathbf{R}^\pi$ . The fifth equivalence holds because the right-hand side is  $\mathbf{0}$  and  $\mathbf{R}^\pi = \mathbf{E}^\pi \mathbf{R}$ . Stacking up Equation (2) for all  $a \in A$ , we obtain Equation (1). ■

**Theorem 1** *IRL algorithms listed in Table 1 are equivalent to computing the MAP estimates with the prior and the likelihood using  $f(\mathcal{X}; \mathbf{R})$  defined as follows:*

$$\begin{aligned} \bullet f_V(\mathcal{X}; \mathbf{R}) &= \hat{V}^E(\mathbf{R}) - V^*(\mathbf{R}) & \bullet f_G(\mathcal{X}; \mathbf{R}) &= \min_i \left[ V_i^{\pi^*}(\mathbf{R}) - \hat{V}_i^E \right] \\ \bullet f_J(\mathcal{X}; \mathbf{R}) &= - \sum_{s,a} \hat{\mu}_E(s) (J(s, a; \mathbf{R}) - \hat{\pi}_E(s, a))^2 & \bullet f_E(\mathcal{X}; \mathbf{R}) &= \log \mathcal{P}_{\text{MaxEnt}}(\mathcal{X} | \mathbf{T}, \mathbf{R}) \end{aligned}$$

where  $\pi^*(\mathbf{R})$  is an optimal policy induced by the reward function  $\mathbf{R}$ ,  $J(s, a; \mathbf{R})$  is a smooth mapping from reward function  $\mathbf{R}$  to a greedy policy such as the soft-max function, and  $\mathcal{P}_{\text{MaxEnt}}$  is the distribution on the behaviour data (trajectory or path) satisfying the principle of maximum entropy.

We prove Theorem 1 by the following lemmas.

**Lemma 1** *The reward function sought by Ng and Russell's IRL algorithm from sampled trajectories [2] is equivalent to the MAP estimate with the uniform prior and the likelihood using  $f_V(\mathcal{X}; \mathbf{R}) = \hat{V}^E(\mathbf{R}) - V^*(\mathbf{R})$ .*

Table 1: IRL algorithms and their equivalent  $f(\mathcal{X}; \mathbf{R})$  and prior for the Bayesian formulation.  $q \in \{1, 2\}$  is for representing  $L_1$  or  $L_2$  slack penalties.

Previous algorithm	$f(\mathcal{X}; \mathbf{R})$	Prior
Ng and Russell’s IRL from sampled trajectories [2]	$f_V$	Uniform
MMP without the loss function [3]	$(f_V)^q$	Gaussian
MWAL [4]	$f_G$	Uniform
Policy matching [1]	$f_J$	Uniform
MaxEnt [5]	$f_E$	Uniform

**Proof** This IRL algorithm seeks the reward function defined by

$$\mathbf{R}_{\text{N\&R}} = \operatorname{argmax}_{\mathbf{R}} \left[ \hat{V}^E(\mathbf{R}) - V^*(\mathbf{R}) \right].$$

The MAP estimate with the uniform prior and the likelihood using  $f_V$  is computed as

$$\begin{aligned} \mathbf{R}_{\text{MAP}} &= \operatorname{argmax}_{\mathbf{R}} P(\mathbf{R}|\mathcal{X}) = \operatorname{argmax}_{\mathbf{R}} \log P(\mathbf{R}|\mathcal{X}) \\ &= \operatorname{argmax}_{\mathbf{R}} [\log P(\mathcal{X}|\mathbf{R}) + \log P(\mathbf{R})] = \operatorname{argmax}_{\mathbf{R}} f_V(\mathcal{X}; \mathbf{R}) \\ &= \operatorname{argmax}_{\mathbf{R}} \left[ \hat{V}^E(\mathbf{R}) - V^*(\mathbf{R}) \right]. \end{aligned}$$

The MAP estimate is thus equivalent to  $\mathbf{R}_{\text{N\&R}}$ . ■

**Lemma 2** *The reward function sought by the MMP algorithm [3] without the loss function is equivalent to the MAP estimate with a Gaussian prior and the likelihood using  $(f_V)^q$  where  $q \in \{1, 2\}$ .*

