

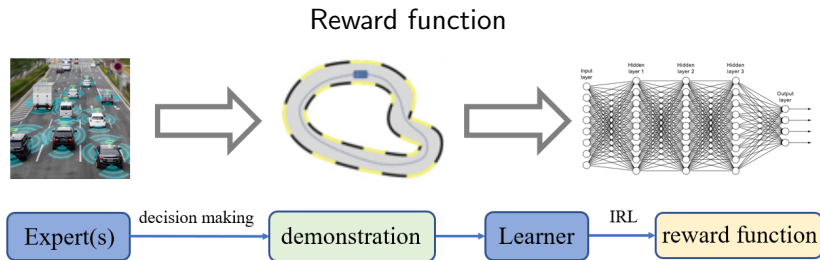
Distributed Inverse Constrained Reinforcement Learning (D-ICRL) for Multi-agent Systems (MASs)

Shicheng Liu & Minghui Zhu

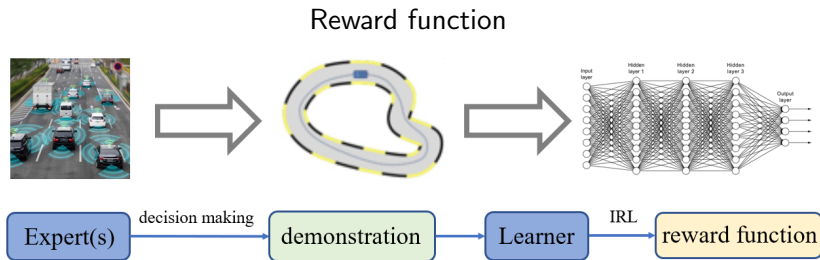
The Pennsylvania State University

Neural Information Processing Systems 2022

Distributed inverse constrained reinforcement learning



Distributed inverse constrained reinforcement learning

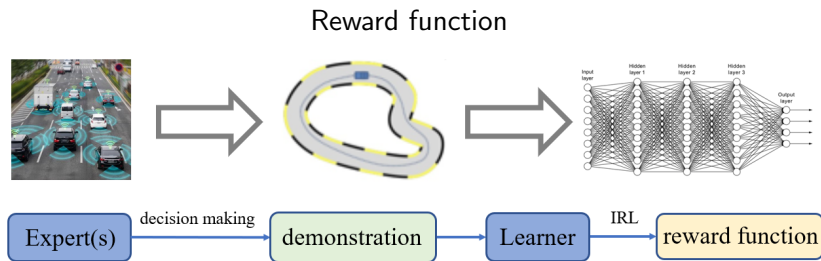


Constraints



Distributed demonstrations

Distributed inverse constrained reinforcement learning



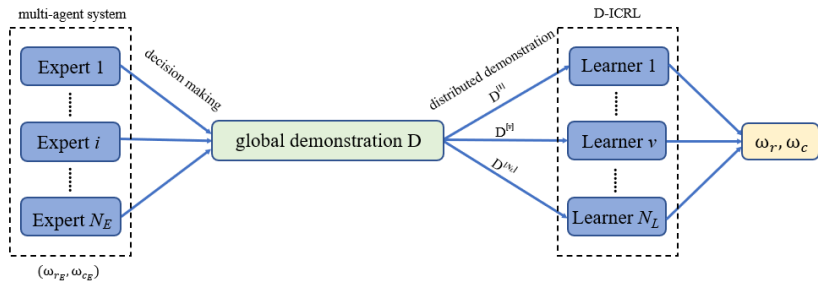
Constraints



Distributed demonstrations

D-ICRL: Solve the three challenges at once.

Model: Multiple experts & multiple learners



N_E cooperative experts: $f_{r_E} = \{r_E^j; r; c_E = \{c_E^j; c; g\}$ $D = fD^{[v]}g_{v=1}^{N_L}$

N_L collaborative learners: $fD^{[v]} = f^j g_{j=1}^{m^{[v]}; r; c; g}$ $f! r; ! c; g$

Distributed bi-level optimization formulation

Distributed bi-level optimization formulation

$$\begin{aligned} \max_{!_c \in \Omega_c} F(!_c; (!_c)) &= \prod_{v=1}^{N_L} F^{[v]}(!_c; (!_c)); && \text{(outer level)} \\ \text{s.t. } (!_c) &= \arg \min \prod_{v=1}^{N_L} m^{[v]} G^{[v]}(; !_c); && \text{(inner level)} \end{aligned}$$

The outer level learns constraints by maximizing the log likelihood $\prod_{v=1}^{N_L} F^{[v]}$ of the demonstrations.

Given a constraint estimate $!_c$, the inner level learns the corresponding reward function and policy by minimizing the dual function $\prod_{v=1}^{N_L} m^{[v]} G^{[v]}$ of maximum causal entropy (MCE).

A perspective of double-loop learning

Double-loop communication: sharing reward and cost function parameters

Inner communication (faster) $W^{[v^0]}(k)$ and $N^{[v]}(k)$.

Outer communication (slower) $W^{[v^0]}(n)$ and $N^{[v]}(n)$.

