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\R $\langle S, A, T, \gamma, \alpha \rangle$, policy π is optimal if and only if reward function **R** satisfies

$$\left[\boldsymbol{I} - (\boldsymbol{I}^A - \gamma \boldsymbol{T})(\boldsymbol{I} - \gamma \boldsymbol{T}^{\pi})^{-1} \boldsymbol{E}^{\pi} \right] \boldsymbol{R} \le \boldsymbol{0}, \tag{1}$$

where \mathbf{E}^{π} 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

$$\Rightarrow \boldsymbol{Q}_{a}^{\pi}(\boldsymbol{R}) \leq \boldsymbol{V}^{\pi}(\boldsymbol{R})$$

$$\Rightarrow \boldsymbol{R}^{a} + \gamma \boldsymbol{T}^{a} \boldsymbol{V}^{\pi}(\boldsymbol{R}) \leq \boldsymbol{R}^{\pi} + \gamma \boldsymbol{T}^{\pi} \boldsymbol{V}^{\pi}(\boldsymbol{R})$$

$$\Rightarrow \boldsymbol{R}^{a} + \gamma \boldsymbol{T}^{a} (\boldsymbol{I} - \gamma \boldsymbol{T}^{\pi})^{-1} \boldsymbol{R}^{\pi} \leq \boldsymbol{R}^{\pi} + \gamma \boldsymbol{T}^{\pi} (\boldsymbol{I} - \gamma \boldsymbol{T}^{\pi})^{-1} \boldsymbol{R}^{\pi}$$

$$\Rightarrow \boldsymbol{R}^{a} - (\boldsymbol{I} - \gamma \boldsymbol{T}^{a}) (\boldsymbol{I} - \gamma \boldsymbol{T}^{\pi})^{-1} \boldsymbol{R}^{\pi} \leq \boldsymbol{R}^{\pi} - (\boldsymbol{I} - \gamma \boldsymbol{T}^{\pi}) (\boldsymbol{I} - \gamma \boldsymbol{T}^{\pi})^{-1} \boldsymbol{R}^{\pi}$$

$$\Rightarrow \boldsymbol{R}^{a} - (\boldsymbol{I} - \gamma \boldsymbol{T}^{a}) (\boldsymbol{I} - \gamma \boldsymbol{T}^{\pi})^{-1} \boldsymbol{R}^{\pi} \leq \boldsymbol{R}^{\pi} - (\boldsymbol{I} - \gamma \boldsymbol{T}^{\pi}) (\boldsymbol{I} - \gamma \boldsymbol{T}^{\pi})^{-1} \boldsymbol{R}^{\pi}$$
(2)

The third equivalence holds by $V^{\pi}(\mathbf{R}) = (\mathbf{I} - \gamma \mathbf{T}^{\pi})^{-1} \mathbf{R}^{\pi}$. The fifth equivalence holds because the right-hand side is 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:*

•
$$f_V(\mathcal{X}; \mathbf{R}) = \hat{V}^E(\mathbf{R}) - V^*(\mathbf{R})$$

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

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}_{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] MMP without the loss function [3] MWAL [4] Policy matching [1] MaxEnt [5]	$\begin{array}{c} f_V \\ (f_V)^q \\ f_G \\ f_J \\ f_E \end{array}$	Uniform Gaussian Uniform Uniform Uniform

Proof This IRL algorithm seeks the reward function defined by

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

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

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

The MAP estimate is thus equivalent to $R_{N\&R}$.

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

$$oldsymbol{R}_{ ext{MMP}} = \operatorname*{argmin}_{oldsymbol{R}} \left[\left(V^*(oldsymbol{R}) - \hat{V}^E(oldsymbol{R})
ight)^q + rac{\lambda}{2} \parallel oldsymbol{R} \parallel_2^2
ight]$$