**Proof** Without the loss function, the MMP algorithm seeks the reward function defined by

$$\mathbf{R}_{\text{MMP}} = \operatorname{argmin}_{\mathbf{R}} \left[ \left( V^*(\mathbf{R}) - \hat{V}^E(\mathbf{R}) \right)^q + \frac{\lambda}{2} \|\mathbf{R}\|_2^2 \right]$$

where  $q \in \{1, 2\}$  denotes  $L_1$  or  $L_2$  slack penalties. The MAP estimate with a Gaussian prior  $\mathcal{N}(0, \sigma^2)$  and the likelihood using  $(f_V)^q$  is computed as

$$\begin{aligned} \mathbf{R}_{\text{MAP}} &= \operatorname{argmax}_{\mathbf{R}} P(\mathbf{R}|\mathcal{X}) = \operatorname{argmax}_{\mathbf{R}} [\log P(\mathcal{X}|\mathbf{R}) + \log P(\mathbf{R})] \\ &= \operatorname{argmax}_{\mathbf{R}} \left[ \beta (f_V(\mathcal{X}; \mathbf{R}))^q - \frac{1}{2\sigma^2} \sum_{s,a} \mathbf{R}(s, a)^2 \right] \\ &= \operatorname{argmax}_{\mathbf{R}} \left[ (f_V(\mathcal{X}; \mathbf{R}))^q - \frac{1}{2\beta\sigma^2} \|\mathbf{R}\|_2^2 \right] \\ &= \operatorname{argmin}_{\mathbf{R}} \left[ \left( V^*(\mathbf{R}) - \hat{V}^E(\mathbf{R}) \right)^q + \frac{1}{2\beta\sigma^2} \|\mathbf{R}\|_2^2 \right]. \end{aligned}$$

If we set  $\lambda = 1/(\beta\sigma^2)$ , the MAP estimate is equivalent to  $\mathbf{R}_{\text{MMP}}$ . ■

**Lemma 3** *When the reward function is linearly parameterized using the weight vector  $\mathbf{w} \geq \mathbf{0}$  such that  $\sum_i w_i = 1$ , the policy sought by the MWAL algorithm [4] is equivalent to an optimal policy on the reward function which is the MAP estimate with the uniform prior and the likelihood using  $f_G(\mathcal{X}; \mathbf{R}) = \min_i [V_i^{\pi^*(\mathbf{R})} - \hat{V}_i^E]$  where  $\pi^*(\mathbf{R})$  is an optimal policy induced by the reward function  $\mathbf{R}$ .*

**Proof** The MWAL algorithm seeks the policy  $\pi_{\text{MWAL}}$  defined by

$$\pi_{\text{MWAL}} = \operatorname{argmax}_{\pi} \min_i \left[ V_i^{\pi} - \hat{V}_i^E \right],$$

with an implicitly computed reward function  $\mathbf{R}_{\text{MWAL}}$  that induces  $\pi_{\text{MWAL}}$  as an optimal policy. Hence, we can rewrite  $\pi_{\text{MWAL}} = \pi^*(\mathbf{R}_{\text{MWAL}})$  where

$$\mathbf{R}_{\text{MWAL}} = \operatorname{argmax}_{\mathbf{R}} \min_i \left[ V_i^{\pi^*(\mathbf{R})} - \hat{V}_i^E \right].$$

The MAP estimate of the reward function with the uniform prior and the likelihood using  $f_G$  is computed as

$$\mathbf{R}_{\text{MAP}} = \operatorname{argmax}_{\mathbf{R}} P(\mathbf{R}|\mathcal{X}) = \operatorname{argmax}_{\mathbf{R}} f_G(\mathcal{X}; \mathbf{R}) = \operatorname{argmax}_{\mathbf{R}} \min_i \left[ V_i^{\pi^*(\mathbf{R})} - \hat{V}_i^E \right].$$

Hence, the optimal policy induced by  $\mathbf{R}_{\text{MAP}}$  is equivalent to  $\pi_{\text{MWAL}}$  since  $\mathbf{R}_{\text{MAP}} = \mathbf{R}_{\text{MWAL}}$ .  $\blacksquare$

**Lemma 4** *The policy sought by the policy matching algorithm [1] is equivalent to an optimal policy on the reward function which is the MAP estimate with the uniform prior and the likelihood using  $f_J(\mathcal{X}; \mathbf{R}) = -\sum_{s,a} \hat{\mu}_E(s) (J(s, a; \mathbf{R}) - \hat{\pi}_E(s, a))^2$ , where  $J(s, a; \mathbf{R})$  is a smooth mapping from reward function  $\mathbf{R}$  to a greedy policy, such as the soft-max function.*