Inner process

Receives $^{[v^0]}(k)$ from neighbor $v^0 \in N^{[v]}(k)$.

$$^{[v]}(k+1) = \frac{1}{P} \sum_{v^0=1}^{N_L} W^{[v^0]}(k) \cdot ^{[v^0]}(k) \quad (k) m^{[v]}_r \quad G^{[v]}(^{[v]}(k); ! c)$$

Runs K iterations

Outer process

Difficulties

Local gradient $\nabla F^V(\theta_c; \pi_c)$ inaccessible.

global gradient $\nabla F(\theta_c; \pi_c)$ inaccessible.

$F(\theta_c; \pi_c)$ non-convex.

Outer process

Difficulties

Local gradient $\nabla F^{[v]}(\mathbf{c}; \mathbf{c}^{[v]})$ inaccessible.

global gradient $\nabla F(\mathbf{c}; \mathbf{c}^{[v]})$ inaccessible.

$F(\mathbf{c}; \mathbf{c}^{[v]})$ non-convex.

Our solutions

Local gradient approximation $\nabla F^{[v]}(\mathbf{c}; \mathbf{c}^{[v]})$.

Global gradient tracking $\mathbf{r}^{[v]}(n) = \sum_{v=1}^{N_L} W^{[v]}(n) \mathbf{r}^{[v]}(n-1)$
 $+ \mathbf{r} F^{[v]}(\mathbf{c}^{[v]}(n); \mathbf{c}^{[v]}(n)) - \mathbf{r} F^{[v]}(\mathbf{c}^{[v]}(n-1); \mathbf{c}^{[v]}(n-1))$.

Successive convex approximation

$\mathbf{c}_c^{[v]}(n) = \text{Project}_c(\mathbf{c}^{[v]}(n) + N_L \mathbf{r}^{[v]}(n))$.

Outer process

Difficulties

Local gradient $\bar{r} F^{[v]}(I_c; (I_c))$ inaccessible.

global gradient $r F(I_c; (I_c))$ inaccessible.

$F(I_c; (I_c))$ non-convex.

Our solutions

Local gradient approximation $\bar{r} F^{[v]}(I_c; \bar{r}^{[v]}(I_c))$.

Global gradient tracking $\bar{r}^{[v]}(n) = \prod_{v=1}^{N_L} \bar{W}^{[v^0]}(n) \bar{r}^{[v^0]}(n-1)$
 $+ \bar{r} F^{[v]}(I_c^{[v]}(n); \bar{r}^{[v]}(I_c^{[v]}(n))) - \bar{r} F^{[v]}(I_c^{[v]}(n-1); \bar{r}^{[v]}(I_c^{[v]}(n-1)))$.

Successive convex approximation

$\bar{r}_c^{[v]}(n) = \text{Project}_{\Omega_c}(I_c^{[v]}(n) + N_L \bar{r}^{[v]}(n))$.

Update rule: $I_c^{[v]}(n+1) = \prod_{v^0=1}^{N_L} (n) \bar{r}_c^{[v^0]}(n) + (1 - \prod_{v^0=1}^{N_L} (n)) I_c^{[v^0]}(n)$

Theoretical guarantee

Convergence rate of inner problem

Suppose $(k) = \frac{1}{k+1}$ where γ is a positive constant, it holds for any learner v and $\epsilon_c \geq \Omega_c$ that

$$J_j^{[v]}(\epsilon_c) - (\epsilon_c) J_j = O\left(\frac{1}{\log K}\right)$$

Asymptotic convergence of outer problem

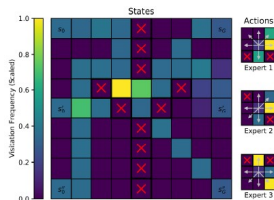
Suppose $(n) \geq (0; 1)$, $\sum_{n=0}^{\infty} (n) = +1$, and $\sum_{n=0}^{\infty} (n)^2 < +1$, it holds for any learner v that

$$\lim_{n \rightarrow \infty} \max_{v; v^0} J_j^{[v]}(\epsilon_c^{[v]}(n)) - (\epsilon_c^{[v]}(n)) J_j = 0;$$

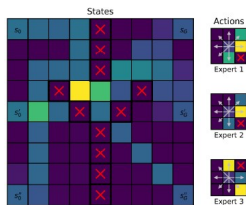
$$\limsup_{n \rightarrow \infty} (r F(\epsilon_c^{[v]}(n); (\epsilon_c^{[v]}(n)))) - (\epsilon_c^{[v]}(n)) = O\left(\frac{1}{\log K}\right)$$

Simulations

Discrete environment

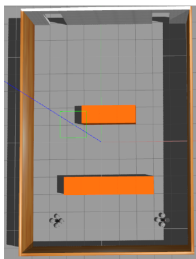


Ground truth environment

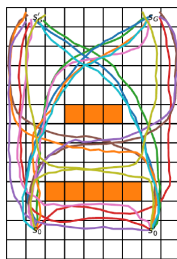


Learned environment

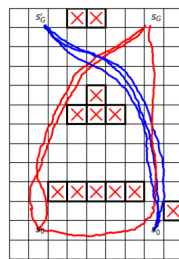
Continuous environment



Simulator environment



Demonstrated trajectories



Learned trajectories

D-ICRL can successfully imitate the experts and recover the constraints.

Conclusion

Solve three challenges at once: Reward function, constraints, and distributed data.

Formulate as a distributed bi-level optimization problem.

D-ICRL: Theoretical framework effective to continuous and discrete environments empirically.