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

$$\begin{split} \boldsymbol{R}_{\text{MAP}} &= \operatorname*{argmax}_{\boldsymbol{R}} P(\boldsymbol{R}|\mathcal{X}) = \operatorname*{argmax}_{\boldsymbol{R}} \left[\log P(\mathcal{X}|\boldsymbol{R}) + \log P(\boldsymbol{R}) \right] \\ &= \operatorname*{argmax}_{\boldsymbol{R}} \left[\beta \left(f_V(\mathcal{X};\boldsymbol{R}) \right)^q - \frac{1}{2\sigma^2} \sum_{s,a} \boldsymbol{R}(s,a)^2 \right] \\ &= \operatorname*{argmax}_{\boldsymbol{R}} \left[\left(f_V(\mathcal{X};\boldsymbol{R}) \right)^q - \frac{1}{2\beta\sigma^2} \parallel \boldsymbol{R} \parallel_2^2 \right] \\ &= \operatorname*{argmin}_{\boldsymbol{R}} \left[\left(V^*(\boldsymbol{R}) - \hat{V}^E(\boldsymbol{R}) \right)^q + \frac{1}{2\beta\sigma^2} \parallel \boldsymbol{R} \parallel_2^2 \right]. \end{split}$$

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

Lemma 3 When the reward function is linearly parameterized using the weight vector $\boldsymbol{w} \geq \boldsymbol{0}$ 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 $f_G(\mathcal{X}; \boldsymbol{R}) = \min_i [V_i^{\pi^*(\boldsymbol{R})} - \hat{V}_i^E]$ where $\pi^*(\boldsymbol{R})$ is an optimal policy induced by the reward function \boldsymbol{R} .

Proof The MWAL algorithm seeks the policy π_{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}_{MWAL} that induces π_{MWAL} as an optimal policy. Hence, we can rewrite $\pi_{MWAL} = \pi^*(\mathbf{R}_{MWAL})$ where

$$m{R}_{ ext{MWAL}} = rgmax_{m{R}}\min_{i} \left[V_{i}^{\pi^{*}(m{R})} - \hat{V}_{i}^{E}
ight].$$

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

$$\boldsymbol{R}_{\text{MAP}} = \operatorname*{argmax}_{\boldsymbol{R}} P(\boldsymbol{R}|\boldsymbol{\mathcal{X}}) = \operatorname*{argmax}_{\boldsymbol{R}} f_{G}(\boldsymbol{\mathcal{X}};\boldsymbol{R}) = \operatorname*{argmax}_{\boldsymbol{R}} \min_{i} \left[V_{i}^{\pi^{*}(\boldsymbol{R})} - \hat{V}_{i}^{E} \right]$$

Hence, the optimal policy induced by R_{MAP} is equivalent to π_{MWAL} since $R_{MAP} = R_{MWAL}$.

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 $\pi_{PM} = J(\mathbf{R}_{PM})$ such that

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

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

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

Hence, $R_{MAP} = R_{PM}$ and the optimal policy induced by R_{MAP} is equivalent to π_{PM} .

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 $f_E(\mathcal{X}; \mathbf{R}) = \log \mathcal{P}_{MaxEnt}(\mathcal{X}|\mathbf{T}, \mathbf{R})$ where \mathcal{P}_{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

$$m{R}_{ ext{MaxEnt}} = rgmax \log \mathcal{P}_{ ext{MaxEnt}}(\mathcal{X}|m{T},m{R})$$

where

$$\mathcal{P}_{\text{MaxEnt}}(\mathcal{X}|\boldsymbol{T},\boldsymbol{R}) = \prod_{m=1}^{M} \mathcal{P}_{\text{MaxEnt}}(\mathcal{X}_{m}|\boldsymbol{T},\boldsymbol{R})$$
$$= \prod_{m=1}^{M} \frac{1}{Z} \exp\left(\sum_{h=1}^{H} \gamma^{h-1}\boldsymbol{R}(s_{h}^{m},a_{h}^{m})\right) \prod_{h=1}^{H-1} \boldsymbol{T}(s_{h}^{m},a_{h}^{m},s_{h+1}^{m}).$$

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

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

The MAP estimate is thus equivalent to R_{MaxEnt} .

Theorem 2 $V^*(R)$ and $Q^*(R)$ are convex.

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

$$V^{*}(\mu \mathbf{R}_{1} + (1-\mu)\mathbf{R}_{2}) = H^{\pi}(\mu \mathbf{R}_{1} + (1-\mu)\mathbf{R}_{2}) = \mu H^{\pi}\mathbf{R}_{1} + (1-\mu)H^{\pi}\mathbf{R}_{2}$$

= $\mu V^{\pi}(\mathbf{R}_{1}) + (1-\mu)V^{\pi}(\mathbf{R}_{2}) \le \mu V^{*}(\mathbf{R}_{1}) + (1-\mu)V^{*}(\mathbf{R}_{2})$

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

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