**Proof** The policy matching algorithm seeks the policy  $\pi_{\text{PM}} = J(\mathbf{R}_{\text{PM}})$  such that

$$\mathbf{R}_{\text{PM}} = \operatorname{argmin}_{\mathbf{R}} \sum_{s,a} \hat{\mu}_E(s) (J(s, a; \mathbf{R}) - \hat{\pi}_E(s, a))^2.$$

The MAP estimate of the reward function with the uniform prior and the likelihood using  $f_J$  is computed as

$$\mathbf{R}_{\text{MAP}} = \operatorname{argmax}_{\mathbf{R}} P(\mathbf{R}|\mathcal{X}) = \operatorname{argmax}_{\mathbf{R}} f_J(\mathcal{X}; \mathbf{R}) = \operatorname{argmin}_{\mathbf{R}} \sum_{s,a} \hat{\mu}_E(s) (J(s, a; \mathbf{R}) - \hat{\pi}_E(s, a))^2.$$

Hence,  $\mathbf{R}_{\text{MAP}} = \mathbf{R}_{\text{PM}}$  and the optimal policy induced by  $\mathbf{R}_{\text{MAP}}$  is equivalent to  $\pi_{\text{PM}}$ .  $\blacksquare$

**Lemma 5** *The reward function sought by the MaxEnt algorithm [5] is equivalent to the MAP estimate with the uniform prior and the likelihood using  $f_E(\mathcal{X}; \mathbf{R}) = \log \mathcal{P}_{\text{MaxEnt}}(\mathcal{X}|\mathbf{T}, \mathbf{R})$  where  $\mathcal{P}_{\text{MaxEnt}}$  is the distribution for the behavior data (trajectory or path) satisfying the principle of maximum entropy.*

**Proof** The MaxEnt algorithm seeks the reward function defined by

$$\mathbf{R}_{\text{MaxEnt}} = \operatorname{argmax}_{\mathbf{R}} \log \mathcal{P}_{\text{MaxEnt}}(\mathcal{X}|\mathbf{T}, \mathbf{R})$$

where

$$\begin{aligned} \mathcal{P}_{\text{MaxEnt}}(\mathcal{X}|\mathbf{T}, \mathbf{R}) &= \prod_{m=1}^M \mathcal{P}_{\text{MaxEnt}}(\mathcal{X}_m|\mathbf{T}, \mathbf{R}) \\ &= \prod_{m=1}^M \frac{1}{Z} \exp \left( \sum_{h=1}^H \gamma^{h-1} \mathbf{R}(s_h^m, a_h^m) \right) \prod_{h=1}^{H-1} \mathbf{T}(s_h^m, a_h^m, s_{h+1}^m). \end{aligned}$$

The MAP estimate with the uniform prior and the likelihood using  $f_E$  is computed as

$$\mathbf{R}_{\text{MAP}} = \operatorname{argmax}_{\mathbf{R}} P(\mathbf{R}|\mathcal{X}) = \operatorname{argmax}_{\mathbf{R}} f_E(\mathcal{X}; \mathbf{R}) = \operatorname{argmax}_{\mathbf{R}} \log \mathcal{P}_{\text{MaxEnt}}(\mathcal{X}|\mathbf{T}, \mathbf{R}).$$

The MAP estimate is thus equivalent to  $\mathbf{R}_{\text{MaxEnt}}$ .  $\blacksquare$

**Theorem 2**  $V^*(\mathbf{R})$  and  $Q^*(\mathbf{R})$  are convex.

**Proof** Let  $C(\pi)$  be the reward optimality region w.r.t.  $\pi$ .  $\mathbf{V}^*(\mathbf{R}) = \mathbf{V}^\pi(\mathbf{R}) = (\mathbf{I} - \gamma\mathbf{T}^\pi)^{-1}\mathbf{E}^\pi\mathbf{R}$  for any  $\mathbf{R} \in C(\pi)$ ,  $\mathbf{V}^*(\mathbf{R})$  is linear w.r.t.  $\mathbf{R}$ . For each and every  $\mathbf{R}_1, \mathbf{R}_2$ , and  $0 \leq \mu \leq 1$ ,

$$\begin{aligned} \mathbf{V}^*(\mu\mathbf{R}_1 + (1 - \mu)\mathbf{R}_2) &= \mathbf{H}^\pi(\mu\mathbf{R}_1 + (1 - \mu)\mathbf{R}_2) = \mu\mathbf{H}^\pi\mathbf{R}_1 + (1 - \mu)\mathbf{H}^\pi\mathbf{R}_2 \\ &= \mu\mathbf{V}^\pi(\mathbf{R}_1) + (1 - \mu)\mathbf{V}^\pi(\mathbf{R}_2) \leq \mu\mathbf{V}^*(\mathbf{R}_1) + (1 - \mu)\mathbf{V}^*(\mathbf{R}_2) \end{aligned}$$

where  $\pi$  is an optimal policy for  $\mu\mathbf{R}_1 + (1 - \mu)\mathbf{R}_2$  and  $\mathbf{H}^\pi = (\mathbf{I} - \gamma\mathbf{T}^\pi)^{-1}\mathbf{E}^\pi$ . Thus,  $\mathbf{V}^*(\mathbf{R})$  is convex. In the same manner, we can also show that  $\mathbf{Q}^*(\mathbf{R})$  is convex using the definition  $\mathbf{Q}^\pi(\mathbf{R}) = \mathbf{R} + \gamma\mathbf{T}\mathbf{E}^\pi\mathbf{Q}^\pi(\mathbf{R})$ . ■

**Theorem 3**  $\mathbf{V}^*(\mathbf{R})$  and  $\mathbf{Q}^*(\mathbf{R})$  are differentiable almost everywhere.

**Proof** Let  $C(\pi)$  be the reward optimality region w.r.t.  $\pi$ . Since  $\mathbf{V}^*(\mathbf{R}) = \mathbf{V}^\pi(\mathbf{R}) = (\mathbf{I} - \gamma\mathbf{T}^\pi)^{-1}\mathbf{E}^\pi\mathbf{R}$  is linear for any  $\mathbf{R} \in C(\pi)$ ,  $\mathbf{V}^*(\mathbf{R})$  is differentiable and  $\nabla_{\mathbf{R}}\mathbf{V}^*(\mathbf{R}) = (\mathbf{I} - \gamma\mathbf{T}^\pi)^{-1}\mathbf{E}^\pi$  when  $\mathbf{R}$  is strictly inside the region. On the boundary,  $\nabla_{\mathbf{R}}\mathbf{V}^\pi(\mathbf{R})$  is a subgradient of  $\mathbf{V}^*(\mathbf{R})$  since the function is convex from Theorem 2 and thus  $\nabla_{\mathbf{R}}\mathbf{V}^\pi(\mathbf{R})(\mathbf{R} - \mathbf{R}') \leq \mathbf{V}^*(\mathbf{R}) - \mathbf{V}^*(\mathbf{R}')$  for any  $\mathbf{R}'$ . In the same manner, we can also show that  $\mathbf{Q}^*(\mathbf{R})$  is differentiable with  $\nabla_{\mathbf{R}}\mathbf{Q}^*(\mathbf{R}) = (\mathbf{I} - \gamma\mathbf{T}\mathbf{E}^\pi)^{-1}$  strictly inside reward optimality regions and  $\nabla_{\mathbf{R}}\mathbf{Q}^\pi(\mathbf{R})$  is a subgradient on the boundaries. ■

## References

- [1] G. Neu and C. Szepesvári. Apprenticeship learning using inverse reinforcement learning and gradient methods. In *Proceedings of UAI*, 2007.
- [2] A. Y. Ng and S. Russell. Algorithms for inverse reinforcement learning. In *Proceedings of ICML*, 2000.
- [3] N. D. Ratliff, J. A. Bagnell, and M. A. Zinkevich. Maximum margin planning. In *Proceedings of ICML*, 2006.
- [4] U. Syed and R. E. Schapire. A game-theoretic approach to apprenticeship learning. In *Proceedings of NIPS*, 2008.
- [5] B. D. Ziebart, A. Maas, J. A. Bagnell, and A. K. Dey. Maximum entropy inverse reinforcement learning. In *Proceedings of AAAI*, 2008.